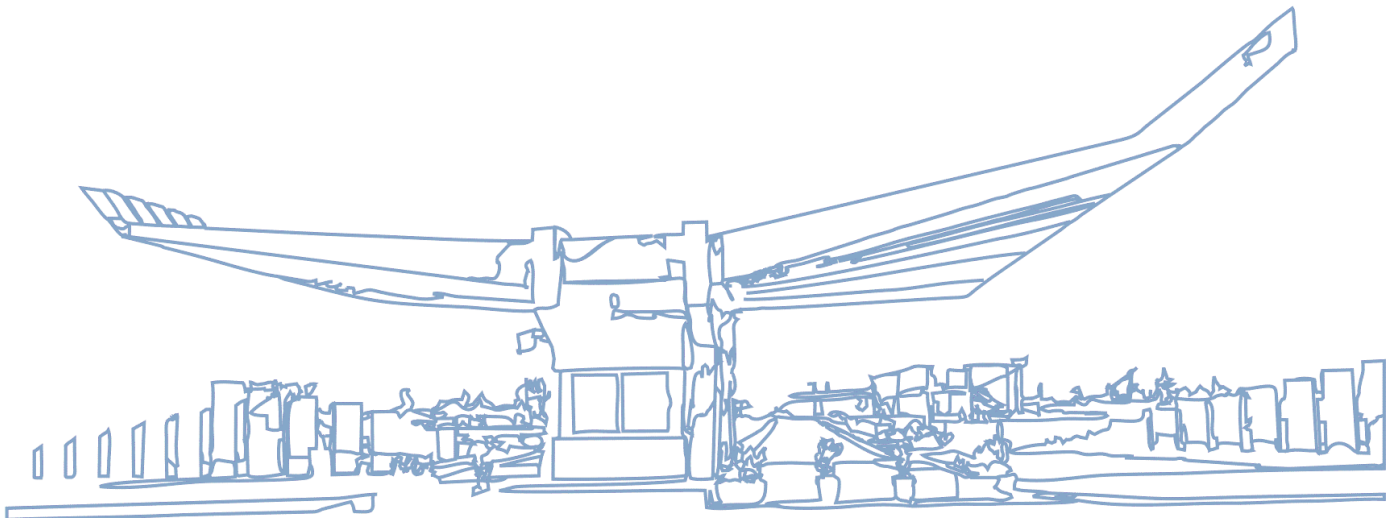


CEN 571 – Data Mining

Assignment 3: MapReduce



PREPARED:
Baftjar TABAKU

26.04.2020
Epoka University
Tirana, ALBANIA

ACCEPTED:
Prof.Dr. Arben Asllani

Assignment tasks and notes

Use the above file (NYSE.csv) as a source and complete these tasks. Consider only records for the last 10 years and with a volume greater than 250,000:

1. Modify an existing or write a new MapReduce code
2. Compile the program as a jar file
3. Upload the input file(s) into the Hadoop AWS cluster
4. Execute the program and display the results

The deliverables for each task are:

- Modified Java code
- List of Linux and HDFS commands that are used to execute the project
- Output results in a text file

Write a Map-Reduce program to find out maximum trading price for each stock

Modify the MapReduce Java code to find out highest price for each stock

Assignment tasks and notes

After modifying the following java code, that also is included and every change is commented, build was successfully.

```

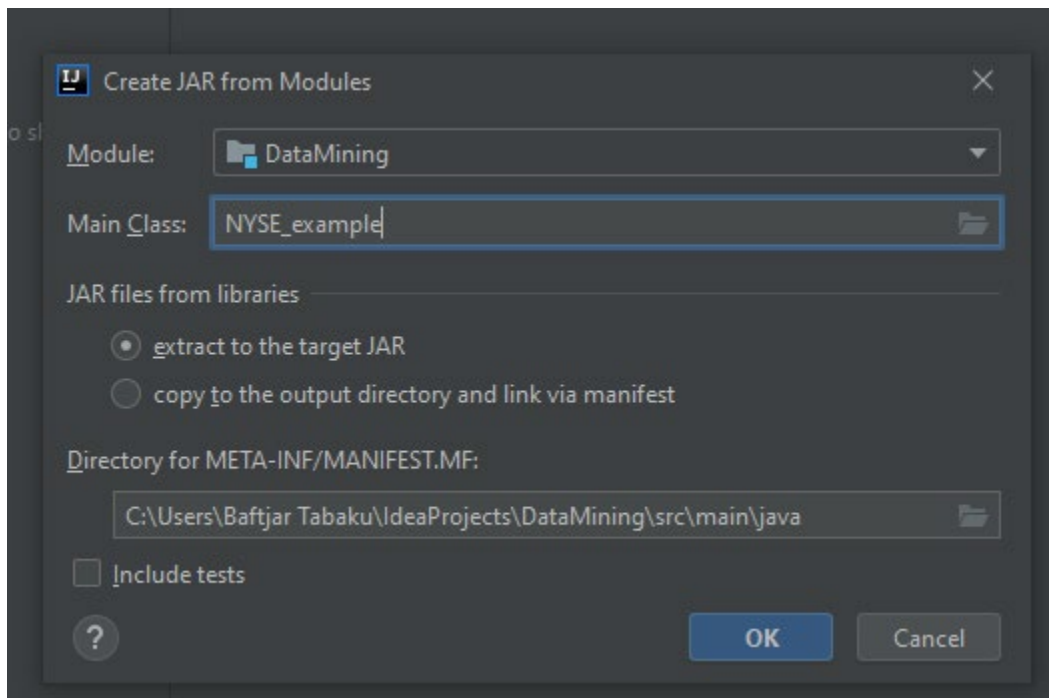
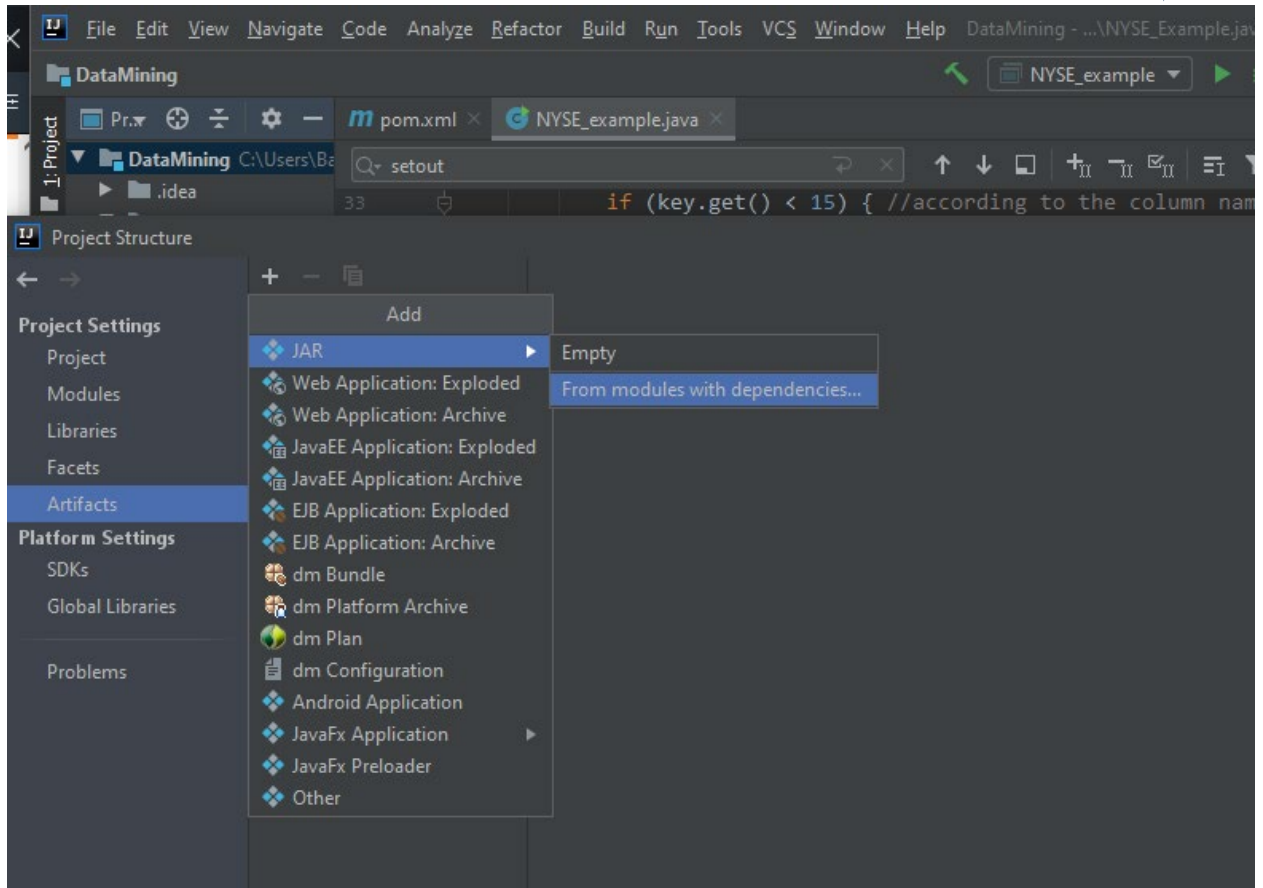
33  if (key.get() < 15) { //according to the column names of csv file
34      return; //byte offset of the file first row must be less than 20
35  }
36  /* we calculate the record year
37   * according to a column we split the year, and we take the last one, and
38   * is 10 last years according to the requirements*/
39  int year = Integer.parseInt(tokens[2].split(" ")[2]);
40
41  if (year > 9) {
42      // we change this ticker.set(tokens[0]); to ticker.set(tokens[1]);
43      ticker.set(tokens[1]);
44      // Changing from trade_price.set(Long.parseLong(tokens[2])); to
45      // trade_price.set(Float.parseFloat(tokens[2])); index from 2 to 6
46      trade_price.set(Float.parseFloat(tokens[6]));
47      output.collect(ticker, trade_price);
48  }
49  }
50  }
51  }
52
53  // we change public static class Reduce extends MapReduceBase implements Reducer
54  public static class Reduce extends MapReduceBase implements Reducer<Text, Float> {
55      @Override
56      public void reduce(Text key, Iterator<FloatWritable> values, OutputCollector<Text, Float> output, ProgressReporter progress) throws IOException {
57          long MaxPrice = 0;
58          float MaxPrice = 0;
59          while (values.hasNext()) {
60              long next_price = values.next().get();
61              float next_price = values.next().get();

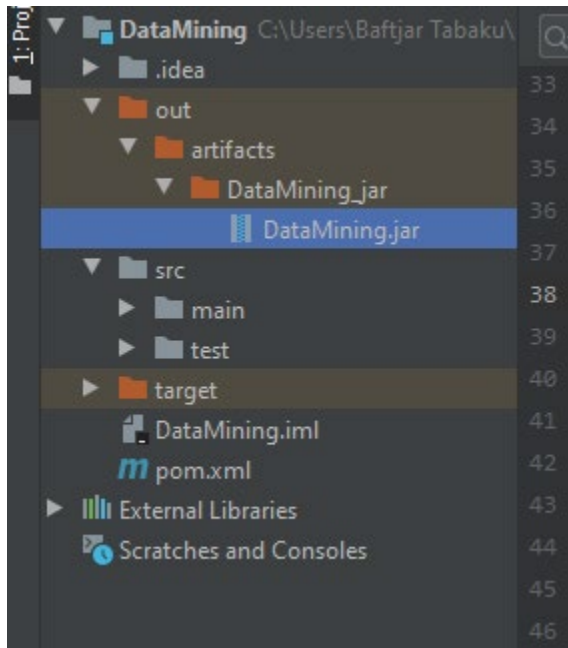
```

Build Messages:

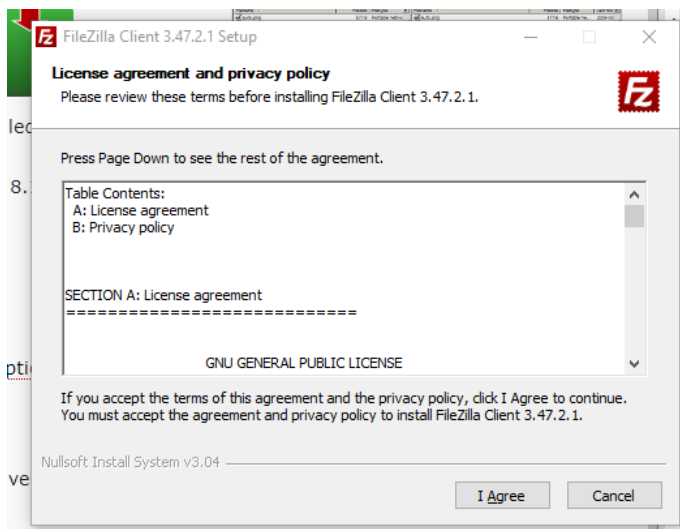
- Information: javac 1.8.0_162 was used to compile java sources
- Information: 4/26/2020 5:15 PM - Build completed successfully with 3 warnings in 2 s 882 ms
- Warning: java: source value 1.5 is obsolete and will be removed in a future release
- Warning: java: target value 1.5 is obsolete and will be removed in a future release
- Warning: java: To suppress warnings about obsolete options, use -Xlint:-options.

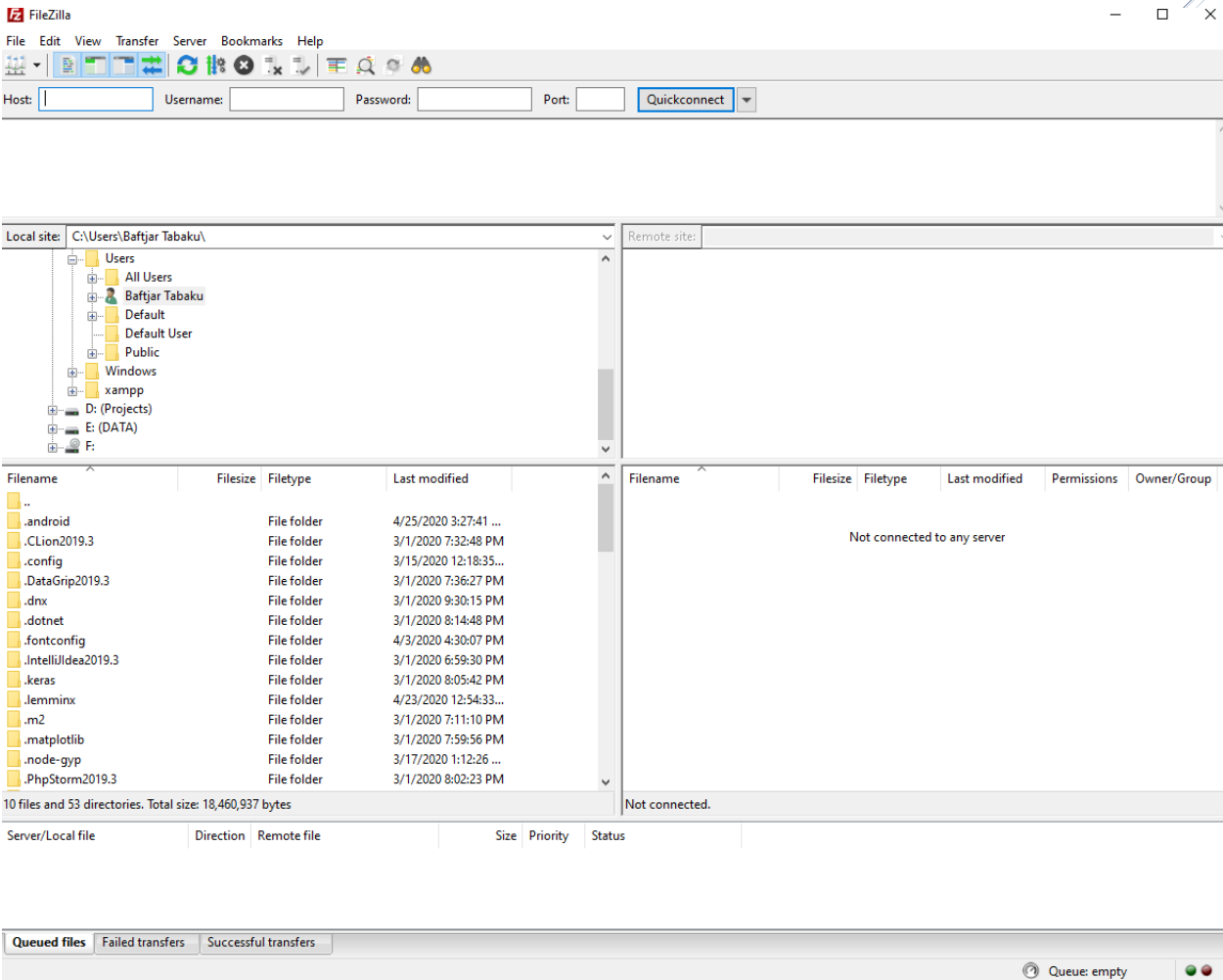
After this compiling the program as a jar file, as shown below





And now, to interact with Hadoop we need to download FileZilla and putty as shown below,





The next step was to install the 'putty' for windows,

Package files

You probably want one of these. They include versions of all the PuTTY utilities.

(Not sure whether you want the 32-bit or the 64-bit version? Read the [FAQ entry](#).)

MSI ('Windows Installer')

32-bit:	putty-0.73-installer.msi	(or by FTP)	(signature)
64-bit:	putty-64bit-0.73-installer.msi	(or by FTP)	(signature)

Unix source archive

.tar.gz:	putty-0.73.tar.gz	(or by FTP)	(signature)
----------	-----------------------------------	-----------------------------	-----------------------------

Alternative

The installer package...

(Not sure whether you want the 32-bit or the 64-bit version? Read the [FAQ entry](#).)

putty.exe (the PuTTY client)

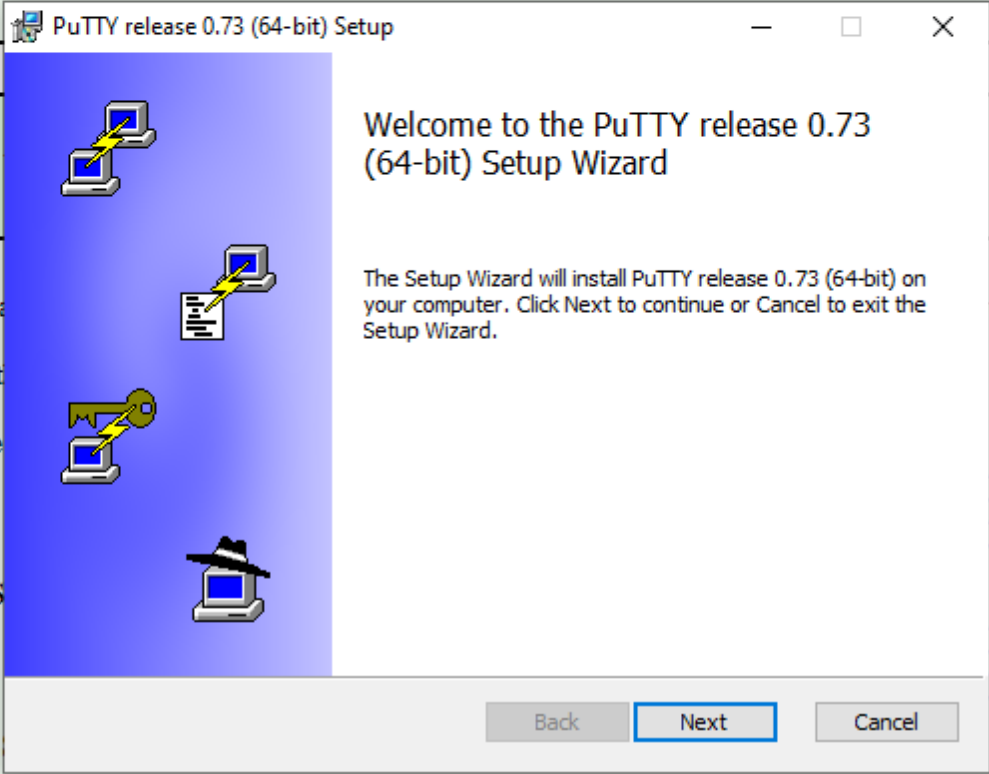
32-bit:	putty-0.73.exe	(or by FTP)	(signature)
64-bit:	putty-64bit-0.73.exe	(or by FTP)	(signature)

pscp.exe (an SFTP client)

32-bit:	pscp-0.73.exe	(or by FTP)	(signature)
64-bit:	pscp-64bit-0.73.exe	(or by FTP)	(signature)

psftp.exe (an SFTP client)

32-bit:	psftp-0.73.exe	(or by FTP)	(signature)
64-bit:	psftp-64bit-0.73.exe	(or by FTP)	(signature)



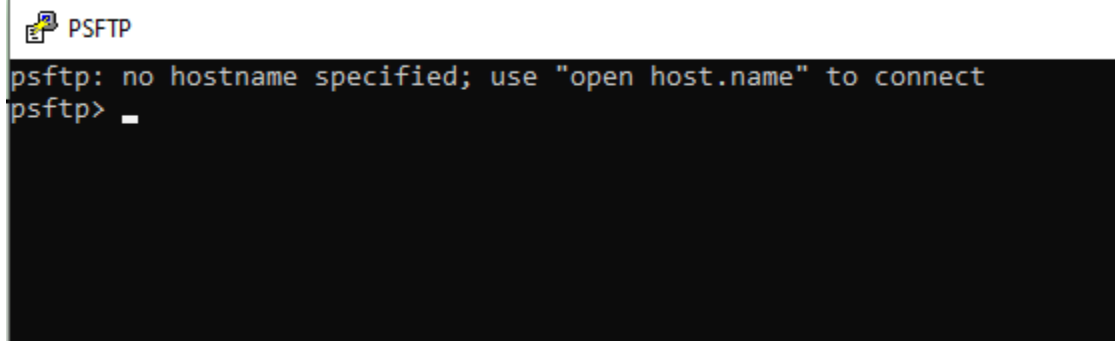
PuTTY release 0.73 (64-bit) Setup

Welcome to the PuTTY release 0.73 (64-bit) Setup Wizard

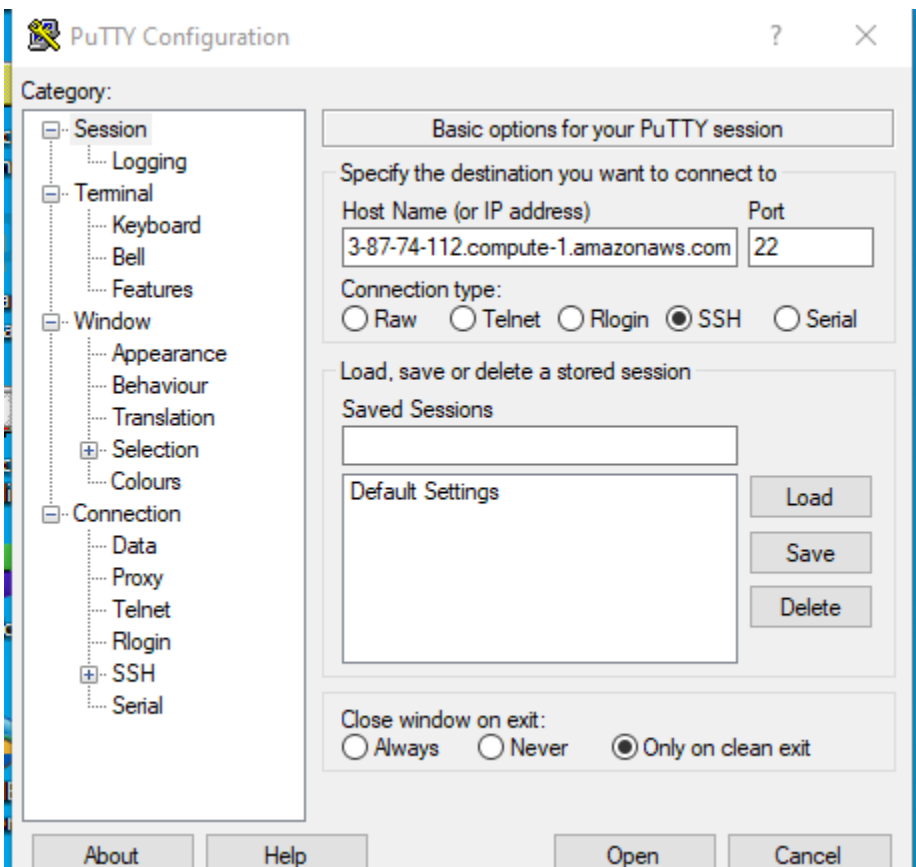
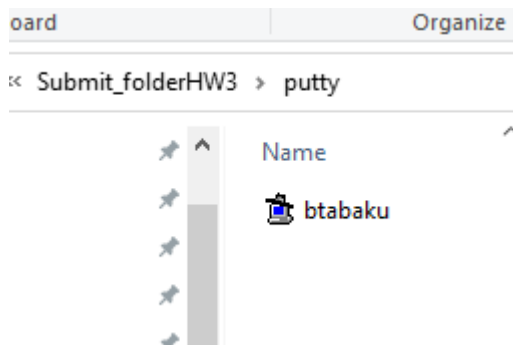
The Setup Wizard will install PuTTY release 0.73 (64-bit) on your computer. Click Next to continue or Cancel to exit the Setup Wizard.

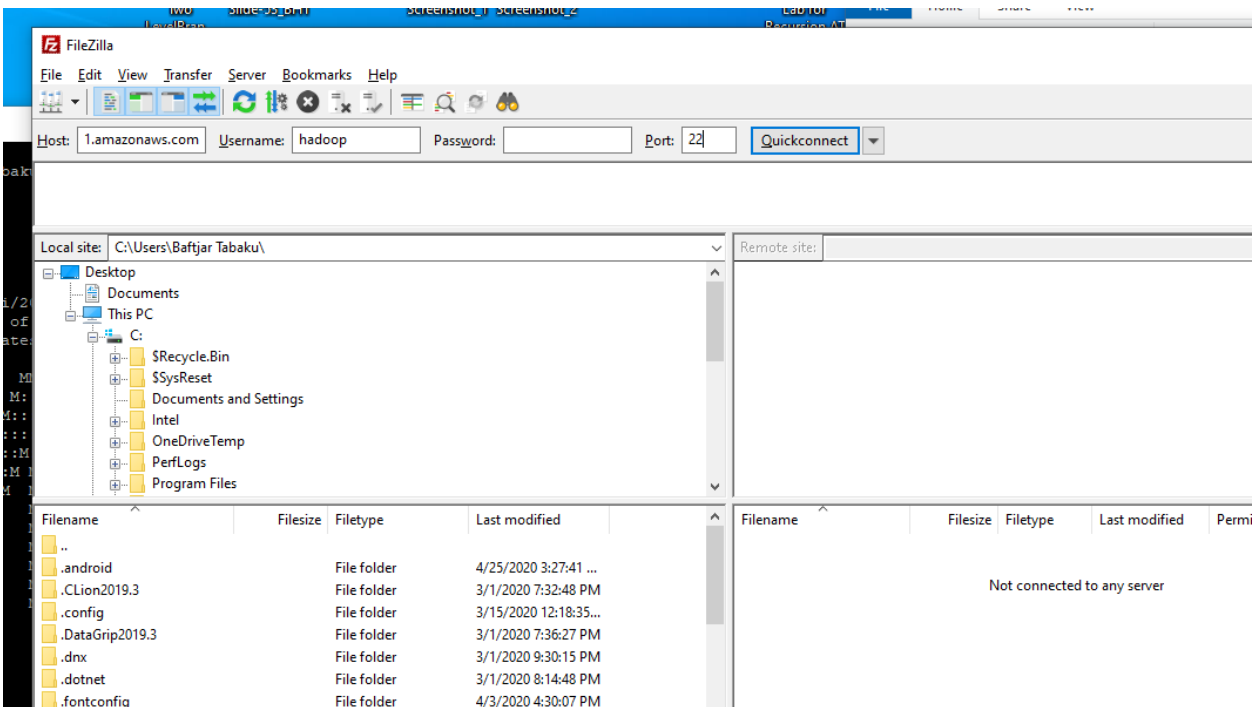
Back Next Cancel

Then starting it



According to the file '.ppk'



[illegible]

sftp://hadoop@ec2-3-87-74-112.compute-1.amazonaws.com - FileZilla

File Edit View Transfer Server Bookmarks Help

Host: sftp://ec2-3-87-74- Username: hadoop Password: Port: Quickconnect

Status: Connecting to ec2-3-87-74-112.compute-1.amazonaws.com...
 Status: Connected to ec2-3-87-74-112.compute-1.amazonaws.com
 Status: Retrieving directory listing...
 Status: Listing directory /home/hadoop
 Status: Directory listing of "/home/hadoop" successful

Local site: ers\Bafjar Tabaku\Desktop\Master_2ndSemester\CEN_571_DATA_MINING\Submit_folderHW3\ Remote site: /home/hadoop

Local site tree:

- CLionProjects
- Contacts
- Cookies
- Desktop
- Internship_LUFTHANSA
- Master_2ndSemester
 - CEN-564-WIRELESS NETWORKS
 - CEN-583-ADVANCED COMPUTER ARCHITECTURE
 - CEN_552_ADVANCED_DATABASE_MANAGEMENT_SYSTEM
 - CEN_571_DATA_MINING
 - Homework03
 - Lessons
 - Submit_folderHW3
- Documents
- Downloads

Remote site tree:

- home
 - hadoop

Filename	Filesize	Filetype	Last modified	Permissions	Owner/Group
..					
.aws		File folder	12/14/2019 5:2...	drwxr-xr-x	hadoop had...
.ssh		File folder	4/26/2020 6:27:...	drwx-----	hadoop had...
.bash_profile	86	BASH_PRO...	12/2/2019 6:32:...	-rw-r--r--	hadoop had...
.bashrc	357	BASHRC File	12/2/2019 6:32:...	-rw-r--r--	hadoop had...

4 files and 1 directory. Total size: 37,655,161 bytes

2 files and 2 directories. Total size: 443 bytes

Server/Local file	Direction	Remote file	Size	Priority	Status
-------------------	-----------	-------------	------	----------	--------

Status: Connecting to ec2-3-87-74-112.compute-1.amazonaws.com...
 Status: Connected to ec2-3-87-74-112.compute-1.amazonaws.com
 Status: Retrieving directory listing...
 Status: Listing directory /home/hadoop
 Status: Directory listing of "/home/hadoop" successful
 Status: Connecting to ec2-3-87-74-112.compute-1.amazonaws.com...
 Status: Connecting to ec2-3-87-74-112.compute-1.amazonaws.com...
 Status: Connected to ec2-3-87-74-112.compute-1.amazonaws.com
 Status: Connected to ec2-3-87-74-112.compute-1.amazonaws.com
 Status: Starting upload of C:\Users\Bafjar Tabaku\Desktop\Master_2ndSemester\CEN_571_DATA_MINING\Submit_folderHW3\NYSE.csv
 Status: Starting upload of C:\Users\Bafjar Tabaku\Desktop\Master_2ndSemester\CEN_571_DATA_MINING\Submit_folderHW3\NYSE_A.jar

Local site: ers\Bafjar Tabaku\Desktop\Master_2ndSemester\CEN_571_DATA_MINING\Submit_folderHW3\ Remote site: /home/hadoop

Local site tree:

- CLionProjects
- Contacts
- Cookies
- Desktop
- Internship_LUFTHANSA
- Master_2ndSemester
 - CEN-564-WIRELESS NETWORKS
 - CEN-583-ADVANCED COMPUTER ARCHITECTURE
 - CEN_552_ADVANCED_DATABASE_MANAGEMENT_SYSTEM
 - CEN_571_DATA_MINING
 - Homework03
 - Lessons
 - Submit_folderHW3
- Documents
- Downloads

Remote site tree:

- home
 - hadoop

Filename	Filesize	Filetype	Last modified	Permissions	Owner/Group
..					
.aws		File folder	12/14/2019 5:2...	drwxr-xr-x	hadoop l
.ssh		File folder	4/26/2020 6:27:...	drwx-----	hadoop l
.bash_profile	86	BASH_PRO...	12/2/2019 6:32:...	-rw-r--r--	hadoop l
.bashrc	357	BASHRC File	12/2/2019 6:32:...	-rw-r--r--	hadoop l

Selected 3 files. Total size: 37,650,999 bytes

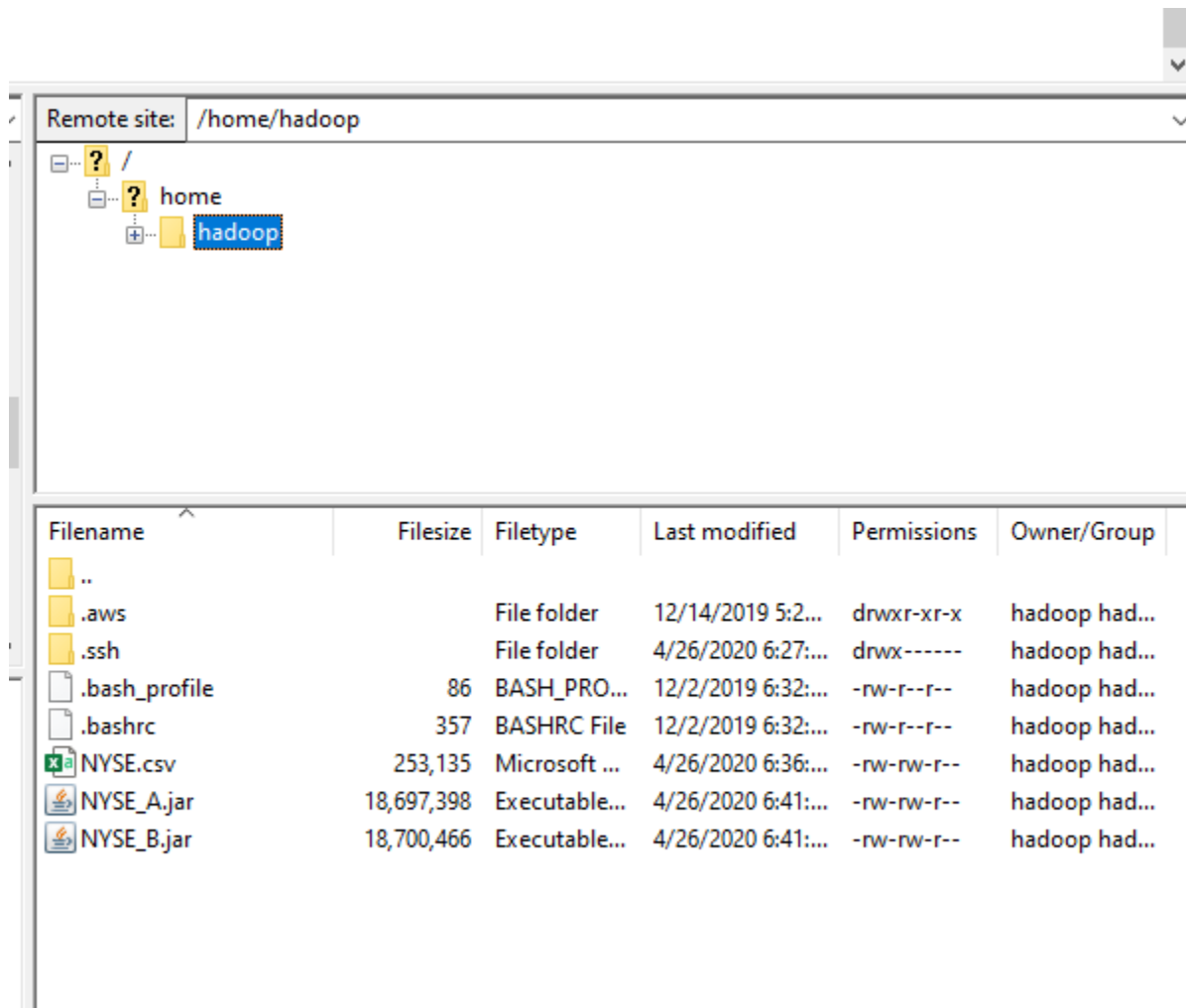
2 files and 2 directories. Total size: 443 bytes

Server/Local file	Direction	Remote file	Size	Priority	Status
sftp://hadoop@ec2-3-87-7...	-->>	/home/hadoop/NYSE.csv	253,135	Normal	Transferring
C:\Users\Bafjar Tabaku\D...	-->>	/home/hadoop/NYSE_A.jar	18,697,398	Normal	Transferring

00:00:02 elapsed 00:00:03 left 64.7% 163,840 bytes (62.2 KiB/s)

Queued files (3) Failed transfers Successful transfers

Queue: 36.0 MiB



Using a simple command to list

```
[hadoop@ip-172-31-22-65 ~]$ ls -a
.  ..  .aws  .bash_profile  .bashrc  NYSE_A.jar  NYSE_B.jar  NYSE.csv  .ssh
[hadoop@ip-172-31-22-65 ~]$
```

Making a folder btabaku16/mr

```
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -mkdir btabaku16
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -mkdir btabaku16/mr
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -copyFromLocal NYSE.csv btabaku16/mr
copyFromLocal: `NYSE.csv': No such file or directory
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -copyFromLocal NYSE.csv btabaku16/mr
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls
Found 1 items
```

```
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -copyFromLocal NYSE.csv btabakul6/mr
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls
Found 1 items
drwxr-xr-x   4 hadoop hadoop          0 2020-04-26 16:43 btabakul6
```

```
Found 1 items
drwxr-xr-x - hadoop hadoop 0 2020-04-26 16:43 btabakul6
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabakul6/mr
Found 1 items
-rw-r--r-- 1 hadoop hadoop 253135 2020-04-26 16:45 btabakul6/mr/NYSE.csv
[hadoop@ip-172-31-22-65 ~]$
```

hadoop@ip-172-31-22-65:~

```

[hadoop@ip-172-31-22-65 ~]$ ls -a
. . . .aws .bash_profile .bashrc NYSE_A.jar NYSE_B.jar NYSE.csv .ssh
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -mkdir btabakul6
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -mkdir btabakul6/mr
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -copyFromLocal BYSE.csv btabakul6/mr
copyFromLocal: `BYSE.csv': No such file or directory
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -copyFromLocal NYSE.csv btabakul6/mr
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls
Found 1 items
drwxr-xr-x - hadoop hadoop 0 2020-04-26 16:43 btabakul6
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabakul6/mr
Found 1 items
-rw-r--r-- 1 hadoop hadoop 253135 2020-04-26 16:45 btabakul6/mr/NYSE.csv
[hadoop@ip-172-31-22-65 ~]$ hadoop jar NYSE_A.jar NYSE_example btabakul6/mr/NYSE.csv btabakul6/mr/outputA
20/04/26 16:53:54 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-22-65.ec2.internal/172.31.2
2.65:8032
20/04/26 16:53:54 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-22-65.ec2.internal/172.31.2
2.65:8032
20/04/26 16:53:54 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Imple
ment the Tool interface and execute your application with ToolRunner to remedy this.
20/04/26 16:53:54 INFO lzo.GPLNativeCodeLoader: Loaded native gpl library
20/04/26 16:53:55 INFO lzo.LzoCodec: Successfully loaded & initialized native-lzo library [hadoop-lzo rev 5f7
88d5e8f90539ee331702c753fa250727128f4]
20/04/26 16:53:55 INFO mapred.FileInputFormat: Total input files to process : 1
20/04/26 16:53:55 INFO mapreduce.JobSubmitter: number of splits:8
20/04/26 16:53:55 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1587918542150_0001
20/04/26 16:53:55 INFO impl.YarnClientImpl: Submitted application application_1587918542150_0001
20/04/26 16:53:55 INFO mapreduce.Job: The url to track the job: http://ip-172-31-22-65.ec2.internal:20888/pro
xy/application_1587918542150_0001/
20/04/26 16:53:55 INFO mapreduce.Job: Running job: job_1587918542150_0001

```

Running the first command and getting the output A, to the folder btabaku16/mr/outputA

```

[hadoop@ip-172-31-22-65 ~]$ hadoop jar NYSE_A.jar NYSE_example btabakul6/mr/NYSE.csv btabakul6/mr/outputA
20/04/26 16:53:54 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-22-65.ec2.internal/172.31.2
2.65:8032
20/04/26 16:53:54 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-22-65.ec2.internal/172.31.2
2.65:8032
20/04/26 16:53:54 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Imple
ment the Tool interface and execute your application with ToolRunner to remedy this.
20/04/26 16:53:54 INFO lzo.GPLNativeCodeLoader: Loaded native gpl library
20/04/26 16:53:55 INFO lzo.LzoCodec: Successfully loaded & initialized native-lzo library [hadoop-lzo rev 5f7
88d5e8f90539ee331702c753fa250727128f4]
20/04/26 16:53:55 INFO mapred.FileInputFormat: Total input files to process : 1
20/04/26 16:53:55 INFO mapreduce.JobSubmitter: number of splits:8
20/04/26 16:53:55 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1587918542150_0001
20/04/26 16:53:55 INFO impl.YarnClientImpl: Submitted application application_1587918542150_0001
20/04/26 16:53:55 INFO mapreduce.Job: The url to track the job: http://ip-172-31-22-65.ec2.internal:20888/pro
xy/application_1587918542150_0001/
20/04/26 16:53:55 INFO mapreduce.Job: Running job: job_1587918542150_0001
20/04/26 16:54:01 INFO mapreduce.Job: Job job_1587918542150_0001 running in uber mode : false
20/04/26 16:54:01 INFO mapreduce.Job: map 0% reduce 0%
20/04/26 16:54:07 INFO mapreduce.Job: map 25% reduce 0%
20/04/26 16:54:10 INFO mapreduce.Job: map 63% reduce 0%
20/04/26 16:54:11 INFO mapreduce.Job: map 88% reduce 0%
20/04/26 16:54:12 INFO mapreduce.Job: map 100% reduce 0%
20/04/26 16:54:15 INFO mapreduce.Job: map 100% reduce 100%
20/04/26 16:54:16 INFO mapreduce.Job: Job job_1587918542150_0001 completed successfully
20/04/26 16:54:17 INFO mapreduce.Job: Counters: 51
File System Counters
FILE: Number of bytes read=20841

```

Which successfully worked and executed.

Then running the second command and getting the output B, to the folder
btabaku16/mr/outputB

```
[hadoop@ip-172-31-22-65 ~]$ hadoop jar NYSE_B.jar NYSE_example_B btabaku16/mr/NYSE.csv btabaku16/mr/outputB
20/04/26 16:59:54 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-22-65.ec2.internal/172.31.22.65:8032
20/04/26 16:59:54 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-22-65.ec2.internal/172.31.22.65:8032
20/04/26 16:59:54 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and e
xecute your application with ToolRunner to remedy this.
20/04/26 16:59:54 INFO lzo.GPLNativeCodeLoader: Loaded native gpl library
20/04/26 16:59:54 INFO lzo.LzoCodec: Successfully loaded & initialized native-lzo library [hadoop-lzo rev 5f788d5e8f90539ee331702c753fa250
727128f4]
20/04/26 16:59:54 INFO mapred.FileInputFormat: Total input files to process : 1
20/04/26 16:59:54 INFO mapreduce.JobSubmitter: number of splits:8
20/04/26 16:59:55 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1587918542150_0002
20/04/26 16:59:55 INFO impl.YarnClientImpl: Submitted application application_1587918542150_0002
20/04/26 16:59:55 INFO mapreduce.Job: The url to track the job: http://ip-172-31-22-65.ec2.internal:20888/proxy/application_1587918542150_
0002/
20/04/26 16:59:55 INFO mapreduce.Job: Running job: job_1587918542150_0002
20/04/26 17:00:00 INFO mapreduce.Job: Job job_1587918542150_0002 running in uber mode : false
20/04/26 17:00:00 INFO mapreduce.Job: map 0% reduce 0%
20/04/26 17:00:06 INFO mapreduce.Job: map 25% reduce 0%
20/04/26 17:00:07 INFO mapreduce.Job: map 38% reduce 0%
20/04/26 17:00:08 INFO mapreduce.Job: map 50% reduce 0%
20/04/26 17:00:09 INFO mapreduce.Job: map 88% reduce 0%
20/04/26 17:00:10 INFO mapreduce.Job: map 100% reduce 0%
20/04/26 17:00:12 INFO mapreduce.Job: map 100% reduce 33%
20/04/26 17:00:13 INFO mapreduce.Job: map 100% reduce 100%
20/04/26 17:00:14 INFO mapreduce.Job: Job job_1587918542150_0002 completed successfully
20/04/26 17:00:14 INFO mapreduce.Job: Counters: 51
File System Counters
FILE: Number of bytes read=20947
FILE: Number of bytes written=1010426
```

And after this, running a list command as shown in the following screenshot below.

```
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabaku16/mr
Found 3 items
-rw-r--r-- 1 hadoop hadoop 253135 2020-04-26 16:45 btabaku16/mr/NYSE.csv
drwxr-xr-x - hadoop hadoop 0 2020-04-26 16:54 btabaku16/mr/outputA
drwxr-xr-x - hadoop hadoop 0 2020-04-26 17:00 btabaku16/mr/outputB
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabaku16/mr/outputA
Found 4 items
-rw-r--r-- 1 hadoop hadoop 0 2020-04-26 16:54 btabaku16/mr/outputA/_SUCCESS
-rw-r--r-- 1 hadoop hadoop 482 2020-04-26 16:54 btabaku16/mr/outputA/part-00000
-rw-r--r-- 1 hadoop hadoop 512 2020-04-26 16:54 btabaku16/mr/outputA/part-00001
-rw-r--r-- 1 hadoop hadoop 502 2020-04-26 16:54 btabaku16/mr/outputA/part-00002
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabaku16/mr/outputB
Found 4 items
-rw-r--r-- 1 hadoop hadoop 0 2020-04-26 17:00 btabaku16/mr/outputB/_SUCCESS
-rw-r--r-- 1 hadoop hadoop 482 2020-04-26 17:00 btabaku16/mr/outputB/part-00000
-rw-r--r-- 1 hadoop hadoop 512 2020-04-26 17:00 btabaku16/mr/outputB/part-00001
-rw-r--r-- 1 hadoop hadoop 501 2020-04-26 17:00 btabaku16/mr/outputB/part-00002
[hadoop@ip-172-31-22-65 ~]$
```


And after this, going to each output folder A and B as shown below and listing their corresponding contents.

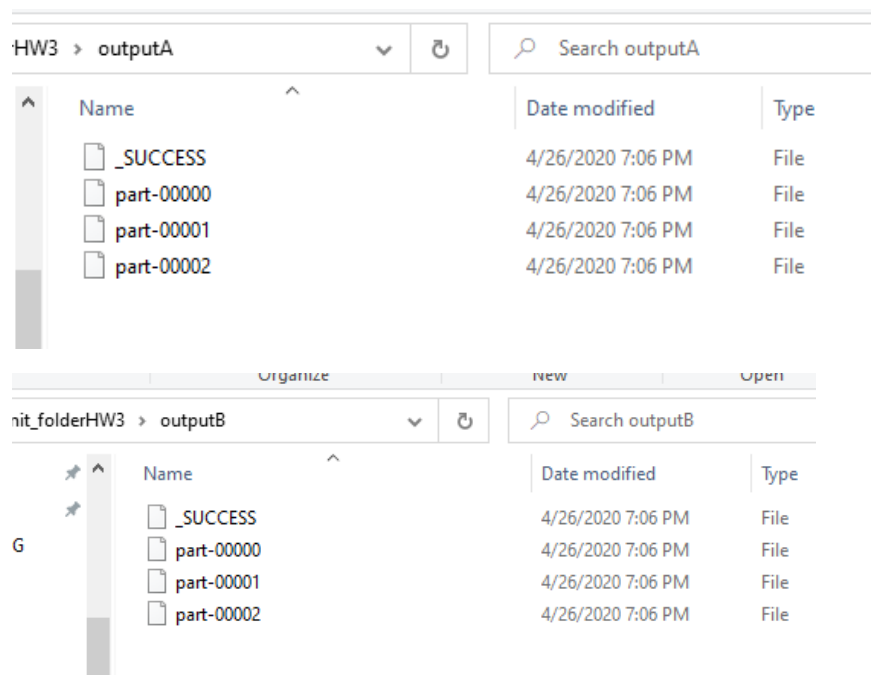
```

Bytes Written=1495
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabakul6/mr
Found 3 items
-rw-r--r-- 1 hadoop hadoop 253135 2020-04-26 16:45 btabakul6/mr/NYSE.csv
drwxr-xr-x - hadoop hadoop 0 2020-04-26 16:54 btabakul6/mr/outputA
drwxr-xr-x - hadoop hadoop 0 2020-04-26 17:00 btabakul6/mr/outputB
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabakul6/mr/outputA
Found 4 items
-rw-r--r-- 1 hadoop hadoop 0 2020-04-26 16:54 btabakul6/mr/outputA/_SUCCESS
-rw-r--r-- 1 hadoop hadoop 482 2020-04-26 16:54 btabakul6/mr/outputA/part-00000
-rw-r--r-- 1 hadoop hadoop 512 2020-04-26 16:54 btabakul6/mr/outputA/part-00001
-rw-r--r-- 1 hadoop hadoop 502 2020-04-26 16:54 btabakul6/mr/outputA/part-00002
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -ls btabakul6/mr/outputB
Found 4 items
-rw-r--r-- 1 hadoop hadoop 0 2020-04-26 17:00 btabakul6/mr/outputB/_SUCCESS
-rw-r--r-- 1 hadoop hadoop 482 2020-04-26 17:00 btabakul6/mr/outputB/part-00000
-rw-r--r-- 1 hadoop hadoop 512 2020-04-26 17:00 btabakul6/mr/outputB/part-00001
-rw-r--r-- 1 hadoop hadoop 501 2020-04-26 17:00 btabakul6/mr/outputB/part-00002
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -copyToLocal btabakul6/mr/outputA ~
[hadoop@ip-172-31-22-65 ~]$ hadoop fs -copyToLocal btabakul6/mr/outputB ~
[hadoop@ip-172-31-22-65 ~]$ ls
NYSE_A.jar NYSE_B.jar NYSE.csv outputA outputB
[hadoop@ip-172-31-22-65 ~]$

```

And everything finished Successfully.

Note: there at the compilation of the jar files, are two classes, the 'NYSE_example.java' that stands for the task 1 and the 'NYSE_example_B.java' that stands for task 2 or B as shown. The output of folders outputA and outputB had like 3 files for each



And all this is splinted into two text files OutputA.txt and OutputB.txt.