Big Data Wrangling With Google Books Ngrams

**Big Data Wrangling With Google Books Ngrams**

Author: Ja`Mone Bridges

Class: Data Science Bootcamp

Institution: BrainStation

Big Data Wrangling With Google Books Ngrams

## Table of Contents

Big Data Wrangling With Google Books Ngrams

Big Data Wrangling With Google Books Ngrams

## Introduction

The purpose of this report is to instruct the user how to use AWS (Amazon Web Services) cloud services to process large datasets and databases. We will spin up (startup) a cluster with three primary services. The first service, Hadoop, allows for large volumes of data to be stored and retrieved at high speed. The second service, PySpark and JupyterHub allows users to perform analysis with RDD (Resilient Distributed Dataset). The third service, S3, is AWS secure cloud data storage that can be public or private. S3 allows access from the web to read or write data.

This report will refer to AWS setup pages and CLI (Command Line Interface) commands. The images of the AWS and CLI are in the appendix. The report will reference the appendix by letters.

Big Data Wrangling With Google Books Ngrams

Instructions

- **Access AWS**: Goto URL https://aws.amazon.com/ and click on "Sign in to the Console". If you need to create a new account, please do so now.

- **Access EMR**: After logging on to AWS, in the upper left is a search bar. Search for EMR and select it. See appendix A.

- **Setup Cluster**: Make sure that Clusters are selected under "EMR on EC2" on the left side panel. Then click on "Create cluster". Refer to appendix B for display.

- **Cluster Configuration**: In the "Name and Applications" field, give the cluster a name. The name can be anything suitable. Under "Amazon EMR release" select "emr-6.10.0". For the Application bundle select Custom. All of the following checkboxes must be selected, "Hue – Livy – Spark – Hadoop – Hive – JupyterHub". The "Operating system option" has "Amazon Linux release". Please, refer to appendix C.

- **Cluster CPU Core Selection**: For "Cluster configuration" select the radio button "Uniform instance groups". Under the "Primary" header and "Choose EC2 instance type" select "m5.xlarge". Under "Core" and "Choose EC2 instance type" select "m5.xlarge". Please, refer to appendix D for a visual.

- **Cluster scaling and provisioning:** Select option "Set cluster size manually". For "Provisioning configuration" type 2 under "Instance(s) size". Under header "Networking" pick options for "Virtual private cloud" and "Subnet". If no options are available, then click on the respective "Create" button on the right. Please, refer to appendix E if a visual guide is desired.

Big Data Wrangling With Google Books Ngrams

- **Cluster Security:** Under "Security configuration and EC2 key pair" look for heading "Amazon EC2 key pair for SSH to the cluster". If a key pair does not exist, then click on the "Create key pair" button to make one. Be sure to select the *.pem option when making and new key pair. NOTE: Please, be aware that the key pair will only work in the AWS region that it is created for. Store the downloaded key *.pem file in a location to be used later to connect to the cluster. Next, browse and select the key pair under "Amazon EC2 key pair for SSH to the cluster". For a visual aid please refer to appendix F.

- **Cluster Identity and Access Management (IAM) roles:** If there is an existing role select it for "Service Role" and "Instance profile", then select "Create cluster" button. Otherwise, click on "Create a service role" and follow the instructions. Once the service role has been created, select and perform the same process for "Instance profile". Now click on the "Create cluster" button. Please refer to append F for a visual aid.

- **EMR and Spark Cluster:** The cluster has now been created and the system is now installing the software and booting the system. It will take several minutes for the process to complete. Under "Status and time" on the right side of the window "Status" will change to "Waiting" when the cluster is ready to be used. Please, see appendix G for a reference.

- **Connect to Hadoop with SSH:** Once the status changes to "Waiting", under "Cluster management" click on "Connect to the Primary node using SSH" and copy the connection link. Open a CLI or terminal with ssh access and paste the connection link. The address of the key pair (*.pem file) will likely require changing. Once the address to the key pair is correct press enter. The window will appear the same as appendix H.

Big Data Wrangling With Google Books Ngrams

- **Large File Transfer to Hadoop:** On the CLI input command "hadoop distcp
  s3://brainstation-dsft/eng_1M_1gram/eng_1M_1gram.csv". Now the data will be
  downloaded and stored directly to the child nodes of Hadoop. For a visual please refer to
  appendix I and J.

- **Setup JupyterHub Secure Port for Localhost:** Open a new CLI or terminal. Type in
  command "ssh -I /Path/To/Key.pem -L 9995:localhost:9443
  hadoop@xxxxxxxxxxxx.compute.amazonaws.com". Replace the x's with the connect
  link portion used to connect to Hadoop via SSH. The window will be like appendix K.

- **Connect to JupyterHub in Browser:** Open a web browser and type in the following
  URL https://localhost:9995. A warning may appear regarding the expiration of the
  certificate, go to advance, and proceed to the site. At the logon page enter username
  "jovyan" and password "jupyter". A successful login will be like appendix N. Click new
  in the upper right corner, then select "PySpark".

- **PySpark:** Please refer to appendix L and M. Type spark in the first line and run to start
  the spark session. Follow the commands in appendix L and M to open the data file on the
  Hadoop server. Then perform a sanity check on the data to ensure that it's loaded. Now,
  check the number of rows and columns. The shape of the dataset is (261,823,225, 5).
  The column names are "token, year, frequency, pages, and books". To run a SQL
  statement on a PySpark dataframe it must be converted to a view. See appendix M In [8]
  for the command. A SQL statement of "SELECT * FROM eng_table WHERE token
  LIKE '%data%' is ran in order to get all the tokens with the work data. The data is stored
  in a new PySpark dataframe. The count of the new dataframe is 24,642. A sanity check is
  also performed on the new dataframe. The new dataframe is then written back to the

Big Data Wrangling With Google Books Ngrams

Hadoop child nodes with command " newdf.write.csv(path, header=True). See appendix

M for the path of how to store the data to Hadoop. PySpark is then stopped with

command spark.stop().

- **Retrieve the Data from Hadoop to the local machine:** In the CLI that is connected to

  Hadoop type the command "hadoop fs -getmerge /user/hadoop/eng_1M_1gram/

  eng_token_data.csv". This will combine all the files into a single file and store it on the

  cluster Linux file system. To gain access to the file for the local machine, it is stored in

  an S3 bucket. Setup IAM with a user and setup the user with a key ID and secret. The

  key ID and secret must be used to access the S3 bucket from the local machine. Search

  the internet on how to setup the IAM account. Store the dataset on the S3 bucket with

  command "s3-dist-cp eng_token_data.csv – dest s3://bucketname.

- **Local Jupyter Notebook:** To access the dataset with python verify that the boto3

  package is installed. If not, you can use the command "! pip install boto3" in python to

  install the package. To use the boto3 package "import boto3" and then set up the

  connection to the S3 bucket. The access key, secret key, and the name of the region the

  S3 bucket is hosted is required for the next command. "s3 = boto3.client('s3',

  aws_access_key_id='', aws_secret_access_key='', region_name=''). The following

  commands will load the dataset into a pandas dataframe. Import pandas first with

  command "import pandas as pd". Access the dataset with command "obj =

  s3.get_object(Bucket='S3BucketName', Key='eng_token_data.csv)". Store the dataset

  in a pandas dataframe with "df = pd.read_csv(obj['Body'])".

Big Data Wrangling With Google Books Ngrams

- **Store the data locally:** Perform a sanity check on the data with 'df.head(10)" and "df.shape".  If the data is present, then store the data with "df.to_csv('path/to/store/file/filename.csv').

- **Terminate the Cluster:** Go to the browser with the cluster console.  Click on "Clusters" on the left-hand panel.  View appendix B to verify the correct page.  Click on the active cluster. On the next page click "terminate".  This will stop any additional charges now that the cluster is no longer required.  The account can now be logged off and all the related windows closed.

Big Data Wrangling With Google Books Ngrams

## Findings

The dataset with tokens of 'data' is a subset of eng_1M_1gram. Hadoop and PySpark were used to analyze and filter the eng_1M_1gram dataset.  Once the dataset was filtered, it was saved to Hadoop.  The file on the Hadoop file system was split into 40 files, one for each partition of PySpark's RDDs.  To consolidate the files into a single file and store it on the file system of the server "hadoop fs -getmerge" command was used.  Then the file was transferred to an S3 bucket so it could be accessed by a local machine.  The file was downloaded in a jupyter notebook using the boto3 package.  The cleaning of the data revealed that there were years missing information amounting to 0.47%.  Because of the small number of missing values, the rows were removed instead of imputed.  The final data shape of the subset was (24,642, 5). The bar graph shows exponential growth in the token 'data' from the 1500's to current times.  A logged scale was used to make the graph readable.  The number of tokens with 'data' was very sparse from 1500 to 1700 and then the token took off with exponential growth.  The likely reason for the initial growth is that the number of published books grew at the same rate.  An analysis of the original dataset would have to be conducted to test this hypothesis.  After 1970 the growth is likely due to the rise of computers in the mainstream.
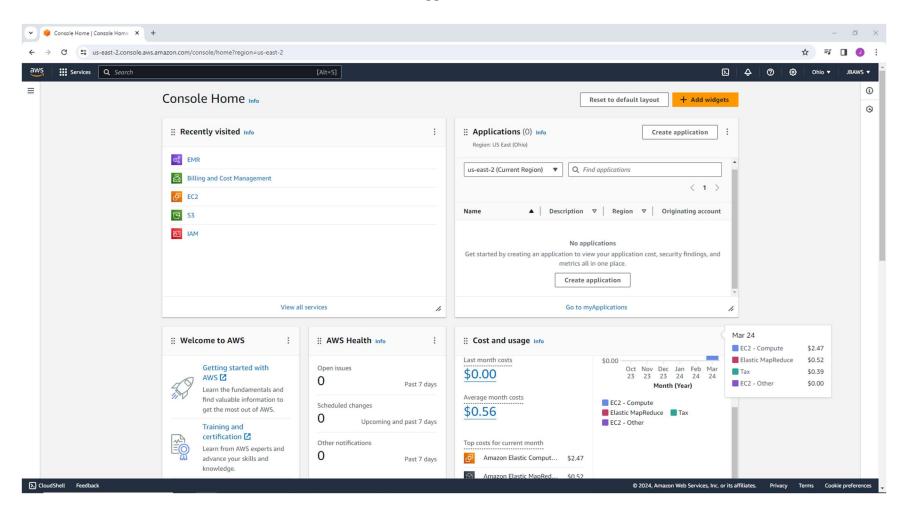
Big Data Wrangling With Google Books Ngrams

## Summary

The Hadoop HDFS stores data through distribution to many Child Nodes. The Head Node breaks up an incoming file up into equal size pieces. The size of the pieces of the file are such that it requires a low amount of memory in both the Head Node and the Child Nodes. Then, the Head Node assigns each piece of data to multiple Child Nodes and transfers the data to those Child Nodes. There are multiple copies of each piece of data that are stored on different Child Nodes. In case of the event that a Child Node fails other Child Nodes have a copy of the same data.
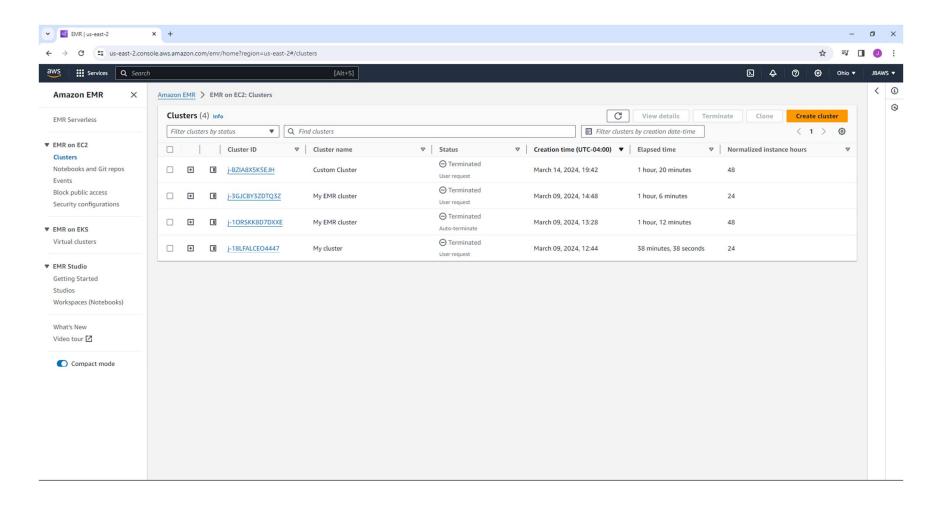
Hadoop and Spark are different systems with different purposes. Hadoop is made to store and retrieve large volumes of data at high speed. It has redundancies (duplicated data) on separate nodes so if a small number of nodes fail the system is still functional. Hadoop was designed for commodity hardware to keep the cost of hardware and maintenance down. It was also designed so that it scales horizontally requiring minimal setup to expand the system. Spark is pandas for very large datasets. Spark runs on clusters allowing it to divide the amount of work and use the servers in the cluster. It accomplishes this by splitting the data into separate groups for each CPU core to process. The system that does this in Spark is called RDD (Resilient Distributed Dataset). Spark also allows users to run SQL queries directly on Spark dataframes. This removes the need to learn all the differences between Spark and pandas commands. Together, Hadoop and Spark are a very powerful combination for Data Scientist working on large dataset. They offer scalable storage and computational power.
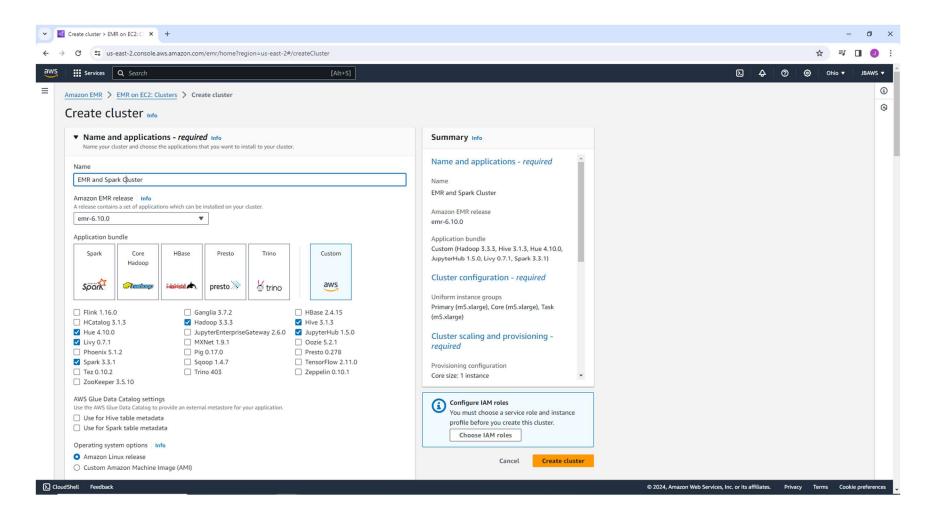
**Big Data Wrangling With Google Books Ngrams**
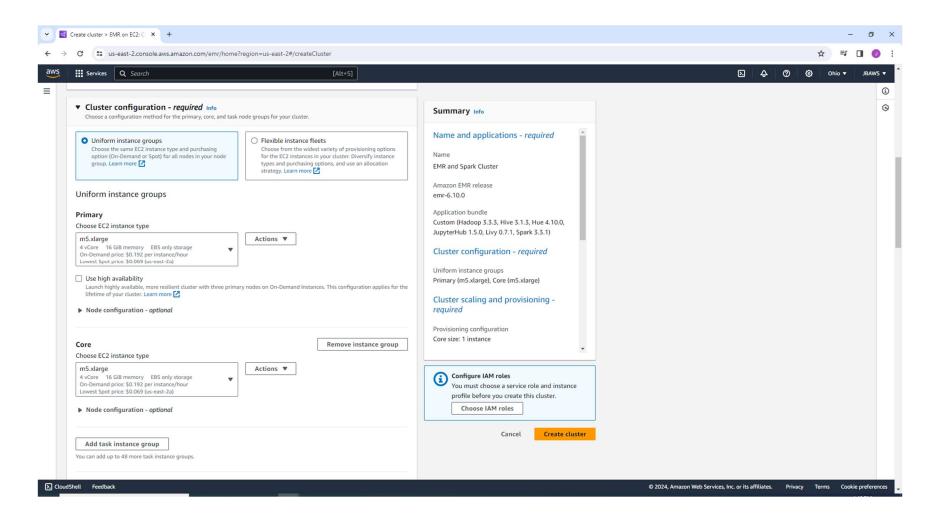
## Appendix

## Appendix A

Big Data Wrangling With Google Books Ngrams

Appendix B

Big Data Wrangling With Google Books Ngrams

Appendix C

Big Data Wrangling With Google Books Ngrams

## Appendix D

Big Data Wrangling With Google Books Ngrams

Appendix E

Big Data Wrangling With Google Books Ngrams

Appendix F

Big Data Wrangling With Google Books Ngrams

Appendix G

Big Data Wrangling With Google Books Ngrams

Appendix H

Big Data Wrangling With Google Books Ngrams

Appendix I



```
hadoop@ip-172-31-5-97:~                                                    —   □   ×

(base) C:\Users\jay\Documents\BrainStation\AWS>ssh -i NewKeyPair.pem hadoop@ec2-3-147-13-110.us-east-2.compute.amazon
aws.com
Last login: Mon Mar 18 18:38:25 2024

       __|  __|_  )
       _|  (     /   Amazon Linux 2 AMI
      ___|\___|___|

https://aws.amazon.com/amazon-linux-2/
64 package(s) needed for security, out of 94 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEEEE MMMMMMMM         MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::::E M:::::::M       M:::::::M R::::::::::::::R
EE:::::EEEEEEEEE:::::E M:::::::M        M:::::::M R:::::RRRRR:::::R
  E::::E        EEEEE M::::::::M       M::::::::M RR::::R    R::::R
  E::::E              M:::::M::::M     M::::M:::::M  R:::R    R::::R
  E:::::EEEEEEEEEE    M::::::M M:::M   M:::M M:::::M  R:::RRRRRR:::::R
  E::::::::::::::E     M::::::M  M:::M:::M  M:::::M    R:::::::::::RR
  E:::::EEEEEEEEEE    M::::::M   M:::::M   M:::::M    R:::RRRRRR::::R
  E::::E              M::::::M    M:::M    M:::::M    R:::R    R::::R
  E::::E        EEEEE M::::::M     MMM     M:::::M    R:::R    R::::R
EE:::::EEEEEEEE::::E M::::::M             M:::::M    R:::R    R::::R
E::::::::::::::::::::E M::::::M             M:::::M RR::::R    R::::R
EEEEEEEEEEEEEEEEEEEEEE MMMMMMM            MMMMMMM RRRRRRR     RRRRRR

[hadoop@ip-172-31-5-97 ~]$ hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram/eng_1M_1gr
am.csv
2024-03-18 18:42:24,529 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteM
issing=false, ignoreFailures=false, overwrite=false, append=false, useDiff=false, useRdiff=false, fromSnapshot=null,
toSnapshot=null, skipCRC=false, blocking=true, numListstatusThreads=0, maxMaps=20, mapBandwidth=0.0, copyStrategy='un
iformsize', preserveStatus=[], atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3://brainstat
ion-dsft/eng_1M_1gram.csv], targetPath=/user/hadoop/eng_1M_1gram/eng_1M_1gram.csv, filtersFile='null', blocksPerChunk
=0, copyBufferSize=8192, verboseLog=false, directWrite=false, useiterator=false}, sourcePaths=[s3://brainstation-dsft
/eng_1M_1gram.csv], targetPathExists=false, preserveRawXattrs=false
2024-03-18 18:42:24,775 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at ip-172-31-5-
97.us-east-2.compute.internal/172.31.5.97:8032
2024-03-18 18:42:24,919 INFO client.AHSProxy: Connecting to Application History server at ip-172-31-5-97.us-east-2.co
mpute.internal/172.31.5.97:10200
2024-03-18 18:42:27,640 INFO tools.SimpleCopyListing: Starting: Building listing using multi threaded approach for s3
://brainstation-dsft/eng_1M_1gram.csv
2024-03-18 18:42:27,642 INFO tools.SimpleCopyListing: Building listing using multi threaded approach for s3://brainst
ation-dsft/eng_1M_1gram.csv: duration 0:00.002s
2024-03-18 18:42:27,773 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
2024-03-18 18:42:27,773 INFO tools.SimpleCopyListing: Build file listing completed.
2024-03-18 18:42:27,774 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort
.mb
2024-03-18 18:42:27,775 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.
sort.factor
2024-03-18 18:42:27,886 INFO tools.DistCp: Number of paths in the copy list: 1
2024-03-18 18:42:27,911 INFO tools.DistCp: Number of paths in the copy list: 1
2024-03-18 18:42:27,936 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at ip-172-31-5-
97.us-east-2.compute.internal/172.31.5.97:8032
2024-03-18 18:42:27,936 INFO client.AHSProxy: Connecting to Application History server at ip-172-31-5-97.us-east-2.co
mpute.internal/172.31.5.97:10200
2024-03-18 18:42:28,022 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/stagi
ng/hadoop/.staging/job_1710782882902_0010
2024-03-18 18:42:28,132 INFO mapreduce.JobSubmitter: number of splits:1
2024-03-18 18:42:28,292 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1710782882902_0010
2024-03-18 18:42:28,292 INFO mapreduce.JobSubmitter: Executing with tokens: []
2024-03-18 18:42:28,468 INFO conf.Configuration: resource-types.xml not found
2024-03-18 18:42:28,468 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
```
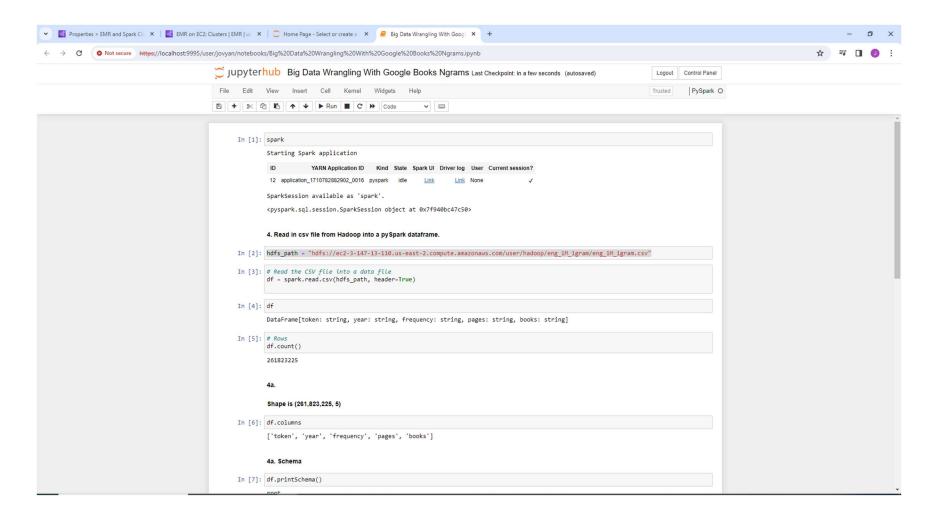
Big Data Wrangling With Google Books Ngrams

Appendix J

```
hadoop@ip-172-31-5-97:~                                                    —    □    ×
2024-03-18 18:42:28,528 INFO impl.YarnClientImpl: Submitted application application_1710782882902_0010
2024-03-18 18:42:28,575 INFO mapreduce.Job: The url to track the job: http://ip-172-31-5-97.us-east-2.compute.interna
l:20888/proxy/application_1710782882902_0010/
2024-03-18 18:42:28,575 INFO tools.DistCp: DistCp job-id: job_1710782882902_0010
2024-03-18 18:42:28,576 INFO mapreduce.Job: Running job: job_1710782882902_0010
2024-03-18 18:42:34,632 INFO mapreduce.Job: Job job_1710782882902_0010 running in uber mode : false
2024-03-18 18:42:34,633 INFO mapreduce.Job:  map 0% reduce 0%
2024-03-18 18:42:50,713 INFO mapreduce.Job:  map 100% reduce 0%
2024-03-18 18:43:37,893 INFO mapreduce.Job: Job job_1710782882902_0010 completed successfully
2024-03-18 18:43:37,977 INFO mapreduce.Job: Counters: 42
        File System Counters
                FILE: Number of bytes read=0
                FILE: Number of bytes written=294915
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=376
                HDFS: Number of bytes written=5292105197
                HDFS: Number of read operations=12
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=5
                HDFS: Number of bytes read erasure-coded=0
                S3: Number of bytes read=5292105197
                S3: Number of bytes written=0
                S3: Number of read operations=0
                S3: Number of large read operations=0
                S3: Number of write operations=0
        Job Counters
                Launched map tasks=1
                Other local map tasks=1
                Total time spent by all maps in occupied slots (ms)=5858592
                Total time spent by all reduces in occupied slots (ms)=0
                Total time spent by all map tasks (ms)=61027
                Total vcore-milliseconds taken by all map tasks=61027
                Total megabyte-milliseconds taken by all map tasks=187474944
        Map-Reduce Framework
                Map input records=1
                Map output records=0
                Input split bytes=137
                Spilled Records=0
                Failed Shuffles=0
                Merged Map outputs=0
                GC time elapsed (ms)=216
                CPU time spent (ms)=59680
                Physical memory (bytes) snapshot=1036165120
                Virtual memory (bytes) snapshot=4434186240
                Total committed heap usage (bytes)=617611264
                Peak Map Physical memory (bytes)=1068023808
                Peak Map Virtual memory (bytes)=4434186240
        File Input Format Counters
                Bytes Read=239
        File Output Format Counters
                Bytes Written=0
        DistCp Counters
                Bandwidth in Bytes=91243193
                Bytes Copied=5292105197
                Bytes Expected=5292105197
                Files Copied=1
[hadoop@ip-172-31-5-97 ~]$
```

Big Data Wrangling With Google Books Ngrams

Appendix K

```
hadoop@ip-172-31-5-97:~                                              —   □   ×

(base) C:\Users\jay\Documents\BrainStation\AWS>ssh -i NewKeyPair.pem -L 9995:localhost:9443 hadoop@ec2-3-147-13-110.u
s-east-2.compute.amazonaws.com
Last login: Mon Mar 18 17:32:48 2024 from d24-150-232-237.home.cgocable.net

       __|  __|_  )
       _|  (     /   Amazon Linux 2 AMI
      ___|\___|___|

https://aws.amazon.com/amazon-linux-2/
64 package(s) needed for security, out of 94 available
Run "sudo yum update" to apply all updates.
Last login: Mon Mar 18 17:32:48 2024 from d24-150-232-237.home.cgocable.net

       __|  __|_  )
       _|  (     /   Amazon Linux 2 AMI
      ___|\___|___|

https://aws.amazon.com/amazon-linux-2/
64 package(s) needed for security, out of 94 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEEE MMMMMMMM            MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::::M          M:::::::M R::::::::::::::R
EE:::::EEEEEEEEE:::E M::::::::M        M::::::::M R:::::RRRRRR:::::R
  E::::E       EEEEE M:::::::::M      M:::::::::M RR::::R      R::::R
  E::::E             M::::::M:::M    M:::M::::::M   R:::R       R::::R
  E:::::EEEEEEEEEE    M::::::M M:::M M:::M M::::::M   R:::RRRRRR:::::R
  E::::::::::::::E    M::::::M  M:::M:::M  M::::::M   R:::::::::::RR
  E:::::EEEEEEEEEE    M::::::M   M:::::M   M::::::M   R:::RRRRRR::::R
  E::::E             M::::::M    M:::M    M::::::M   R:::R       R::::R
  E::::E       EEEEE M::::::M     MMM     M::::::M   R:::R       R::::R
EE:::::EEEEEEEE::::E M::::::M             M::::::M   R:::R       R::::R
E::::::::::::::::::E M::::::M             M::::::M RR::::R       R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM             MMMMMMM RRRRRRR       RRRRRR

[hadoop@ip-172-31-5-97 ~]$
```

# Big Data Wrangling With Google Books Ngrams

## Appendix L

# Big Data Wrangling With Google Books Ngrams

## Appendix M



```
|-- token: string (nullable = true)
|-- year: string (nullable = true)
|-- frequency: string (nullable = true)
|-- pages: string (nullable = true)
|-- books: string (nullable = true)
```

**4b. Make a new datafrom from a sql qurey with column token containing the word 'data'.**

```
In [8]:  # Register dataframe as a view
         df.createOrReplaceTempView("eng_table")

         # Run sql Query
         newdf = spark.sql("SELECT * FROM eng_table WHERE token LIKE '%data%'")
```

**4b. Describe the new dataset.**

```
In [9]:  newdf.count()

         24642
```

```
In [10]: newdf.columns

         ['token', 'year', 'frequency', 'pages', 'books']
```

```
In [11]: newdf.head(5)

         [Row(token='laticaudata', year='1800', frequency='1', pages='1', books='1'), Row(token='laticaudata', year='1823', frequency
         ='2', pages='2', books='2'), Row(token='laticaudata', year='1827', frequency='1', pages='1', books='1'), Row(token='laticaudat
         a', year='1843', frequency='2', pages='2', books='1'), Row(token='laticaudata', year='1844', frequency='6', pages='6', books
         ='4')]
```

```
In [12]: newdf.printSchema()

         root
          |-- token: string (nullable = true)
          |-- year: string (nullable = true)
          |-- frequency: string (nullable = true)
          |-- pages: string (nullable = true)
          |-- books: string (nullable = true)
```

**4c. Write new dataframe to HDFS**

```
In [13]: hdfs_write_path = "hdfs://ec2-3-147-13-110.us-east-2.compute.amazonaws.com/user/hadoop/eng_1M_1gram/eng_token_data.csv"
```

```
In [14]: newdf.write.csv(hdfs_write_path, header=True)
```

```
In [15]: spark.stop()
```

Big Data Wrangling With Google Books Ngrams

Appendix N