LAB - 1

X. Difference b/w data and big data

| Traditional data | Big data |
|---|---|
| a. Small amount of data that can be collected and analysed using traditional methods easily. | a. Big amount of data that cannot processing and analysed easily using traditional methods. |
| b. Usually structured data. Stored in databases spread sheets. | b. Includes semi-structured and unstructured data too |
| c. Manual data collection | c. Automatic data collection |
| d. Slow, gradual analysis | d. Fast, instant analysis |
| e. Typically processed in batches | e. Developed and proceed in real time. |

Applications of Big data

ⓐ Tracking customers spending habit, shopping
ⓑ Recommendation
ⓒ. Smart traffic system
ⓓ Auto driving car
ⓔ Secure car traffic system
ⓕ Virtual personal assistant tool
ⓖ IOT Sensors
ⓗ Education sector.

Database available for big data

* AWS Dynamo DB
× Azur cosmos DB
× Amazon keyspaces
× Amazon Document DB
× Amazon Redshift

× Write SQL queries to do th following

ⓐ Create + insert 2 values in the table
Create table department (
   empID INTEGER primary key,
   name VARCHAR (20) Not NULL,
   loc_name VARCHAR (15) Not NULL,
)

Create table employee (
   empID INTEGER primary key,
   name varchar (20) not null,
   jobname varchar (20),
   manager id integer.

hire_data data.
salary varchar
commission varchar.
dep - id integer
)

Insert into emplayee values ( 1. Ram, SDI , 1
25 - 01 - 2024 , 140000,
50000 , 4 ):
Insert into employee value ( 2, Raj , SDE , 2
25 - 01 - 2014 , 20000
50000 , 3 ):

① Find Salaries of all employee
Select salary name
from employee:

② Find those employees with hire data lik-
February 22, 1991
Select name
From employee
where job - name = 'Analyst hire - data like'
'February 22, 1991':

④ Calculate avg salary of employees who
work as analyst
Select avg (salary)
from employee
where job - name = 'Analyst'

③ Find those employee who are either clerk or Manager

```
Select emp_name
From employees
where job_name in ('clerk', 'Manager')
```

④ Search for all the employees with an annual salary b/w 24000 and 50000

```
Select salary, emp_name from employees
where salary BETWEEN 24000 AND 50000;
```

## LAB-02

I. Perform the following DB operations using MangoDB.

1. Create database "student" with attributes Rollno, age, contact no, email-id
   db. Create Collection ("student").

2. Insert appropriate values
   db. Student . insertMany (
   { "Roll No" : 10 ; "Age" : 20 , "Contact No :
   " 7676526898 ", "Email ID" : " Jai @gmail.com

   { "Roll No" : 11 , "Age" : 21 , "Contact No:
   " 610564758 ", "Email ID " : Ray @gmail.com

   acknowledged : true
   inserted Ids : { '0' : objectId ('      ')
                    '1' : object ('      ) }

3. Write query to update email of student with Roll no
   db. Student. update ( { ROll No : 103 ,
   { $set : { email ID :" Jam @gmail.com "}

   acknowledged : True
   inserted ID : null
   matched Count : 1
   modified count : 1
   upsorted count : 0 }

4. Replace the student name from "ABC" to "JEM" of roll = 11

// insert apt data

db. student. insert ({ Roll No: 11, Age: 21, Nv
ABC, contact No: 90076814, email JD:
ABC@gmail.com }):

// replace
db. student. replace ({ Roll No: 11 Age: 2
Name: JEM, contact No: 90076814, @mai
Jem@gmail.com }):

Ⅱ  Perform the following DB operations using
Mango DB

1  Create collection by name customers
Cust_id, Acc_bal, Acc-Typ
db. create collection ("customer")

2. Insert 5 values
db. customer. insert Many ([
{ cust_id: 1, balance: 200, Type: s}
{ cust_id: 2, balance: 1000, Type: z }
{ cust_id: 8, balance: 300, type: z }
{ cust_id: 3, balance: 500, type: s
// view
db. customer find ()

3. Display those records whose total balance
1200 of type z for each id.

db. customer. aggregate (

```
{ $match : { type : 23}.
{ $group : { _id : " $cust_id ".
  totaccbal : { $ sum : " $balance " }}}
{ $match : { totAccbal : { $gt : 1200}}}) :
  _id : 2 , totAccbal : 1300
```

6. Determine min and max account balance for each . 1

```
db . customers . aggregate (
  { $ group: { _id : " $ cust_id " ,
    minAccbal : { $min : " $balance" }
    maxAcc bal : { $max : " $balance" }}}
```

LAB - 03

Perform the following DB operations using Cassandra.

1. Create a keyspace by name employee

```
CREATE KEYSPACE "employee"
WITH replication = { 'class' : 'SimpleStradegy' ,
  ' replication _factor' : 2 } ;
```

Use 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

```
CREATE TABLE emp_info (
    emp_id int    PRIMARY KEY
    emp - Name,    emp - designation
    emp _doj , emp_ sal , emp _dept _name );
```

3. Insert the values into table in batch

```
BEGIN BATCH
INSERT INTO emp - info (empid, emp_name
emp_ designation, emp-doj , emp- sal , emp- dept_name
value ( +234 , 'Hiya' , 'Manager', 20-01-24 , 1000
 'Biotech ' );
```

4. Update Employee name and department of emp 1234

~~code~~ UPDATE emp -info
&ET emp _name = 'Hiru', emp -deptnam
'Chemixtry'
where emp id = 1234 ;

## LAB - 04

Execution of HDFS Commands for interaction w Hadoop Environment

1. mkdir        2. ls        3. get        4. Put
5. copy fromlocal    6. Copy tolocal    7

1. 1 To create directories
→ hadoop fs - mkdir /newDataFlair

2. To enlist the files and directories present in HDFS
→ hadoop fs - ls/

3. To copy localfile of the localfile system to the Hadoop filesystem
→ hadoop fs - put ~ /localfile / first.text

4 To copy the first file present in the local file system to the newdataFlair directory of Hadoop
→ hadoop fs - copyFromLocal ~ /first / newDataFlair / copy text

→ hadoop fs - cat /newDataFlair / copytext

6. To copy the file of the hadoop file system to the local file system

→ hadoop fs - get /first ~/ copy from hadoop

13

7. Hadoop HDFS get command Example hadoop fs - getfacl /abl

   This apache hadoop command shows the access control lists of files and directories

8. Cat , hdfs dfs - cat /abc/wc. txt

9. mv HDFS mv command Example hadoop fs - mv /abc/FFF .h

10. CP , hadoop fs - cp /cse /LLL hadoop fs - ls /LLL

Wordcount MapReduce program

Terminal
→ start - all. sh
→ jps
→ sudo snap install eclipse--classic
→ su hadoop

1. Open eclipse create new java project

2. Click right on the wordcount project → build configure build path → libraries → add jar files from (share / common and share / mapreduce)

3. Right click on the wordcount add closser called WCDriver, WCMapper, WCReducer

4. Go to Home (etc / Hadoop / . xml and modify it

5. Fight Click on the wordcount → export → jar file → Create jar file wordcount.jar

6. Open terminal
→ hadoop fs - mkdir / rqo
→ hadoop fs - copy From Local / home / hadoop / Desktop / sample .txt / rqn / test. txt
→ hadoop jar / home / hadoop / Desktop / word-count.jar WCDriver / rqn / test.-txt / output
→ hadoop . output
→ hadoop - cat (output / part -0000

x. Create Map Reduce program to
@ Find average temperature for each year
from NCDC data set

Mapper :-

```
import sys

for line in syn.stdin:
    line = line.strip()
    parts = line.split ("\t")
    date = parts [0]
    temperature = parts [1]
    year = date [:4]


try:
    temperature = float (temperature)
    print (f" {year} \t {temperature}")
except value error:
    continue
```

Reducer :-

```
import sys

current_year = None
Sum_temp = 0
count_temp = 0

for line in sys. stdin:
    line = line . strip ()
```

```
year, temperature = line.split ("t")

try :
    temperature = float (temperature)
except ValueError:
    continue

if current - year = year:
    sum - temp + = temperature
    count temp + = 1
else :
        if current - year:
        average - temp = sum temp / count temp

    current - year = year
    sum_temp = temperature
    count temp = 1

if current - year:
    average - temp = sum - temp / count temp
    print (f" { current - year } \t { average -
    temp } ")
```

x. create MapReduce program to
Q. Find the mean max temperature for
every month, give

Mapper :

import sys

```
for line in sys. stdin :
    line = line. strip()
    parts = line. strip()
    date = parts[0]
    temperature = parts[1]
    year_month = date[7]
```

```
try :
    temperature = float (temperature)
    print (f" {year_month} \t {temperature}")
except valueError :
    continue
```

Reducer :

import sys

```
current_month = None
sum_temp = 0
count_temp = 0
```

```
for line in sys.stdin:
    line = line. strip()
```

```
month, temperature = line split ("\t")


try:
    temperature = float (temperature)
except valueError :
    continue


if current_month == month :
    sum_temp + = temperature
    count_temp + = 1
else :
    if current_month :
        average_temp = sum_temp / count_temp
        print { current _month average temp}


current_month = month
sum_temp = temperature
count_temp = 1


if current_month :
    average_temp = sum_temp / count_temp
    print ( current_month average_temp)
```

X. Create a Map Reduce program to sort the content in an alphabetic order listing only top 10 maximum occurrences of words.

```
import sys

current_word = None
current_count = 0

for line in sys.stdin:
    line = line.strip()
    word, count = line.split('\t', 1)

try:
    count = int(count)
except valueError:
    continue

if current_word == word:
    current_count += count
else:
    if current_word:
        current_word = word
        current_word = count

if current_word == word:
    print( {current_word} \t {current_word})
```

```
for line in sys.stdin:
    print (line.strip())


word_counts = [ ]


for line in sys.stdin:
    line = line.strip()
    word, count = line.split('\t', 1)


try:
    count = int (count)
    word_counts.append ((word, count))
except ValueError:
    continue


top_10_words = nlargest (10, word_counts,
    key = itemgetter (1))


top_10_words.sort (key = itemgetter (0))


for word, count in top_10_words:
    print (f'{word}\t{count}')
```