**Hadoop Mapreduce**

**Case study: Analysis of the sale of certain products in the world**
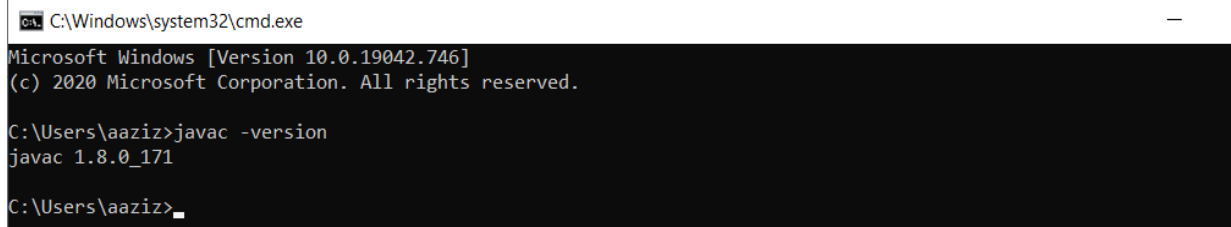
In this project, I will use Hadoop with MapReduce. The input data that we are going to use is a sale products dataset calls SalesJa2009.csv that we can find on kaggle website by clicking on this link: https://www.kaggle.com/jensroderus/salesjan2009csv. It contains Sales related information like product name, price, payment mode, city, country of client etc. The goal is to **find out the Number of Products Sold in Each Country.**

| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Transaction_c | Product | Price | Payment_Typ | Name | City | State | Country | Account_Crea | Last_Login | Latitude | Longitude |
| 2 | ######### | Product1 | 1200 | Mastercard | carolina | Basildon | England | United Kingdc | ######### | ######### | 51.5 | -1.1166667 |
| 3 | ######### | Product1 | 1200 | Visa | Betina | Parkville | MO | United States | ######### | ######### | 39.195 | -94.68194 |
| 4 | ######### | Product1 | 1200 | Mastercard | Federica e An | Astoria | OR | United States | ######### | ######### | 46.18806 | -123.83 |
| 5 | ######### | Product1 | 1200 | Visa | Gouya | Echuca | Victoria | Australia | ######### | ######### | -36.1333333 | 144.75 |
| 6 | ######### | Product2 | 3600 | Visa | Gerd W | Cahaba Heigh | AL | United States | ######### | ######### | 33.52056 | -86.8025 |
| 7 | ######### | Product1 | 1200 | Visa | LAURENCE | Mickleton | NJ | United States | ######### | ######### | 39.79 | -75.23806 |
| 8 | ######### | Product1 | 1200 | Mastercard | Fleur | Peoria | IL | United States | ######### | ######### | 40.69361 | -89.58889 |
| 9 | ######### | Product1 | 1200 | Mastercard | adam | Martin | TN | United States | ######### | ######### | 36.34333 | -88.85028 |
| 10 | ######### | Product1 | 1200 | Mastercard | Renee Elisabe | Tel Aviv | Tel Aviv | Israel | ######### | ######### | 32.0666667 | 34.7666667 |
| 11 | ######### | Product1 | 1200 | Visa | Aidan | Chatou | Ile-de-France | France | ######### | ######### | 48.8833333 | 2.15 |
| 12 | ######### | Product1 | 1200 | Diners | Stacy | New York | NY | United States | ######### | ######### | 40.71417 | -74.00639 |
| 13 | ######### | Product1 | 1200 | Amex | Heidi | Eindhoven | Noord-Braban | Netherlands | ######### | ######### | 51.45 | 5.4666667 |
| 14 | ######### | Product1 | 1200 | Mastercard | Sean | Shavano Park | TX | United States | ######### | ######### | 29.42389 | -98.49333 |
| 15 | ######### | Product1 | 1200 | Visa | Georgia | Eagle | ID | United States | ######### | ######### | 43.69556 | -116.35306 |
| 16 | ######### | Product1 | 1200 | Visa | Richard | Riverside | NJ | United States | ######### | ######### | 40.03222 | -74.95778 |
| 17 | ######### | Product1 | 1200 | Diners | Leanne | Julianstown | Meath | Ireland | ######### | ######### | 53.6772222 | -6.3191667 |
| 18 | ######### | Product1 | 1200 | Visa | Janet | Ottawa | Ontario | Canada | ######### | ######### | 45.4166667 | -75.7 |
| 19 | ######### | Product1 | 1200 | Diners | barbara | Hyderabad | Andhra Prade: | India | ######### | ######### | 17.3833333 | 78.4666667 |
| 20 | ######### | Product2 | 3600 | Visa | Sabine | London | England | United Kingdc | ######### | ######### | 51.52721 | 0.14559 |
| 21 | ######### | Product1 | 1200 | Diners | Hani | Salt Lake City | UT | United States | ######### | ######### | 40.76083 | -111.89028 |
| 22 | ######### | Product1 | 1200 | Visa | Jeremy | Manchester | England | United Kingdc | ######### | ######### | 53.5 | -2.2166667 |
| 23 | ######### | Product1 | 1200 | Diners | Janis | Ballynora | Cork | Ireland | ######### | ######### | 51.8630556 | -8.58 |
| 24 | ######### | Product1 | 1200 | Mastercard | Nicola | Roodepoort | Gauteng | South Africa | ######### | ######### | -26.1666667 | 27.8666667 |
| 25 | ######### | Product1 | 1200 | Visa | asuman | Chula Vista | CA | United States | ######### | ######### | 32.64 | -117.08333 |
| 26 | ######### | Product1 | 1200 | Mastercard | Lena | Kuopio | Ita-Suomen L | Finland | ######### | ######### | 62.9 | 27.6833333 |
| 27 | ######### | Product1 | 1200 | Visa | Lisa | Sugar Land | TX | United States | ######### | ######### | 29.61944 | -95.63472 |
| 28 | ######### | Product1 | 1200 | Diners | Bryan Kerrene | New York | NY | United States | ######### | ######### | 40.71417 | -74.00639 |
| 29 | ######### | Product1 | 1200 | Visa | chris | London | England | United Kingdc | ######### | ######### | 51.52721 | 0.14559 |
| 30 | ######### | Product1 | 1200 | Visa | Maxine | Morton | IL | United States | ######### | ######### | 40.61278 | -89.45917 |
| 31 | ######### | Product1 | 1200 | Visa | Family | Los Gatos | CA | United States | ######### | ######### | 37.22667 | -121.97361 |
| 32 | ######### | Product1 | 1200 | Mastercard | Katherine | New York | NY | United States | ######### | ######### | 40.71417 | -74.00639 |

SalesJan2009

First, we will start by installing the hadoop framework in our Windows 10 machine, then configure it before starting to write our java mapreduce code for our case study.

## I- Installation of Hadoop

We assume that java version 1.8 is already installed on our computer. In our case, version 1.8.0_171 is installed in our machine.
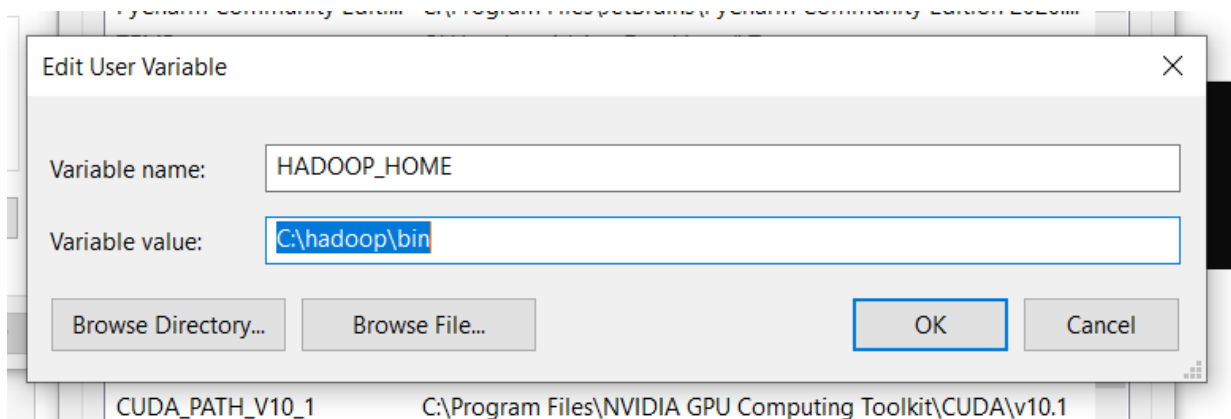


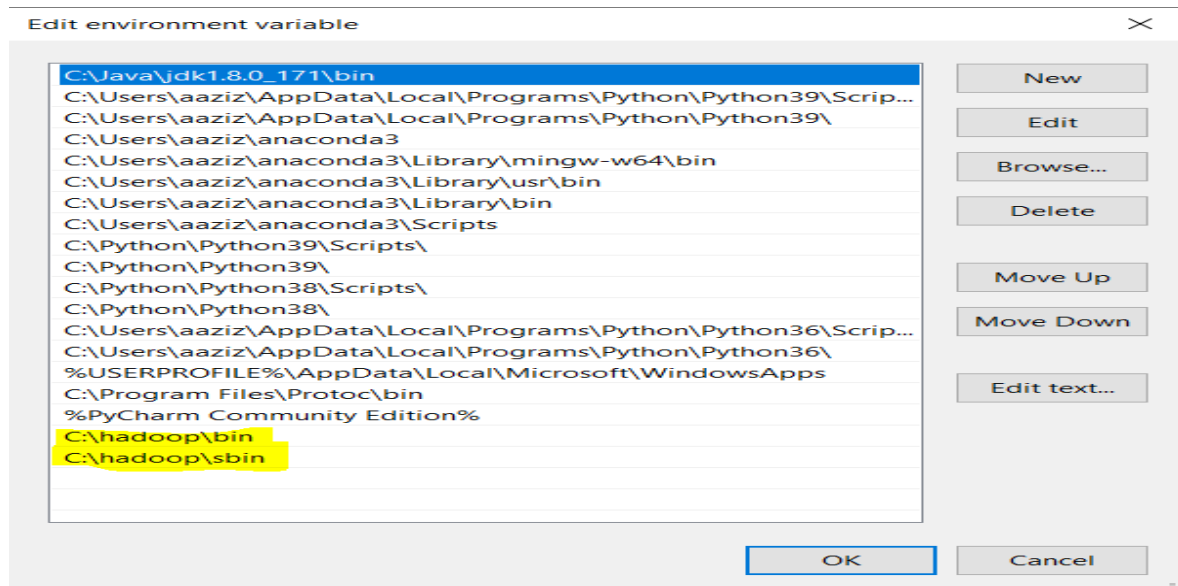To install Hadoop in our windows machine, we need to follow the following steps:

**Step 1:** Download Hadoop binary package from apache website (https://hadoop.apache.org/releases.html)

In our case we have downloaded hadoop version 2.7.6

**Step 2:** unpack the package, copy the folder in our C directory and configure the path in our environment variable.

**Step 3:** modification of core-site.xml, mapred-site.xml, hdfs-site.xml and yarn-site.xml which are in C: //hadoop/etc/hadoop.

```
new 1.txt ☒   hdfs-site.xml ☒   yarn-site.xml ☒   core-site.xml ☒   mapred-site.xml ☒

 1   <?xml version="1.0"?>
 2   <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
 3   <!--
 4      Licensed under the Apache License, Version 2.0 (the "Licer
 5      you may not use this file except in compliance with the Li
 6      You may obtain a copy of the License at
 7
 8        http://www.apache.org/licenses/LICENSE-2.0
 9
10      Unless required by applicable law or agreed to in writing,
11      distributed under the License is distributed on an "AS IS"
12      WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either expre
13      See the License for the specific language governing permis
14      limitations under the License. See accompanying LICENSE fi
15   -->
16
17   <!-- Put site-specific property overrides in this file. -->
18
19   <configuration>
20       <property>
21           <name>mapreduce.framework.name</name>
22           <value>yarn</value>
23       </property>
24   </configuration>
25
```

Before modifying hdfs-site.xml and yarn-site.xml, we must create a folder in C: //hadoop calls data and create two more folders (datanode and namenode).

```
new 1.txt    hdfs-site.xml    yarn-site.xml    core-site.xml    mapred-site.xml
 1   <?xml version="1.0" encoding="UTF-8"?>
 2   <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
 3   <!--
 4      Licensed under the Apache License, Version 2.0 (the "Lice
 5      you may not use this file except in compliance with the L
 6      You may obtain a copy of the License at
 7
 8         http://www.apache.org/licenses/LICENSE-2.0
 9
10      Unless required by applicable law or agreed to in writing
11      distributed under the License is distributed on an "AS IS
12      WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either expr
13      See the License for the specific language governing permi
14      limitations under the License. See accompanying LICENSE f
15   -->
16
17   <!-- Put site-specific property overrides in this file. -->
18
19   <configuration>
20       <property>
21           <name>dfs.replication</name>
22           <value>1</value>
23       </property>
24       <property>
25           <name>dfs.namenode.name.dir</name>
26           <value>C:\hadoop\data\namenode</value>
27       </property>
28       <property>
29           <name>dfs.datanode.data.dir</name>
30           <value>C:\hadoop\data\datanode</value>
31       </property>
32   </configuration>
33
```

```
new 1.txt    hdfs-site.xml    yarn-site.xml    core-site.xml    mapred-site.xml
 7         http://www.apache.org/licenses/LICENSE-2.0
 8
 9      Unless required by applicable law or agreed to in writing, software
10      distributed under the License is distributed on an "AS IS" BASIS,
11      WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
12      See the License for the specific language governing permissions and
13      limitations under the License. See accompanying LICENSE file.
14   -->
15   <configuration>
16   <!-- Site specific YARN configuration properties -->
17       <property>
18           <name>yarn.nodemanager.aux-services</name>
19           <value>mapreduce_shuffle</value>
20       </property>
21       <property>
22           <name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
23           <value>org.apache.hadoop.mapred.ShuffleHandler</value>
24       </property>
25   </configuration>
26
```

**Step 4:** Format the namenode and start all daemons

To format the namenode, we will use the following command: **hdfs namenode -format**

Our namenode has started successfully. We are now going to start our daemons which are: namenode, datanode, resourcemanager and nodemanager. We will use the following command: **start-all.cmd**



The **jps** command allows us to check the demons that are started on our machine.

We can see that all of our demons are cast. Now that we're done with that, let's see our cluster and our namenode in localhost. The address page for our is **localhost:8088** and for our namenode page is **localhost:50070**

Now we are going to set up a folder name Abdel_Aziz_KAMO_MEGNA (which is my name) in our hadoop space by using command line: **hadoop fs -mkdir /Abdel_Aziz_KAMO_MEGNA**

## Browse Directory

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| drwxr-xr-x | aaziz | supergroup | 0 B | 4/29/2021, 1:54:05 AM | 0 | 0 B | Abdel_Aziz_KAMO_MEGNA |

Hadoop, 2018.

## II-    Map, Reduce and Driver code for our case study

### 1- SalesMapper.java

```
SalesReducer.java    SalesDriver.java    SalesMapper.java ⊠

 1 package sales;
 2
 3 import java.io.IOException;
 4
 5 import org.apache.hadoop.io.IntWritable;
 6 import org.apache.hadoop.io.LongWritable;
 7 import org.apache.hadoop.io.Text;
 8 import org.apache.hadoop.mapred.*;
 9
10 public class SalesMapper extends MapReduceBase implements Mapper<LongWritable, Text, Text, IntWritable> {
11     private final static IntWritable one = new IntWritable(1);
12
13     public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
14
15         String valueString = value.toString();
16         String[] SingleCountryData = valueString.split(",");
17         output.collect(new Text(SingleCountryData[7]), one);
18     }
19 }
20
```

## 2- SalesReduce.java

SalesReducer.java    SalesDriver.java    SalesMapper.java

```java
1 package sales;
2
3 import java.io.IOException;
4 import java.util.*;
5
6 import org.apache.hadoop.io.IntWritable;
7 import org.apache.hadoop.io.Text;
8 import org.apache.hadoop.mapred.*;
9
10 public class SalesReducer extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
11
12     public void reduce(Text t_key, Iterator<IntWritable> values, OutputCollector<Text,IntWritable> output, Reporter reporter) throws IOException {
13         Text key = t_key;
14         int frequencyForCountry = 0;
15         while (values.hasNext()) {
16             // replace type of value with the actual type of our value
17             IntWritable value = (IntWritable) values.next();
18             frequencyForCountry += value.get();
19
20         }
21         output.collect(key, new IntWritable(frequencyForCountry));
22     }
23 }
24
```
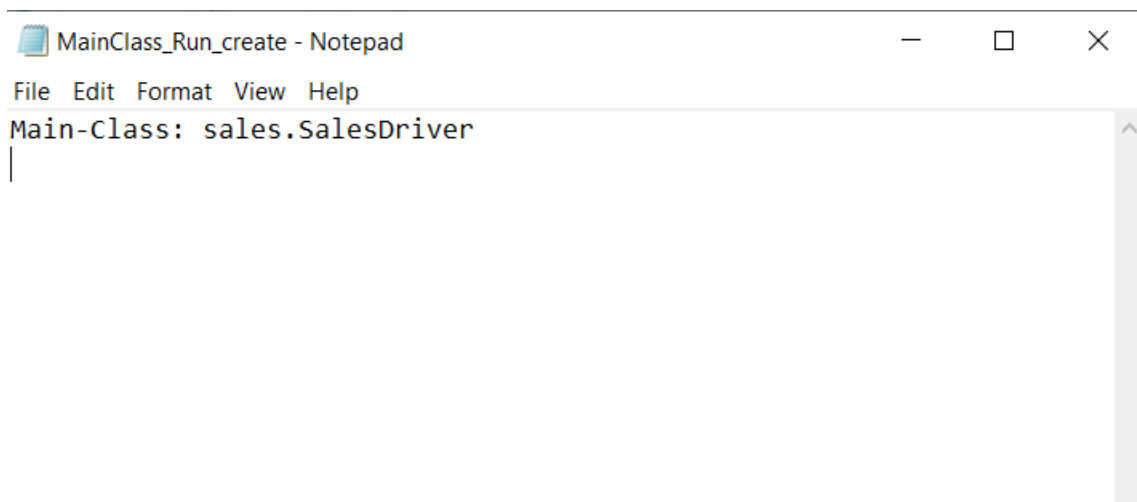
## 3- SalesDriver.java

SalesReducer.java    SalesDriver.java    SalesMapper.java

```java
1 package sales;
2
3 import org.apache.hadoop.fs.Path;
4 import org.apache.hadoop.io.*;
5 import org.apache.hadoop.mapred.*;
6
7 public class SalesDriver {
8     public static void main(String[] args) {
9         JobClient my_client = new JobClient();
10         // Create a configuration object for the job
11         JobConf job_conf = new JobConf(SalesDriver.class);
12
13         // Set a name of the Job
14         job_conf.setJobName("SalePerCountry");
15
16         // Specify data type of output key and value
17         job_conf.setOutputKeyClass(Text.class);
18         job_conf.setOutputValueClass(IntWritable.class);
19
20         // Specify names of Mapper and Reducer Class
21         job_conf.setMapperClass(sales.SalesMapper.class);
22         job_conf.setReducerClass(sales.SalesReducer.class);
23
24         // Specify formats of the data type of Input and output
25         job_conf.setInputFormat(TextInputFormat.class);
26         job_conf.setOutputFormat(TextOutputFormat.class);
27
28         // Set input and output directories using command line arguments,
29         //arg[0] = name of input directory on HDFS, and arg[1] =  name of output directory to be created to store the output file.
30
31         FileInputFormat.setInputPaths(job_conf, new Path(args[0]));
32         FileOutputFormat.setOutputPath(job_conf, new Path(args[1]));
33
34         my_client.setConf(job_conf);
35         try {
36             // Run the job
37             JobClient.runJob(job_conf);
38         } catch (Exception e) {
39             e.printStackTrace();
40         }
41     }
42 }
```

Once we have compile of our java files, the files: **SalesMapper.class, SalesReduce.class and SaleDriver.class** are automatically created.



| | | | | |
|---|---|---|---|---|
| SalesDriver.class | ⊘ | 4/26/2021 2:52 PM | CLASS File | 2 KB |
| SalesDriver | ⊘ | 4/26/2021 2:52 PM | JAVA File | 2 KB |
| SalesMapper.class | ⊘ | 4/26/2021 11:48 AM | CLASS File | 3 KB |
| SalesMapper | ⊘ | 4/26/2021 11:47 AM | JAVA File | 1 KB |
| SalesReducer.class | ⊘ | 4/26/2021 11:50 AM | CLASS File | 3 KB |
| SalesReducer | ⊘ | 4/26/2021 11:50 AM | JAVA File | 1 KB |

### 4- Specification of our main class

We now need to specify our main class. This will be done by creating new text file (MainClass_run_create) and we will specify it inside.


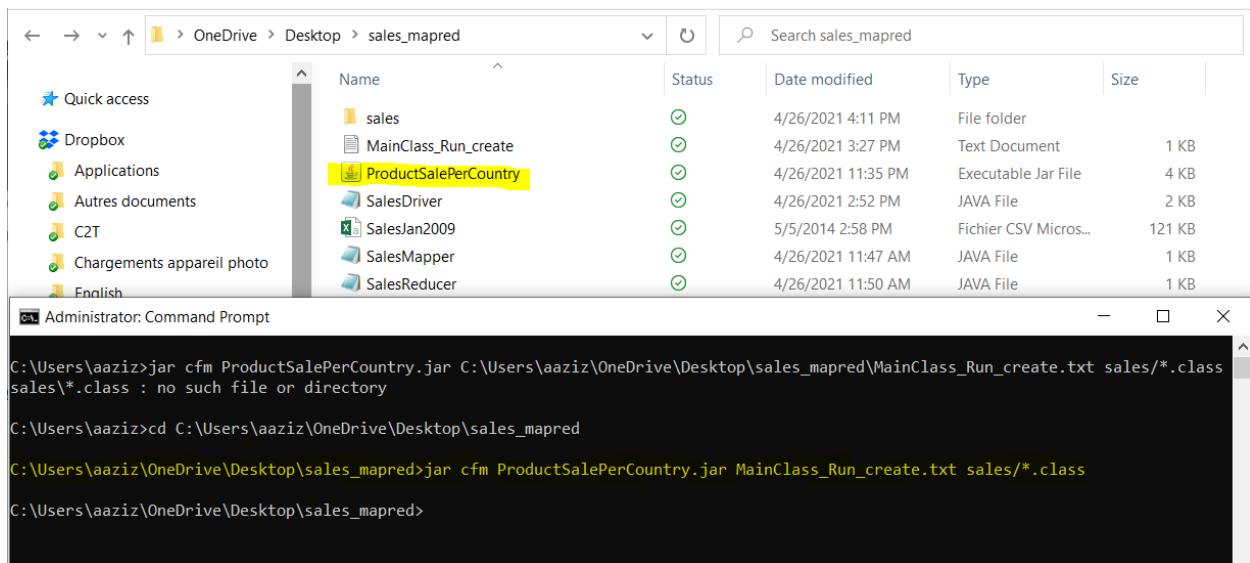
**Sales.SaleDriver** is the name of our main class
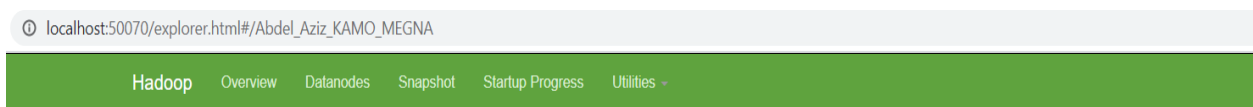
**Note:** we have to hit enter key at end of this line.

### 5- Create jar file for hadoop execution

This part consists of the creation of the jar file which will be executed by hadoop, this file contains the different classes of our application in which our main class is specified.

## 6- Creation of our input and output directories in hadoop space

- hadoop fs -mkdir /Abdel_Aziz_KAMO_MEGNA/input
- hadoop fs -mkdir /Abdel_Aziz_KAMO_MEGNA/output



### Browse Directory

/Abdel_Aziz_KAMO_MEGNA                                                                                     Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|------------|-------|-------|------|---------------|-------------|------------|------|
| drwxr-xr-x | aaziz | supergroup | 0 B | 4/30/2021, 4:36:00 PM | 0 | 0 B | input |
| drwxr-xr-x | aaziz | supergroup | 0 B | 4/30/2021, 4:38:33 PM | 0 | 0 B | output |

Hadoop, 2018.

## 7- Copy of our dataset file in our input directory

- hadoop fs -put C:\Users\aaziz\OneDrive\Desktop\sales_mapred\SalesJan2009.csv \Abdel_Aziz_KAMO_MEGNA\input



## 8- Running our application

To run our application and put the result into our output folder, we have to execute this command: **hadoop jar ProductSalePerCountry.jar \Abdel_Aziz_KAMO_MEGNA\input Abdel_Aziz KAMO_MEGNA\output**

localhost:50070/explorer.html#/user/aaziz/Abdel_Aziz_KAMO_MEGNA/output

**Hadoop**    Overview    Datanodes    Snapshot    Startup Progress    Utilities ▾

# Browse Directory

/user/aaziz/Abdel_Aziz_KAMO_MEGNA/output     Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | aaziz | supergroup | 0 B | 4/30/2021, 5:01:49 PM | 1 | 128 MB | _SUCCESS |
| -rw-r--r-- | aaziz | supergroup | 661 B | 4/30/2021, 5:01:48 PM | 1 | 128 MB | part-00000 |

Hadoop, 2018.

**Select Administrator: Command Prompt**

```
C:\Users\aaziz>hadoop jar C:\Users\aaziz\OneDrive\Desktop\sales_mapred\ProductSalePerCountry.jar \Abdel_Aziz_KAMO_MEGNA\inp
ut Abdel_Aziz_KAMO_MEGNA\output
21/04/30 17:00:42 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
21/04/30 17:00:42 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
21/04/30 17:00:44 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool
interface and execute your application with ToolRunner to remedy this.
21/04/30 17:00:44 INFO mapred.FileInputFormat: Total input paths to process : 1
21/04/30 17:00:45 INFO mapreduce.JobSubmitter: number of splits:2
21/04/30 17:00:45 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619788973438_0001
21/04/30 17:00:46 INFO impl.YarnClientImpl: Submitted application application_1619788973438_0001
21/04/30 17:00:46 INFO mapreduce.Job: The url to track the job: http://DESKTOP-SN0MSI7:8088/proxy/application_1619788973438
_0001/
21/04/30 17:00:46 INFO mapreduce.Job: Running job: job_1619788973438_0001
21/04/30 17:01:11 INFO mapreduce.Job: Job job_1619788973438_0001 running in uber mode : false
21/04/30 17:01:11 INFO mapreduce.Job:  map 0% reduce 0%
21/04/30 17:01:28 INFO mapreduce.Job:  map 100% reduce 0%
21/04/30 17:01:50 INFO mapreduce.Job:  map 100% reduce 100%
21/04/30 17:01:53 INFO mapreduce.Job: Job job_1619788973438_0001 completed successfully
21/04/30 17:01:53 INFO mapreduce.Job: Counters: 49
        File System Counters
                FILE: Number of bytes read=17747
                FILE: Number of bytes written=406500
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=127591
                HDFS: Number of bytes written=661
```

**Select Administrator: Command Prompt**

```
                Map output materialized bytes=17753
                Input split bytes=236
                Combine input records=0
                Combine output records=0
                Reduce input groups=58
                Reduce shuffle bytes=17753
                Reduce input records=999
                Reduce output records=58
                Spilled Records=1998
                Shuffled Maps =2
                Failed Shuffles=0
                Merged Map outputs=2
                GC time elapsed (ms)=180
                CPU time spent (ms)=4074
                Physical memory (bytes) snapshot=622325760
                Virtual memory (bytes) snapshot=834326528
                Total committed heap usage (bytes)=450363392
        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=127355
        File Output Format Counters
                Bytes Written=661

C:\Users\aaziz>
```

Once done, we can view and download our output result file to our output directory on hadoop.

**File information - part-00000** ✕

Download

**Block information --** Block 0 ⌄

Block ID: 1073741832

Block Pool ID: BP-20745571-192.168.1.39-1619648186150

Generation Stamp: 1008

Size: 661

Availability:

- DESKTOP-SN0MSI7

Close

## 9- Output result

```
part-00000 - Notepad
File  Edit  Format  View  Help
Argentina          1
Australia          38
Austria 7
Bahrain 1
Belgium 8
Bermuda 1
Brazil   5
Bulgaria           1
CO       1
Canada   76
Cayman Isls        1
China    1
Costa Rica         1
Country 1
Czech Republic     3
Denmark 15
Dominican Republic        1
Finland 2
France   27
Germany 25
Greece   1
Guatemala          1
Hong Kong          1
Hungary 3
Iceland 1
India    2
Ireland 49
Israel   1
Italy    15
Japan    2
Jersey   1
Kuwait   1
Latvia   1
Luxembourg         1
Malaysia           1
Malta    2
Mauritius          1
Moldova 1
Monaco   2
Netherlands        22
New Zealand        6
Norway   16
Philippines        2

Poland   2
Romania 1
Russia   1
South Africa       5
South Korea        1
Spain    12
Sweden   13
Switzerland        36
Thailand           2
The Bahamas        2
Turkey   6
Ukraine 1
United Arab Emirates      6
United Kingdom  100
United States   462
```