

Bharathwin MA
205229105

Exercise 06: Earth Quack data analytics using MapReduce (Java & Python)

This exercise's MapReduce process is doing Earth Quack data analysis. This analysis is used to find maximum magnitude earth quack in each region. In this exercise students try to create Mapper and Reducer process using Java and Python.

Prerequisites

Ensure that Hadoop is installed, configured and is running. More details:

Single Node Setup for first-time users.

Cluster Setup for large, distributed clusters.

Inputs and Outputs

- i. **Input file should be in : /earth/in/**

WaData.txt

Copy the content text from earth.csv, Which is attached in Google classroom.

- ii. **Output file should be in /earth/out/**

Step 1

Create and Compile EarthQuack.java and create an EarthQuack.jar:

(i) Create EarthQuack.java project.

```
import org.apache.hadoop.fs.Path; import
org.apache.hadoop.io.DoubleWritable; import
org.apache.hadoop.io.Text; import
org.apache.hadoop.mapreduce.Job;

import org.apache.hadoop.mapreduce.lib.input.FileInputFormat; import
org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```

public class EarthQuake {

    public static void main(String[] args) throws Exception {
    if (args.length != 2) {

        System.err.println("Usage: hadoopex <input path> <output path>");
        System.exit(-1);
    }

    // Create the job specification object
    Job job = new Job();
    job.setJarByClass(EarthQuake.class);
    job.setJobName("Earthquake Measurment");

    // Setup input and output paths
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    // Set the Mapper and Reducer classes
    job.setMapperClass(EarthQuakeMapper.class);
    job.setReducerClass(EarthQuakeReducer.class);

    // Specify the type of output keys and values
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(DoubleWritable.class);

    // Wait for the job to finish before terminating    // Wait for the job to finish before
terminating
    System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}

```

(ii) Create EarthquakeMapper.java project.

```

import org.apache.hadoop.io.DoubleWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text; import
org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class EarthquakeMapper extends
    Mapper<LongWritable, Text, Text, DoubleWritable>
{

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

```

```

        String[] line = value.toString().split(",", 12);

        // Ignore invalid lines
        if (line.length != 12) {
            System.out.println("- " +
                line.length);          return;
        }

String[] line = value.toString().split(",", 12);

        // Ignore invalid lines
        if (line.length != 12) {
            System.out.println("- " + line.length);
            return;
        }

        // The output `key` is the name of the region
        String outputKey = line[1];

        // The output `value` is the magnitude of the earthquake
        double outputValue = Double.parseDouble(line[9]);

        // Record the output in the Context object
        context.write(new Text(outputKey), new DoubleWritable(outputValue));

    }
}

```

(iii) Create EarthquakeMapper.java project.

```

import org.apache.hadoop.io.DoubleWritable;

import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException; import

org.apache.hadoop.io.Text;

public class EarthquakeReducer extends

    Reducer<Text, DoubleWritable, Text, DoubleWritable>

{

    @Override

```

```

    public void reduce(Text key, Iterable<DoubleWritable> values,
        Context context) throws IOException, InterruptedException {

        // Standard algorithm for finding the max value

        double maxMagnitude = Double.MIN_VALUE;
        for (DoubleWritable value : values) {
            maxMagnitude = Math.max(maxMagnitude, value.get());
        }

        context.write(key, new DoubleWritable(maxMagnitude));

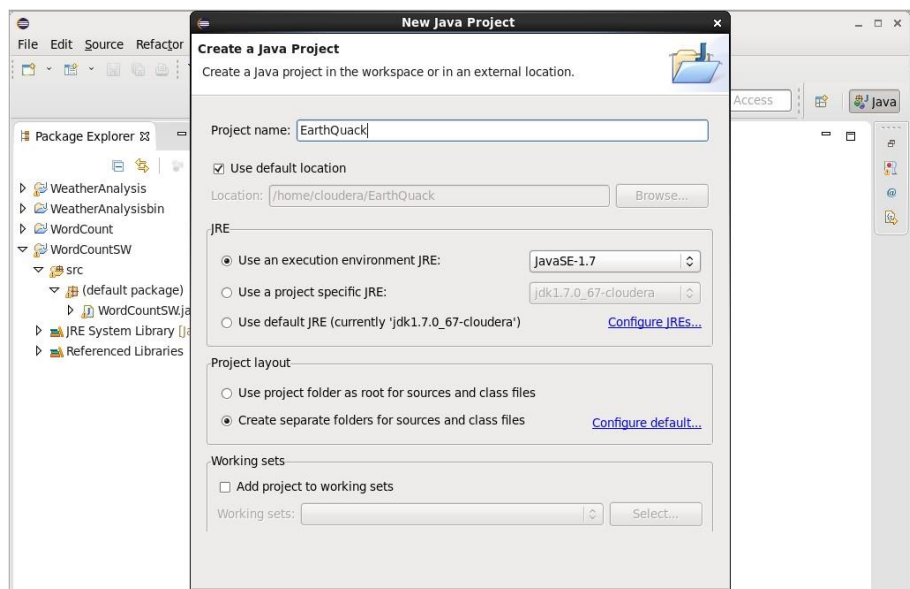
    }

    context.write(key, new DoubleWritable(maxMagnitude));

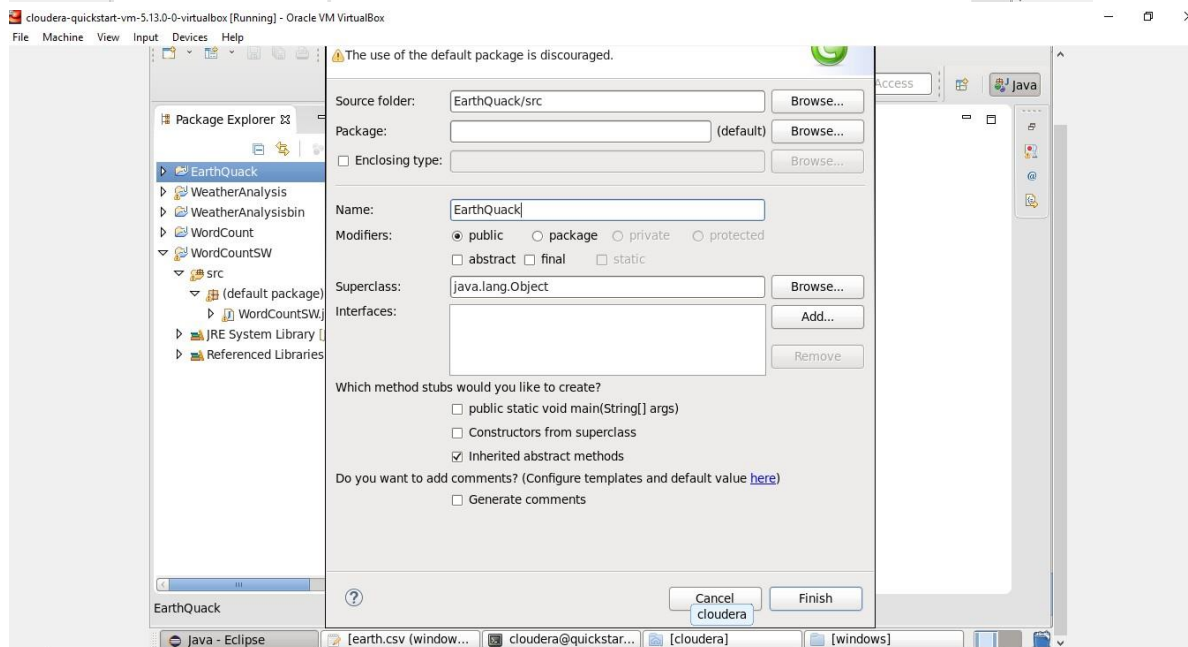
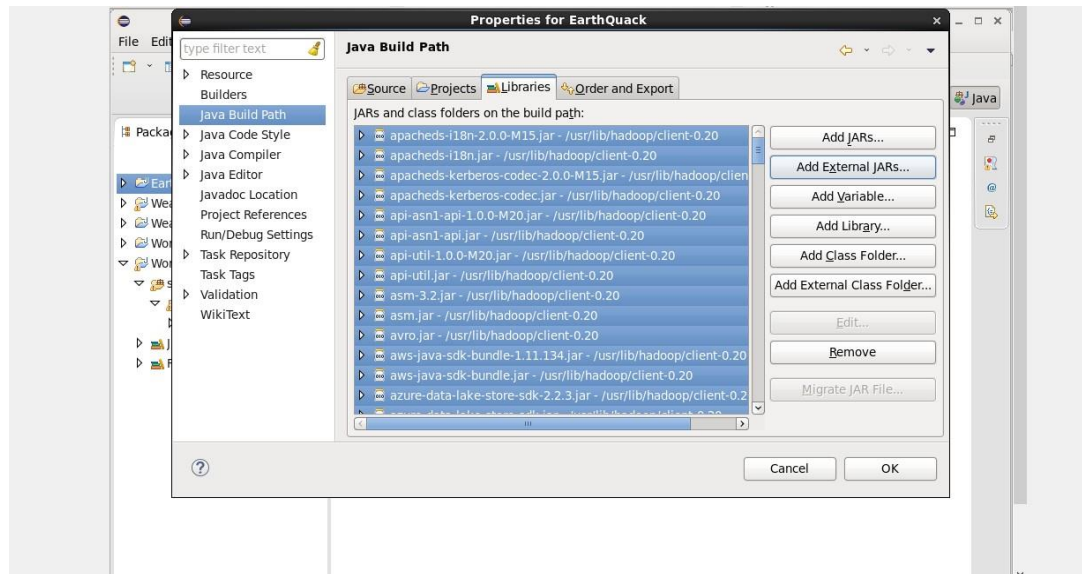
}
}

```

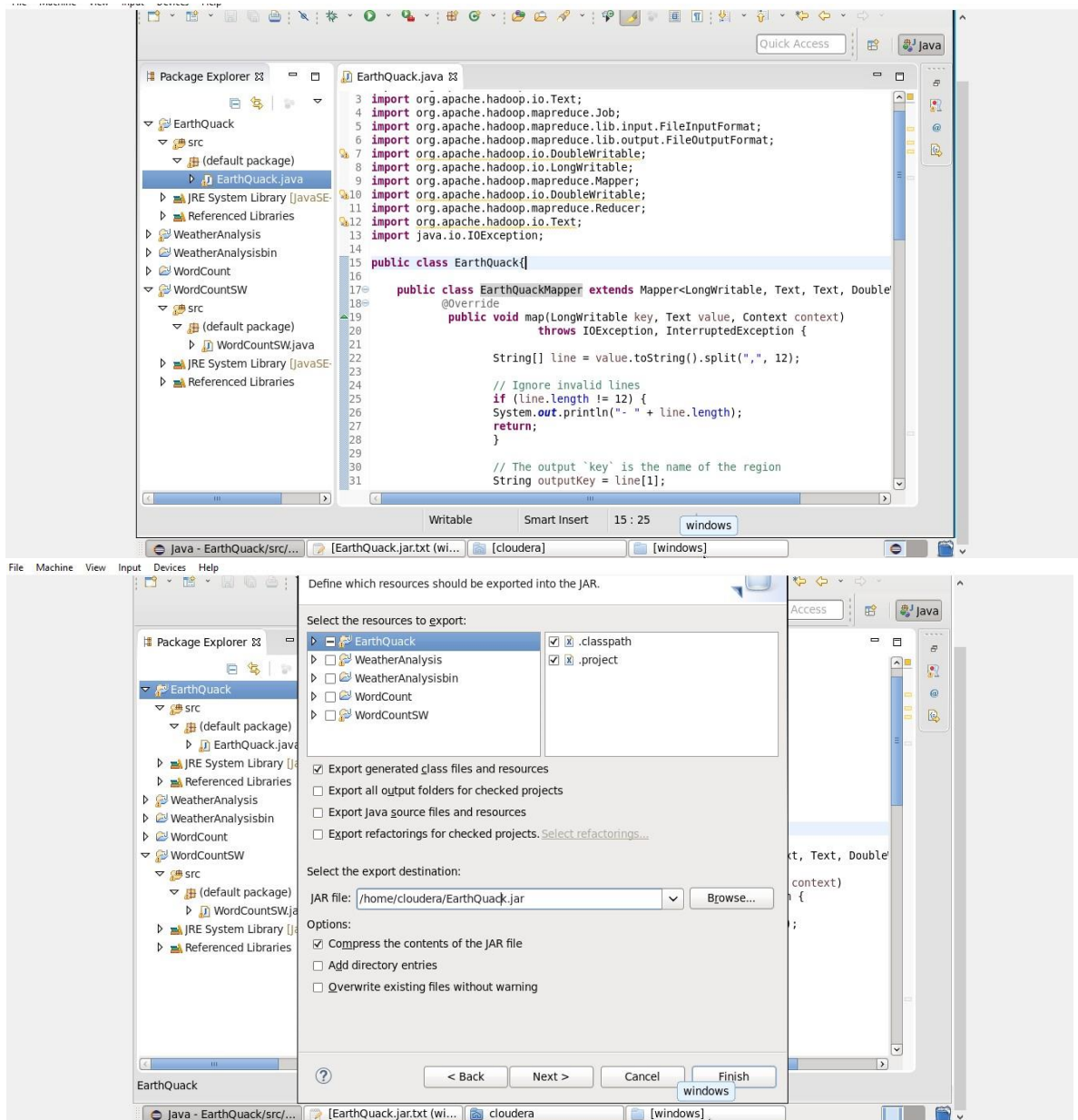
(iv) Import external .jar files



(v)



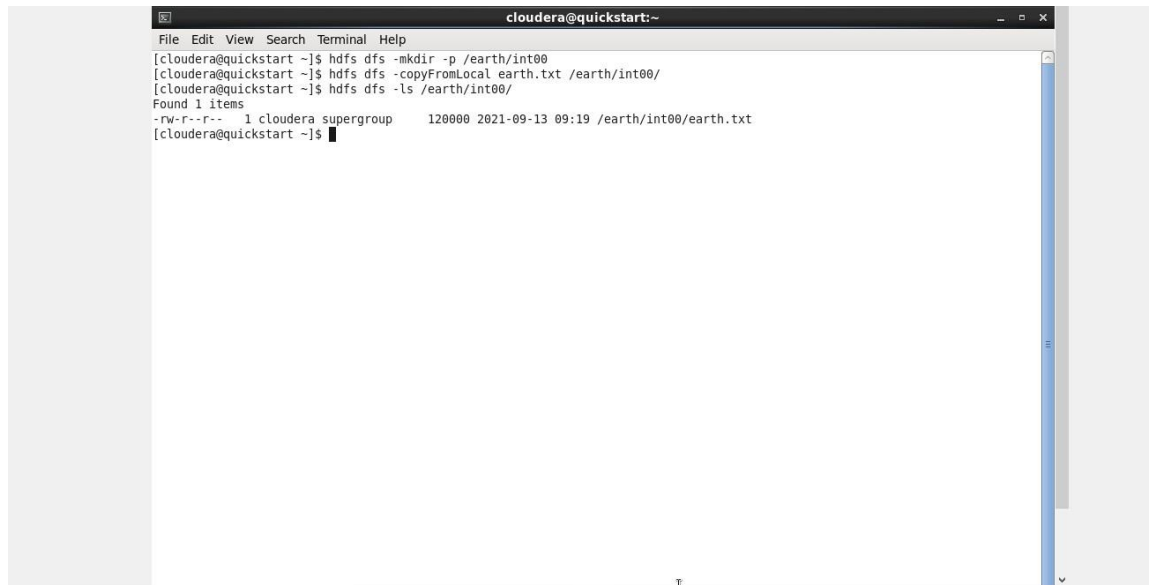
(vi) Create EarthQuake.jar file



Step 2

Create following folders in HDFS:

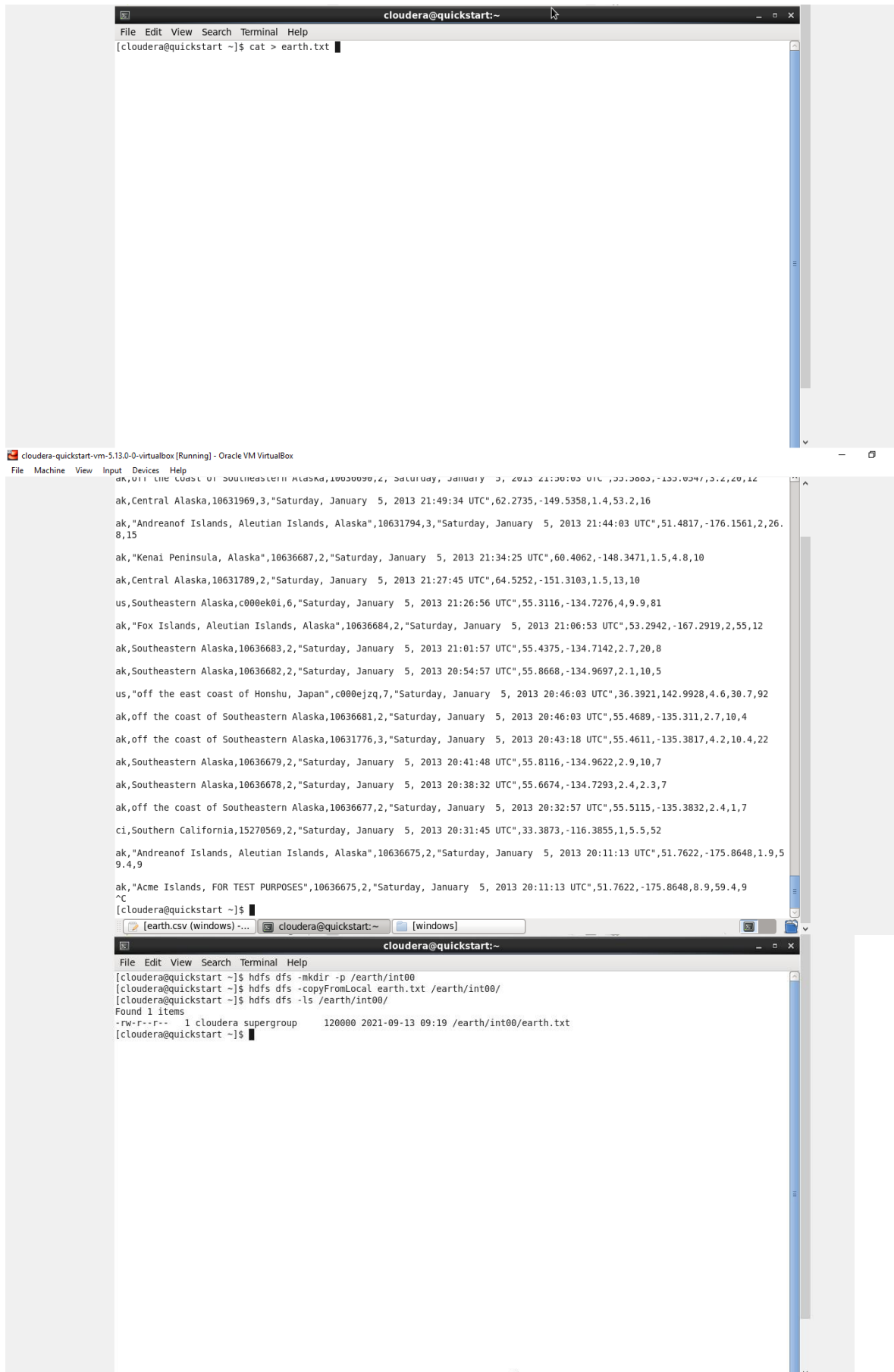
- `/earth/in` - input directory in HDFS
- `/earth/out` - output directory in HDFS

A terminal window titled 'cloudera@quickstart:~' with a menu bar (File, Edit, View, Search, Terminal, Help). The terminal shows the following commands and output:

```
[cloudera@quickstart ~]$ hdfs dfs -mkdir -p /earth/int00
[cloudera@quickstart ~]$ hdfs dfs -copyFromLocal earth.txt /earth/int00/
[cloudera@quickstart ~]$ hdfs dfs -ls /earth/int00/
Found 1 items
-rw-r--r-- 1 cloudera supergroup 120000 2021-09-13 09:19 /earth/int00/earth.txt
[cloudera@quickstart ~]$
```

Step 3

Create and copy earth.txt-files into input folder:




```
[cloudera@quickstart ~]$ hdfs dfs -ls /earth/in00/
```

Found 1 items

```
-rw-r--r-- 1 cloudera supergroup 12054 2021-08-26 15:48 /earth/in/earth.txt
```

Step 4

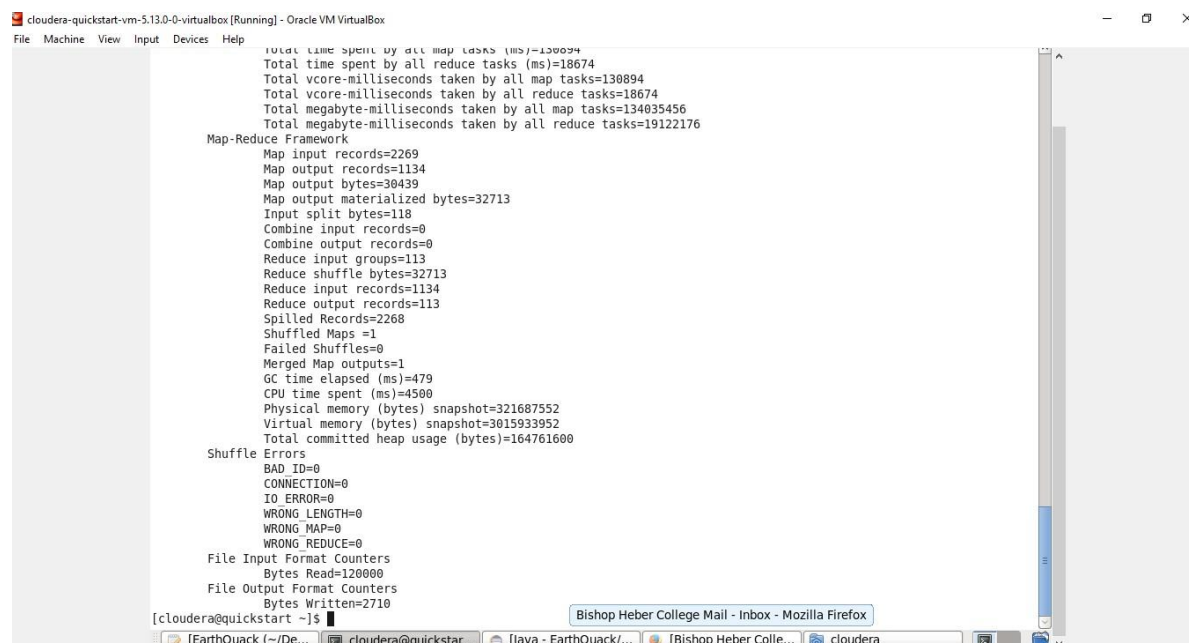
Run the MapReduce application :

Java

```
[cloudera@quickstart ~]$ hadoop jar EarthQuake.jar EarthQuake /earth/in/earth.txt /earth/out/
```

Python

```
[cloudera@quickstart ~]$ hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streaming/hadoopstreaming-2.6.0-mr1-cdh5.13.0.jar -file /home/cloudera/map.py /home/cloudera/reduce.py mapper "python map.py" -reducer "python reduce.py" -input /earth/in/earth.txt -output /earth/out Show MapReduce Framework
```



The screenshot shows a terminal window titled "cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox". The terminal displays the output of a MapReduce job. The output includes various statistics such as "Total time spent by all map tasks (ms)=130094", "Total time spent by all reduce tasks (ms)=18674", and "Map-Reduce Framework" details like "Map input records=2269", "Map output records=1134", and "Reduce input groups=113". It also shows "Shuffle Errors" and "File Input Format Counters". The terminal window has a menu bar with "File", "Machine", "View", "Input", "Devices", and "Help". The bottom of the window shows a taskbar with several open applications, including "EarthQuack", "cloudera@quickstar...", "Java - EarthQuack/...", "Bishop Heber Colle...", and "cloudera".

```
cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Total time spent by all map tasks (ms)=130094
Total time spent by all reduce tasks (ms)=18674
Total vcore-milliseconds taken by all map tasks=130094
Total vcore-milliseconds taken by all reduce tasks=18674
Total megabyte-milliseconds taken by all map tasks=134035456
Total megabyte-milliseconds taken by all reduce tasks=19122176
Map-Reduce Framework
Map input records=2269
Map output records=1134
Map output bytes=30439
Map output materialized bytes=32713
Input split bytes=118
Combine input records=0
Combine output records=0
Reduce input groups=113
Reduce shuffle bytes=32713
Reduce input records=1134
Reduce output records=113
Spilled Records=2268
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=479
CPU time spent (ms)=4500
Physical memory (bytes) snapshot=321687552
Virtual memory (bytes) snapshot=3015933952
Total committed heap usage (bytes)=164761600
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=120000
File Output Format Counters
Bytes Written=2710
[cloudera@quickstart ~]$
```

Step 5

Output:

```
[cloudera@quickstart ~]$ hdfs dfs -ls /earth/out/
```

Found 2 items

```
-rw-r--r--  1 cloudera supergroup      0 2021-08-26 15:50 /weather/out00/_SUCCESS
-rw-r--r--  1 cloudera supergroup  228 2021-08-26 15:50 /weather/out00/part-r-00000
```

cloudera-quickstart-vm-5.13.0-0-virtualbox (Running) - Oracle VM VirtualBox

File Machine View Input Devices Help

cloudera@quickstart:~

File Edit View Search Terminal Help

```
[cloudera@quickstart ~]$ hdfs dfs -ls /earth/out00/
Found 2 items
-rw-r--r-- 1 cloudera supergroup 0 2021-09-13 09:39 /earth/out00/ SUCCESS
-rw-r--r-- 1 cloudera supergroup 2710 2021-09-13 09:39 /earth/out00/part-r-00000
[cloudera@quickstart ~]$ hdfs dfs -cat /earth/out00/part-r-00000
"Acme Islands 4.9E-324
"Andaman Islands 92.3832
"Andreasof Islands 52.8796
"Anguilla region 4.9E-324
"Antofagasta 4.9E-324
"Arumachal Pradesh 94.3088
"Babuyan Islands region 121.2571
"Baja California 4.9E-324
"British Columbia 4.9E-324
"Channel Islands region 4.9E-324
"Fox Islands 54.115
"Greater Los Angeles area 4.9E-324
"Gulf of Santa Catalina 4.9E-324
"Halmahera 127.4821
"Hawaii region 4.9E-324
"Island of Hawaii 4.9E-324
"Izu Islands 141.5995
"Jujuy 4.9E-324
"Kenai Peninsula 4.9E-324
"Kodiak Island region 4.9E-324
"Lassen Peak area 4.9E-324
"Mona Passage 4.9E-324
"New Britain region 152.7111
"New Guinea 141.9851
"Newberry Caldera area 4.9E-324
"Oaxaca 4.9E-324
"Oklahoma City urban area 4.9E-324
"Olympic Peninsula 4.9E-324
"Papua 138.859
"Puget Sound region 4.9E-324
"Rat Islands 51.8203
"Salta 4.9E-324
```

cloudera-quickstart-vm-5.13.0-0-virtualbox (Running) - Oracle VM VirtualBox

File Machine View Input Devices Help

```
cloudera@quickstart:~$ hdfs dfs -ls /earth/out00/
Found 2 items
-rw-r--r-- 1 cloudera supergroup 0 2021-09-13 09:39 /earth/out00/ SUCCESS
-rw-r--r-- 1 cloudera supergroup 2710 2021-09-13 09:39 /earth/out00/part-r-00000
[cloudera@quickstart ~]$ hdfs dfs -cat /earth/out00/part-r-00000
"Acme Islands 4.9E-324
"Andaman Islands 92.3832
"Andreasof Islands 52.8796
"Anguilla region 4.9E-324
"Antofagasta 4.9E-324
"Arumachal Pradesh 94.3088
"Babuyan Islands region 121.2571
"Baja California 4.9E-324
"British Columbia 4.9E-324
"Channel Islands region 4.9E-324
"Fox Islands 54.115
"Greater Los Angeles area 4.9E-324
"Gulf of Santa Catalina 4.9E-324
"Halmahera 127.4821
"Hawaii region 4.9E-324
"Island of Hawaii 4.9E-324
"Izu Islands 141.5995
"Jujuy 4.9E-324
"Kenai Peninsula 4.9E-324
"Kodiak Island region 4.9E-324
"Lassen Peak area 4.9E-324
"Mona Passage 4.9E-324
"New Britain region 152.7111
"New Guinea 141.9851
"Newberry Caldera area 4.9E-324
"Oaxaca 4.9E-324
"Oklahoma City urban area 4.9E-324
"Olympic Peninsula 4.9E-324
"Papua 138.859
"Puget Sound region 4.9E-324
"Rat Islands 51.8203
"Salta 4.9E-324
```

cloudera-quickstart-vm-5.13.0-0-virtualbox (Running) - Oracle VM VirtualBox

File Machine View Input Devices Help

Automatically generated by Colaboratory.

```
Original file is located at
https://colab.research.google.com/drive/18CYB7k35XYbAIuwc0AlWmzJEqS_9D0a
```

```
"""
```

```
import re
import sys
```

```
for line in sys.stdin:
    val = line.strip()
    (reg, mag) = (val[1], val[9])
    print "%s\t%s" % (reg, mag)
```

```
[cloudera@quickstart Desktop]$ cat reduce.py
# -*- coding: utf-8 -*-
"""max_temperature_reduce.ipynb
```

Automatically generated by Colaboratory.

```
Original file is located at
https://colab.research.google.com/drive/18CYB7k35XYbAIuwc0AlWmzJEqS_9D0a
```

```
"""
```

```
import sys
```

```
(last key, max val) = (None, -sys.maxint)
for line in sys.stdin:
    (key, val) = line.strip().split("\t")
    if last key and last key != key:
        print "%s\t%s" % (last key, max val)
        (last key, max val) = (key, int(val))
    else:
        (last key, max val) = (key, max(max val, int(val)))
```

```
if last key:
    print "%s\t%s" % (last key, max val)
```

```
[cloudera@quickstart Desktop]$ hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streming/hadoop-streaming-2.6.0-mrl-cdh5.13
```

cloudera@quickstart:~