

# VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



## LAB REPORT

on

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

*Submitted by:*

**PRAJWAL C(1BM22CS198)**

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

*in partial fulfillment for the award of the degree of*  
**BACHELOR OF ENGINEERING**  
*in*  
**COMPUTER SCIENCE AND ENGINEERING**



**B.M.S. COLLEGE OF ENGINEERING**

(Autonomous Institution under VTU)

**BENGALURU-560019**

**March 2024 - June 2024**

**B. M. S. College of Engineering,  
Bull Temple Road, Bangalore 560019**  
(Affiliated To Visvesvaraya Technological University, Belgaum)  
**Department of Computer Science and Engineering**



**CERTIFICATE**

This is to certify that the Lab work entitled "**Big Data Analytics**" carried out by PRAJWAL C(1BM22CS198), who is bonafide student of **B. M. S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum during the year 2024. The Lab report has been approved as it satisfies the academic requirements in respect of **Big Data Analytics – (23CS6PCBDA)** work prescribed for the said degree.

**Vikranth B.M. Dr.**

Associate Professor  
Department of CSE  
BMSCE, Bengaluru

**Kavitha sooda**

Professor and Head  
Department of CSE  
BMSCE, Bengaluru

## Table Of Contents

<b>Sl.no</b>	<b>Program details</b>	<b>Pg no</b>
1	MongoDB- CRUD Operations Demonstration (Practice and Self Study)	1-8
2	Perform the DB operations using Cassandra.	9-13
3	Perform the DB operations using Cassandra	14-16
4	Execution of HDFS Commands for interaction with Hadoop Environment. (Minimum 10 commands to be executed)	17-19
5	Implement Wordcount program on Hadoop framework	20-23
6	a)Create a MapReduce program to find average temperature for each year from the NCDC data set. b) find the mean max temperature for every month.	24-30
7	For a given Text file, Create a Map Reduce program to sort the content in an alphabetic order listing only top 10 maximum occurrences of words.	31-34
8	Write a Scala program to print numbers from 1 to 100 using a for loop.	35
9	Using RDD and FlatMap count how many times each word appears in a file and write out a list of words whose count is strictly greater than 4 using Spark.	36-37
10	Write a simple streaming program in Spark to receive text data streams on a particular port, perform basic text cleaning (like white space removal, stop words removal, lemmatization, etc.), and print the cleaned text on the screen. (Open Ended Question).	38-39

Github Link - <https://github.com/Pra-jwalC/BDA-LAB>

## **Course Outcomes**

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

**CO2:** Analyse data analytic techniques for a given problem.

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

# Program 1

## MongoDB- CRUD Operations Demonstration (Practice and Self Study)

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

### Observation:

1 AB-02

papergrid  
Date: / /

```
1) Perform following DB operations using MongoDB
   db.Customer.insertMany([
     { Cust_id: 1, Acc_Bal: 1500, Acc_Type: '2' },
     { Cust_id: 2, Acc_Bal: 900, Acc_Type: '2' },
     { Cust_id: 3, Acc_Bal: 1000, Acc_Type: '1' },
     { Cust_id: 4, Acc_Bal: 1800, Acc_Type: '2' },
     { Cust_id: 5, Acc_Bal: 1900, Acc_Type: '2' }
   ]);

   acknowledged: true

2) db.Customer.find({ Acc_Bal: { $gt: 1200 },
   Acc_Type: '2' });

   [
     {
       _id: ObjectId('62cffbf5f5eef600'),
       Cust_id: 1,
       Acc_Bal: 1500,
       Acc_Type: '2'
     },
     {
       _id: ObjectId('62cffbf5f5eef601'),
       Cust_id: 4,
       Acc_Bal: 1800,
       Acc_Type: '2'
     }
   ];

   acknowledged: true

3) db.Customer.find({ Acc_Bal: { $gt: 1200 },
   Acc_Type: '2' });

   [
     {
       _id: ObjectId('62cffbf5f5eef600'),
       Cust_id: 1,
       Acc_Bal: 1500,
       Acc_Type: '2'
     }
   ];

   acknowledged: true

4) db.Customer.aggregate([
   {
     $group: {
       _id: '$Cust_id',
       min_balance: { $min: '$Acc_Bal' }
     }
   }
]);
```

```
max_balance: { $max: "f Acc-Bal" }
```

```
}
```

```
}
```

```
});
```

```
{ { id: 1, min_balance: 1800, max_balance: 1800 },
```

```
{ id: 2, min_balance: 2000, max_balance: 2000 },
```

```
{ id: 1, min_balance: 1800, max_balance: 1800 },
```

```
{ id: 2, min_balance: 900, max_balance: 900 },
```

```
{ id: 5, min_balance: 1200, max_balance: 1200 }
```

### Product

2) use e-commerce

switched to db ecommerce

```
db.createCollection("Product")
```

```
db.products.insertMany({
```

```
    Product_id: "P001",
```

```
    name: "Laptop",
```

```
    category: "Electronics",
```

```
    price: 999.99,
```

```
    quantity: 50,
```

```
    description: "High-end gaming laptop",
```

```
},
```

```
    Product_id: "P0.02",
```

```
    name: "Headphones",
```

```
    category: "Electronics",
```

```
    price: 199.99,
```

```
    quantity: 100,
```

```
    description: "Noise-cancelling Headphones"
```

```
},
```

```
})
```

## Code with Output:

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

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

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

```
[ { _id: 1, fruits: [ 'grapes', 'banana', 'apple' ] },
  { _id: 2,
    fruits: [ 'grapes', 'mango', 'cherry' ],
    price: [ { grapes: 80, mango: 100, cherry: 200 } ]
  },
  { _id: 3, fruits: [ 'banana', 'mango' ] }
]
MyDB> db.food.find().pretty();
[
  { _id: 1, fruits: [ 'grapes', 'banana', 'apple' ] },
  { _id: 2,
    fruits: [ 'grapes', 'mango', 'cherry' ],
    price: [ { grapes: 80, mango: 100, cherry: 200 } ]
  },
  { _id: 3, fruits: [ 'banana', 'mango' ] }
]
MyDB> db.createCollection("customer");
{ ok: 1 }
MyDB> var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
Uncaught:
SyntaxError: Unexpected token, expected "," (1:144)
> 1 | var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
|   ^
z |
MyDB> var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
MyDB> db.customer.insert(val);
{
  acknowledged: true,
  insertedIds: [
    '0' : ObjectId("683405cb126af1945c47587e"),
    '1' : ObjectId("683405cb126af1945c47587f"),
    '2' : ObjectId("683405cb126af1945c475872"),
    '3' : ObjectId("683405cb126af1945c475873")
  ]
}
MyDB> db.customer.aggregate({$group:{_id:'$custid',totalbal:{$sum:'$accbal'}}});
[ { _id: 1, totalbal: 200 }, { _id: 2, totalbal: 400 } ]
MyDB> db.customers.aggregate( { $match:{AcctType:"S"}},{$group : { _id : "$custID",TotAccBal :
Uncaught:
SyntaxError: Unexpected character '''. (1:43)
> z | db.Customers.aggregate ( {$match:{AcctType:"S"}},{$group : { _id : "$custID",TotAccBal :
```

```

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

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

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

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

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

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

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

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

```

## Program 2

Perform the following DB operations using Cassandra.

a) Create a keyspace by name Employee

b) Create a column family by name

Employee-Info with attributes Emp\_Id Primary Key, Emp\_Name, Designation, Date\_of\_Joining, Salary, Dept\_Name

c) Insert the values into the table in batch

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

e) Sort the details of Employee records based on salary

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

g) Update the altered table to add project names.

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

**Observation:**

```
    { acknowledged : true
      User
        db.createCollection("users")
        db.user.insertOne({
          username: "john_doe",
          password: "hashedpassword123",
          email: "john.doe@example.com",
          phone_no: "123-456-7890",
          shipping_address: {
            street: "123 Main St",
            city: "Somewhere",
            state: "CA",
            postal_code: "90001",
            country: "USA"
          }
        })
        created_at: new Date(),
        updated_at: new Date()
      }
      acknowledged : true
      Cart
        db.createCollection("cart")
        db.cart.insertOne({
          user_id: ObjectID("user_id_here"),
          products: [
            {
              product_id: ObjectId("product_id_here"),
              quantity: 5,
              price_at_time: 999.99
            }
          ]
        })
    }
```

Total price: 1999.99,  
Created at: new Date(),  
Update at: new Date()  
})

### Order

db.createCollection("order")  
db.order.insertOne({  
 user\_id: ObjectId("user\_id here"),  
 order\_status: "Pending",  
 shipping\_address: {  
 street: "123 Main St",  
 City: "Somewhere",  
 state: "CA",  
 Postal\_code: "90001",  
 Country: "USA"},  
 products: ["P1", {  
 product\_id: ObjectId("product\_id here"),  
 quantity: 1,  
 price\_at\_line: 999.99  
 }],  
 total\_price: 999.99,  
 created\_at: new Date(),  
 update\_at: new Date()  
})

a) db.products.find({})

{  
 id: ObjectId("67 ---"),  
 product\_id: "P001",  
 Name: "laptop",  
 Category: "Electronics",  
 Price: 999.99,

quantity: 50,  
description: 'High-end gaming laptop',  
image: ['mag1.jpg', 'mag2.jpg'],  
{  
    id: ObjectID('67 - -'),  
    product\_id: 'P002',  
    name: "Headphones",  
    category: 'Electronics',  
    price: 199.99,  
    quantity: 100,  
    description: 'Noise-cancelling headphones'  
}

b) db.products.find({ category: "Electronics",  
    {  
        id: ObjectID('67 - -'),  
        product\_id: 'P001',  
        name: 'Laptop',  
        category: 'Electronics',  
        price: 999.99,  
        quantity: 50,  
        description: 'High-end gaming laptop',  
    }  
    {  
        id: ObjectID('67 - -'),  
        product\_id: 'P003',  
        name: "Headphones",  
        category: 'Electronics',  
        price: 199.99,  
        quantity: 100,  
        description: 'Noise-cancelling headphones'  
    }

## Code with Output:

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

student    system      system_distributed  system_traces  system_virtual_schema
students   system_auth  system_schema       system_views

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

library_book  st_info  students_info  userlogin

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

CREATE TABLE students.st_info (
    rollno int PRIMARY KEY,
    doj timestamp,
    name text,
    percent double
) WITH additional_write_policy = '99p'
    AND bloom_filter_fp_chance = 0.01
    AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}
    AND cdc = false
    AND comment = ''
    AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy', 'max_threshold': '32', 'min_threshold': '4'}
    AND compression = {'chunk_length_in_kb': '16', 'class': 'org.apache.cassandra.io.compress.LZ4Compressor'}
    AND memtable = 'default'
    AND crc_check_chance = 1.0
    AND default_time_to_live = 0
    AND extensions = {}
    AND gc_grace_seconds = 864000
    AND max_index_interval = 2048
    AND memtable_flush_period_in_ms = 0
    AND min_index_interval = 128
    AND read_repair = 'BLOCKING'
    AND speculative_retry = '99p';
cqlsh:students> begin batch
... insert into st_info(rollno,name,doj,percent)
```

```
cqlsh:students> select * from st_info;
+-----+-----+-----+
| rollno | doj           | name    | percent |
+-----+-----+-----+
| 1      | 2010-02-28 18:30:00.000000+0000 | preeti  | 90      |
| 2      | 2010-03-19 18:30:00.000000+0000 | prajwal | 89      |
| 4      | 2010-04-22 18:30:00.000000+0000 | rachana | 90      |
+-----+-----+-----+
(3 rows)
cqlsh:students> select * from st_info where rollno in(1,2);
+-----+-----+-----+
| rollno | doj           | name    | percent |
+-----+-----+-----+
| 1      | 2010-02-28 18:30:00.000000+0000 | preeti  | 90      |
| 2      | 2010-03-19 18:30:00.000000+0000 | prajwal | 89      |
+-----+-----+-----+
(2 rows)
cqlsh:students> select * from st_info where name="preeti";
SyntaxException: line 1:42 no viable alternative at input ';' (...* from st_info where name=["preet]i";)
cqlsh:students> create index on st_info(name);
cqlsh:students> select * from st_info where name="preeti";
SyntaxException: line 1:42 no viable alternative at input ';' (...* from st_info where name=["preet]i";)
cqlsh:students> select * from st_info where name='preeti';
+-----+-----+-----+
| rollno | doj           | name    | percent |
+-----+-----+-----+
| 1      | 2010-02-28 18:30:00.000000+0000 | preeti  | 90      |
+-----+-----+-----+
(1 rows)
cqlsh:students> select rollno,name,percent from st_info limit 2;
+-----+-----+-----+
| rollno | name    | percent |
+-----+-----+-----+
| 1      | preeti  | 90      |
| 2      | prajwal | 89      |
+-----+-----+-----+
(2 rows)
cqlsh:students> slect rollno as usn from st_info;
SyntaxException: line 1:0 no viable alternative at input 'slect' ([slect]...)
cqlsh:students> select rollno as usn from st_info;
+-----+
| usn |
+-----+
| 1   |
+-----+
```

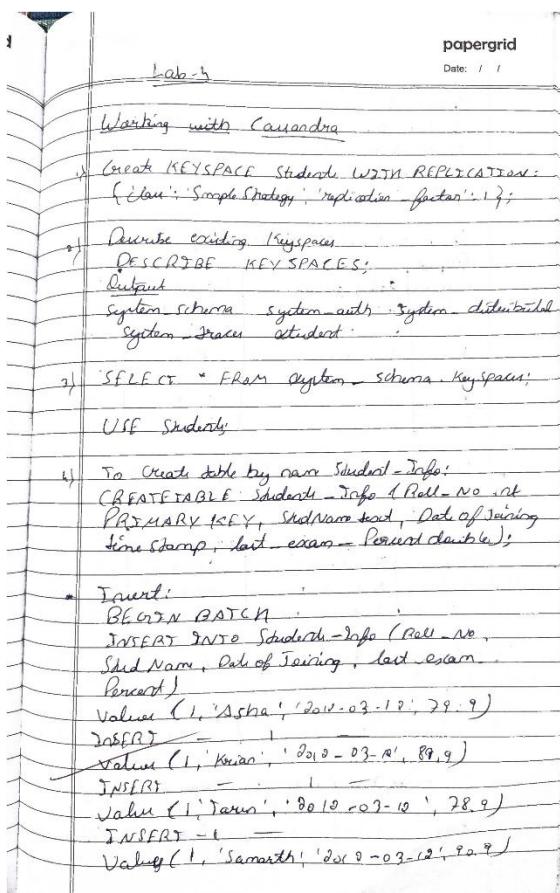
```
+-----+
| usn |
+-----+
| 1   |
| 2   |
| 4   |
+-----+
(3 rows)
cqlsh:students> create table library(c_val counter,book_name varchar,stud_name varchar,primary key(book_name,stud_name));
cqlsh:students> update library set c_val=c_val+1 where book_name='BDA' and stud_name='preeti';
cqlsh:students> create table userlogin(id int primary key,pass text);
AlreadyExists: Table 'students.userlogin' already exists
cqlsh:students> create table login(id int primary key,pass text);
cqlsh:students> insert into login(id,pass) values(1,'infy')using ttl 30;
cqlsh:students> select ttl(pass) from login where id=1;
+-----+
| ttl(pass) |
+-----+
| 3          |
+-----+
```

# Program 3

Perform the following DB operations using Cassandra.

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

**Observation:**



Date: / /

INSERT

Value (1, 'Smitha', '2010-03-12', 62.9)

INSERT

Value (1, 'Rohan', '2010-03-12', 56.9)

APPLY BATCH;

SELECT \* FROM Student-Info;

Output

roll_no	date of joining	last exam period	Stud Name
5	2010-03-11 18:30:00	67.9	Smitha
1	2010-03-11	79.9	Asha
2	2010-03-11	89.9	Kiran
4	2010-03-11	90.9	Samarth
6	2010-03-11	56.9	Rohan
3	2010-03-11	78.9	Tarun

SELECT \* FROM Student-Info WHERE

Roll-No IN (1, 2, 3);

Output

roll_no	date of joining	last exam period	Stud Name
1	2010-03-11	79.9	Asha
2	2010-03-11	89.9	Kiran
3	2010-03-11	78.9	Tarun

Select \* from Student-Info where StudName = 'Asha';

Output

roll_no	date of joining	last exam period	Stud Name
1	2010-03-11	79.9	Asha

To Create INDEX on StudName column of the

Student-Info column family

CREATE INDEX ON Student-Info (StudName);

## Code with Output:

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

libraries    students      system_distributed   system_views
library      system       system_schema        system_virtual_schema
student     system_auth  system_traces

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

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

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

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

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

(1 rows)
cqlsh:libraries> □
```

## Program 4

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

Observation:

papergrid  
Date: / /

Lab 6: Hadoop Exercise

- Start Hadoop  
\$ start-all.sh
- Create directory  
hdfs dfs -mkdir /bda/hadoop
- List contents  
hadoop fs -ls /
- put - copy from Local to HDFS  
hdfs dfs -copyFromLocal /home/Desktop/bda-local.txt /bda/hadoop/file.txt
- cat - display file content  
hdfs dfs -cat /bda/hadoop/file.txt  
"Hello, BMSC"
- hdfs dfs -cat /bda/hadoop/file.txt  
get (download from HDFS)
- hdfs dfs -get /bda/hadoop/file.txt  
/home/Desktop/downloaded-file.txt
- get merge (multiple files)  
hdfs dfs -get merge /bda/hadoop/file.txt

-> get fact

hadoop fs -get fact : (bda - hadoop)

-> copy to Local

hdfs dfs -copyToLocal /abc /www  
(home/Desktop)

-> mv (move file / directory)

hadoop fs -mv /bda - hadoop /abc

get fact output:

# file : /bda - hadoop

# owner: hadoop

# group: supergroup

User: rwx

group: r-x

other: r-x

-> cp (copy a file from one directory to another directory)

hadoop fs -cp /home /Desktop  
/bda - hadoop

Q11/11/15

## Code with Output:

```
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [bmscecse-HP-Elite-Tower-600-G9-Desktop-PC]
Starting resourcemanager
Starting nodemanagers
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -mkdir /bda_hadoop
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -ls /
Found 4 items
drwxr-xr-x  - hadoop supergroup          0 2025-04-15 15:07 /abc
drwxr-xr-x  - hadoop supergroup          0 2025-05-26 14:13 /bda_hadoop
drwxr-xr-x  - hadoop supergroup          0 2025-05-22 16:32 /pqr
drwxr-xr-x  - hadoop supergroup          0 2025-05-20 16:36 /rgs
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -put /home/hadoop/sample.txt /bda_hadoop/file.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/sample.txt /bda_hadoop/local.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -cat /bda_hadoop/file.txt
hi how are you
how is your job
how is your family
how is your brother
how is your sister
eof
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -get /bda_hadoop/file.txt /home/hadoop/sample.txt
get: '/home/hadoop/sample.txt': File exists
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -get /bda_hadoop/file.txt /home/hadoop/get.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop fs -getfacl /bda_hadoop/
# file: /bda_hadoop
# owner: hadoop
# group: supergroup
user::rwx
group::r-x
other::r-x

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

## Program 5

### Implement Wordcount program on Hadoop framework

#### Observation:

**papergrid**  
Date: / /

Lab6 - Hadoop Exercise

- Start Hadoop  
\$ start-all.sh
- Create directory  
hdfs dfs -mkdir /bda-hadoop
- List contents  
hadoop fs -ls /
- put - copy from Local to HDFS  
hdfs dfs -copyFromLocal /home/Downloads/bda-local.txt /bda-hadoop/file.txt
- cat - display file contents  
hdfs dfs -cat /bda-hadoop/file.txt  
"Hello, BMSC"
- hdfs dfs -cat /bda-hadoop/file.txt  
get (download from HDFS)
- hdfs dfs -get /bda-hadoop/file.txt  
/home/Desktop/downloaded-file.txt
- get many (multiple files)  
hdfs dfs -get many /bda-hadoop/file.txt

-> get fact

hadoop fs -get fact : (bda - hadoop)

im  
ii

-> copy to Local

hadoop fs -copyToLocal /abc /local  
(home / desktop)

pu  
rr

-> mv (move file / directory)

hadoop fs -mv /bda - hadoop /abc

f

get fact output:

# file : bda - hadoop

# owner : hadoop

# group : supergroup

User : wr

group : r-x

other : r-x

-> cp (copy a file from one directory to another directory)

hadoop fs -cp /home /Desktop  
/bda - hadoop

07/15/15

## Code with Output:

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

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

## Program 6

From the following link extract the weather data

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

a) Create a MapReduce program to find average temperature for each year from the NCDC data set.

b) find the mean max temperature for every month

Observation:

Lab 7  
WC Mapper Java

papergrid  
Date: / /

```
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reporter;

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

## WCReducer.java

```
import java.io.IOException;
import java.util.Iterator;
import org.apache.hadoop.io.IntWritable;
import java.io.IOException;
import java.util.Iterator;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.Reporter;
```

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

LG

WC Reducer.java

```

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

```

public class WCReducer extends Configuration implements

Tool

public int run(String args[]) throws IOException

{

if (args.length < 2)

{

System.out.println("Please give valid input");

return -1;

}

JobConf conf = new JobConf(WCReducer.class);

FileInputFormat.setInputPath(conf,

new Path(args[0]));

FileOutputFormat.setOutputPath(conf,

new Path(args[1]));

conf.setMapperClass(WCMapper.class);

conf.setReducerClass(WCReducer.class);

conf.setMapOutputKeyClass(Text.class);

conf.setMapOutputValueClass(IntWritable.class);

conf.setOutputKeyClass(IntWritable.class);

conf.setOutputValueClass(IntWritable.class);

JobClient.runJob(conf);

## Code with Output:

### a) `Average temperature

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

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

## b) Maximum temperature

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

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

## Program 7

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

Observation:

Lab 7:  
WC Mapper.java  
papergrid  
Date: / /

```
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reporter;

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

### WC Reducer.java

```

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

```

public class WCReducer extends Configuration implements

Tool

public int run(String args[]) throws IOException

{

if (args.length < 2)

{

System.out.println("Please give valid input");

return -1;

}

JobConf conf = new JobConf(WCReducer.class);

FileInputFormat.setInputPath(conf,

new Path(args[0]));

FileOutputFormat.setOutputPath(conf,

new Path(args[1]));

conf.setMapperClass(WCMapper.class);

conf.setReducerClass(WCReducer.class);

conf.setMapOutputKeyClass(Text.class);

conf.setMapOutputValueClass(Text.class);

conf.setOutputKeyClass(IntWritable.class);

conf.setOutputValueClass(IntWritable.class);

JobClient.runJob(conf);

## Code with Output:

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

```

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

```

## Program 8

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

## **Observation:**

papergrid

Date: / /

Lab 9.

Write a Scala Program to print Numbers from 1 to 100 using for loop

object PrintNumbers {

```
def main(args: Array[String]): Unit = {
    for (i ← 1 to 100) {
        print(i)
    }
}
```

} Print Numbers, main (Array ())

q) Create MapReduce Program to sort the content in an alphabetic.

Check out list of words where count is strictly greater than 6 using spark

```
val data = sc.textFile("sparkdata.txt")
data.collect()
val splitData = data.flatMap(line =>
    line.split(" "))
splitData.collect()
val mapData = splitData.map(record => (record,
    1))
mapData.collect()
val reduceData = mapData.reduceByKey(_ + _)
reduceData.collect()
```

Output

spark: 5

is: 6

## Code with Output:

## Program 9

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

Observation:

papergrid  
Date: / /

Lab 9 -

1) Write a Scala Program to print Numbers from 1 to 100 using for loop

```
object PrintNumbers {
    def main(args: Array[String]): Unit = {
        for (i <= 1 to 100) {
            printIn(i)
        }
    }
}
```

Print Numbers.main(Array())

2) Create MapReduce Program to sort the content in an alphabat

Check out list of words whose count is strictly greater than 4 using spark

```
val data = sc.textFile("sparkdata.txt")
data.collect()
val splitData = data.flatMap(line =>
    line.split(" "))
splitData.collect()
val prepData = splitData.map((word => (word, 1)))
prepData.collect()
val reduceData = prepData.reduceByKey(_ + _)
reduceData.collect()
```

Output  
spark: 5  
is: 6

## Code with Output:

## Program 10

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

### Observation:

Lab 8  
Mean max using python  
memas\_mapper.py  
import sys  
MISSING = 999  
  
for line in sys.stdin:  
 line = line.strip()  
 if len(line) < 93:  
 continue  
 month = line[19:21]  
 if line[87] == '+':  
 temperature = int(line[88:92])  
 else:  
 temperature = int(line[87:92])  
 quality = line[92:93]  
 if temperature != MISSING and quality in  
 ['0', '1', '4', '5', '9']:  
 print(f'{month} {temperature}')  
  
mean\_max\_reducer.py  
  
import sys  
current\_month = None  
temp = []  
  
def emit\_avg(month, temp\_list):  
 max\_temp = 0  
 total = 0  
 count = 0  
 day = 0  
 for t in temp\_list:  
 if t > max\_temp:  
 max\_temp = t

```

count += 1
if count == 3:
    total += max_temp
    max_temp = 0
    count = 0
days += 1
if days > 0:
    average = total / days
    print(f'month {t} ({average})')
for line in sys.stdin:
    line = line.strip()
    if not line:
        continue
    month, temp = line.split(' ')
    temp = int(temp)
    if current_month == month:
        temps.append(temp)
    else:
        emit_average(current_month, temps)
        current_month = month
        temps = [temp]
if current_month:
    emit_average(current_month, temps)

```