

Working with MongoDB

I) Creating DATABASE In MongoDB

```
use myDB;
```

```
db;
```

```
show dbs;
```

(CREATE, READ, UPDATE, DELETE)

II) CRUD (CREATE, READ, UPDATE, DELETE) Operations

```
→ db.createCollection("Student")
```

```
→ db.Student.drop()
```

```
→ db.Student.insert({ _id: 1, StudName: "Michellejintan",  
                      Grade: "VII", Hobbies: "internet surfing" })
```

```
→ db.Student.update({ _id: 3, StudName: "Aryan",  
                      Grade: "VIII" }, { $set: { Hobbies: "Skating" } })
```

```
→ db.Student.find({ StudName: "Aryan",  
                     Grade: "VIII" }).pretty();
```

FIND METHOD

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

```
→ db.Student.find({ _id: 1, Grade: 1 })
```

```
→ db.Student.find({ Grade: { $eq: "VII" } })
```

```
Pretty();
```

```
→ db.Student.find({ Hobbies: { $in: ["Chess",  
                               "Skating"] } }).pretty();
```

```
→ db.Student.find({ StudName: /IM/ }).pretty();
```

```
Pretty();
```

db. Student. find ({ StudName: /e/ }), pretty()
db. Student. count()
db. Student. find(). sort(' StudName:-1'), pretty()

III import data from a csv file

CSV file "Sample.txt"

- Import file into the MongoDB Collection, "SampleJSon".
The collection is in the dataset "test".

mongorestore --db Student --collection airlines --type csv
/home/1hduru/Desktop/airline.csv

IV Export data to a csv file

mongoexport --host localhost --db Student --collection airlines
--csv --out /home/1hduru/Desktop/output.txt -fields
"Year", "Quarter"

V Save method.

db. Student. Save ({ StudName: "Vamei", Grade: 1 })

add a new field

db. Students. Update ({ _id: 4 }, { \$set: { location: "Nehru Nagar" } })

remove a field
db. Students. Update ({ _id: 4 }, { \$unset: { location: "Nehru Nagar" } })

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

db. Student . find ({ Grade: "VII" }, pretty());

Find doc Student name ends with Sia with student document

db. Student . find ({ Stud / name : / s / 1 / 3 }, pretty());

db. Student . update ({ Stud / name : / s / 3 }, pretty());

db. Student . update ({ Stud / name : / s / 3 }, pretty());

→ set a Particular field value to null

db. Student . update ({ id: 3 }, { \$set: { location: null } })

→ Count the number of doc

or one - file db. Student . count ({ Grade: "VII" })

retrieve first 3 documents

db. Student . find ({ Grade: "VII" }). limit (3), pretty();

db. Student . find ({ Grade: "VII" }). limit (3), pretty();

and : begin and end

begin and end

(begin) between : 0

: 1

: 2

: 3

: 4

: 5

: 6

: 7

: 8

: 9

: 10

: 11

: 12

: 13

: 14

: 15

: 16

: 17

: 18

: 19

: 20

docs: [{ _id: "5a0310a310a310a3" }, limit, remainder, db]

docs: [{ _id: "5a0310a310a310a3" }, limit, remainder, db]

docs: [{ _id: "5a0310a310a310a3" }, limit, remainder, db]

: 1

: 2

: 3

: 4

: 5

: 6

: 7

: 8

: 9

: 10

: 11

: 12

: 13

: 14

: 15

: 16

: 17

: 18

: 19

: 20

: 21

: 22

: 23

: 24

: 25

: 26

: 27

: 28

: 29

: 30

: 31

: 32

: 33

: 34

: 35

: 36

: 37

: 38

: 39

: 40

: 41

: 42

: 43

: 44

: 45

: 46

: 47

: 48

: 49

: 50

: 51

: 52

: 53

: 54

: 55

: 56

: 57

: 58

: 59

: 60

: 61

: 62

: 63

: 64

: 65

: 66

: 67

: 68

: 69

: 70

: 71

: 72

: 73

: 74

: 75

: 76

: 77

: 78

: 79

: 80

: 81

: 82

: 83

: 84

: 85

: 86

: 87

: 88

: 89

: 90

: 91

: 92

: 93

: 94

: 95

: 96

: 97

: 98

: 99

: 100

: 101

: 102

: 103

: 104

: 105

: 106

: 107

: 108

: 109

: 110

: 111

: 112

: 113

: 114

: 115

: 116

: 117

: 118

: 119

: 120

: 121

: 122

: 123

: 124

: 125

: 126

: 127

: 128

: 129

: 130

: 131

: 132

: 133

: 134

: 135

: 136

: 137

: 138

: 139

: 140

: 141

: 142

: 143

: 144

: 145

: 146

: 147

: 148

: 149

: 150

: 151

: 152

: 153

: 154

: 155

: 156

: 157

: 158

: 159

: 160

: 161

: 162

: 163

: 164

: 165

: 166

: 167

: 168

: 169

: 170

: 171

: 172

: 173

: 174

: 175

: 176

: 177

: 178

: 179

: 180

: 181

: 182

: 183

: 184

: 185

: 186

: 187

: 188

: 189

: 190

: 191

: 192

: 193

: 194

: 195

: 196

: 197

: 198

: 199

: 200

: 201

: 202

: 203

: 204

: 205

: 206

: 207

: 208

: 209

: 210

: 211

: 212

: 213

: 214

: 215

: 216

: 217

: 218

: 219

: 220

: 221

: 222

: 223

: 224

: 225

: 226

: 227

: 228

: 229

: 230

: 231

: 232

: 233

: 234

: 235

: 236

: 237

: 238

: 239

: 240

: 241

: 242

: 243

: 244

: 245

: 246

: 247

: 248

: 249

MongoDB Lab Exercises

1) Exercise : Collection Database.

① db.createCollection("customers")

Output: {ok: 1}

② db.customers.insertMany([{"Cust-id": 1, "Acc-Bal": 1500, "Acc-Type": "Z"},

{Cust-id: 2, Acc-Bal: 900, Acc-Type: "A"},

{Cust-id: 3, Acc-Bal: 1300, Acc-Type: "Z"},

{Cust-id: 4, Acc-Bal: 700, Acc-Type: "B"},

{Cust-id: 5, Acc-Bal: 1600, Acc-Type: "Z"}])

Output

Acknowledged: true,

inserted: 5

'0': ObjectId("62cfd111"),

'1': "-4-",

'2': "-1-",

'3': "-1-",

'4': "-1-",

'5': "-1-",

③ }

③ db.customers.find({Acc-Bal: {\$gt: 1200}},

Output { Acc-Type: "Z" })

-Id: ObjectId("62cfd111"),

Cust-id: 1,

Acc-Bal: 1500,

Acc-Type: "Z"

}

① db. customers . aggregate []

{

\$groups {

- id : "\$cust_id",

Min_Acc_Bal : { \$min : "\$acc_bal" } ,

Max_Acc_Bal : { \$max : "\$acc_bal" } }

} }

Output:

[{

atmos

{ - id : "5", Min_Acc_Bal : 1600, Max_Acc_Bal : 1600 },

{ - id : "5", Min_Acc_Bal : 1600, Max_Acc_Bal : 1600 }

{ - id : "4", Min_Acc_Bal : 700, Max_Acc_Bal : 700 }

{ - id : "3", Min_Acc_Bal : 600, Max_Acc_Bal : 600 }

② Exercise

Product Database []

db.createCollection("Products")

db. Products.insertMany []

{ - id: ObjectId ("507f1f1"), name: "Laptop",
Category: "Electronics", Price: 800, Quantity: 10 },

{ - id: ObjectId ("507f1f2"), name: "SmartPhone",
Category: "Electronics", Price: 500, Quantity: 5 },

{ - id: ObjectId ("507f1f3"), name: "Shoes",
Category: "Fashion", Price: 50, Quantity: 20 },

{ - id: ObjectId ("507f1f4"), name: "Tablet",
Category: "Electronics", Price: 300, Quantity: 8 },

{ - id: ObjectId ("507f1f5"), name: "T-shirt",
Category: "Fashion", Price: 20, Quantity: 15 }

} }



Output

acknowledged: true,

inserted_ids: [

'0': ObjectId("4507f111"),

'1': ObjectId("507f12"),

'2': ObjectId("507f13")

ii) db.createCollection("Orders")

lok: 1 }

db.Orders.insertMany([

{

_id: ObjectId("607f121"),

user_id: "123abc",

items: [{product_id: ObjectId("507f112"),

Quantity: 1, Price: 500 }],

total_Price: 500

{

_id: ObjectId("607f1022"),

user_id: "123abc",

items: [{product_id: ObjectId("507f13")},

Quantity: 3, Price: 450 }],

total_Price: 1350

{

_id: ObjectId("607f123"),

user_id: "11156def",

items: [{product_id: ObjectId("507f111"),

Quantity: 2, Price: 1600 }],

total_Price: 3200

}
])

'o' : ObjectId ('607021')

"": object_id ("60742")

'2': Object 2d (NGC 7134)

2. Object 20 40 10

۳

db.createCollection ("wants")

$$\{0\} = \{1\}$$

def (inertOneLL) kmt, stewart, sh

un-*id* : "789ghi"

items := [

2 {"ProductId": "507f11", "Quantity": 1},
2 {"ProductId": "507f2", "Quantity": 2}] })

6P

acknowledge : true ,

inverted ID: Observed ("Gridbooks")

elb. Products - find 1)

L-1d - Object ID (45014143), b14 - detail - sb

Name : 'Bob Bob' (Bobde) is playing.

Category : 'Electronics', 'radio': what

Anice L-800

Quantity : 10 { Pkt : 100g

i.d.: object sd ($50 \text{ ft } 5''$),

Name : 'T-shirt' (T恤)

"Fact'ion", *Confederacy*

price : 20,

Quantity = 15 L

Output: Objectid

db. Products . find ({ Category: "Electronics" })

[
 - id : Objectid ("507f1")
 - name : "Laptop",

Category : "Electronics",

Price : 800,

Quantity : 10 ,

db. Products . find ({ Quantity : 10 })

[
 - id : Objectid ("507f1")
 - name : "Laptop",

Category : "Electronics",

Price : 800 ,

Quantity : 10 ,

db. Products . find (, Sort ({ Price: 1 }))

{
 - id : Objectid ("507f1")
 - name : "Fashion",

Price : 20

Quantity : 15 ,

db. Products . find ({ Price : { \$lt: 100 } })

{
 - id : Objectid ("507f1")
 - name : "Shoes",

Category : "Fashion",

Price : 20 ,

Quantity : 10 ,



db. - Conts - aggregate (f)

{ \$match: { user_id: "7899hi7y3" }

{ \$ unwind: "\$items" }

{ \$ lookup: {

from: "products"

as: ProdDetails

}, {

})

OP:

{ id: ObjId ("6071") }

user_id: "123abc"

items: [

{ prod_id: ObjId ("5071") }

}, {

price: 500

}, {

db. Orders - aggregate (L)

{ \$match: { user_id: "23abc" } },

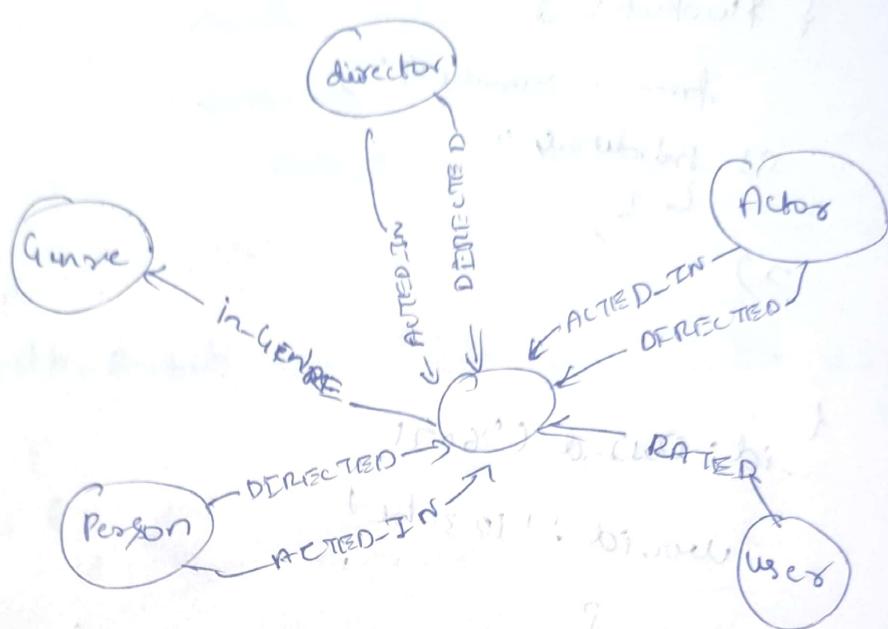
{ \$group: { _id: "\$user_id", total_qty: }

{ \$group: { _id: "\$user_id", total_price: } }]

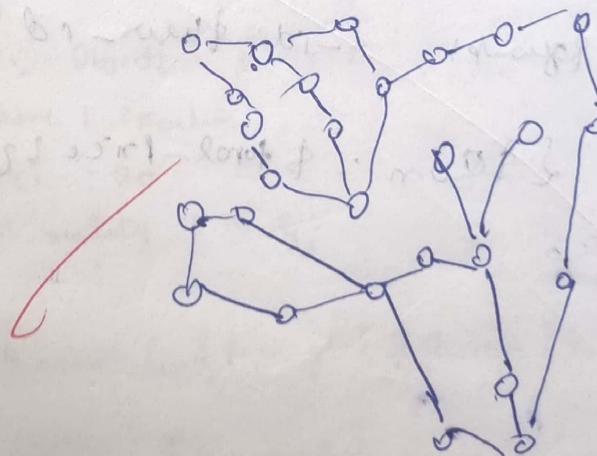
11/3/25

Lab 3 (Neo4j)

1) Coll. db-Schema, Visualization



2) MATCH (n) return n limit 100



1/4/2025

Lab 4

Exploring operations on NOSQL - Cassandra.

Create keySpace:

→ CREATE KEYSPACE Student WITH REPLICATION = { 'class' : 'SimpleStrategy', 'replication_factor' : 1 } ;
keySpace Created Successfully

→ DESCRIBE KEYSPACES;

Student System-auth System-Schema System-views
System System-distributed system-traces system-normal-schema

→ SELECT * FROM System-Schema.keyspaces;

→ SELECT * FROM System-Schema.keyspaces;

KeySpace-name durable_writes replication
True { 'class' : 'SimpleStrategy', 'replication_factor' : 1 }
System-auth off-the-shelf apache-cassandra.
impl strategy!, 'replication_factor' : 1 } color.S

(6 rows)

→ USE Student;

CREATE TABLE Student-Info (roll_no int PRIMARY
KEY, std_name text, DateOfJoining timestamp,
lost_exam_percent double);

→ DESCRIBE TABLES;

Student-Info user login



Scanned with OKEN Scanner

→ DESCRIBE TABLE Student - Info;

/* Description of Table */

BEGIN BATCH

INSERT INTO Student - Info (Roll-No, Stud Name,
Date of Joining, last_exam_Percent
VALUES (1, 'Asha', '2012-03-12', 79.9)

;

INSERT INTO Student - Info (Roll-No, Stud Name,
Date of Joining, last_exam_Percent)
VALUES (1, 'Rohan', '2012-03-12', 56.9)

→ SELECT * FROM Student - Info;

→ SELECT * FROM Student - Info WHERE Roll-No
IN (1, 2, 3);

→ SELECT * FROM Student - Info WHERE
Stud Name = 'Asha';

→ CREATE INDEX X ON STUDENTS - Info (Stud Name);

→ select from Student - Info where Studname = 'Asha';

→ Select Roll-No AS "USN" from Student - Info;

→ Update Student - Info Set Studname = 'David Shein'
where Roll No = 2;

→ DECIES IS FROM Student_info WHERE Roll_no = 2;

ALTER TABLE Student_info ADD hobbies SET <text>;
UPDATE Student_info
SET hobbies = hobbies + 'Chess', 'Table Tennis'
WHERE Roll_no = 1;

List collection:

ALTER TABLE STUDENTS_INFO Add long wt <text>;
ALTER TABLE Stud_info
Update Stud_info
SET long = [Hindi, English] - but 2
where roll no = 1;

Counter Table

create table library_book {
book_name varchar
stud_name varchar
};

Time To Live

create Table user_info (
User_id
Password
);

RAJ25
NAT



Lab - 5
Cassandra Exercise.

1) Create a keybase by name library

library WITH REPLICATION { 'class': 'SimpleStrategy',
'replication_factor': 1 };

2) USE Library ;

CREATE TABLE Library-Info

Stud-ID int PRIMARY KEY,

Counter-value Counter,

Stud-Name text,

Book-Name text,

Book-ID int,

Date-of-issue timestamp);

3) BEGIN BATCH

INSERT INTO Library-Info (Stud-ID, Counter-value,
Stud-Name, Book-Name, Book-ID, Date-of-issue)
VALUES (111, 1, 'Alice', 'Introduction to Cassandra',
101, '2025-04-10 10:00:00+0000');

INSERT INTO Library-Info (Stud-ID, Counter-value,
Stud-Name, Book-Name, Book-ID, Date-of-issue)
VALUES (112, 1, 'Bob', 'BDA', 205, '2025-04-11
14:30:00+0000');

INSERT INTO Library-Info (Stud-ID, Counter-value,
Stud-Name, Book-Name, Book-ID, Date-of-issue)
VALUES (113, 4, 1 David, 'NOSQL', 206,
'2025-06-11', '14:30+00+0000');

4) DESCRIBBLE TABLE Library-Info

Output

```
CREATE TABLE Library-Info (
    Stud-ID int PRIMARY KEY,
    Book-ID int,
    Book-Name text,
    Counter-value Counter,
    Date-of-issue timestamp,
    Stud-Name text);
```

5) Increase counter val for Student with Stud-ID = 112.

UPDATE Library-Info SET Counter-value = Counter-value + 1 WHERE Stud-ID = 112;

SELECT Stud-ID, Counter-value, Stud-Name,
 Book-Name, Book-Name FROM Library-Info
 WHERE Stud-ID = 112;

Output

Stud-ID	Counter-value	book-name	Stud-name
112	2	BDA	Bob

CREATE TABLE Library-Book-Tracking (

```
Stud-ID int,
Book-name text,
Issue-number int,
Date-of-issue timestamp,
PRIMARY KEY (Stud-ID, Book-Name),
```

SELECT * FROM Library-Info WHERE Stud-ID = 112 AND Book-Name = 'BDA'

Stud-ID	Book-ID	Book-Name	Counter-value	Date-of-issue
112	205	BDA	2	2025-04-15
112	205	BDA	2	2025-04-15

(2 rows)

b) Export

COPY Library.Library-Info TO '/path/to/your/import-data.csv'
with header = TRUE;

"SELECT Stud-ID, Counter-value, Stud-Name,
Book-Name, Book-ID FROM Library.Library-Info"

Y Stud-ID, Counter-value, Stud-Name, Book-Name,
Book-ID, Date-of-issue

201, 1, "Eve", "DS", 305, "2025-04-15"
202, 2, "Frank", "Algorithms", 308, "2025-03-15"

COPY Library.Library-Info (Stud-ID, Counter-val,
Stud-Name, Book-Name, Book-ID,
From '/path/to/your/import-data.csv',
With header = TRUE;

Output Processed 2 rows in 0.*** seconds

lab-6
Hadoop

- \$ start-all.sh
- \$ hdfs dfs -mkdir /bda-hadoop
- hdfs dfs -ls /
└── band | items
 drwxr-xr-x - hadoop subgraph
- hdfs dfs -ls /bda-hadoop /
└── band | items
 -rw-r--r-- 1 hadoop subgraph
- hdfs - dfs - cat /bda-hadoop/file.txt.
- hdfs - get | abc | wk.txt /home/hadoop
/Downloads /annual.income.csv.
ls /home/hadoop /Downloads /
- 6) → hdfs - dfs - copyToLocal /abc /annual.csv /
home /hadoop /Desktop /
ls /home /hadoop /Desktop /
- hdfs dfs -cat | abc | annual.csv /
- hadoop subgraph
- OP: drwxr-xr-x - x

→ hadoop fs -cp / core / 1uu.

hadoop fs -ls / LLL

DP: bound + items

draw X r - rr - x - hadoop Subgroups

0 XYXY - MM-DD HH:MM/Lucky

15/A/125

(1) Disjointed subgroups -> 216

(2) overlapping subgroups -> 216

(3) hadoop -> 216

(4) 216 -> 216

(5) 216 -> 216

(6) 216 -> 216

(7) 216 -> 216

(8) 216 -> 216

(9) 216 -> 216

(10) 216 -> 216

(11) 216 -> 216

January 25

Scala

Scala

val x = 5 + 2

answ: Int = 42

Scala > 0.5 * 2.0

null: Double = 21.0

Scala > "Hello," + null

null: String = Hello, null

null: Double = 42.0

Scala > val answer = 5 + 2

answer: Int = 42

Scala > val answer = 0.5 * answer

null: Double = 21.0

Scala > var counter = 0

counter: Int = 0

Scala > counter = 1

counter: Int = 1

Scala > val greeting: String = null

greeting: String = null

Scala > val greeting: Any = "Hello"

greeting: Any = Hello

Scala > val xmax, ymax = 100

xmax: Int = 100

ymax: Int = 100

Scala > var greeting, message: String = null

greeting: String = null

message: String = null

Scala > `toString()`
res0: String = @ 2a2c2ue

Scala > `toString(10)`

res1: String = 10

Scala > `"Hello".intake("World")`

res2: String = lo

Scala > `1 to 10`

res3: Scala.collection.immutable.Range = Range(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)

Inclusive = Range(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)

Scala > `Scala.math.Sort(2)`

res4: Double = 1.41421

Scala > `import Scala.math`

import Scala.math

~~def~~
~~def~~



TOP N

Create a MapReduce Program to find average
temp for each year from NCDC data set package
Sample job:

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.Writable
```

```
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
import org.apache.hadoop.mapreduce.lib.output.MultipleOutputs;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

```
import org.apache.hadoop.util.GenericOptionsParser;
```

```
public class TOP N & main (String[] args) throws
```

```
public static void
```

```
Exception {
```

```
Configuration conf = new Configuration();
```

```
String [] otherArgs = GenericOptionsParser.  
getRemainingArgs();
```

```
if (otherArgs.length != 2) {
```

```
System.out.println ("Usage: TOPN<in>
```

```
<out> "1");
```

```
System.exit(2);
```

3

```
Job job = Job.getinstance (conf);
job.setJobName ("TOP N");
```

- job. Set Job By class (TopN + class)
- job. Set Mapper class (Top N Mapper class).
- job. Set Reducer class (Top N Reducer class).
- job. Set Output key class (Text class);
- job. Set Output value class (IntWritable, key);
- file input format.add input path(job, new Path(otherArgs[1]));
- System.exit(0); wait for workstation (true);
- c Static class TopNMapper extends Mapper<IntWritable, IntWritable, > {

Top N Reduce? May.

```

    package Sambler TopN;
    import java.io.IOException;
    import java.util.HashMap;
    import java.util.List;
    import org.apache.hadoop.io.IntWritable;
    import org.apache.hadoop.io.Text;
    import org.apache.hadoop.mapred.Reduce;
    import org.apache.hadoop.mapred.lib.MiniWritables;
    import util.List;

    public class TopN Reducer extends Reducer<Text, IntWritable, Text, IntWritable> {
        private Map<Text, IntWritable> countMap = new HashMap<Text, IntWritable>();
        public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException, InterruptedException {
            int sum = 0;
            for (IntWritable val : values) {
                sum += val.get();
            }
            countMap.put(key, new IntWritable(sum));
        }
        protected void writeKeyValue(Text key, IntWritable value) {
            System.out.println(key + " : " + value);
        }
    }

```

b) mean max temperature for every month

Package main;

Public class MeanMaxDriver

Public static void main (String args) throws

Exception {

if (args.length != 2) {

System.exit(1);

Job job = new Job(args[0]);

Job.setJobByUser (meanMaxDriver, day);

fileInputStream add InputPath (job, new Path (args[0]))

System.exit (job.isSuccessful () ? 0 : 1);

(*lose character array)

} student will write box reading

Mean Max Number does.

Package meanmax;

Public class MeanMaxNumber extends

number < long writable, Text, Text, intwritable > {

private Map < Text, intwritable > countMap = new

HashMap <>;

public void reduce (Text key, Iterable < intwritable >)

values, Reduces < Text, intwritable, Text >

intwritable > context Context);

throws IOException, InterruptedException Exception {

int sum = 0;

for (int writable & Val : values)

sum += val.get();

this.countMap.put (new Text (key),

new intwritable (sum));

}

Protected void cleanMob (Reduces Text, intwritable
Text, int writable > context context)

exception 2

mob < Text, intwritable > Sarcoid Mob =
new Util . Sort By Value (this . countmob);

int Counter = 0;

for (Text key : Sarcoid mob . key set ()){

if ((Counter + tno) >= 20)

break;

Context . write (key, Sarcoid mob . get (key))

}

}

class * producer mob also {

void write (Text, intwritable, int writable)

set of p;

002

- (Q) Ques 8
Topic: ArrayList
- ① Write a Scala Program to Print number from 1 to 100 using for loop
- > Scala abt install scala [3.2.2]. Will run.
- > nano PrintNumbers.scala
- ```
object PrintNumbers {
 def main(args: Array[String]): Unit = {
 for (i ← 1 to 100) {
 print(i)
 }
 }
}
```
- > Scala PrintNumbers. Scala

> Scala PrintNumbers

Output:

|    |                      |
|----|----------------------|
| 1  | Numbers Printed from |
| 2  | 1 to 100             |
| 3  | 4                    |
| 5  | 6                    |
| 7  | 8                    |
| 9  | 10                   |
| 11 | 12                   |
| 13 | 14                   |
| 15 | 16                   |
| 17 | 18                   |
| 19 | 20                   |
| 21 | 22                   |
| 23 | 24                   |
| 25 | 26                   |
| 27 | 28                   |
| 29 | 30                   |
| 31 | 32                   |
| 33 | 34                   |
| 35 | 36                   |
| 37 | 38                   |
| 39 | 40                   |
| 41 | 42                   |
| 43 | 44                   |
| 45 | 46                   |
| 47 | 48                   |
| 49 | 50                   |
| 51 | 52                   |
| 53 | 54                   |
| 55 | 56                   |
| 57 | 58                   |
| 59 | 60                   |
| 61 | 62                   |
| 63 | 64                   |
| 65 | 66                   |
| 67 | 68                   |
| 69 | 70                   |
| 71 | 72                   |
| 73 | 74                   |
| 75 | 76                   |
| 77 | 78                   |
| 79 | 80                   |
| 81 | 82                   |
| 83 | 84                   |
| 85 | 86                   |
| 87 | 88                   |
| 89 | 90                   |
| 91 | 92                   |
| 93 | 94                   |
| 95 | 96                   |
| 97 | 98                   |
| 99 | 100                  |

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

> choose "hello word" when spark in hello world [In spark in post] & hello spark spark  
spark" > input.txt

spark - shell

> val rdd = Seq[TextFile]("input.txt")  
> val words = rdd.flatMap(\_.lines).map(\_.split(" ")).  
val words = words.map(word => word.toLowerCase)  
val pairs = words.map(word => word.toLowerCase->1)  
val count = pairs.reduceByKey(\_+\_).map{case (word, count) => (word, count)}  
val filtered = count.filter{case (word, count) => count > 4}  
filtered.collect().foreach{case (word, count) => println(s"\$word : \$count")}

Output:

Spark: 6 different & nonempty words

( "had", "said", "and", "the", "in", "it" ) words  
( "is", "it", "a", "of", "to" ) words  
( "the", "a", "in", "it" ) words

~~Program in~~

3) Write a Simple Streaming ~~Program in~~

3) Map-reduce Program to Sort the Content in a alphabetical order listing only top 10 max occurrences of words (Top10k Prog).

WordCount MapReduce.java

```
import java.io.IOException;
import java.util.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.Reducer;
```

Public class WordCountReducer Extends Reducer<Text, IntWritable, Text, IntWritable>

```
private Map<Text, IntWritable> CountMap = new
HashMap<Text, IntWritable>();
```

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

```
(throws IOException, InterruptedException)
```

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

```
CountMap.put(new Text(key), new IntWritable
(sum));
```

Method void map(Text key, IntWritable value, Context context) throws  
IOException, InterruptedException

CountMap. entrySet(). Stream()

```
+ Sorted<Entry<Text, IntWritable>> l;
int cmp = e2.getValue(). compareTo(e1.getValue());
if (cmp == 0) return e1.getKey(). toString();
Comparable e1.getKey(). toString());
```

return cmp;

}

```
+ limit(10)
+ for each ee -> l.
```



Connect . write (get key (key), get value (val));  
} catch (IOException | InterruptedException ex)  
{  
 d  
 e  
 f  
 g  
 h  
 i  
 j  
 k  
 l  
 m  
 n  
 o  
 p  
 q  
 r  
 s  
 t  
 u  
 v  
 w  
 x  
 y  
 z  
 {  
 throw new RuntimeException (ex);  
 }  
}

Word Count Driver.java

import org.apache.hadoop.conf.Configuration;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.\*;

public class WordCountDriver {

public static void main (String [] args)

throws Exception {

Configuration conf = new Configuration();

Job job = Job . get instance (conf, "Job 10  
Word Count");

job . setJarByClass (WordCountDriver . class);

job . setMapperClass (WordCountMapper . class);

job . setReduceClass (WordCountReducer . class);

job . setOutputKeyClass (Text . class);

FileInputFormat.addInputPath (job, new Path

(args [0]));

System . exit (job . waitForCompletion (true))

} else



|           |               |                                             |
|-----------|---------------|---------------------------------------------|
| OP:       | <u>hadoop</u> | Want to filter, don't want to have any log. |
|           | and           | Want to filter, don't want to have any log. |
|           | an            |                                             |
|           | are           | (most) Val would go to down.                |
|           | big           | big + not red. Brown (or brown) will fit.   |
| data      | 2             | (brown) dictionary                          |
| framework | #             |                                             |
| hadoop    | 3             | (brown) dictionary                          |
| is        | 1             | (not) dictionary                            |
| wrong     |               | (not) dictionary                            |
| amount    | \$            |                                             |

Ex: Textmate

#### a) (Open Ended)

Write a simple Streaming Program in Spark to receive (text) data streams on a particular space (remove stop words, remove punctuation, lemmatization etc.) and print cleaned text on screen.

Terminal P:

Ex: Textmate

```

> spark-shell
> import org.apache.spark.streaming
> import org.apache.spark.streaming.StreamingContext
> val ssc = new StreamingContext(ssc, second(5))
val stopwords = Set("a", "an", "the", "is", "are",
 "and", "or", "to", "in")
val lines = ssc.socketTextStream("localhost", 9999)

```



val cleaned = lines . flatmats (unlocal host area)  
val cleaned = lines . flatmats (-split ("1111"))  
, mat (-to lower left + trim)  
. filter (word => word, nonEmpty, (all) NotEmpty  
contains (word ?))  
cleaned . print ()  
SS c. start ()  
SS c. await Termination  
Terminal B

→ next area (host area) (a few seconds later)  
→ Terminal A: loading word recognition  
Terminal B: (run me by prog) screen  
> start submit word count by keyboard  
> start submit word count by keyboard

### Terminal B

The quick brown fox jumps over the lazy dog  
The quick brown fox jumps over the lazy dog  
The quick brown fox jumps over the lazy dog  
The quick brown fox jumps over the lazy dog  
The quick brown fox jumps over the lazy dog

### Terminal A

The "at", "no", "o" are standard for  
Ouch  
brown  
fox  
jumps  
over  
the  
lazy  
dog