

Big Data Wrangling With Google Books Ngrams

The scope of data processing and analysis is outlined in steps and screenshots for the questions below. This file showcases screenshots of steps necessary for completing the tasks: 1 to 6. Parts of step 4 are in the jupyterhub notebook with PySpark. Parts of step 6-8 are located in the second jupyter notebook on my local machine.


1. Spin up a new EMR cluster on AWS for using Spark and EMR notebooks - follow the same instructions as for the Spark Lab.


Name and applications [Info](#)


Name

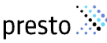
Amazon EMR release [Info](#)
A release contains a set of applications which can be installed on your cluster.


Application bundle


Spark


Core Hadoop


HBase


Presto


Trino


Custom


Applications included in bundle
Hadoop 3.3.3 with Hive 3.1.3, Hue 4.10.0, Pig 0.17.0 and Tez 0.10.2

AWS Glue Data Catalog settings
Use the AWS Glue Data Catalog to provide an external metastore for your application.
☐ Use for Hive table metadata

Operating system options [Info](#)
☒ Amazon Linux release
☐ Custom Amazon Machine Image (AMI)

Cluster configuration [Info](#)

Choose a configuration method for the primary, core, and task node groups for your cluster.

☒ Instance groups

Choose one instance type per node group

☐ Instance fleets

Choose any combination of instance types within each node group

Instance groups

Primary

Choose EC2 instance type

m5.xlarge

4 vCore 16 GiB memory EBS only storage
On-Demand price: \$0.192 per instance/hour
Lowest Spot price: \$0.053 (us-east-2b)

Actions ▼

☐ Use multiple primary nodes

To improve cluster availability, use 3 primary nodes with the same configuration and bootstrap actions. You can not use multiple primary nodes with instance fleets.

► Node configuration - optional

Core

Remove instance group

Choose EC2 instance type

m5.xlarge

4 vCore 16 GiB memory EBS only storage
On-Demand price: \$0.192 per instance/hour
Lowest Spot price: \$0.053 (us-east-2b)

Actions ▼

► Node configuration - optional

Add task instance group

You can add up to 48 more task instance groups.

► EBS root volume - optional

Cluster termination [Info](#)

☐ Manually terminate cluster
 ☐ Automatically terminate cluster after last step ends
 ☒ Automatically terminate cluster after idle time (Recommended)

Idle time

Enter the time until your cluster terminates.

0 days ▼

04:00:00

Choose a time that is greater than 1 minute (00:01:00) and less than 7 days. The time is in hh:mm:ss (24-hour) format.

☐ **Use termination protection**
 Protect your EC2 instances from accidental termination.

Security configuration and EC2 key pair - optional [Info](#)

Security configuration

Select your cluster encryption, authentication, authorization, and instance metadata service settings.

Amazon EC2 key pair for SSH to the cluster [Info](#)

2. Connect to the head node of the cluster using SSH.
 - Use your file path for your pemkey on your local machine. To connect to the Primary Node of the cluster, and access JupyterHub in a browser window, edit the following bash command:
 - `ssh -i mykey.pem -L 9995:localhost:9443 hadoop@xxxxxxxxxxxxx.compute.amazonaws.com`
 - Replace the 'xxxxxxxxxx' with your 'Primary node public DNS' which you can find on the overview page of your cluster.
 - If the connection is successful, you should be greeted by the EMR ASCII banner:
 - Here is how I connected to my SSH cluster: `~/desktop/pemkey/aws_cloud.pem -L 9995:localhost:9443 hadoop@ec2-3-131-97-241.us-east-2.compute.amazonaws.com`

My_emr_cluster

Updated 1 minute ago



Actions ▾

▼ Summary

Cluster info

Cluster ID
j-WMC80NAWA6TJCluster configuration
Instance groupsCapacity
1 Primary 2 Core 0 Task

Applications

Amazon EMR version
emr-6.10.0Installed applications
Hadoop 3.3.3, Hive 3.1.3, Hue 4.10.0,
JupyterHub 1.5.0, Livy 0.7.1, Spark 3.3.1

Cluster management

Log destination in Amazon S3
Logging not configuredPersistent application UIs
[Spark History Server](#)
[YARN timeline server](#)
[Tez UI](#)Primary node public DNS
[ec2-3-131-97-241.us-east-2.compute.amazonaws.com](#)
[Connect to the Primary Node using SSH](#)

Status and time

Status
● **Waiting**
Creation time
August 22, 2023, 11:00 (UTC-04:00)
Elapsed time
1 day, 15 hours

Connect to the primary node using SSH



You can connect to the Amazon EMR primary node using SSH to perform actions like running interactive queries, examining log files, submit Linux commands, and view web interfaces hosted on Amazon EMR clusters. [Learn more](#)

Windows

Mac/Linux

1. Open a terminal window. On Mac OS X, choose Applications > Utilities > Terminal. On other Linux distributions, terminal is typically found at Applications > Accessories > Terminal.
2. To establish a connection to the primary node, enter the following command. Replace `~/aws_cloud.pem` with the location and filename of the private key file (.pem) that you used to launch the cluster.

```
ssh -i ~/aws_cloud.pem hadoop@ec2-3-131-97-241.us-east-2.compute.amazonaws.com
```



3. Enter yes to dismiss the security warning.

[View web interfaces hosted on Amazon EMR clusters](#)

Close

```
Last login: Tue Aug 22 11:04:43 on ttys004
(base) markbenhain@Marks-MBP ~ % ssh -i ~/desktop/pemkey/aws_cloud.pem -L 9995:localhost:9443 hadoop@ec2-3-131-97-241.us-east-2.compute.amazonaws.com
Last login: Tue Aug 22 15:20:43 2023
```

```
--| ( _ )
--| ( _ ) / Amazon Linux 2 AMI
--| \_/_/_/_|
```

```
https://aws.amazon.com/amazon-linux-2/
14 package(s) needed for security, out of 16 available
Run "sudo yum update" to apply all updates.
```

```
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::::M M::::::::M R:::::::::R
EE::::::::EEEEEEEE::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 M:::M R:::R R:::R
E:::EEEEEEEEEE M:::M M:::M M:::M M:::M R:::RRRRR:::R
E::::::::::::E M:::M M:::M M:::M M:::M R:::RRRRR:::R
E:::EEEEEEEEEE M:::M M:::M M:::M M:::M R:::RRRRR:::R
E:::E M:::M M:::M M:::M M:::M R:::R R:::R
E:::E EEEEE M:::M M:::M M:::M M:::M R:::R R:::R
EE:::EEEEEEEE::E M:::M M:::M M:::M M:::M R:::R
E::::::::::::E M:::M M:::M M:::M M:::M R:::R
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRR RRRRRR
```

```
[hadoop@ip-172-31-28-187 ~]$ hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram/
```

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

- Connect to the Brainstation public S3 bucket to access the dataset within the data folder by typing this code: `hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram/`

-

```

EEEEEEEEEEEEEEEEEEEE MMMMMMM      MMMMMMM RRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::::M      M::::::::M R::::::::::::R
EE::::::::EEEEEEEE::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::::EEEEEEEE M::::M M::M M::M M::M      R::RRRRR:::R
        E::::::::::::E M::::M M::M::M M::M      R:::::::::RR
          EEEEEEEEEEE M::::M M::M::M M::M      R::RRRRR:::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-172-31-28-187 ~]$ hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram/
```

- Using pyspark, read the data you copied into HDFS in Step 3. You may either use Jupyterhub on EMR (the default user and password are `jovyan` and `jupyter`) or work from pyspark in the terminal if you prefer.

- Use this link to access jupyterhub in your browser: <https://localhost:9995/>
- Make sure your cluster is in the waiting stage and rerun the ssh command
- Login to Jupyterhub using the login info then get to this page and create a notebook

Sign in

Username:

Password:

Sign In

jupyterhub

Logout Control Panel

Files Running Clusters

Select items to perform actions on them.

Upload New ↻

<input type="checkbox"/>	0	/	Name ↓	Last Modified	File size
<input type="checkbox"/>			Mark_Benaim_Unit_4_Big_Data_Wrangling_With_Google_Books_Ngrams_Part_1.ipynb	Running 13 minutes ago	12 kB
<input type="checkbox"/>			jupyterhub-proxy.pid	2 days ago	2 B
<input type="checkbox"/>			jupyterhub.sqlite	2 minutes ago	102 kB
<input type="checkbox"/>			jupyterhub_cookie_secret	2 days ago	65 B

5. Collect the contents of the directory into a single file on the local drive of the head node using **getmerge** and move this file into a S3 bucket in your account.

- Input code in bash:
- Note: filtered_1gram is the name of my file. Use your file name that you gave from question 4.
- You'll also need to create your own s3 bucket in your AWS account under the S3 page.
- Use the same region as your cluster. Name your bucket
- Create the bucket

```
[hadoop@ip-172-31-28-187 ~]$ hadoop fs -getmerge /user/hadoop/filtered_1gram/ ~/filtered_1gram.csv
```

Amazon S3 > Buckets > Create bucket

Create bucket [Info](#)

Buckets are containers for data stored in S3. [Learn more](#)

General configuration

Bucket name

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

AWS Region

Copy settings from existing bucket - optional
Only the bucket settings in the following configuration are copied.

Object Ownership [Info](#)

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

☒ **ACLs disabled (recommended)**
All objects in this bucket are owned by this account.

☐ **ACLs enabled**
Objects in this bucket can be owned by other AWS

Default encryption [Info](#)

Server-side encryption is automatically applied to new objects stored in this bucket.

Encryption type [Info](#)

☒ **Server-side encryption with Amazon S3 managed keys (SSE-S3)**

☐ Server-side encryption with AWS Key Management Service keys (SSE-KMS)

☐ Dual-layer server-side encryption with AWS Key Management Service keys (DSSE-KMS)
Secure your objects with two separate layers of encryption. For details on pricing, see [DSSE-KMS pricing](#) on the **Storage** tab of the [Amazon S3 pricing page](#).

Bucket Key
Using an S3 Bucket Key for SSE-KMS reduces encryption costs by lowering calls to AWS KMS. S3 Bucket Keys aren't supported for DSSE-KMS. [Learn more](#)

☐ Disable

☒ **Enable**

► **Advanced settings**

After creating the bucket, you can upload files and folders to the bucket, and configure additional bucket settings.

6.
 - Use 'aws configure' in bash and place your codes and default region:
 - First, get your codes from your aws account
 - Check your s3 bucket to see if your file is there by clicking on the name of your bucket
 - The file should be shown under objects
 - Go to your IAM dashboard and create a user and name it
 - Afterwards, go to the drop down add permissions and click create inline policy
 - Click S3 and all S3 actions then click next then create policy
 - Go to Security and Recommendations on your IAM dashboard and click manage access keys
 - Under access keys, click create access key then check I understand and click create access key
 - Download your access key and secret key
 - Now you can use 'aws configure' in bash and place your codes and default region:


```
[[hadoop@ip-172-31-28-187 ~]$ aws configure
AWS Access Key ID [*****LTUB]: AKIA6KCSJI30LJHELTUB
AWS Secret Access Key [*****hcVi]: FxPAXJgzRU3I3Yff0iffbMYDy10gsf/RhNv0hcVi
Default region name [us-east-2]: us-east-2
Default output format [None]:
```


[Amazon S3](#) > Buckets


Account snapshot
Storage lens provides visibility into storage usage and activity trends. [Learn more](#)

[View Storage Lens dashboard](#)

Buckets (1) [Info](#)
Buckets are containers for data stored in S3. [Learn more](#)

 [Copy content](#) [Empty](#) [Delete](#) [Create bucket](#)

< 1 > 

	Name ▲	AWS Region ▼	Access ▼	Creation date ▼
	markbigdatabucket	US East (Ohio) us-east-2	Bucket and objects not public	August 22, 2023, 14:09:24 (UTC-04:00)

markbigdatabucket [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](#)

[Refresh](#) [Copy S3 URI](#) [Copy URL](#) [Download](#) [Open](#) [Delete](#) [Actions](#)

[Create folder](#) [Upload](#)

< 1 > [Settings](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	filtered_1gram.csv	csv	August 24, 2023, 00:54:38 (UTC-04:00)	7.2 KB	Standard

Identity and Access Management (IAM)

Dashboard

Access management

User groups
Users
Roles
Policies
Identity providers
Account settings

Access reports

[IAM](#) > Dashboard

IAM Dashboard

Security recommendations 2

[Add MFA](#)

[Deactivate or delete access keys for root user](#)

[Manage access keys](#)

IAM resources

Resources in this AWS Account

AWS Account

Account ID
[983720085212](#)

Account Alias
[Create](#)

Sign-in URL for IAM users in this account
<https://983720085212.signin.aws.amazon.com/console>

[IAM](#) > Users


Users (1) [Info](#)

An IAM user is an identity with long-term credentials that is used to interact with AWS in an account.



< 1 > [Settings](#)

<input type="checkbox"/>	User name	Path	Group	Last activity	MFA	Password
<input type="checkbox"/>	mbenheim	/	0	.	-	-

Summary

ARN  arn:aws:iam::983720085212:user/mbenhaim	Console access Disabled	Access key 1 Create access key
Created August 22, 2023, 15:35 (UTC-04:00)	Last console sign-in -	


Permissions policies (1)
Permissions are defined by policies attached to the user directly or through groups.

  Remove

Add permissions ▲
Add permissions
Create inline policy

Filter by Type
All types ▼


☐

Policy name 


▲

Type

▼

Attached via 

☐

 [S3](#)

Customer inline


Inline

Access keys (1)

Use access keys to send programmatic calls to AWS from the AWS CLI, AWS Tools for PowerShell, AWS SDKs, or direct AWS API calls. You can have a maximum of two access keys (active or inactive) at a time. [Learn more](#)

Actions ▾

Create access key

	Access key ID	Created on	Access key last used	Region last used	Service last used	Status
<input type="radio"/>	AKIA6KCSJI3OLJHELTUB	Yesterday	19 minutes ago	us-east-2	s3	 Active

Alternatives to root user access keys [Info](#)



Root user access keys are not recommended

We don't recommend that you create root user access keys. Because you can't specify the root user in a permissions policy, you can't limit its permissions, which is a best practice.

Instead, use alternatives such as an IAM role or a user in IAM Identity Center, which provide temporary rather than long-term credentials. [Learn More](#)

If your use case requires an access key, create an IAM user with an access key and apply least privilege permissions for that user. [Learn More](#)

Continue to create access key?

- ☒ I understand creating a root access key is not a best practice, but I still want to create one.

Cancel

Create access key

Note: The answers corresponding to questions 1-8 are also located in within the 2 notebooks:

- Mark_Benheim_Unit_4_Big_Data_Wrangling_With_Google_Books_Ngrams_Part_1
- Mark_Benheim_Unit_4_Big_Data_Wrangling_With_Google_Books_Ngrams_Part_2