Ministry of Higher Education and Research
Higher School of Computer Science 08 May 1945 - Sidi Bel Abbes

Second Year Second Cycle - Artificial Intelligence and Data Science

# LAB 01 : HADOOP

Presented By : FELLAH Abdelnour.

Date: February 17, 2024

# 1 INTRODUCTION

In this Big Data Lab, we will explore Apache Hadoop and MapReduce basics.where we'll understand managing vast data efficiently. In today's digital era, scalable solutions like Hadoop are essential for processing massive datasets. MapReduce complements Hadoop by providing a powerful processing model. We'll dive into setting up clusters, designing MapReduce jobs.

# 2 PART 01 : SETUP

**1- Downloading the spark hadoop image :**

Downloading the spark hadoop docker image can be acheived by either searching for the image and downloading it using docker desktop or running the following unix command :

```
# docker pull liliasfaxi/spark-hadoop:hv-2.7.2
```

**2- Creating the bridge network and the docker containers :**
to create the bridge network that connects the three containers,we can use the following command :

```
#  docker network create --driver=bridge hadoop
```

and to create the three container with the specified ports,we run the following command :

```
#  docker run -itd --net=hadoop -p 50070:50070 -p 8088:8088 -p 7077:7077 -p 16010:16010  --name
hadoop-master --hostname hadoop-master  liliasfaxi/spark-hadoop:hv-2.7.2

#  docker run -itd -p 8040:8042 --net=hadoop  --name hadoop-slave1
--hostname hadoop-slave1  liliasfaxi/spark-hadoop:hv-2.7.2

#  docker run -itd -p 8041:8042 --net=hadoop  --name hadoop-slave2
--hostname hadoop-slave2  liliasfaxi/spark-hadoop:hv-2.7.2
```

**3- Entering the master container :**

This can be done using the following commands :
```
# docker exec -it hadoop-master bash
root@hadoop-master:#./start-hadoop.sh
```

# 3 PART 02 : FIRST STEPS WITH HADOOP :

**1- Create a directory named input :**

```
root@hadoop-master:#hadoop fs -mkdir -p input
```

**2- Load the puchases file into recently created input directoty :**

```
root@hadoop-master:#hadoop fs -put purchases.txt input
```

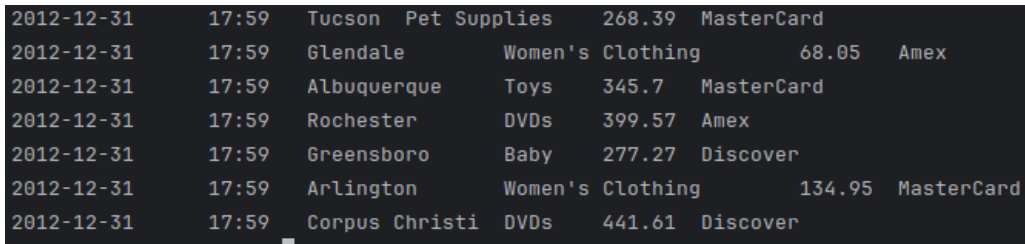**3- show the content of the directory input :**

```
root@hadoop-master:#hadoop fs -ls input
```



Figure 1: The display of the content of the directory input

**4- show the tail of the purchases.txt file :**

```
root@hadoop-master:# hadoop fs -tail input/purchases.txt
```

```
2012-12-31    17:59    Tucson  Pet Supplies    268.39  MasterCard
2012-12-31    17:59    Glendale        Women's Clothing       68.05    Amex
2012-12-31    17:59    Albuquerque     Toys    345.7   MasterCard
2012-12-31    17:59    Rochester       DVDs    399.57  Amex
2012-12-31    17:59    Greensboro      Baby    277.27  Discover
2012-12-31    17:59    Arlington       Women's Clothing      134.95   MasterCard
2012-12-31    17:59    Corpus Christi  DVDs    441.61  Discover
```

Figure 2: The display of the couple last rows of the file purchases.txt

# 4 PART 03 : MAP REDUCE

A map reduce job is mainly composed of tow main parts,the mappers and the reducers,the mappers are responsible about returning the key-value pairs,the reducers that are responsible about sorting and aggregating the results.

We are going to test a MapReduce program using a very simple example, WordCount that calculates the number of occurences for each unique word in the text by following the following steps :

- First we create a new maven project,make sure you include the right java target (8) and the necessary dependencies,see the pom.xml file bellow for more details.

- Create three classes Main class (driver),the WordsCountMap class and the WordsCountReducer class.

- We then implement the code for each class (see below for more details).

```xml
<?xml version="1.0" encoding="UTF-8"?>
<project xmlns="http://maven.apache.org/POM/4.0.0"
        xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
        xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org/xsd
        /maven-4.0.0.xsd">
    <modelVersion>4.0.0</modelVersion>

    <groupId>org.example</groupId>
    <artifactId>TP2-BDT</artifactId>
    <version>1.0-SNAPSHOT</version>

    <properties>
        <maven.compiler.source>8</maven.compiler.source>
        <maven.compiler.target>8</maven.compiler.target>
        <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
    </properties>
    <dependencies>
        <dependency>
            <groupId>org.apache.hadoop</groupId>
            <artifactId>hadoop-common</artifactId>
            <version>2.7.2</version>
        </dependency>
        <!-- https://mvnrepository.com/artifact/org.apache.hadoop/hadoop-mapreduce-client-co
        <dependency>
            <groupId>org.apache.hadoop</groupId>
            <artifactId>hadoop-mapreduce-client-core</artifactId>
            <version>2.7.2</version>
        </dependency>
        <!-- https://mvnrepository.com/artifact/org.apache.hadoop/hadoop-hdfs -->
        <dependency>
            <groupId>org.apache.hadoop</groupId>
            <artifactId>hadoop-hdfs</artifactId>
            <version>2.7.2</version>
        </dependency>
        <dependency>
            <groupId>org.apache.hadoop</groupId>
            <artifactId>hadoop-mapreduce-client-common</artifactId>
            <version>2.7.2</version>
        </dependency>
    </dependencies>
</project>
```

Figure 3: the content of the pom.xml file.

```
1  package org.example;
2
3
4  import org.apache.hadoop.fs.Path;
5  import org.apache.hadoop.mapreduce.Job;
6  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
7  import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
8  import org.apache.hadoop.conf.Configuration;
9  import org.apache.hadoop.util.GenericOptionsParser;
10 import org.apache.hadoop.io.Text;
11 import org.apache.hadoop.io.IntWritable;
12
13 import java.io.IOException;
14
15 public class Main {
16     public static void main(String[] args) throws IOException, InterruptedException,
           ClassNotFoundException {
17
18         Configuration conf = new Configuration();
19         String[] remainingArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
20         Job job = Job.getInstance(conf, "Words Counter v1.0");
21
22         job.setJarByClass(Main.class);
23         job.setMapperClass(WordsCountMap.class);
24         job.setReducerClass(WordsCountReduce.class);
25
26         job.setOutputKeyClass(Text.class);
27         job.setOutputValueClass(IntWritable.class);
28
29         FileInputFormat.addInputPath(job, new Path(remainingArgs[0]));
30         FileOutputFormat.setOutputPath(job, new Path(remainingArgs[1]));
31
32         if (job.waitForCompletion(true)) {
33             System.exit(0);
34         }
35
36         System.exit(-1);
37
38     }
39 }
```

Figure 4: The Main class code

```
1  package org.example;
2
3  import org.apache.hadoop.mapreduce.Job;
4  import org.apache.hadoop.io.Text;
5  import org.apache.hadoop.io.IntWritable;
6  import java.util.StringTokenizer;
7  import org.apache.hadoop.mapreduce.Mapper;
8  import java.io.IOException;
9
10 public class WordsCountMap extends Mapper<Object, Text, Text, IntWritable> {
11
12     private static final IntWritable ONE = new IntWritable(1);
13
14     @Override
15     protected void map(Object key, Text value, Mapper<Object, Text, Text, IntWritable>.
           Context context) throws IOException, InterruptedException {
16
17         // create a word tokenizer
18         StringTokenizer tokenizer = new StringTokenizer(value.toString());
19
20         while (tokenizer.hasMoreTokens()) {
21             Text word = new Text(tokenizer.nextToken());
22             context.write(word, ONE);
23         }
24     }
25 }
```

Figure 5: The code for the WordsCountMap class

```
1  package org.example;
2
3  import org.apache.hadoop.io.Text;
4  import org.apache.hadoop.io.IntWritable;
5  import org.apache.hadoop.mapreduce.Reducer;
6  import java.util.Iterator;
7  import java.io.IOException;
8
9  public class WordsCountReduce extends Reducer<Text,IntWritable,Text,Text> {
10
11     @Override
12     protected void reduce(Text key, Iterable<IntWritable> values, Reducer<Text,
          IntWritable, Text, Text>.Context context) throws IOException, InterruptedException
           {
13
14        Iterator<IntWritable> i = values.iterator();
15        int count = 0;
16
17        while (i.hasNext()) {
18           count += i.next().get();
19        }
20
21        context.write(key, new Text(count + " occurences"));
22
23     }
24  }
```

Figure 6: The code for the WordsCountReduce class

Now that we've implemented the different required class we now generate a .jar file using maven.

We then copy the generated .jar file to the master node with the following command :

```
# docker cp [file].jar hadoop-master:/root/[file].jar
```

In our case the command is :

```
$ docker cp target/TP2-BDT-1.0-SNAPSHOT.jar hadoop-master:/root/TP-2BDT-1.0-SNAPSHOT.jar
```

Now we can use the run the .jar file and give it as arguments the main class,the input folder and the output folder,the command is :

```
root@hadoop-master:#hadoop jar [file].jar [Main-CLass] [input-folder] [output-folder]
```
In our case the command is :

```
root@hadoop-master:~# hadoop jar TP-2BDT-1.0-SNAPSHOT.jar org.example.Main input output
```

the output of this command yields some statistics and information about the performed job.

Figure 7: The oputput of the terminal after running the .jar file

to display the last rows of the generated file we run the command :

```
root@hadoop-master:#hadoop fs -tail output/part-r-00000
```



Figure 8: A sample of the content of the output file.

# 5 Part 04 : Python version

to implement MapReduce using python,follow the steps bellow :

Step 01 : create tow python files (mapper.py and reducer.py),with the following code :

```python
#!/usr/bin/env python
"""mapper.py"""

import sys

# input comes from STDIN (standard input)
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # split the line into words
    words = line.split()
    # increase counters
    for word in words:
        # write the results to STDOUT (standard output);
        # what we output here will be the input for the
        # Reduce step, i.e. the input for reducer.py
        #
        # tab-delimited; the trivial word count is 1
        print("%s\t%s" % (word, 1))
```

Figure 9: The code for the mapper.py file

```python
#!/usr/bin/env python
"""reducer.py"""

import sys

current_word = None
current_count = 0
word = None

# input comes from STDIN
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()

    # parse the input we got from mapper.py
    word, count = line.split("\t", 1)

    # convert count (currently a string) to int
    try:
        count = int(count)
    except ValueError:
        # count was not a number, so silently
        # ignore/discard this line
        continue

    # this IF-switch only works because Hadoop sorts map output
    # by key (here: word) before it is passed to the reducer
    if current_word == word:
        current_count += count
    else:
        if current_word:
            # write result to STDOUT
            print("%s\t%s" % (current_word, current_count))
        current_count = count
        current_word = word

# do not forget to output the last word if needed!
if current_word == word:
    print("%s\t%s" % (current_word, current_count))
```

Figure 10: The code for the reducer.py file

Step 02 : Download hadoop-streaming JAR 2.7.3 from this link.

Step 03 : Move the downloaded jar file,and the tow python programs to the master container,using the following commands :

```
# docker cp /path/to/hadoop-streaming-2.7.3.jar hadoop-master:/root/hadoop-streaming-2.7.3.jar

# docker cp /path/to/mapper.py hadoop-master:/root/mapper.py

# docker cp /path/to/reducer.py hadoop-master:/root/reducer.py
```

Step 04: Change the permissions of the mapper.py and reducer.py so they are executable by the user.
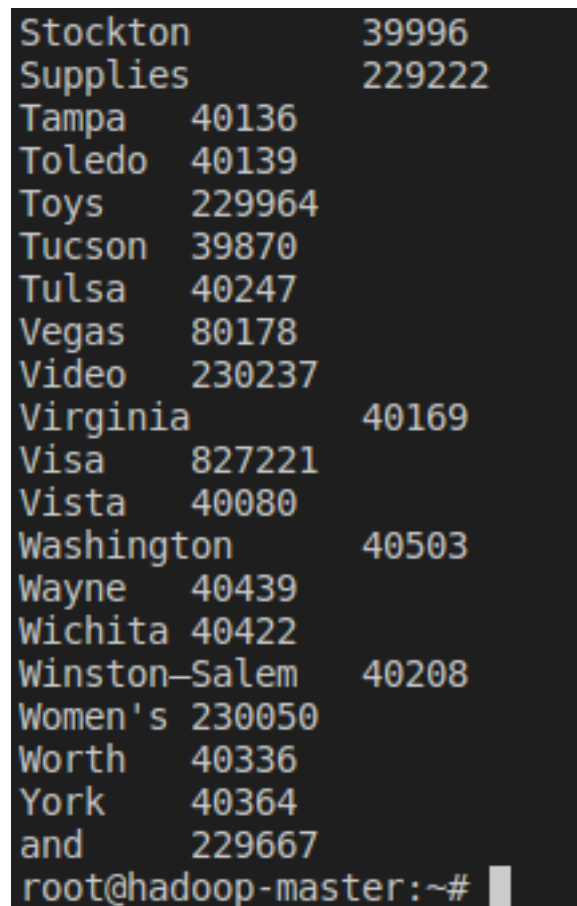
```
root@hadoop-master:#chmod 744 mapper.py reducer.py
```

Step 05 : start the MapReduce job with the following command.

```
root@hadoop-master:#hadoop jar hadoop-streaming-2.7.3.jar
-file ./mapper.py -mapper "python3 mapper.py"
-file ./reducer.py -reducer "python3 reducer.py"
-input input/purchases.txt  -output py-output
```
Now you can check out the result of this command :

```
root@hadoop-master:#hadoop fs -tail py-output/part-00000
```



Figure 11: A sample of the content of the output file (python version).

# 6 CONCLUSION

Our Big Data Lab journey has been enlightening and empowering. We've explored Hadoop and MapReduce hands-on, gaining insights into distributed computing's transformative power. With Hadoop and MapReduce, we can handle complex data tasks confidently, driving positive change. As we conclude, let's carry forward our newfound skills to make an impact. In business, science, or public services, our expertise positions us at the forefront of the data-driven revolution. We're ready to shape the future of big data analytics with innovation and meaningful contributions.