

Project #6: Hadoop MapReduce

Systems Programming
Department of Computer Science and Engineering
Sogang University



Due: June 27 (Thu), 11:59PM (KST)



Goal

1

Goal

The goal of this project is to improve your understanding basic of Hadoop and MapReduce.

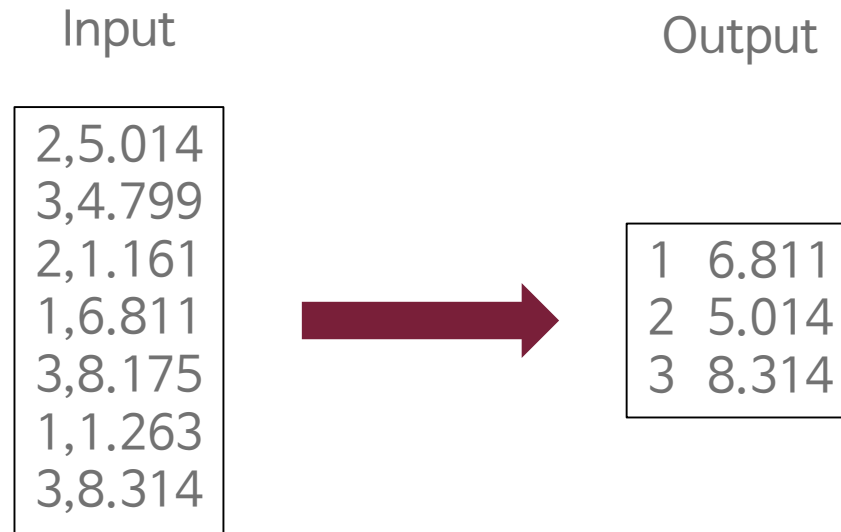
1. Write all the code in Python3.
2. Do with Azure HDInsight service.
3. Use MapReduce technique for processing of data.
4. Solve the group by max problem for a large size of data.

2

Group by max

This problem is to find the maximum value in each group. The Input has a one integer and one real number each line, each number is separated by a comma. A integer number part represents a group number. A real number part means a value. You must find the maximum value for each group for given data.

Here is an example of this problem:



3

How to generate input.data and upload it to the server

Run the given generate.py file in linux server (cspiro, or azure cloud shell) with the command : “python3 generate.py 1000000000 16”

```
Bash
@Azure: ~$ python3 generate.py 1000000000 16
0
10000
20000
30000
40000
50000
60000
70000
80000
90000
100000
110000
120000
```

... ..

```
Bash
99880000
99890000
99900000
99910000
99920000
99930000
99940000
99950000
99960000
99970000
99980000
99990000
Done!
@Azure: ~$
```

The dataset ‘input.data’ to be given consists of 100 million real numbers from 0 to 10000. And there are 16 group keys.

3

How to generate input.data and upload it to the server

Upload the input.data file in your local directory to the hdfs environment with the command : “hdfs dfs -put input.data /example/data/input.data”

```
Bash
sshuser@hn0-cluste:~$ ls
generate.py input.data mapper.py reducer.py
sshuser@hn0-cluste:~$ hdfs dfs -put input.data /example/data/input.data
sshuser@hn0-cluste:~$ hdfs dfs -ls /example/data
Found 8 items
-rw-r--r-- 1 root supergroup 66 2019-06-09 13:40 /example/data/fruits.txt
drwxr-xr-x - root supergroup 0 2019-06-09 13:40 /example/data/gutenberg
-rw-r--r-- 1 sshuser supergroup 2071324728 2019-06-09 13:54 /example/data/input.data
-rw-r--r-- 1 root supergroup 77 2019-06-09 13:40 /example/data/people.json
drwxr-xr-x - root supergroup 0 2019-06-09 13:40 /example/data/people.parquet
drwxr-xr-x - root supergroup 0 2019-06-09 13:40 /example/data/people.seq
-rw-r--r-- 1 root supergroup 97884 2019-06-09 13:40 /example/data/sample.log
-rw-r--r-- 1 root supergroup 62 2019-06-09 13:40 /example/data/yellowthings.txt
sshuser@hn0-cluste:~$
```

Check that the file is uploaded well into the hdfs environment with the command : hdfs dfs -ls /example/data

4

Hadoop File System command

```
hdfs dfs -put <local_src> . . . <dst>
```

Copy single src, or multiple srcs from local file system to the destination file system.

```
hdfs dfs -put input.data /example/data/input.data
```

```
hdfs dfs -get <src> <local_dst>
```

Copy files to the local file system.

```
hdfs dfs -put /example/data/output.data output.data
```

```
hdfs dfs -ls <args>
```

Files within a directory are order by filename by default.

```
hdfs dfs -ls /example/wordcountout/
```

```
hdfs dfs -text <src>
```

Takes a source file and outputs the file in text format.

```
hdfs dfs -text /example/wordcountout/part-00000
```

```
hdfs dfs -rm [-R] URI [URI . . .]
```

Delete files specified as args. The -R option deletes the directory and any content under it recursively.

```
hdfs dfs -rm -R /example/wordcountout*
```


5

Requirements

You need to resolve the group by max for the input.data you created yourself.

Generate 10000000 numbers for 16 groups using the given code.

You must include your student ID in output file name when use Hadoop MapReduce.

```
yarn jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar ₩  
-files mapper.py,reducer.py-mapper mapper.py -reducer reducer.py ₩  
-input /example/data/input.data ₩  
-output /example/result[your Student ID]
```

If your student ID is 20161234, then output file is “result2016124”

Submission

1

Things

(1) Python codes

- Each python code file named “mapper.py”, “reducer.py”

(2) A document file

- This document file should describe how you implemented your programs.
- You should insert a data flow diagram that shows how MapReduce worked.
- A sample document will be posted on cyber campus.

(3) Capture images

- Two capture images for execution and output about MapReduce.
- “20161234_1.png” is captured image about with execution command.
- “20161234_2.png” is captured image about MapReduce result.
- The numeric part should be **your student ID**.

2 Example of a captured image “20161234_1.png”

```

Bash
sshuser@hn0-cluste:~$ yarn jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar -files mapper.py, reducer
.py -mapper mapper.py -reducer reducer.py -input /example/data/input.data -output /example/data/result20161234
packageJobJar: [] [/usr/hdp/2.6.5.3008-11/hadoop-mapreduce/hadoop-streaming-2.7.3.2.6.5.3008-11.jar] /tmp/streamjob47
43914934645318502.jar tmpDir=null
19/06/09 14:04:38 INFO client.AHSPProxy: Connecting to Application History server at headnodehost/10.0.0.21:10200
19/06/09 14:04:38 INFO client.AHSPProxy: Connecting to Application History server at headnodehost/10.0.0.21:10200
19/06/09 14:04:39 INFO client.RequestHedgingRMFailoverProxyProvider: Looking for the active RM in [rm1, rm2]...
19/06/09 14:04:39 INFO client.RequestHedgingRMFailoverProxyProvider: Found active RM [rm2]
19/06/09 14:04:40 INFO lzo.GPLNativeCodeLoader: Loaded native gpl library
19/06/09 14:04:40 INFO lzo.LzoCodec: Successfully loaded & initialized native-lzo library [hadoop-lzo rev b5efb3e531b
c1558201462b8ab15bb412ffa6b89]
19/06/09 14:04:40 INFO mapred.FileInputFormat: Total input paths to process : 1
19/06/09 14:04:40 INFO mapreduce.JobSubmitter: number of splits:4
19/06/09 14:04:40 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1560086904039_0002
19/06/09 14:04:41 INFO impl.YarnClientImpl: Submitted application application_1560086904039_0002
19/06/09 14:04:41 INFO mapreduce.Job: The url to track the job: http://hn1-cluste.3oquisiz15gu3pixh02cs3blzc.psx.inte
rnal.cloudapp.net:8088/proxy/application_1560086904039_0002/
19/06/09 14:04:41 INFO mapreduce.Job: Running job: job_1560086904039_0002
19/06/09 14:05:00 INFO mapreduce.Job: Job job_1560086904039_0002 running in uber mode : false
19/06/09 14:05:00 INFO mapreduce.Job: map 0% reduce 0%
19/06/09 14:05:27 INFO mapreduce.Job: map 5% reduce 0%
19/06/09 14:05:30 INFO mapreduce.Job: map 9% reduce 0%
  
```

2 Example of a captured image “20161234_2.png”

```
Bash
sshuser@hn0-cluste:~$ hdfs dfs -text /example/data/result20161234/part-00000
19/06/09 14:16:31 INFO Izo.GPLNativeCodeLoader: Loaded native gpl library
19/06/09 14:16:31 INFO Izo.LzoCodec: Successfully loaded & initialized native-lzo library [hadoop-lzo rev b5efb3e531b
c1558201462b8ab15bb412ffa6b89]
11      9997.47050915
10      9996.36036315
13      9993.13887596
12      9989.68431495
15      9999.11036256
14      9994.30161411
1       9997.27526803
0       9994.8446836
3       9992.5913585
2       9991.17793128
5       9992.72677359
4       9996.35257724
7       9997.67399153
6       9997.73439591
9       9998.20677729
8       9994.16280426
sshuser@hn0-cluste:~$
```

3

Instructions

- Make a directory named “sp20161234_proj6”. The numeric part should be **your student ID**.
- Put all the files in the directory, and compress the directory itself using tar or zip.
- When you make a tar file, do NOT use the z option (which makes a gz compressed file.)

Example:

```
sp20161234_proj6/  
    document.docx  
    20161234_1.png  
    20161234_2.png  
    mapper.py  
    reducer.py
```

3

Instructions

The file for submission

sp20161234_proj6.tar or sp20161234_proj5.zip

Upload this file on the cyber campus.

Late Submission

No late submissions accepted for this project