**LABORATORY RECORD**

**On**

**BIG DATA ANALYTICS**

**B.E (IT) – VII Sem**

**By**

**N.Durga Sai Lakshmi(160117737006)**

**Under the guidance of**

**Sri.K.Gangadhar Rao**

**Assistant Professor,**

**Dept. of IT,CBIT.**

**DEPARTMENT OF INFORMATION TECHNOLOGY**



**CHAITANYA BHARATHI INSTITUTE OF TECHNOLOGY (A)**

**(Affiliated to Osmania University; Accredited by NBA(AICTE) and NAAC(UGC), ISO Certified 9001:2015)**

**KOKAPET(V),GANDIPET(M), RR District HYDERABAD -75**

**Website: www.cbit.ac.in**

**2020-2021**

**Problem Statement-1:** Write a Map-reduce application to find number of occurrences of each word from the given dataset.

**Description:**

In the MapReduce word count example, we find out the frequency of each word. Here, the role of Mapper is to map the keys to the existing values and the role of Reducer is to aggregate the keys of common values. So, everything is represented in the form of a Key-value pair.

**Procedure:**

$cd hadoop-3.2.1/

$cd sbin/

$cd start-all.sh

$jps

$cd..

$cd wc.py/

$hadoop fs -mkdir -p /wordcount/p1

$hadoop fs -copyFromLocal word.txt /wordcount/p1

$hadoop jar /home/hduser/hadoop-3.2.1/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar -file wcmap.py -mapper wcmap.py -file wcred.py -reducer wcred.py -input /wordcount/p1 -output /wordcount/p1/output

$hadoop fs -cat /wordcount/p1/output/part-00000

**Code:**

**wcmap.py**

```
import sys

# input comes from STDIN (standard input)

for line in sys.stdin:

        # remove leading and trailing whitespace
```

```
        line = line.strip()

        # split the line into words

        words = line.split()

        # increase counters

        for word in words:

        # write the results to STDOUT (standard output);

        # what we output here will be the input for the

        # Reduce step, i.e. the input for reducer.py

        #

        # tab-delimited; the trivial word count is 1

                print '%s\t%s' % (word, 1)
```

**wcred.py**

```
from operator import itemgetter
import sys

current_word = None
current_count = 0
word = None

# input comes from STDIN
for line in sys.stdin:
        # remove leading and trailing whitespace
        line = line.strip()

        # parse the input we got from mapper.py
        word, count = line.split('\t', 1)

        # convert count (currently a string) to int
        try:
```

```
count = int(count)
except ValueError:
# count was not a number, so silently
# ignore/discard this line
continue

# this IF-switch only works because Hadoop sorts map output
# by key (here: word) before it is passed to the reducer
if current_word == word:
current_count += count
else:
if current_word:
# write result to STDOUT
print '%s\t%s' % (current_word, current_count)
current_count = count
current_word = word
if current_word == word:
print '%s\t%s' % (current_word, current_count)
```

**INPUT :**

**OUTPUT:**



**Problem Statement-2:** Write a 3map-reduce application to predict maximum stock price for given dataset.

**Description:**

We are trying to find out the maximum closing price of each stock symbol. This means that we have to group the records by symbol so that we can calculate the maximum closing price by symbol. So we will output Stock Symbol as the key and close price as the value for each record. We now know what is going to be the Map's input and what is going to be the maps output.

**Procedure:**

$cd hadoop-3.2.1/

$cd sbin/

$cd start-all.sh

$jps

$cd..

$cd stock.py/

$hadoop fs -mkdir /stocksip

$hadoop fs -copyFromLocal stocks.txt  /stocksip

$hadoop jar /home/hduser/hadoop-3.2.1/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar -file
stock-mapper.py  -mapper stock-mapper.py  -file stock-reducer.py -reducer stock-reducer.py
-input /stocksip/stocks.txt -output /stocksout/output

$hadoop fs –cat stocksout/output/part-00000

**Code:**

**Stock-mapper.py**

```
import sys

for line in sys.stdin:
        line = line.strip()
        data = line.split(",")
        stock, price = data[1], data[6]
        print("%s\t%s"%(stock,price))
```

**stock-reducer.py**

```
import sys

max_price = 9999
max_stock = None

for line in sys.stdin:
        line = line.strip()
        stock,price = line.split("\t",1)
```

```
if max_stock and max_stock!=stock:
        if max_price > price:
                max_price = price
                max_stock = stock


    else:
            max_stock, max_price = stock,max(max_price,price)


if max_stock:
    if max_price > price:
                    max_price = price
                    max_stock = stock


print("%s\t%s"%(max_stock,max_price))
```

OUTPUT:

**Problem Statement-3:** Write a map-reduce application to find maximum temperature for a given year from the NCDC weather dataset.

**Description:**

MapReduce is based on set of key value pairs. So first we have to decide on the types for the key/value pairs for the input.

**Map Phase:** The input for Map phase is set of weather data files as shown in snap shot. The types of input key value pairs are *LongWritable* and Text and the types of output key value pairs are *Text* and *IntWritable*. Each Map task extracts the temperature data from the given year file. The output of the map phase is set of key value pairs. Set of keys are the years. Values are the temperature of each year.

**Reduce Phase:** Reduce phase takes all the values associated with a particular key. That is all the temperature values belong to a particular year is fed to a same reducer. Then each reducer finds the highest recorded temperature for each year. The types of output key value pairs in Map phase is same for the types of input key value pairs in reduce phase (*Text* and *IntWritable*). The types of output key value pairs in reduce phase is too *Text* and *IntWritable*.

**Procedure:**

$cd hadoop-3.2.1/

$cd sbin/

$cd start-all.sh

$jps

$cd..

$cd temp.py/

$ hadoop fs -mkdir /NCDCWeatherData

$hadoop fs -copyFromLocal NCDCWeatherData/*  /NCDCWeatherData

$hadoop jar /home/hduser/hadoop-3.2.1/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar -file temp-mapper.py  -mapper temp-mapper.py  -file temp-reducer.py -reducer temp-reducer.py -input /NCDCWeatherData -output /weatherout/output

$ hadoop fs –cat weatherout /output/part-00000

**Code:**
**temp-mapper.py**

```
import re
import sys
for line in sys.stdin:
val = line.strip()
(year, temp, q) = (val[15:19], val[87:92], val[92:93])
if (temp != "+9999" and re.match("[01459]", q)):
print "%st%s" % (year, temp)
```

**temp-reducer.py**

```
import sys
(last_key, max_val) = (None, 0)
for line in sys.stdin:
(key, val) = line.strip().split("t")
if last_key and last_key != key:
print "%st%s" % (last_key, max_val)
(last_key, max_val) = (key, int(val))
else:
(last_key, max_val) = (key, max(max_val, int(val)))
if last_key:
print "%st%s" % (last_key, max_val)
```

**Expected Output:**
**1901** year has maximum temperature.

**Problem Statement-4:** Write a map-reduce application to find how many times a particular page has been accessed (use from the Apache Web Server log data).

**Description:**
In today's world the usage of internet has become very high and using all the logs from web server we can actually predict the customer moods in buying the product or can analyze the interests of Customer.

**Procedure:**

$cd hadoop-3.2.1/

$cd sbin/

$cd start-all.sh

$jps

$cd..

$cd web-log.py/

$hadoop fs -mkdir -p /weblogip

$hadoop fs -copyFromLocal test_access_log  /weblogip
$hadoop jar /home/hduser/hadoop-3.2.1/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar -file my_mapper_by_page.py  -mapper my_mapper_by_page.py  -file my_reducer.py -reducer my_reducer.py -input /weblogip/test_access_log -output /weblogout/output

$ hadoop fs –cat weblogout /output/part-00000

**Code:**
**my_mapper_by_page.py**

```
import sys

for line in sys.stdin:
        data = line.strip().split(" ")
        if len(data) == 10:
        page = data[6]
        print page
```

**my_reducer.py**

```
import sys

number = 0
oldKey = None

for line in sys.stdin:
        thisKey = line

    if oldKey and oldKey != thisKey:
        print oldKey, "\t", number
```

```
        oldKey = thisKey;
        number = 0

        oldKey = thisKey
        number += 1


if oldKey != None:
        print oldKey, "\t", number
```

OUTPUT:

**Problem Statement-5:** Write a pig script to find max temp from the given dataset.

**Description:**

The Apache Pig MAX function is used to find out the maximum of the numeric values or chararrays in a single-column bag. It requires a preceding GROUP ALL statement for global maximums and a GROUP BY statement for group maximums. However, it ignores the NULL values.

**Procedure:**

$cd hadoop-3.2.1/

$cd sbin/

$cd start-all.sh

$jps

$cd..

$pig –x local

**Code:**

grunt>records = LOAD 'sample.txt' AS (year:chararray, temperature:int, quality:int);



grunt>filtered_records = FILTER records BY temperature != 9999 AND (quality == 0 OR quality == 1 OR quality == 4 OR quality == 5 OR quality == 9);

grunt>grouped_records = GROUP filtered_records BY year;



grunt>max_temp = FOREACH grouped_records GENERATE group,
MAX(filtered_records.temperature);

grunt>DUMP max_temp;

**Sample.txt**

| 1950 | 0 | 1 |
| 1950 | 22 | 1 |
| 1950 | -11 | 1 |
| 1949 | 111 | 1 |
| 1949 | 78 | 1 |

**Problem Statement-6:** Write Pig script to implement Word Count Job.

**Description:** Pig is a high-level programming language useful for analyzing large data sets. Pig was a result of development effort at Yahoo!

In a MapReduce framework, programs need to be translated into a series of Map and Reduce stages. However, this is not a programming model which data analysts are familiar with. So, in order to bridge this gap, an abstraction called Pig was built on top of Hadoop.

**Procedure:**

$cd hadoop-3.2.1/

$cd sbin/

$cd start-all.sh

$jps

$cd..

$pig –x local

**Code:**

grunt>input1 = load 'sample.txt' as (line);
Output



grunt>words = foreach input generate flatten(TOKENIZE(line)) as word;

grunt>grpd = group words by word;



grunt>cntd = foreach grpd generate group, COUNT(words);

grunt>dump cntd;

**Sample.txt**

Mary had a little lamb
its fleece was white as snow
and everywhere that Mary went
the lamb was sure to go

**Problem Statement-7:** Find the Number of Products Sold in Each Country for the given dataset using pig framework.

**Description:**

Apache Pig enables people to focus more on analyzing bulk data sets and to spend less time writing Map-Reduce programs. Similar to Pigs, who eat anything, the Apache Pig programming language is designed to work upon any kind of data.

**Procedure:**

$cd hadoop-3.2.1/

$cd sbin/

$cd start-all.sh

$jps

$cd..

$pig –x local

**Code:**

grunt>salesTable = LOAD '/SalesJan2009.csv' USING PigStorage(',') AS
(Transaction_date:chararray,Product:chararray,Price:chararray,Payment_Type:chararray,Name:c
hararray,City:chararray,State:chararray,Country:chararray,Account_Created:chararray,Last_Logi
n:chararray,Latitude:chararray,Longitude:chararray);



grunt>GroupByCountry = GROUP salesTable BY Country;

grunt>CountByCountry = FOREACH GroupByCountry GENERATE
CONCAT((chararray)$0,CONCAT(':',(chararray)COUNT($1)));



grunt>STORE CountByCountry INTO 'pig_output_sales' USING PigStorage('\t');

**Problem Statement-8:** Illustrate the concept of bucketing and partitioning and bucketing in hive.

**Description:**

**Partitioning –** Apache Hive organizes tables into partitions for grouping the same type of data together based on a column or partition key. Each table in the hive can have one or more partition keys to identify a particular partition. Using partition we can make it faster to do queries on slices of the data.

The Hive command for Partitioning is:

CREATE TABLE table_name (column1 data_type, column2 data_type) PARTITIONED BY (partition1 data_type, partition2 data_type,….);

**Bucketing –** In Hive Tables or partitions are subdivided into buckets based on the hash function of a column in the table to give extra structure to the data that may be used for more efficient queries.

The Hive command for Bucketing is:

CREATE TABLE table_name PARTITIONED BY (partition1 data_type, partition2 data_type,….) CLUSTERED BY (column_name1, column_name2, …) SORTED BY (column_name [ASC|DESC], …)] INTO num_buckets BUCKETS;

**Code:**

set hive.exec.dynamic.partition.mode=nonstrict;

Create table students (id int, name string, year int, dept string) row format delimited fields terminated by „‚"; Load data inpath „hdfs://localhost:54310/lab10/data" into table students;

Create table parteddepartment (id int, name string, year int) PARTITIONED by (dept string) row format delimited fields terminated by „‚";

insert overwrite table parteddepartment PARTITION(dept) SELECT id,name,year,dept from students;

Create table samplebucket (id int, name string, year int) clustered by (name) into 2 buckets row format delimited fields terminated by „‚";

From parteddepartment insert overwrite table samplebucket select id,name,year;

```
hive> create table parteddepartment(id int,name string,year int) PARTITIONED BY(dept string);
OK
Time taken: 0.167 seconds
hive> show tables;
OK
parteddepartment
students
Time taken: 0.126 seconds, Fetched: 2 row(s)
hive>
```

```
Loaded : 3/3 partitions.
        Time taken to load dynamic partitions: 0.619 seconds
        Time taken for adding to write entity : 0.001 seconds
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1   Cumulative CPU: 2.77 sec   HDFS Read: 4272 HDFS Write: 272 SUCCESS
Total MapReduce CPU Time Spent: 2 seconds 770 msec
OK
Time taken: 24.837 seconds
```

```
hive> show partitions parteddepartment;
OK
dept=cse
dept=eee
dept=it
Time taken: 0.174 seconds, Fetched: 3 row(s)
hive>
```

**Problem Statement-9:** From the given Aadhaar dataset , write the hive queries for the following.

a. Count the number of Identities generated in each state
b. Count the number of Identities generated by each Enrolment Agency
c. For how many states the Aadhaar data exists?
d. What is the % of Aadhaar being approved per state?

**Description:**
Hive provides SQL type querying language for the ETL purpose on top of Hadoop file system.

Hive Query language (HiveQL) provides SQL type environment in Hive to work with tables, databases, queries.

We can have a different type of Clauses associated with Hive to perform different type data manipulations and querying. For better connectivity with different nodes outside the environment. HIVE provides JDBC connectivity as well.

Hive queries provides the following features:

- Data modeling such as Creation of databases, tables, etc.
- ETL functionalities such as Extraction, Transformation, and Loading data into tables
- Joins to merge different data tables
- User specific custom scripts for ease of code
- Faster querying tool on top of Hadoop

**Code:**

create table aad(register string,enrolment_agency string,state string,district string,sub_district string,pincode int,gender string,age int,aad_generated int,enroll_rejected int,res_providing_email int,res_providing_number int)row format delimited fields terminated by ',' stored as textfile;

load data local inpath '/home/ak/Desktop/adata.txt' overwrite into table aad;

a. select state,count(*) from aad group by state;

b. select enrolment_agency,count(*) from aad group by enrolment_agency;

c. select count(distinct state) from aad;

d. select state,((sum(aad_generated)-sum(enroll_rejected))/(sum(aad_generated))*100) from aad group by state;

**OUTPUTS:**
**a)**

**b)**



```
Transline Technologies P Ltd      119
Transmoovers India      6
Twinstar Industries Ltd.      4005
UID e-Seva Society Ahmedabad      560
UIDAI-EA      19
UMC Technologies Pvt. Ltd      29
UNITED DATA SERVICES PRIVATE LIMITED    3
UT Computers Educational & Welfare Soc  188
UT of Daman and Diu      50
United Telecoms Ltd      1272
United Telecoms e-Services Pvt Ltd      27
Urmila Info solution      162
Utility Forms Pvt Ltd    3004
VAP INFOSOLUTIONS      414
VEETECHNOLOGIES PVT. LTD      2428
VIKALP MULTIMEDIA      1
VIRGO SOFTECH LIMITED      670
VISESH INFOTECNICS LIMITED      235
VISION COMPTECH INTEGRATOR LTD   985
Vakrangee Softwares Limited      3454
Vayam technologies Ltd  415
Vedavaag Systems Limited      4654
Viesa Technologies      7
Virinchi Technologies Ltd      615
WEBEL   1
WEBEL TECHNOLOGY LIMITED      461
Wedha Communication Pvt Ltd      131
Wipro Ltd      6175
Women and Child Development      39
Yash Ornaments Pvt. Ltd  405
Yashi Informatics LLP   6
Yuvaan Infotech 608
Zephyr System Pvt.Ltd.  5949
e-Seva Society  Chhotaudepur      315
e-Seva Society UID Dang 108
eCentric solutions pvt ltd      21
Time taken: 16.327 seconds, Fetched: 326 row(s)
hive>
```

**c)**



```
OK
Barddhaman      4276      2859
North 24 Parganas      3772      3121
South 24 Parganas      3630      2448
Bhagalpur      3543      1744
Patna   3485      1766
Nadia   3460      1673
Murshidabad      3018      1399
Gaya    2915      1487
Kolkata 2678      1388
Katihar 2622      1352
Time taken: 40.263 seconds, Fetched: 10 row(s)
hive> select count(distinct state) from aad;
Query ID = hdoop_20210102195453_ae73de31-9a69-4d3a-9a9a-0fab7688e2df
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1609586600747_0005, Tracking URL = http://user-HP-Laptop-15-bs1xx:8088/proxy/application_1609586600747_0005/
Kill Command = /home/hdoop/hadoop-3.2.1/bin/mapred job  -kill job_1609586600747_0005
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-01-02 19:54:59,341 Stage-1 map = 0%,  reduce = 0%
2021-01-02 19:55:04,474 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 2.5 sec
2021-01-02 19:55:08,585 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 3.92 sec
MapReduce Total cumulative CPU time: 3 seconds 920 msec
Ended Job = job_1609586600747_0005
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1   Cumulative CPU: 3.92 sec   HDFS Read: 46493467 HDFS Write: 102 SUCCESS
Total MapReduce CPU Time Spent: 3 seconds 920 msec
OK
38
Time taken: 16.293 seconds, Fetched: 1 row(s)
hive>
```

**d)**



```
Arunachal Pradesh        91.34720700985761
Assam    99.12854030501089
Bihar    95.56661152348916
Chandigarh       96.52509652509652
Chhattisgarh     94.24591156874621
Dadra and Nagar Haveli  97.85714285714285
Daman and Diu    98.09523809523809
Delhi    92.98599572751009
Goa      89.11739502999143
Gujarat 92.82516358627024
Haryana 92.18106995884774
Himachal Pradesh         92.04912734324499
Jammu and Kashmir        83.46839546191248
Jharkhand        95.16619375760033
Karnataka        92.68872697834446
Kerala  96.44720332827049
Lakshadweep      75.0
Madhya Pradesh  93.41354456040243
Maharashtra      93.03047728579644
Manipur 85.63869992441421
Meghalaya        99.27797833935018
Mizoram 97.05367096671445
Nagaland         92.11009174311927
Odisha  91.2220877791222
Others  100.0
Puducherry       87.95180722891565
Punjab 93.529050107593
Rajasthan        94.8572150619156
Sikkim  96.0
State    NULL
Tamil Nadu       96.32753578574726
Telangana        94.53965723395775
Tripura 93.17180616740089
Uttar Pradesh    94.90589493769696
Uttarakhand      92.7799198608906
West Bengal      94.68561563289714
Time taken: 16.497 seconds, Fetched: 38 row(s)
hive>
```

**Problem Statement-10:** Implement the word count job in Scala.

**Description:**
Scala is a general-purpose programming language. It supports object oriented, functional and imperative programming approaches. It is a strong static type language. In scala, everything is an object whether it is a function or a number. It does not have a concept of primitive data.

**Procedure:**
hdfs dfs -mkdir /spark
hdfs dfs -put /home/kgr/sparkdata.txt /spark
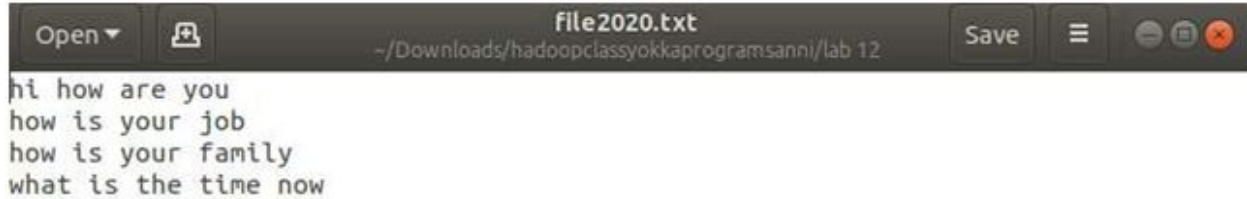
**Code:**
```scala
scala> val data=sc.textFile("sparkdata.txt");
val splitdata = data.flatMap(line => line.split(""));
val mapdata = splitdata.map(word => (word,1));
val reducedata = mapdata.reduceByKey(_+_);
reducedata.collect;
```

For dataset:

```
val text = sc.textFile("mytextfile.txt")
val counts = text.flatMap(line => line.split(" ")).map(word => (word,1)).reduceByKey(_+_)
counts.collect
```

**Input file:**

| Open ▾ | 📁 | file2020.txt<br>~/Downloads/hadoopclassyokkaprogramsanni/lab 12 | Save | ≡ | ⊖ ⊡ ⊗ |

```
hi how are you
how is your job
how is your family
what is the time now
```

**Outputs:**

```
scala> var sampleFile = sc.textFile("hdfs://localhost:54310/lab12/file2020.txt");
sampleFile: org.apache.spark.rdd.RDD[String] = hdfs://localhost:54310/lab12/file2020.txt MapPartitionsRDD[12] at textFile at <console>:24

scala> var wCount = sampleFile.flatMap(line => line.split(" "))
wCount: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[13] at flatMap at <console>:25

scala> wCount.collect
res4: Array[String] = Array(hi, how, are, you, how, is, your, job, how, is, your, family, what, is, the, time, now)

scala> var mapOp=wCount.map(w => (w,1))
mapOp: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[14] at map at <console>:25

scala> mapOp.collect
res5: Array[(String, Int)] = Array((hi,1), (how,1), (are,1), (you,1), (how,1), (is,1), (your,1), (job,1), (how,1), (is,1), (your,1), (family,1), (what,1), (is,1), (the,1), (time,1), (now,1))

scala> var reduceOp=mapOp.reduceByKey(_+_)
reduceOp: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[15] at reduceByKey at <console>:25

scala> reduceOp.collect
res6: Array[(String, Int)] = Array((are,1), (is,3), (family,1), (how,3), (what,1), (now,1), (job,1), (you,1), (hi,1), (time,1), (your,2), (the,1))
```