

# Lab 3 Spark,

By:

<b>Mohamed Said</b>	<b>57</b>
<b>Ahmed El Bawab</b>	<b>8</b>
<b>Khalil Ismail</b>	<b>23</b>

**Screen Shots:**

The screenshot shows a presentation slide titled "Word Count" with the following text: "Now, you will create the Word Count application and run it as a Spark job on some files on HDFS." Below this, there is a list of steps: "You can download the code for WordCount example from [here](#).", "Extract the downloaded zip and make it your current directory.", "Compile the code using mvn clean compile.", "Create a jar file using mvn package -DskipTests.", "You are ready to run the application using java -jar wordcount.jar [HDFS\\_INPUT\\_FILE\\_ON\\_HDFS](#) command.", "Copy the output to a file." A terminal window is overlaid on the slide, showing the command "spark-shell" and the prompt "said@said:~\$".

Department: \_\_\_\_\_ Date: 29/3/2020

### Word Count

Now, you will create the Word Count application and run it as a Spark job on some files on HDFS.

- You can download the code for WordCount example from [here](#).
- Extract the downloaded zip and make it your current directory.
- Compile the code using `mvn clean compile`.
- Create a jar file using `mvn package -DskipTests`.
- You are ready to run the application using `java -jar wordcount.jar HDFS\_INPUT\_FILE\_ON\_HDFS` command.
- Copy the output to a file.

### Notes

- You can work in groups of 2 or 3.

Good Luck

Terminal output: said@said:~\$ spark-shell

Applications Places System 66% 40 B/s 56 B/s Terminal File Edit View Search Terminal Help 13°C en 8:42, 22 مار 2020

lap3

Previous Next 2 (2 of 2) 200%

Departement Due: 20/03/2020

Word

Now,

HDFS

```
said@said:~  
sai@said:~$ spark-shell  
20/03/22 08:41:57 WARN Utils: Your hostname, said resolves to a loopback address: 127.0.1.1; using 192.168.1.2 instead (on interface wlp3s0)  
20/03/22 08:41:57 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address  
20/03/22 08:41:58 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties  
Setting default log level to "WARN".  
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).  
Spark context Web UI available at http://192.168.1.2:4040  
Spark context available as 'sc' (master = local[*], app id = local-1584859328880).  
Spark session available as 'spark'.  
Welcome to  
Databricks version 3.0.0-preview2  
Using Scala version 2.12.10 (OpenJDK 64-Bit Server VM, Java 1.8.0_242)  
Type in expressions to have them evaluated.  
Type :help for more information.  
scala>
```

Notes

- You can work in groups of 2 or 3.

Good Luck

[sai@said: ~/Desktop/... [IMP Hadoop] [Apache Hadoop 3.2.1... lap3 sai@said: ~]

Applications Places System 80% 94 B/s 117 B/s Terminal File Edit View Search Terminal Tabs Help 13°C en 8:42, 22 مار 2020

lap3

Previous Next 2 (2 of 2) 200%

Departement Due: 20/03/2020

Word

Now,

HDFS

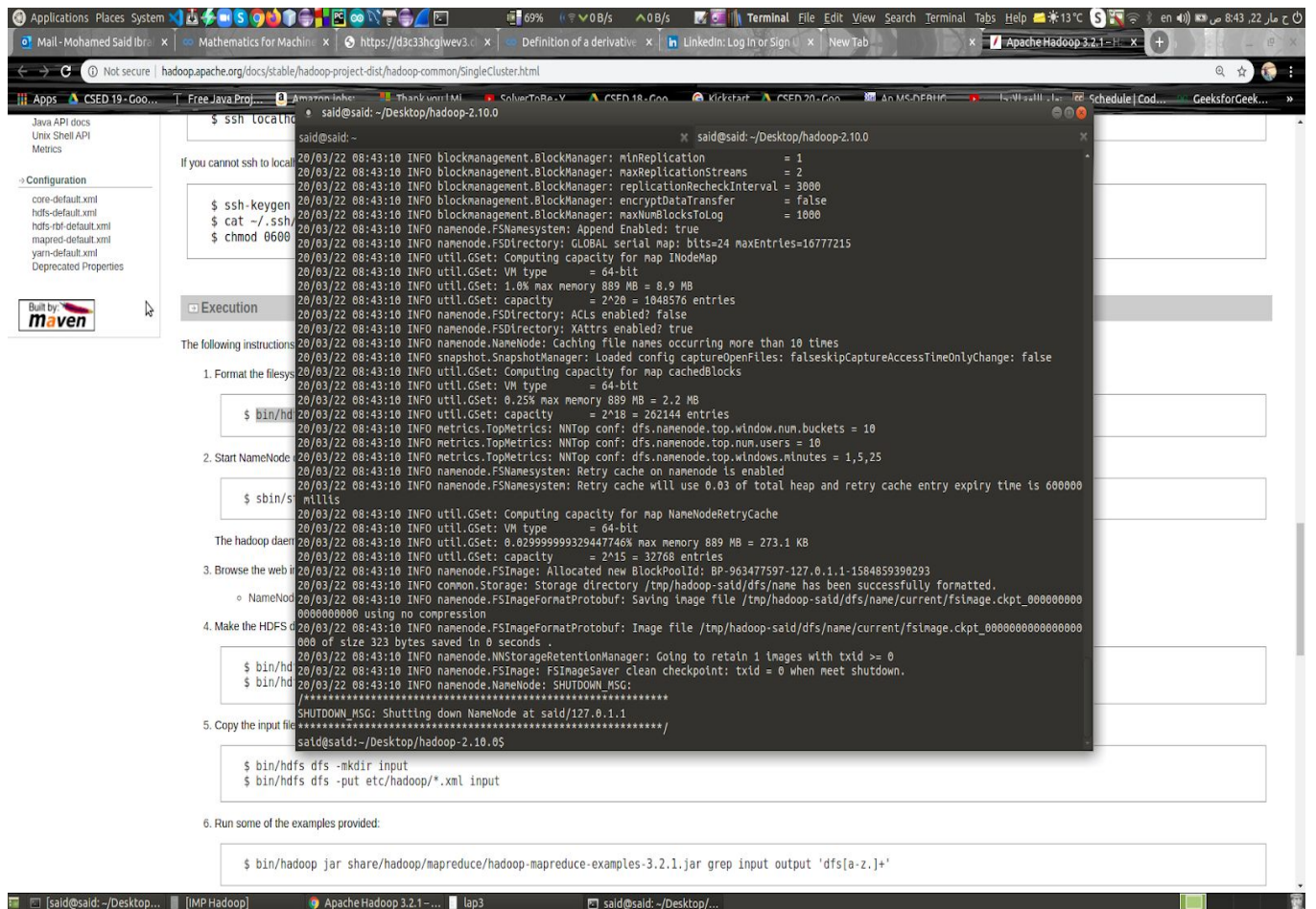
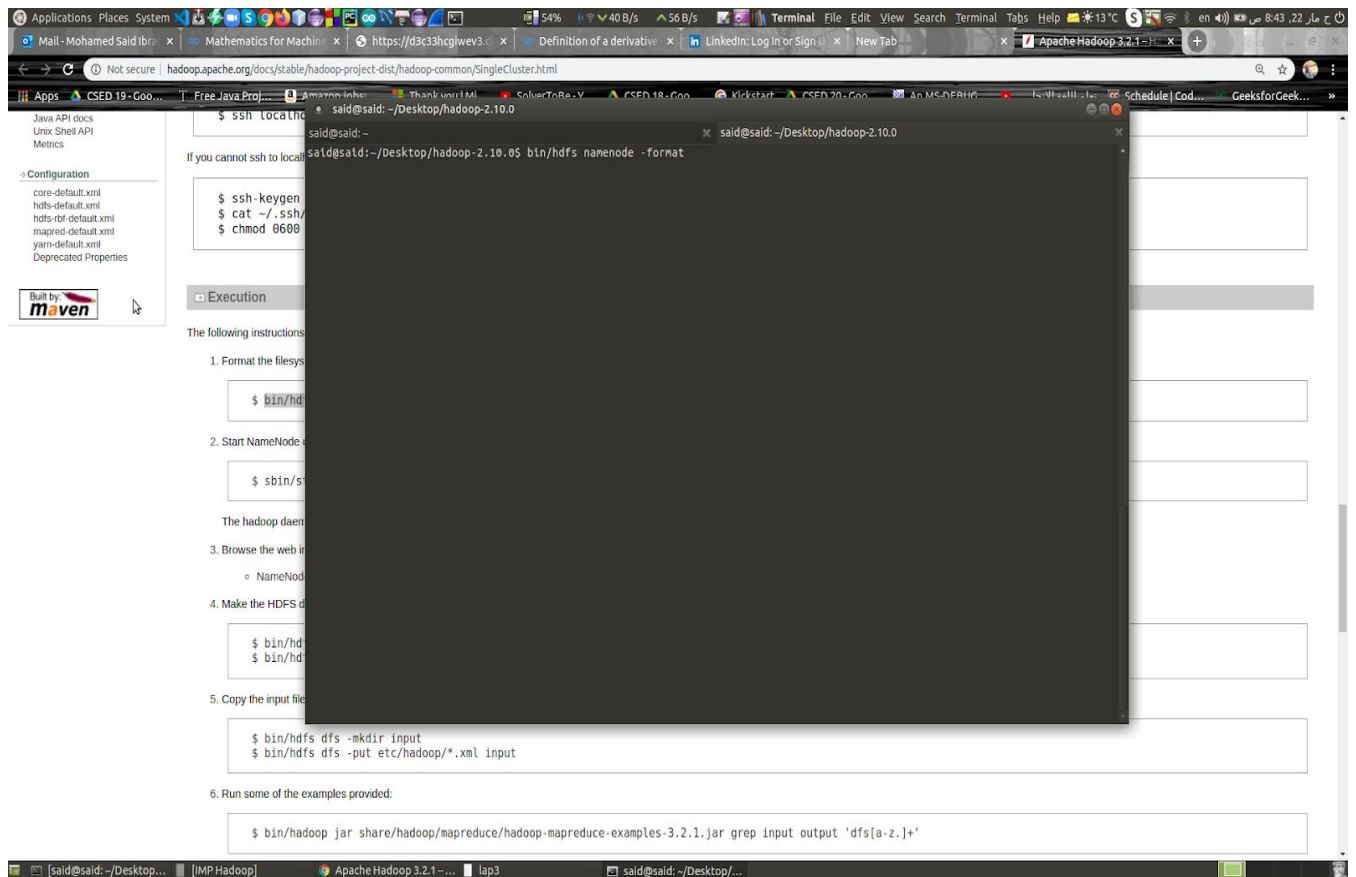
```
said@said:~/Desktop/hadoop-2.10.0  
sai@said: ~  
sai@said:~/Desktop/hadoop-2.10.0
```

Notes

- You can work in groups of 2 or 3.

Good Luck

[sai@said: ~/Desktop/... [IMP Hadoop] [Apache Hadoop 3.2.1... lap3 sai@said: ~/Desktop/...]







Applications Places System 24% 22 B/s 33 B/s Terminal File Edit View Search Terminal Tabs Help 13°C en 8:45, 22 ج ٢٠

Mail - Mohamed Said x Mathematics for Mac x https://d3c3hcg1wex x Definition of a derivat x LinkedIn: Log in or Sig x New Tab x Apache Hadoop 3.2.1 x Namenode Informa: x

localhost:50070/dfshealth.html#tab-overview

Apps CSED 19 - Goo... Free Java Proj... Amazon Inhe... Thank you! M... SolverToBe... CSED 18 - Goo... Kickstart CSED 20 - Goo... An MS-DOS... In-Vi-sual... Schedule | Cod... GeeksforGeek...

Hado said@said: ~

```

said@said: ~/Desktop/hadoop-2.10.0$ sbin/start-dfs.sh
Starting namenodes on [localhost]
said@localhost's password:
localhost: starting namenode, logging to /home/said/Desktop/hadoop-2.10.0/logs/hadoop-said-namenode-said.out
said@localhost's password:
localhost: starting datanode, logging to /home/said/Desktop/hadoop-2.10.0/logs/hadoop-said-datanode-said.out
Starting secondary namenodes [0.0.0.0]
said@0.0.0.0's password:
0.0.0.0: Permission denied, please try again.
said@0.0.0.0's password:
0.0.0.0: Permission denied, please try again.
said@0.0.0.0's password:
said@said: ~/Desktop/hadoop-2.10.0$ 0.0.0.0: Permission denied (publickey,password).
said@said: ~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -mkdir /user

```

DFS Remaining:	28.2 GB (19.84%)
Block Pool Used:	24 KB (0%)
DataNodes usages% (Min/Median/Max/stdDev):	0.00% / 0.00% / 0.00% / 0.00%
Live Nodes	1 (Decommissioned: 0, In Maintenance: 0)
Dead Nodes	0 (Decommissioned: 0, In Maintenance: 0)

[said@said: ~/Desktop... [IMP Hadoop] Namenode Informa... [lap3] said@said: ~/Desktop/...

Applications Places System 65% 17 B/s 72 B/s Terminal File Edit View Search Terminal Tabs Help 13°C en 8:45, 22 ج ٢٠

Mail - Mohamed Said x Mathematics for Mac x https://d3c3hcg1wex x Definition of a derivat x LinkedIn: Log in or Sig x New Tab x Apache Hadoop 3.2.1 x Namenode Informa: x

localhost:50070/dfshealth.html#tab-overview

Apps CSED 19 - Goo... Free Java Proj... Amazon Inhe... Thank you! M... SolverToBe... CSED 18 - Goo... Kickstart CSED 20 - Goo... An MS-DOS... In-Vi-sual... Schedule | Cod... GeeksforGeek...

Hado said@said: ~

```

said@said: ~/Desktop/hadoop-2.10.0$ sbin/start-dfs.sh
Starting namenodes on [localhost]
said@localhost's password:
localhost: starting namenode, logging to /home/said/Desktop/hadoop-2.10.0/logs/hadoop-said-namenode-said.out
said@localhost's password:
localhost: starting datanode, logging to /home/said/Desktop/hadoop-2.10.0/logs/hadoop-said-datanode-said.out
Starting secondary namenodes [0.0.0.0]
said@0.0.0.0's password:
0.0.0.0: Permission denied, please try again.
said@0.0.0.0's password:
0.0.0.0: Permission denied, please try again.
said@0.0.0.0's password:
said@said: ~/Desktop/hadoop-2.10.0$ 0.0.0.0: Permission denied (publickey,password).
said@said: ~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -mkdir /user
said@said: ~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -mkdir /user/said
said@said: ~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -copyFromLocal /home/said/Desktop/files2/ /user/said/

```

DFS Remaining:	28.2 GB (19.84%)
Block Pool Used:	24 KB (0%)
DataNodes usages% (Min/Median/Max/stdDev):	0.00% / 0.00% / 0.00% / 0.00%
Live Nodes	1 (Decommissioned: 0, In Maintenance: 0)
Dead Nodes	0 (Decommissioned: 0, In Maintenance: 0)

[said@said: ~/Desktop... [IMP Hadoop] Namenode Informa... [lap3] said@said: ~/Desktop/...

Applications Places System 64% 0.0 B/s 0.0 B/s Terminal File Edit View Search Terminal Tabs Help 13°C 8:46, 22 مار 2020

Mail - Mohamed Said x Mathematics for Mac x https://d3c33hgjwv... x Definition of a deriva x LinkedIn: Log In or Si x New Tab x Apache Hadoop 3.2.1 x Namenode Informa... x

localhost:50070/dfshealth.html#tab-overview

Apps CSED 19- Goo... Free Java Proj... Jammam-Inher... Thank you! M... SolverToBe... CSED 18- Goo... Kickstart... CSED 20- Goo... An MCS-DEPH... لا والله... Schedule | Cod... GeeksforGeek...

Terminal

```
said@said: ~/Desktop/WordCount/target
said@said: ~/Desktop/hadoop-2.10.0
said@said: ~/Desktop/WordCount/target

said@said:~/Desktop/hadoop-2.10.0$ cd ..
said@said:~/Desktop$ cd WordCount/
said@said:~/Desktop/WordCount$ ls
pom.xml src target
said@said:~/Desktop/WordCount$ cd target/
said@said:~/Desktop/WordCount/target$ ls
classes lib naven-archiver wordcount-1.0.jar
said@said:~/Desktop/WordCount/target$
```

DFS Remaining: 28.2 GB (19.84%)

Block Pool Used: 24 KB (0%)

DataNodes usages% (Min/Median/Max/stdDev): 0.00% / 0.00% / 0.00% / 0.00%

Live Nodes 1 (Decommissioned: 0, In Maintenance: 0)

Dead Nodes 0 (Decommissioned: 0, In Maintenance: 0)

[said@said: ~/Desktop/... [IMP Hadoop] Namenode Informa... [lap3] said@said: ~/Desktop/...

Applications Places System 53% 50 B/s 42 B/s Terminal File Edit View Search Terminal Tabs Help 13°C 8:49, 22 مار 2020

Mail - Mohamed Said x Mathematics for Mac x https://d3c33hgjwv... x Definition of a deriva x LinkedIn: Log In or Si x New Tab x Apache Hadoop 3.2.1 x Browsing HDFS x

localhost:50070/explorer.html#/user/said/input

Apps CSED 19- Goo... Free Java Proj... Jammam-Inher... Thank you! M... SolverToBe... CSED 18- Goo... Kickstart... CSED 20- Goo... An MCS-DEPH... لا والله... Schedule | Cod... GeeksforGeek...

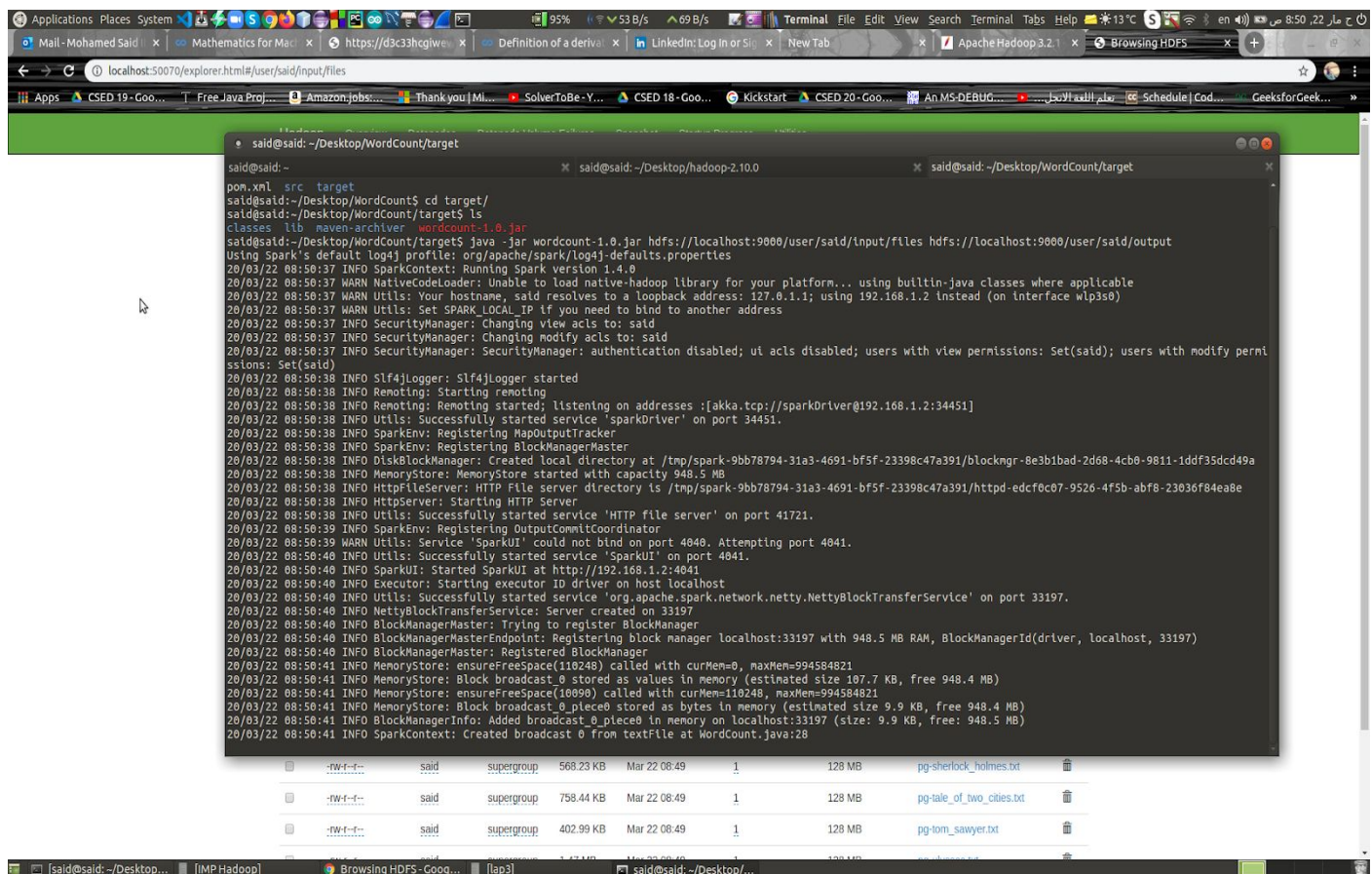
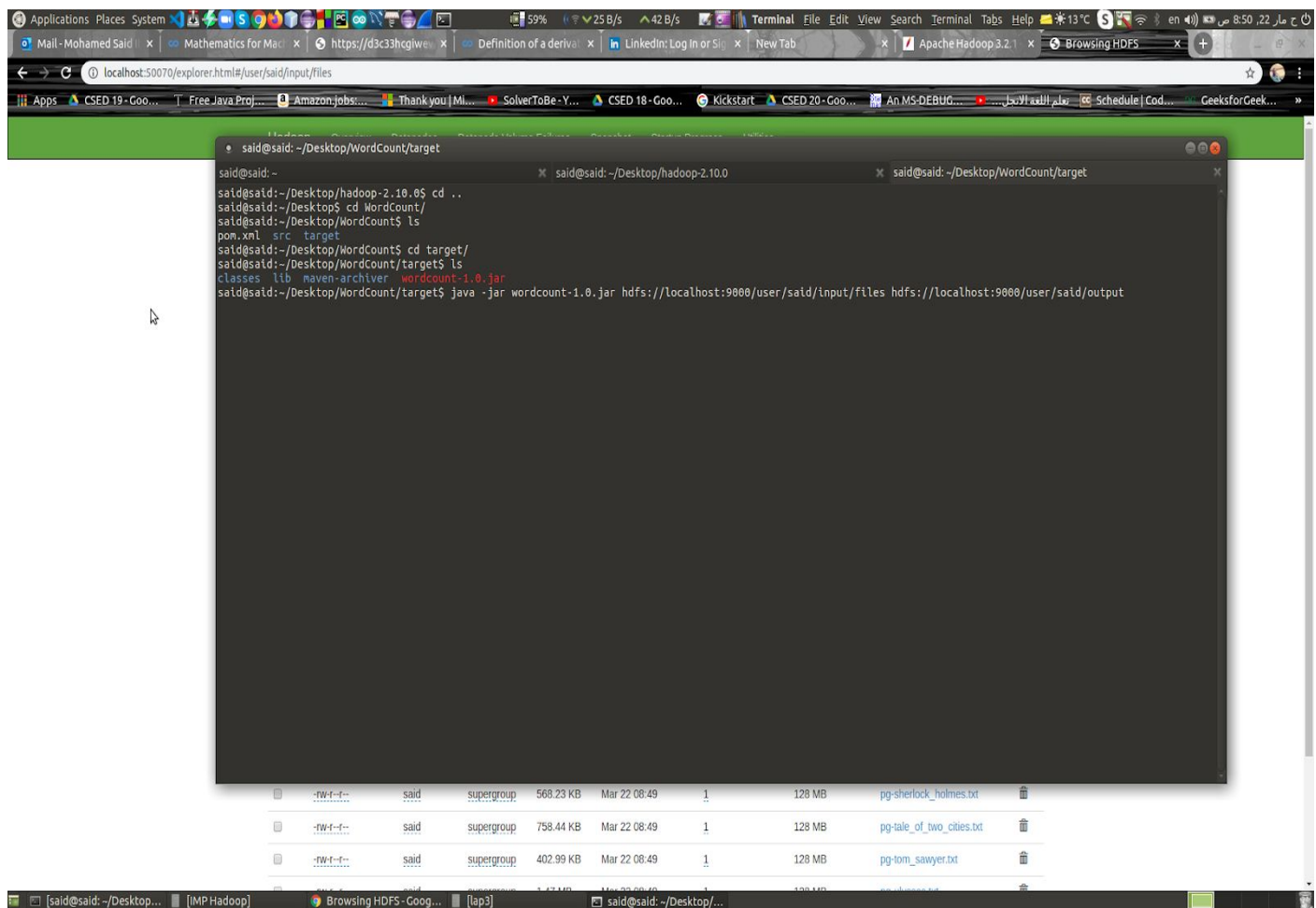
Terminal

```
said@said: ~/Desktop/hadoop-2.10.0
said@said: ~/Desktop/hadoop-2.10.0
said@said: ~/Desktop/WordCount/target

said@said:~/Desktop/hadoop-2.10.0$ bin/start-dfs.sh
Starting namenodes on [localhost]
said@localhost's password:
localhost: starting namenode, logging to /home/said/Desktop/hadoop-2.10.0/logs/hadoop-said-namenode-said.out
said@localhost's password:
localhost: starting datanode, logging to /home/said/Desktop/hadoop-2.10.0/logs/hadoop-said-datanode-said.out
Starting secondary namenodes [0.0.0.0]
said@0.0.0.0's password:
0.0.0.0: Permission denied, please try again.
said@0.0.0.0's password:
0.0.0.0: Permission denied, please try again.
said@0.0.0.0's password:
said@0.0.0.0's password:
said@said:~/Desktop/hadoop-2.10.0$ 0.0.0.0: Permission denied (publickey,password).

said@said:~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -mkdir /user
said@said:~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -mkdir /user/said
said@said:~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -copyFromLocal /home/said/Desktop/files2/ /user/said/
said@said:~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -mkdir /user/said/input
said@said:~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -copyFromLocal /home/said/Desktop/files/ /user/said/input
copyFromLocal: 'user/said/input': No such file or directory: 'hdfs://localhost:9000/user/said/user/said/input'
said@said:~/Desktop/hadoop-2.10.0$ bin/hdfs dfs -copyFromLocal /home/said/Desktop/files/ input
said@said:~/Desktop/hadoop-2.10.0$
```

[said@said: ~/Desktop/... [IMP Hadoop] Browsing HDFS - Goo... [lap3] said@said: ~/Desktop/...





Applications Places System 26% 0 B/s Terminal File Edit View Search Terminal Tabs Help 13°C en 8:50:22

Mail - Mohamed Said Mathematics for Mac https://d3c3hcgjwv... Definition of a deriva LinkedIn: Log in or Si New Tab Apache Hadoop 3.2.1 Browsing HDFS

localhost:50070/explorer.html#/user/said/input/files

Apps CSED 19 - Goo Free Java Proj Amazon.jobs Thank you | M... SolverToBe - Y... CSED 18 - Goo Kickstart CSED 20 - Goo An MS-DEBUG... علم اللغة الانجليز Schedule | Cod... GeeksforGeek...

Terminal

```
said@said: ~/Desktop/WordCount/target
k 202003220850_0001_m_000013
20/03/22 08:50:48 INFO SparkHadoopMapRedUtil: attempt 202003220850_0001_m_000013_29: Committed
20/03/22 08:50:48 INFO Executor: Finished task 13.0 in stage 1.0 (TID 29): 886 bytes result sent to driver
20/03/22 08:50:48 INFO TaskSetManager: Starting task 14.0 in stage 1.0 (TID 30, localhost, PROCESS_LOCAL, 1165 bytes)
20/03/22 08:50:48 INFO TaskSetManager: Finished task 13.0 in stage 1.0 (TID 29) in 84 ms on localhost (14/16)
20/03/22 08:50:48 INFO Executor: Running task 14.0 in stage 1.0 (TID 30)
20/03/22 08:50:48 INFO ShuffleBlockFetcherIterator: Getting 16 non-empty blocks out of 16 blocks
20/03/22 08:50:48 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
20/03/22 08:50:48 INFO FileOutputCommitter: Saved output of task 'attempt 202003220850_0001_m_000014_30' to hdfs://localhost:9000/user/said/output/_temporary/0/tas
k 202003220850_0001_m_000014
20/03/22 08:50:48 INFO SparkHadoopMapRedUtil: attempt 202003220850_0001_m_000014_30: Committed
20/03/22 08:50:48 INFO Executor: Finished task 14.0 in stage 1.0 (TID 30): 886 bytes result sent to driver
20/03/22 08:50:48 INFO TaskSetManager: Starting task 15.0 in stage 1.0 (TID 31, localhost, PROCESS_LOCAL, 1165 bytes)
20/03/22 08:50:48 INFO TaskSetManager: Finished task 14.0 in stage 1.0 (TID 30) in 93 ms on localhost (15/16)
20/03/22 08:50:48 INFO Executor: Running task 15.0 in stage 1.0 (TID 31)
20/03/22 08:50:48 INFO ShuffleBlockFetcherIterator: Getting 16 non-empty blocks out of 16 blocks
20/03/22 08:50:48 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
20/03/22 08:50:48 INFO FileOutputCommitter: Saved output of task 'attempt 202003220850_0001_m_000015_31' to hdfs://localhost:9000/user/said/output/_temporary/0/tas
k 202003220850_0001_m_000015
20/03/22 08:50:48 INFO SparkHadoopMapRedUtil: attempt 202003220850_0001_m_000015_31: Committed
20/03/22 08:50:48 INFO Executor: Finished task 15.0 in stage 1.0 (TID 31): 886 bytes result sent to driver
20/03/22 08:50:48 INFO TaskSetManager: Finished task 15.0 in stage 1.0 (TID 31) in 91 ms on localhost (16/16)
20/03/22 08:50:48 INFO DAGScheduler: ResultsStage 1 (saveAsTextFile at WordCount.java:52) finished in 3.294 s
20/03/22 08:50:48 INFO TaskSchedulerImpl: Removed TaskSet 1.0, whose tasks have all completed, from pool
20/03/22 08:50:48 INFO DAGScheduler: Job 0 finished: saveAsTextFile at WordCount.java:52, took 6.396979 s
20/03/22 08:50:49 INFO SparkContext: Invoking stop() from shutdown hook
20/03/22 08:50:49 INFO SparkUI: Stopped Spark web UI at http://192.168.1.2:4041
20/03/22 08:50:49 INFO DAGScheduler: Stopping DAGScheduler
20/03/22 08:50:49 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
20/03/22 08:50:49 INFO Utils: path = /tmp/spark-9bb78794-31a3-4691-bf5f-23398c47a391/blockmgr-8e3b1bad-2d68-4c6b-9811-1ddf35dcd49a, already present as root for del
etion.
20/03/22 08:50:49 INFO MemoryStore: MemoryStore cleared
20/03/22 08:50:49 INFO BlockManager: BlockManager stopped
20/03/22 08:50:49 INFO BlockManagerMaster: BlockManagerMaster stopped
20/03/22 08:50:49 INFO SparkContext: Successfully stopped SparkContext
20/03/22 08:50:49 INFO Utils: Shutdown hook called
20/03/22 08:50:49 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
20/03/22 08:50:49 INFO Utils: Deleting directory /tmp/spark-9bb78794-31a3-4691-bf5f-23398c47a391
20/03/22 08:50:49 INFO RemoteActorRefProvider$RemotingTerminator: Shutting down remote daemon.
20/03/22 08:50:49 INFO RemoteActorRefProvider$RemotingTerminator: Remote daemon shut down; proceeding with flushing remote transports.
20/03/22 08:50:49 INFO RemoteActorRefProvider$RemotingTerminator: Remoting shut down.
said@said: ~/Desktop/WordCount/target$
```

ls -l

-rw-r--r--	said	supergroup	568.23 KB	Mar 22 08:49	1	128 MB	pg-sherlock_holmes.txt			
-rw-r--r--	said	supergroup	758.44 KB	Mar 22 08:49	1	128 MB	pg-tale_of_two_cities.txt			
-rw-r--r--	said	supergroup	402.99 KB	Mar 22 08:49	1	128 MB	pg-tom_sawyer.txt			

Applications Places System 70% 27 B/s 24 B/s 13°C en 8:51:22

Mail - Mohamed Said Mathematics for Mac https://d3c3hcgjwv... Definition of a deriva LinkedIn: Log in or Si New Tab Apache Hadoop 3.2.1 Browsing HDFS

localhost:50070/explorer.html#/user/said/

Apps CSED 19 - Goo Free Java Proj Amazon.jobs Thank you | M... SolverToBe - Y... CSED 18 - Goo Kickstart CSED 20 - Goo An MS-DEBUG... علم اللغة الانجليز Schedule | Cod... GeeksforGeek...

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

## Browse Directory

/user/said/ Go

Show 25 entries Search:

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
	drwxr-xr-x	said	supergroup	0 B	Mar 22 08:45	0	0 B	files2	
	drwxr-xr-x	said	supergroup	0 B	Mar 22 08:49	0	0 B	input	
	drwxr-xr-x	said	supergroup	0 B	Mar 22 08:50	0	0 B	output	

Showing 1 to 3 of 3 entries

Previous 1 Next

Hadoop, 2019.



Applications Places System 70% 95 B/s 67 B/s 13°C en 8:51 22 مار 2019

Mail - Mohamed Said Mathematics for Mac https://d3c33hcgwe... Definition of a deriva... LinkedIn: Log in or Si... New Tab Apache Hadoop 3.2.1 Browsing HDFS

localhost:50070/explorer.html#/user/said/output

user/said/output Go

Show 25 entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	said	supergroup	0 B	Mar 22 08:50	3	128 MB	_SUCCESS
-rw-r--r--	said	supergroup	133.67 KB	Mar 22 08:50	3	128 MB	part-00000
-rw-r--r--	said	supergroup	133.72 KB	Mar 22 08:50	3	128 MB	part-00001
-rw-r--r--	said	supergroup	137.17 KB	Mar 22 08:50	3	128 MB	part-00002
-rw-r--r--	said	supergroup	135.62 KB	Mar 22 08:50	3	128 MB	part-00003
-rw-r--r--	said	supergroup	135.52 KB	Mar 22 08:50	3	128 MB	part-00004
-rw-r--r--	said	supergroup	136.06 KB	Mar 22 08:50	3	128 MB	part-00005
-rw-r--r--	said	supergroup	134.29 KB	Mar 22 08:50	3	128 MB	part-00006
-rw-r--r--	said	supergroup	135.92 KB	Mar 22 08:50	3	128 MB	part-00007
-rw-r--r--	said	supergroup	136.31 KB	Mar 22 08:50	3	128 MB	part-00008
-rw-r--r--	said	supergroup	135.79 KB	Mar 22 08:50	3	128 MB	part-00009
-rw-r--r--	said	supergroup	135.95 KB	Mar 22 08:50	3	128 MB	part-00010
-rw-r--r--	said	supergroup	133.96 KB	Mar 22 08:50	3	128 MB	part-00011
-rw-r--r--	said	supergroup	134.27 KB	Mar 22 08:50	3	128 MB	part-00012
-rw-r--r--	said	supergroup	135.47 KB	Mar 22 08:50	3	128 MB	part-00013
-rw-r--r--	said	supergroup	133.38 KB	Mar 22 08:50	3	128 MB	part-00014
-rw-r--r--	said	supergroup	134.39 KB	Mar 22 08:50	3	128 MB	part-00015

Showing 1 to 17 of 17 entries

Previous 1 Next

Hadoop, 2019.

[said@said: ~/Desktop/...] [IMP Hadoop] Browsing HDFS - Goog... [lap3] said@said: ~/Desktop/...

part-00001 (~Downloads) - Pluma

Open Save Undo Cut Copy Paste Find

part-00001

```
1 (she.,1)
2 (seriousness,,2)
3 (small,,2)
4 (desolated,,2)
5 (massproduct,,1)
6 (monologue,,4)
7 (bardoor,,1)
8 (obediently,,3)
9 (maiden,,2)
10 (parlour,,10)
11 (making,,4)
12 (reinforced,,3)
13 ("Gimme,,4)
14 (unintentionally,,1)
15 ("Nothing,,8)
16 (_Sonnez,,1)
17 (onerous,,2)
18 (peremptory,,1)
19 (practise,,2)
20 (toll,,1)
21 (speeding,,1)
22 (scale,,1)
23 (kink,,1)
24 (negro,,3)
25 (deceitver,,1)
26 (revision,,1)
27 (Sandycovel,,1)
28 (if-what,,1)
29 (boldest,,4)
30 (chewing,,5)
31 (really,,29)
32 (fairness,,1)
33 (gourme,,1)
34 (interlaced,,1)
35 (HABES,,1)
36 (tillac,,9)
37 (harshly,,1)
38 (pouff,,1)
39 (cabin-scuttle,,1)
40 (urgent,,12)
41 (tolls-one,,1)
42 (milk,,17)
43 (shape,,86)
44 (Domini,,1)
45 (Piles,,2)
46 (Elephant,,4)
47 (anuck,,1)
48 (propitiate-when,,1)
49 (Hacking,,1)
50 (Sarsfield,,1)
51 (tan,,10)
52 (seducer,,1)
53 (landmark,,1)
```

Plain Text Tab Width: 4 Ln 1, Col 1 INS

[said@said: ~/Desktop/...] [IMP Hadoop] Browsing HDFS - Goog... [lap3] said@said: ~/Desktop/... part-00001 (~Downloads)

