# IEMS5730 Homework2

## Q1 Pig Setup



**P1 Pig install**



**P2 Pig environment setting**



**P3 Connect to pig**

**P4 Combine two tables into one table and upload to hdfs**

## Pig Script:

A = load '/bigrams/table' as (bigram: chararray, year: int, match_count: int, volume_count: int);
describe A;
grouped = group A by bigram;
avgoccurence = foreach grouped generate group, AVG(A.match_count) as avgoccur;
sorted = order avgoccurence by avgoccur desc;
top_20 = limit sorted 20;
describe top_20;
dump top_20;
Store top_20 into '/bigrams/top_20';

hadoop fs -cat /bigrams/top_20

**P5 Pig result**


```
2021-03-02 11:09:59,178 [main] INFO  org.apache.pig.Main - Pig script completed in 45 minutes, 27 seconds and 784 milliseconds
(2727784 ms)
```

**P6 Pig running time(45minutes, 27seconds and 784milliseconds)**

# Q2 Hive Setup


```
unset i
unset -f pathmunge
export JAVA_HOME=/usr/lib/jvm/java
export PATH=$JAVA_HOME/bin:$PATH
export HADOOP_HOME=/usr/hadoop/hadoop-2.10.1
export PATH=$PATH:$HADOOP_HOME/sbin:$HADOOP_HOME/bin
HADOOP_STREAM=$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-2.10.1.jar

export HADOOP_STREAM

export PIG_HOME=/usr/local/pig-0.17.0
export PIG_CLASSPATH=/usr/local/hadoop-2.10.1/conf
export PATH=$PATH:/usr/hadoop/hadoop-2.10.1/bin:$PIG_HOME/bin
export HIVE_HOME=/usr/hadoop/apache-hive-1.2.2-bin
export HIVE_CONF=${HIVE_HOME}/conf
```

**P7 Hive install and environment setting**


```
[root@master conf]# hadoop fs -mkdir -p /user/hive/warehouse
[root@master conf]# hadoop fs -chmod 777 /user/hive/warehouse
```

**P8 Grant permissions**

```
hive> load data inpath '/bigrams/test' into table bigrams;
Loading data to table default.bigrams
Table default.bigrams stats: [numFiles=1, totalSize=94]
OK
Time taken: 1.306 seconds
hive> select* from bigrams;
OK
circumvallate    1978    335     NULL
circumvallate    1979    261     95
asds    1968    234    23
   > ad    1987    111    57
```

**P9 Hive test**

## Hive scripts:

create table bigrams(bigram string, year int, match_count int, volume int) row format delimited fields terminated by '\t';
show tables;
LOAD DATA INPATH '/bigrams/table' INTO TABLE bigrams;
SELECT bigram,AVG(match_count) as avg FROM bigrams GROUP BY bigram order by avg desc limit 20;

```
and       2.593207744E7
and_CONJ        2.5906234451764707E7
a         1.6665890811764706E7
a_DET     1.6645121127058823E7
as        6179734.075294117
be        5629591.52
be_VERB   5621156.232941177
as_ADP    5360443.872941176
by        5294067.04
by_ADP    5272951.997647059
are       4298564.341176471
are_VERB        4298561.303529412
at        3676050.1529411767
at_ADP    3670625.785882353
an        2979272.7411764706
an_DET    2977977.8870588234
but       2471102.4964705883
but_CONJ        2468978.0564705883
all       2189962.722352941
all_DET 2161257.294117647
Time taken: 257.885 seconds, Fetched: 20 row(s)
```

**P10 Hive result(same as pig)**

```
Stage-Stage-1: Map: 12  Reduce: 12   Cumulative CPU: 240.05 sec   HDFS Read: 3070095634 HDFS Write: 90818230 SUCCESS
Stage-Stage-2: Map: 1  Reduce: 1   Cumulative CPU: 16.87 sec   HDFS Read: 90825662 HDFS Write: 466 SUCCESS
Total MapReduce CPU Time Spent: 4 minutes 16 seconds 920 msec
```
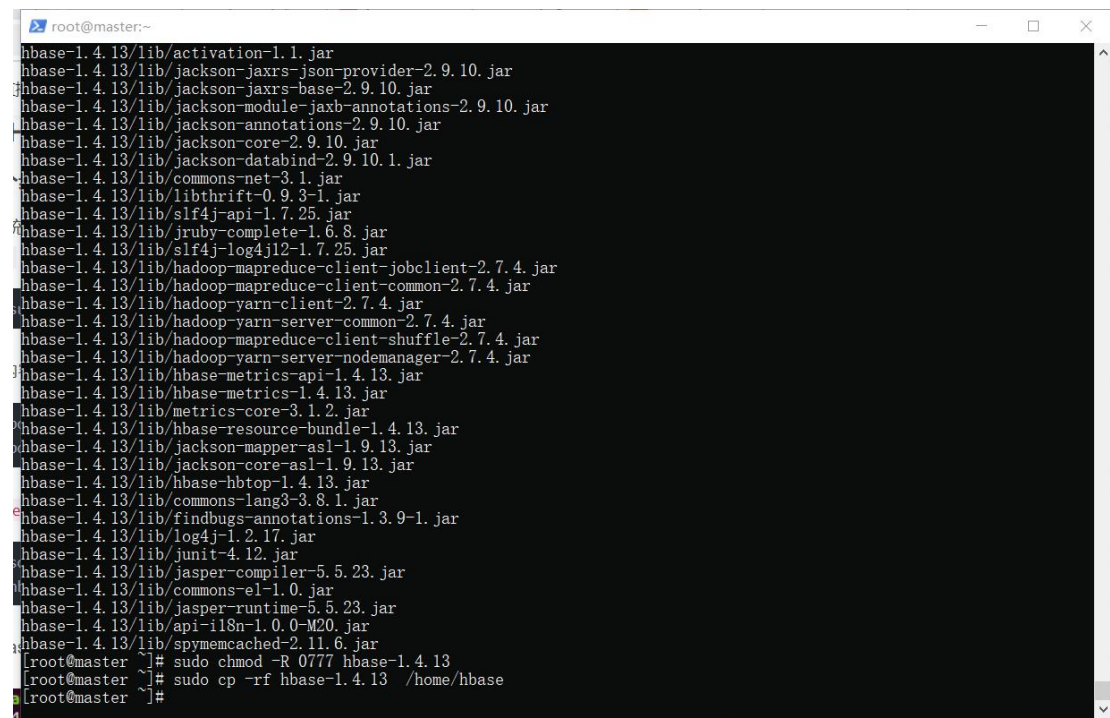
**P11 Hive running time(4 minutes 16 seconds 920 msec)**

In P6 and P11, we can compare the running time of pig and hive. From my results, the hive running time is much more less than pig's. The pig spent 45 minutes while hive spent only about 4 minutes. Put aside the time complexity of my program and the system delay, hive can run faster than pig to get the same result.

## Q3: Hbase Setup



**P12 Hbase setup**

**P13 Environment configuration in master and slaves**

**P14 Hbase cluster started successfully**

## ImportTsv script:

hbase org.apache.hadoop.hbase.mapreduce.ImportTsv '-Dimporttsv.separator= '
-Dimporttsv.bulk.output=/output/hfile
-Dimporttsv.columns=HBASE_ROW_KEY,cf:a test    /bigrams/table2

### All Applications

| Apps Pending | Apps Running | Apps Completed | Containers Running | Used Resources | Total Resources |
|---|---|---|---|---|---|
| 0 | 1 | 1 | 11 | <memory:12288, vCores:11> | <memory:24576, vCores:24> |

| Decommissioning Nodes | Decommissioned Nodes | Lost Nodes | Unhealthy Nodes | Rebooted Nod |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |

| Scheduling Resource Type | Minimum Allocation | Maximum Allocation |
|---|---|---|
| [<name=memory-mb default-unit=Mi type=COUNTABLE>, <name=vcores default-unit= type=COUNTABLE>] | <memory:1024, vCores:1> | <memory:8192, vCores:4> |

| | User | Name | Application Type | Queue | Application Priority | StartTime | LaunchTime | FinishTime | State | FinalStatus | Running Containers | Allocated CPU VCores | Allocated Memory MB | Allocated GPUs | Reserved CPU VCores | Reserved Memory MB | Reserved GPUs | % of Queue |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0003 | root | importtsv_test | MAPREDUCE | default | 0 | Tue Mar 2 09:12:08 +0800 2021 | Tue Mar 2 09:24:36 +0800 2021 | N/A | RUNNING | UNDEFINED | 11 | 11 | 12288 | -1 | 0 | 0 | -1 | 50.0 |

**P15 ImportTsv mission on yarn**

**Hbase shell script:**
**create 'table2', {NAME => 'bigram'}, {NAME => 'year'}, {NAME => 'match_count'},{NAME => 'volume_count'}**
**describe 'table2'**
**put 'table2', 'ierg4330', '2019', '100','4'**
**scan 'table2', FILTER=>"ValueFilter(=, 'binary:1671')"**

# Bonus: Setup Pig on hadoop over Kubernetes



```
and      2.593207744E7
and_CONJ      2.5906234451764707E7
a        1.6665890811764706E7
a_DET    1.6645121127058823E7
as       6179734.075294117
be       5629591.52
be_VERB  5621156.232941177
as_ADP   5360443.872941176
by       5294067.04
by_ADP   5272951.997647059
are      4298564.341176471
are_VERB      4298561.303529412
at       3676050.1529411767
at_ADP   3670625.785882353
an       2979272.7411764706
an_DET   2977977.8870588234
but      2471102.4964705883
but_CONJ      2468978.0564705883
all      2189962.722352941
```

**P16 Pig result**

```
HadoopVersion   PigVersion    UserId  StartedAt          FinishedAt
2.10.1   0.17.0   root    2021-03-02 18:26:28    2021-03-02 18:34:30
```

**P17 Running time**

In this part, the running time of the same script is 8 minutes and 2 seconds, compared to the running time in Q1(45 minutes). It's much less in the pig via Kubernetes.