

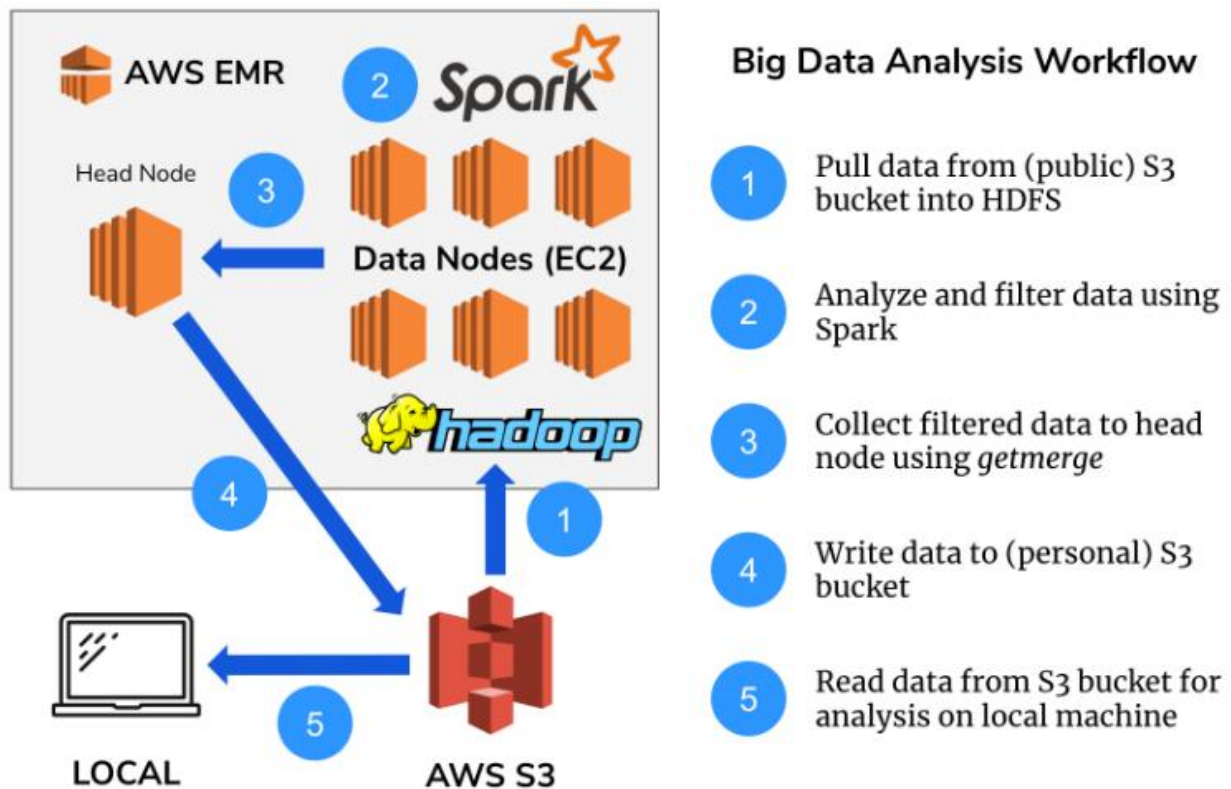
Harsh Sharma

Big Data assignment (Unit 4) Report

5 December 2022

Submitted To: Borna Ghotbi, Fallyn Thompson, Arati Santhanakrishnan

Prepare a professional report to summarize the findings and be sure to include an appendix with screenshots of the steps completed for Questions 1 and 2.



The following screenshots/ snippets are the replication of the above flow process.

Question 1 Creating a cluster

[Create cluster](#) [View details](#) [Clone](#) [Terminate](#)

Filter: All clusters 3 clusters (all loaded)

	Name	ID	Status	Creation time (UTC-8)	Elaps
<input type="checkbox"/>	bigdataassgn	j-13TMCCCKGI49QX	Waiting Cluster ready	2022-12-02 13:16 (UTC-8)	1 day

If you want to create a cluster, you can click on the create cluster. For this assignment I created **bigdataassgn** as my cluster.

Question 2 Connect to the head node of the cluster using SSH.

Step 1: Open an SSH Tunnel to the Amazon EMR Master Node - [Learn more](#)

Windows

Mac / Linux

1. Open a terminal window. On Mac OS X, choose Applications > Utilities > Terminal. On other Linux distributions, terminal is found at Applications > Accessories > Terminal.
2. To establish an SSH tunnel with the master node using dynamic port forwarding, type the following command. Replace ~/spark_nov.pem with the location and filename of the private key file (.pem) used to launch the cluster.

```
ssh -i ~/spark_nov.pem -ND 8157 hadoop@ec2-3-144-102-165.us-east-2.compute.amazonaws.com
```

```
user@LAPTOP-JH75U61V MINGW64 ~/desktop
$ ssh -i spark_nov.pem hadoop@ec2-3-144-107-186.us-east-2.compute.amazonaws.com
Last login: Wed Nov 30 01:21:45 2022

 _ _ _ _ _
 _ | ( _ /   Amazon Linux 2 AMI
 _|_|_|_|_|_|

https://aws.amazon.com/amazon-linux-2/
79 package(s) needed for security, out of 113 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRRRRRRRRR
E::::::::::::::::::E M::::::::M          M::::::::M R:::::::::R
EE::::::::::::::::::E M::::::::M          M::::::::M R:::::::::R
E::::E      EEEEE M::::::::M          M::::::::M RR::::R      R::::R
E::::E      M::::::::M M::M M::M M::M M::M M::M R::RRRRR::::R
E::::::::::::::::::E M::::::::M M::M M::M M::M M::M R:::::::::RR
E::::E      M::::::::M M::M M::M M::M M::M R::RRRRR::::R
E::::E      M::::::::M M::M M::M M::M M::M R::R      R::::R
E::::E      EEEEE M::::::::M MMM M::M M::M R::R      R::::R
EE::::::::::::::::::E M::::::::M          M::::::::M R::R      R::::R
E::::::::::::::::::E M::::::::M          M::::::::M RR::::R      R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRR      RRRRRR

[hadoop@ip-172-31-10-230 ~]$ |
```

From the snippet, I have connected to server on my local windows machine.

Question 3 Copy the data folder from the S3 bucket *directly* into a directory on the Hadoop File System (HDFS) named `/user/hadoop/eng_1M_1gram`

1. Checking if I have anything in the hdf5.

```

EEEEEEEEEEEEEEEEEEEE MMMMMMMM                      MMMMMMMM RRRRRRRR                      RRRRRR

[hadoop@ip-172-31-15-205 ~]$ hadoop fs -ls
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -lsr
lsr: DEPRECATED: Please use 'ls -R' instead.
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -ls -R
[hadoop@ip-172-31-15-205 ~]$ client_loop: send disconnect: Connection reset by peer
hh(base)

```

From above, it can be observed that there is nothing in the hdf5.

2. Create a directory and do a sanity check

```
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -mkdir /user/hadoop/eng_1M_1gram
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -ls
Found 1 items
drwxr-xr-x  - hadoop hadoop          0 2022-12-02 23:56 eng_1M_1gram
```

So now I have eng_1M_1gram directory in hdfs.

3. Copying the .csv file from the S3 brainstation bucket DIRECTLY to the HDFS using ***distcp***.

```
[hadoop@ip-172-31-15-205 ~]$ hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram/eng_1M_1gram.csv
2022-12-03 00:00:34,547 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, userDiff=false, fromSnapshot=null, toSnapshot=null, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=0.0, copyStrategy='uni
SIZE', atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3://brainstation-dsft/eng_1M_1gram.csv], targetPath=/user/hadoop/eng_1M_1gram/er
l', blocksPerChunk=0, copyBuffer=Size=8192, verboseLog=false, directWrite=false}, sourcePaths=[s3://brainstation-dsft/eng_1M_1gram.csv], targetPathExists=false,
2022-12-03 00:00:34,821 INFO client.RWProxy: Connecting to ResourceManager at ip-172-31-15-205.us-east-2.compute.internal/172.31.15.205:8032
2022-12-03 00:00:34,998 INFO client.AHSProxy: Connecting to Application History server at ip-172-31-15-205.us-east-2.compute.internal/172.31.15.205:10200
2022-12-03 00:00:38,517 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
2022-12-03 00:00:38,517 INFO tools.SimpleCopyListing: Build file listing completed
```

- Sanity Check to see if the data is really copied or not.

```
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -ls
Found 1 items
drwxr-xr-x   - hadoop hadoop          0 2022-12-03 00:01 eng_1M_1gram
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -lsr
lsr: DEPRECATED: Please use 'ls -R' instead.
drwxr-xr-x   - hadoop hadoop          0 2022-12-03 00:01 eng_1M_1gram
-rw-r--r--   1 hadoop hadoop 5292105197 2022-12-03 00:01 eng_1M_1gram/eng_1M_1gram.csv
[hadoop@ip-172-31-15-205 ~]$ pwd
/home/hadoop
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -ls /user/hadoop/eng_1M_1gram/
Found 1 items
-rw-r--r--   1 hadoop hadoop 5292105197 2022-12-03 00:01 /user/hadoop/eng_1M_1gram/eng_1M_1gram.csv
```

You can see that eng_1M_1gram.csv file is sitting in eng_1M_1gram directory in HDFS.

5. Make a new dataframe and writing it into the directory.

```
3]: df2.write.csv("/user/hadoop/eng_1M_1gram/datatoken.csv", header=True)
```

► Spark Job Progress

After doing this I wanted it to appear in the directory in HDFS.

```
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -ls
Found 2 items
drwxr-xr-x - hadoop hadoop 0 2022-12-03 00:01 eng_1M_1gram
drwxr-xr-x - livy hadoop 0 2022-12-03 18:16 eng_1M_1gram2
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -lsr
lsr: DEPRECATED: Please use 'ls -R' instead.
drwxr-xr-x - hadoop hadoop 0 2022-12-03 00:01 eng_1M_1gram
-rw-r--r-- 1 hadoop hadoop 5292105197 2022-12-03 00:01 eng_1M_1gram/eng_1M_1gram.csv
drwxr-xr-x - livy hadoop 0 2022-12-03 18:16 eng_1M_1gram2
drwxr-xr-x - livy hadoop 0 2022-12-03 18:17 eng_1M_1gram2/datatoken.csv
-rw-r--r-- 1 livy hadoop 0 2022-12-03 18:17 eng_1M_1gram2/datatoken.csv/_SUCCESS
-rw-r--r-- 1 livy hadoop 33 2022-12-03 18:16 eng_1M_1gram2/datatoken.csv/part-00000-88ec2ae9-7637-400c-91f3-9f67125cfc84-c000.csv
-rw-r--r-- 1 livy hadoop 7305 2022-12-03 18:16 eng_1M_1gram2/datatoken.csv/part-00023-88ec2ae9-7637-400c-91f3-9f67125cfc84-c000.csv
```

The datatoken.csv file was saved as a directory itself as it can be observed that memory occupied by it is 0. However, 2 separate csv are created as can be seen in last 2 lines.

6. Next Step: Merging these 2 csvs' on the local drive of the head node.

```
[hadoop@ip-172-31-15-205 ~]$ hadoop fs -getmerge /user/hadoop/eng_1M_1gram2/datatoken.csv/part-00000-88ec2ae9-7637-400c-91f3-9f67125cfc84-c000.csv /user/hadoop/eng_1M_1gram2/datatoken.csv/part-00023-88ec2ae9-7637-400c-91f3-9f67125cfc84-c000.csv mergedfile.csv
[hadoop@ip-172-31-15-205 ~]$ cat mergedfile.csv
token,year,frequency,pages,books
token,year,frequency,pages,books
data,1584,16,14,1
data,1614,3,2,1
data,1627,1,1,1
data,1631,22,18,1
data,1637,1,1,1
data,1638,2,2,1
data,1640,1,1,1
data,1642,1,1,1
data,1644,4,4,1
data,1647,1,1,1
data,1651,1,1,1
data,1674,1,1,1
```

Also after doing the sanity check you can see (using cat mergedfile.csv) that the 2 csvs' are now combined/merged.

However, the name of the columns also appear 2 times as can be seen from above.

7. Let's move this into the S3 bucket (forbigdata3) and then try to open it without spark in jupyter notebook with pandas.

```
ERROR: Invalid argument type
[hadoop@ip-172-31-15-205 ~]$ aws s3 cp mergedfile.csv s3://forbigdata3
upload: ./mergedfile.csv to s3://forbigdata3/mergedfile.csv
[hadoop@ip-172-31-15-205 ~]$ aws configure
```

forbigdata3 Info

Objects

Properties

Permissions

Metrics

Management

Access Points

Objects (1)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)



Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Find objects by prefix

< 1 > ⚙

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	mergedfile.csv	csv	December 3, 2022, 13:09:45 (UTC-08:00)	7.2 KB	Standard

Now you can open this mergedfile.csv into jupyter notebook using Pandas.

Question 8: Compare Hadoop and Spark as distributed file systems.

- What are the advantages/ differences between Hadoop and Spark? List two advantages for each.
- Explain how the HDFS stores the data.

1. Advantages: Hadoop

- It is very highly scalable storage platform because it can store and distribute very large datasets across hundreds of servers that work in parallel.
- Hadoop is fault tolerant. When data is sent to a individual node, that data is also replicated to other nodes in the cluster, which means that in the event of failure, there is another copy available to use.

Advantages: Spark

- a) Spark has all the advantages as Hadoop has. But Spark offers outstanding streaming with the latest methods for writing and maintaining streaming code. By spark streaming, you can use a variety of applications such as stream mining, network optimization, real time scoring of analytical models and more.
- b) Spark provides python addition, also called PySpark which let programmers to interface with the Spark framework and learn how to manipulate data at a large scale.

Difference between Spark and Hadoop.

Unlike Hadoop, Spark does not write to disk but keeps data in memory and so is typically much faster for performing data engineering, analysis, and modeling tasks within Big Data environments and with large data sets.

Spark runs 100 times faster in memory and ten times faster on disk than Hadoop MapReduce.

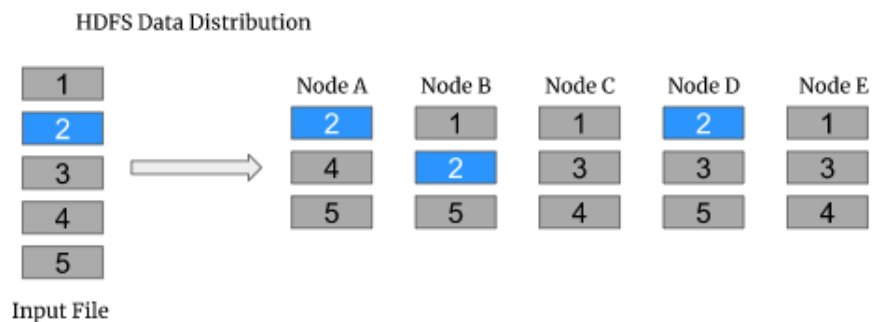
Question 8.b Explain how HDFS store the data.

Answer:

In HDFS, data is split into blocks typically 64MB or 128MB. And these blocks are split into machines at load time. Also blocks are replicated across multiple machines as shown below.

HDFS - Data Replication

Default replication is 3-fold



NameNode keeps track of which blocks make up a file and where they are stored.

