

# Lab 3

**Roll No.:** J019-Ayush Hendre, J031-Rohit Mittal

**Aim:** Word Count Using Map Reduce

**Objectives:**

- 1.To run Hive command.
2. Copy Data file from Local to HDFS.
3. Generate a Word count query.
4. Display Word count of the file

**Codes:**

```
//Map Reduce in
```

```
HIVE hive
```

```
CREATE TABLE FILES (line STRING);
```

```
LOAD DATA INPATH 'data1.txt' OVERWRITE INTO TABLE FILES;
```

```
CREATE TABLE word_count AS
```

```
SELECT w.word, count(1) AS count from
```

```
(SELECT explode(split(line, ' ')) as word from
```

```
FILES) w GROUP BY w.word
```

```
ORDER BY w.word;
```

```
SELECT * FROM word_count ;
```

```
[cloudera@quickstart hive1]$ hadoop fs -put data.txt data1.txt
[cloudera@quickstart hive1]$ hive
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties
WARNING: Hive CLI is deprecated and migration to Beeline is recommended.
hive> LOAD DATA INPATH 'data1.txt' OVERWRITE INTO TABLE FILES;
Loading data to table default.files
chgrp: changing ownership of 'hdfs://quickstart.cloudera:8020/user/hive/warehouse/files/data1.txt':
does not belong to supergroup
Table default.files stats: [numFiles=1, numRows=0, totalSize=50, rawDataSize=0]
nk
```

In order to change the average load •or a reducer (in bytes ) :

```
set hive.exec.seducers.bytes.per.reducer=<number>
```

In order to limit the maximum number of seducers:

```
set hive.exec.seducers.max:<number>
```

In order to set a constant number of seducers:

```
set mapreduce.job.reduces=<numben>
```

Starting Job: job\_1614416156655\_B082, Tracking URL: http://quickstart.cloudera:8088/

proxyfapplication\_1614416156655\_B B2/

Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job\_1614416156655\_9682

MapReduce job Information for Stage-2: number of mappers : 1; number of reducers : 1 2621—

2021-02-27 02:09:52, 727 Stage-2 map = 8g, reduce = a'g

2021-02-27 02:09:52, 727 Stage-2 map = 1001, reduce = B •, Curru1at1 ve CPU 1.5 sec

2021-02-27 02:09:52, 727 Stage-2 map = 1B0g, reduce = been, Cumu1 ative CPU 5.84 sec

MapReduce Total cuoulat1ve CPU t:ime: 5 seconds 48 msec

Ended Job = job\_1614416156655\_BB6Z

Moving data to: hdfs://qul c kstart. cloudera:ae2a/user/ hive/warehouse/word\_count

Table default.word\_count stats: [numFiles=1, numRows=7, totalSize=54, rawDataSize=47]

MapReduce Jobs Launched:

Stage-Stage-1: Map: 1 Reduce: 1 Cumulative cPu: 9.74 etc HDFS Read: 7589 HDFS Write: 262

SUCCESS Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 5.04 sec HDFS Read: 4886 HDFS

Write: 128 SUCCESS Total MapReduce CPU Time Spent: 8 seconds 7ae msec

h1ve> CREATE TABLE word\_count AS

```
> SELECT w.xord, count(1) AS count from
```

```
>> (SELECT explode(sp1lt(11ne, ' ')) AS word FROM FILES) u
```

```
> GROUP BY w.word
```

```
> ORDER BY x.xord;
```

query ID = c louder a\_2e21e227820808\_9b7e 516d -essb-<wg -sfc6 -56dra2c78bbc

Total jobs = 2

Launching Job 1 out of 2

Number of reduce tasks not specified. Estimated from input data size: 1 In order

to change the average load -for a reducer ( in bytes ) :

```
set hive.exec.seducers.bytes.per.reducer=<number>
```

In order to limit the maximum number of seducers:

```
set hive.exec.seducers.max=<number>
```

In order to set a constant number of seducers :

```
set mapreduce.job.reduces=<nuober>
```

Starting Job = job\_1614416156655\_0001, Tracking URL: http://

quickstart.cloudera:8088/proxy/application\_1614416156655\_ee01/

Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job\_1614416156655\_ee01

MapReduce job Information -For Stage-1: number of mappers : 1; number of seducers : 1 2021-

2021-02-27 02:09:52, 727 Stage-1 map = US, reduce = BE

2021-02-27 02:09:52, 727 Stage-1 map = 1001, reduce = 0X•, Cumu1at1ve CPU 2. e2 sec

2021-02-27 02:09:52, 727 Stage-1 map = 1BBX, reduce = 18Bg, Cumu1at1ve CPU 3.74 sec

MapReduce Total cumu1at1ve CPU t:ime: 3 seconds 74ig Insec

Ended Job = job\_1614416156655\_0001

Launching Job 2 out of 2

Number of reduce tasks determined at compile time: 1

h1ve SELE \* FRON

```
> CT word_count
```

OK

Th1s 2

a 2

h1ve 1

is 2

spar 1

k

tutorial 1

tutorial. 1

Time taken: e.ea2 seconds, Fetched: 7 row(s)