

Lab 3

Roll No.: J061-Anirudh VM, J069-Shubh Lilani

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 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=<number>
```

Starting Job: job_1614416156655_B082, Tracking URL:

http://quickstart.cloudera:8088/proxy/application_1614416156655_B082/

Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1614416156655_9682

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

02—27 B2:B9:4e,727 Stage-2 map = 8g, reduce = a'g

2021-62-27 02:09:52, B72 Stage-2 map = 1001, reduce = B*, Cumulative CPU 1.5 sec

2621-B2-27 BZ:GB:BB,B57 Stage-2 map = 1B0g, reduce = been, Cumulative CPU 5.84 sec

MapReduce Total cumulative CPU time: 5 seconds 48 msec

Ended Job = job_1614416156655_BB6Z

Moving data to: hdfs://quickstart.cloudera:8020/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 74 msec

hive> CREATE TABLE word_count AS

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

```
» (SELECT explode(sp11t(11ne, ' ')) AS word FROM FILES) u
```

```
» GROUP BY w.word
```

```
» ORDER BY x.xord;
```

query ID = c_louder_a_2e21e227820808_9b7e516d-essb-cwg-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=<number>
```

Starting Job = job_1614416156658_0001, Tracking URL:

http://quickstart.cloudera:8088/proxy/application_1614416156658_0001/

Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1614416156655_ee01

Job Information for Stage-1: number of mappers: 1; number of reducers: 1 2021-e2-

27 ez:e9:83,631 Stage-1 map = US, reduce = BE

2021-82-27 e2:e9:15,283 Stage-1 map = 1001, reduce = 0X*, Cumulative CPU 2.02 sec

2B21-82-27 02:B9:27,523 Stage-1 map = 1BBX, reduce = 18Bg, Cumulative CPU 3.74 sec

MapReduce Total cumulative CPU time: 3 seconds 74 msec

Ended Job = job_1614416156658_0001

Launching Job 2 out of 2

Number of reduce tasks determined at compile time: 1

hive> SELECT * FROM word_count

OK

Th1s 2

a 2

hive 1

is 2

spark 1

tutorial 1

tutorial. 1

Time taken: 0.02 seconds, Fetched: 7 row(s)