

Experiment - 5

Write a program of word count using Map-Reduce

Creating the path files/directories

hdfs dfs -mkdir -p /user/process

Create a file wordcount.txt and copy this file in this directory.

hdfs dfs -put /netjs/MapReduce/wordcount.txt /user/process

Map Reduce java code

```

import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

public class WordCount {

public static class MyMapper extends Mapper<LongWritable, Text,
Text, IntWritable> {

private Text Word = new Text();

public void map (LongWritable key, Text value, Context context)
throws IOException, InterruptedException {

String [] stringArr = value.toString().split("\\s+");

for (String str : stringArr) {
word.set(str);
content.write

}

>

>

public static class MyReducer extends Reducer<Text, IntWritable,
Text, IntWritable> {

private IntWritable result = new IntWritable();

public void reduce (Text key, Iterable<IntWritable> values,
Context context)

throws IOException, InterruptedException {

int sum = 0;

for (IntWritable val : values) {

sum += val.get();

>

result.set(sum);

content.write(key, result);

>

>

```

public static void main (String [] args) throws Exception {
    Configuration conf = new Configuration ();
    Job job = Job . get Instante (conf, "wc");
    job . set Jar By Class (wordCount . class);
    job . set Mapper Class (MyMapper . class);
    job . set Reducer class (MyReducer . class);
    job . set Output Key Class (Text . class);
    job . set Output Value Class (IntWritable . class);
    File Input Format . add InputPath (job, new Path (args [0]));
    File Output Format . set OutputPath (job, new Path (args [1]));
    system . exit (job . wait For Completion (true) ? 0 : 1);
}
>
>

```

Required jars for hadoop MapReduce code:-

These are found under inside

share / hadoop directory

hadoop - common - 2.9.0 . jar
 hadoop - hdfs - 2.9.0 . jar
 hadoop - hdfs - client - 2.9.0 . jar
 hadoop - mapreduce - client - core - 2.9.0 . jar
 hadoop - mapreduce - common . 2.9.0 . jar
 hadoop - mapreduce - client - jobclient . 2.9.0 . jar
 hadoop - mapreduce - client - hs - 2.9.0 . jar
 hadoop - mapreduce - client - opp - 2.9.0 . jar
 Commons - io - 2.4 . jar

Creating jar of your wordcount mapreduce code
once the code is compiled, create a jar file, in eclipse IDE.

Select Export - Java - jar file

To run the map reduce code -

go to directory

bin/hadoop jar /net.js/MapReduce/wordcount.jar org.net.js.WordCount /user/process/user/out

Jar Path file as net.js/MapReduce/wordcount.jar org.net.js.WordCount
org.net.js.WordCount is fully qualified path to your jar
file to java program class

/user/process - path to ⁱⁿoutput directory

/user/out - path to output directory

Experiment

Implementation of Matrix multiplication using map reduce.

Before writing the code let's first create Matrices and put them in HDFs.

Create two files M1 and M2 and put the matrix values
 (separate columns with spaces and rows with a line break!)

Put the above files to HDFs at location

/user/cloudera/matrices/

hdfs dfs -mkdir /user/cloudera/matrices

hdfs dfs -put /path/to/M1 /user/cloudera/matrices/

hdfs dfs -put /path/to/M2 /user/cloudera/matrices/

Creating two program Mapper.py and Reducer.py

defining the dimension of matrix

`#!/usr/lib/python`

`import sys`

`m-r = 2`

`m-c = 3`

`n-r = 3`

`n-c = 2`

`i = 0`

Read each line ie a row from stdin and Split then to separate elements. Map int to each element as we

read element as string from stdin.

for line in sys.stdin: el = map(int, line.split())

The mapper will first read the first matrix and then the second differentiate them we can keep a count i of the line number we are reading and first $m - \sigma$ line will bring the first matrix.

if ($i < m - \sigma$):

for j in range(len(el)):

for k in range(n - c):

print " $\{0\} \text{t} \{2\} \text{t} \{3\}$ ", format(i, k, j, el[j])

else :

for j in range(len(el)):

for k in range(m - \sigma):

print " $\{0\} \text{t} \{1\} \text{t} \{2\} \text{t} \{3\}$ ".format(k, j, i - m - \sigma, el[j])

$i = i + 1$

pointing the key values

$\{0\} \{1\} \{2\}$ are the part of key and $\{3\}$ is the value

$\{0\} \{1\} \{2\}$ actually represent the position of element from $A \otimes B$
to $A^* B$

$\{0\}$ is the row position of the element.

$\{1\}$ is the column position of the element

$\{2\}$ is the position of the element in addition.

We can see that A's element is repeated B's number of column times ie 2 and B's element is repeated A's number of row times ie 2.

i is used to iterate through each row.

j is used to iterate through each column.

k is used to iterate through each duplicate produced.

for each element in matrix A:

element remains in same row, therefore $\langle 0 \rangle = i$

element is duplicated and distributed to each column, therefore column pos in $A \times B$ = Duplication order of element ie $\langle 1 \rangle = k$

and $\langle 2 \rangle = j$

for each element in matrix B:-

Each element remain in the same column, so $\langle 1 \rangle = j$

$\langle 0 \rangle = k$

and $\langle 2 \rangle = i - m + j$

Reducer code

```
#!/usr/bin/python
```

```
import sys
```

```
m-j=2
```

```
m-c=3
```

```
n-j=3
```

```
n-c=2, matrix = [ ]
```

for row in range(m-j):

```
j=[ ]
```

for col in range(m-c):

S = 0

for el in range(m - c):

mul = 1

for num in range(2):

line = sys.stdin.readline()

n = map(int, line.split('\t'))[-1]

mul * n

S += mul

s.append(s)

matrix.append(s)

print('In'.join([str(x) for x in matrix]))

Running the Map-Reduce job on hadoop

\$ chmod +x /Desktop/mo/matrix-mul/Mapper.py \$ chmod +x

\$ chmod +x /Desktop/mo/matrix-mul/Reducer.py

\$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar

> -input /user/cloudera/matrices/

-output /user/cloudera/mat-output

mapper /Desktop/mo/matrix-mul/Mapper.py |

reducer /Desktop/mo/matrix-mul/Reducer.py

get output

hdfs dfs -cat /user/cloudera/mat-output/*

Experiment -

Write a program to of map reduce for weather dataset.

The aim of program to find the max temp. recorded for each year of NCDC data.

Map reduce is based on set of key value pair. 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 the types of input key value pairs are long writable and text and the types of output key value pairs are text and int writable.

Reduce Phase :- Reduce phase takes all the values associated with a particular key. That is all the temp. values belong to a particular year is fed to a same reducer. Then each reducer finds the highest recorded temp.

files that we will create three java classes:-

- HighestMapper.java
- HighestReducer.java
- highestDriver.java

HighestMapper.java

```
import java.io.IOException;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapred.*;
```

```
public class HighestMapper extends MapReduceBase implements
Mapper<LongWritable, Text, Text, IntWritable>
```

<

```
public static final int MISSING = 9999;
```

```
public void map(LongWritable key, Text value, OutputCollector<
<Text, IntWritable> output, Reporter reporter) throws IOException
```

<

```
String line = value.toString();
```

```
String year = line.substring(15, 19);
```

```
int temperature;
```

```
if (line.charAt(87) == '+')
```

```
temperature = Integer.parseInt(line.substring(88, 92));
```

```
else
```

```
temperature = Integer.parseInt(line.substring(87, 92));
```

```
String quality = line.substring(92, 93);
```

```
if (temperature != MISSING && quality.matches("[0-459]"))
```

```
output.collect(new Text(year), new IntWritable(temperature));
```

>

>

Highest Reducer.java

```
import java.io.IOException;
```

```
import java.util.Iterator;
```

```
import org.apache.hadoop.io.*;
```

```
import org.apache.hadoop.mapred.*;
```

```

public class HighestReducer extends MapReduceBase implements
Reducer<Text, IntWritable, Text, IntWritable>
{
    public void reduce(Text key, Iterator<IntWritable> values,
                      OutputCollector<Text, IntWritable> output, Reporter reporter)
    throws IOException
    {
        int maxTemp = 0;
        while (values.hasNext())
        {
            int current = values.next().get();
            if (maxTemp < current)
                maxTemp = current;
        }
        output.collect(key, new IntWritable(maxTemp / 10));
    }
}

```

HighestDriver.java

```

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapred.*;
import org.apache.hadoop.util.*;

```

```

public class HighestDriver extends Configured implements Tool {

```

```

public int sum (String [] args) throws Exception {
    JobConf conf = new JobConf (getConf (), HighestDriver.class);
    conf.setJobName ("Highest Driver");
    conf.setOutputKeyClass (Text.class);
    conf.setOutputValueClass (IntWritable.class);
    conf.setMapperClass (HighestMapper.class);
    conf.setReducerClass (HighestReducer.class);
    Path inp = new Path (args [1]);
    Path out = new Path (args [2]);
    FileInputFormat.addInputPath (conf, inp);
    FileOutputFormat.setOutputPath (conf, out);
    JobClient.runJob (conf);
    return 0;
}

```

public static void main (String [] args) throws Exception

```

{
    int res = ToolRunner.run (new Configuration (), new
        HighestDriver (), args ());
    System.exit (res);
}

```

Build jar file → Add jar file using
 /usr/local/hadoop/share/hadoop/common
 /usr/local/hadoop/share/hadoop/mapreduce

Copy input files from local system to hdfs

hadoop fs - copy from local source-path (local System) Destination-path(HDFS)

hd user @ MKP: /usr/local/hadoop/sbin \$ hadoop fs - copy from
local /home/mkp/Desktop/igo1/Highest Temp/input

copy input files from local system to HDFS

Hadoop fs - copy From local source-path (local system) Destination - Path

Executing Jar file

hadoop jar jar-file-path main-class-name input-file
directory-path output-directory-path

Experiment

Install and run Pig then write Pig Latin script to sort group, join, project and filter your data.

Downloading apache PIG from

<https://download.apache.org/Pig/>

Pig releases → Please make sure you're downloading from the archives.

Extract the pig directory files as they are in zip format into

"E:/hadoop-env"

Setting up Environmental Variables

Add new environmental variable to /as

name - PIG_HOME

value E:/hadoop-env/pig-0.17.0

Add a new path variable as

%PIG_HOME% %\bin

Starting Apache Pig

Note that Hadoop must be running in background.

open cmd as administrator and run

pig - version

we received an exception, to fix this error need to change pig.cmd file located in the "pig-0.17.0\bin" directory by changing the Hadoop-Bin-PATH value from "%Hadoop-Home%\bin" to "%Hadoop-home%\libexec".

Now let's try to run the 'pig-version'

1. Writing a Pig script

Create and open an Apache Pig script file in an editor ie:- gedit

Command sudo gedit /home/user/output.pig

Writing some PIG commands inside output pig

A = Load "/user/information.txt"

using pigStorage('') as (FName:chararray, LName:chararray,
MobileNo:chararray, City:chararray, Profession:chararray);

B = ForEach A generate FName, Mobile, Profession;

Dump B;

Execute the apache PIG script

Command pig /home/user/output.pig

Experiment :-

Install and run hive then use hive to ~~or~~ create, alter and drop database, tables, views, functions and indexes.

Downloading microsoft SDK → downloading microsoft SDK v.7.1 to get the tools, compilers, headers and libraries that are necessary to run hadoop.

Cygwin : Download and install Unix command line tool Cygwin to run the unix commands on windows as per your windows.

Maven Download and install maven 3.1.1 . The installation of apache Maven is simple as extracting the files and adding bin directory to path environmental variable.

Step 1 :- Start windows 7 pc and open a new command prompt

Run → mvn -v to check installation.

Download protocol buffers 2.5.0 and extract to a folder in c drive

Type the below commands prompts to check environmental variables.

Download source and build a windows package. Download hadoop sources tarball hadoop-2.6.4-src.tar.gz and extract to a folder having short path.

use C:\hdcp to avoid runtime problem due to maximum path length limitation in windows.

Extract hadoop-2.6.0.tar.gz to a folder. Add Environment variable Hadoop-home = "C:\hdcp\hadoop-2.6.4-src" and edit

19. Open in browser

HIVE Install

1. Open cmd and write

java -version

2. Download any binary for hive from any source available on

3. Paste the HIVE in C drive, eg c:\hdparclic-hive-2.1.0-bin

4. Set environmental variables named 'HIVE_HOME' \bin

5. Install and start MySQL, if you have not.

6. download mysql-connector-java-5.0.5.jar

7. Use MySQL command line and execute

\$ Mysql -u root -p

Enter password

mysql > create Database metastore db;

8. Create a user [hive user] in MySQL database using root

user. Let's take hive user as 'userhive' and 'hivepassword' as 'hivepwd'.

mysql > CREATE USER 'userhive'@'%' IDENTIFIED BY 'hivepassword'

mysql > GRANT all on *.* to 'userhive'@localhost identified by
hivepassword.

mysql > flush privileges;

9. Go to %hive-Home%\conf folder and copy "hive-default.xml.template" file. Now rename the copied file as "hive-site.xml"

The template file has formatting needed for "hive-site.xml"

10. Edit `hive-site.xml` in `%HIVE_HOME%/.conf` directory and add the configuration.

Database commands in HIVE

`Create (Database) [If not exists] db-name`

[Comment database-comment]

[Location hdfs-path]

[With DB Properties (property-name = property-value, --);]

Example.

`hive > create database if not exists firstDB comment`

'this is first demo' location

location '/user/hive/Warehouse/newdb'

with DBProperties ('createby'='abhay', 'createfor'='dezyse');

Drop database

`Drop(Database) [IF Exist] database-name [Restrict|cascade]`

`hive > drop database if exist firstDB cascade;`

Alteration

`Alter (Database) database-name SET DBProperties
(property-name = property-value, --);`

`Alter (Database) database-name SET OWNER [USER|ROLE] user-to
role;`

Describing table

Describe [Extended /Formatted] [db-name.] table-name
[. col-name (.Field-name)]

hive > describe extended college.college-students;