# Big Data Wrangling With Google Books Ngrams

Author: Tilova Shahrin

1. On AWS Management Console, click on the search bar and type 'EMR'. Click on EMR.
   a. Cloning cluster from spark lab.
   b. My_spark_cluster_2 is the cluster we will clone from, since it has all the features we need from spark lab.

   c.


   d. Click on 'Clone Cluster'



   e.

2. Connect to the head node of the cluster using SSH.
   a. In your terminal, make sure you are in the same directory as your .pem file.
   b. Follow these instructions under 'Connect to the primary node using SSH'. Paste the command onto the terminal. You can find your primary node public DNS in your cluster page, see the following command and image.

c. *ssh -i aws_cloud_tilova.pem -L 9995:localhost:9443*

*hadoop@ec2-18-118-27-195.us-east-2.compute.amazonaws.com*



d.

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`.

   a. Add this command onto the terminal: *hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram.*

   b. You can double check if your file exists in terminal using *hadoop fs -ls*

   ```
   [[hadoop@ip-172-31-7-10 ~]$ hadoop fs -ls
   Found 1 items
   -rw-r--r--   1 hadoop hdfsadmingroup 5292105197 2024-04-01 02:41 eng_1M_1gram.csv
   ```
   c.

4. In your cluster page, click on the Applications tab and check your port for JupyterHub.



   a. Type exit in your terminal and apply this command (follow pem file). *ssh -i aws_cloud_tilova.pem -L 9995:localhost:9443*

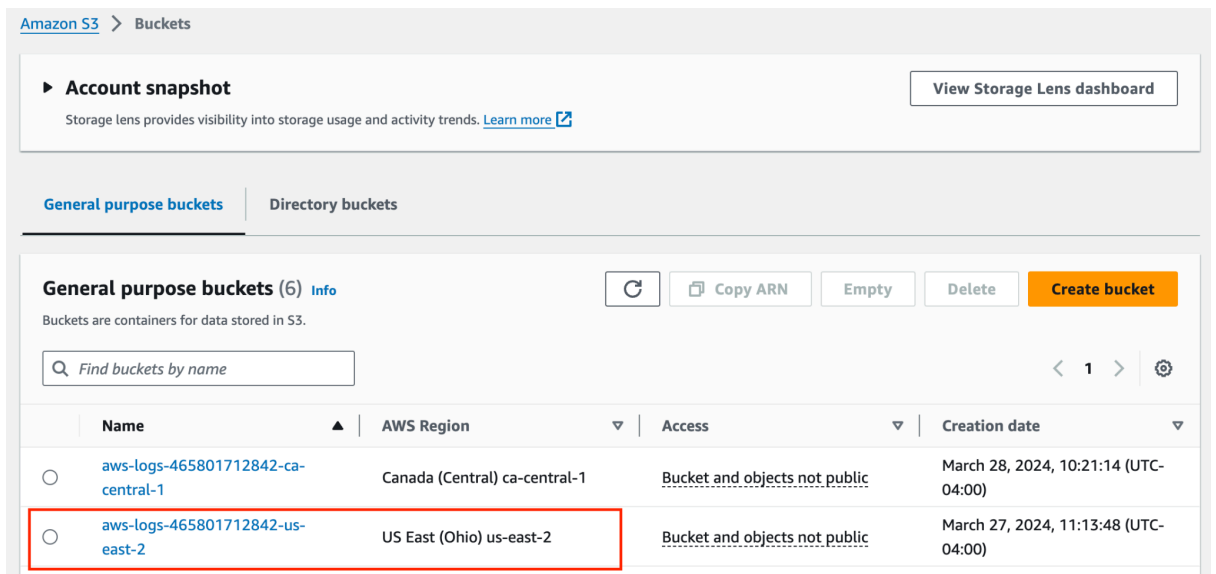   *hadoop@ec2-18-118-27-195.us-east-2.compute.amazonaws.com*

b. Make sure the port matches your jupyter hub port. Once it's connected click on the link and proceed to 'Advanced' button and click on the link.

c. Once you enter the login page, use jovyan as username and jupyter as password.

d. When starting a new notebook, make sure it is a PySpark kernel.



e.

Appendix F

f. Continue onto the jupyter notebook.

5. On the terminal, start to merge the filtered data frame using this command.

a. *hadoop fs -getmerge hdfs://ec2-18-118-205-216.us-east-2.compute.amazonaws.com:8020/user/hadoop/eng _1M_1gram file_name_filtered_data.csv*

b. Apply the csv file onto the s3 bucket using this command.
*aws s3 mv file_name_filtered_data.csv s3://aws-logs-465801712842-us-east-2/eng_1M_1gram.csv*

c. Then, check if the csv file is in the correct s3 bucket using this command.
*aws s3 ls s3://aws-logs-465801712842-us-east-2*

d. You should be able to see your csv file in your aws s3 bucket. Click on the appropriate bucket.

e.

f. Then click on the file you've saved and click the download button to open your csv file.



g.

6. In the same page from the image shown above, grab your S3 URI and copy it so you can read csv from your local jupyter notebook. Check Jupyter Notebook - 'Big Data Local - Tilova Shahrin.ipynb'

7. Check Jupyter Notebook - 'Big Data Local - Tilova Shahrin.ipynb'

8. Hadoop integrates with external libraries to provide machine learning capabilities. Spark has built-in machine learning libraries. You would need Spark to run jupyter hub, which runs much faster than Hadoop. Hadoop runs at a lower cost since it relies on any disk storage type for data processing. Spark uses RAM for in memory processing so it costs more than Hadoop.