

PROYECTO 1 – PARTE 2: MANUAL DE REPRODUCIBILIDAD

ALEJANDRO ARANGO GIRALDO

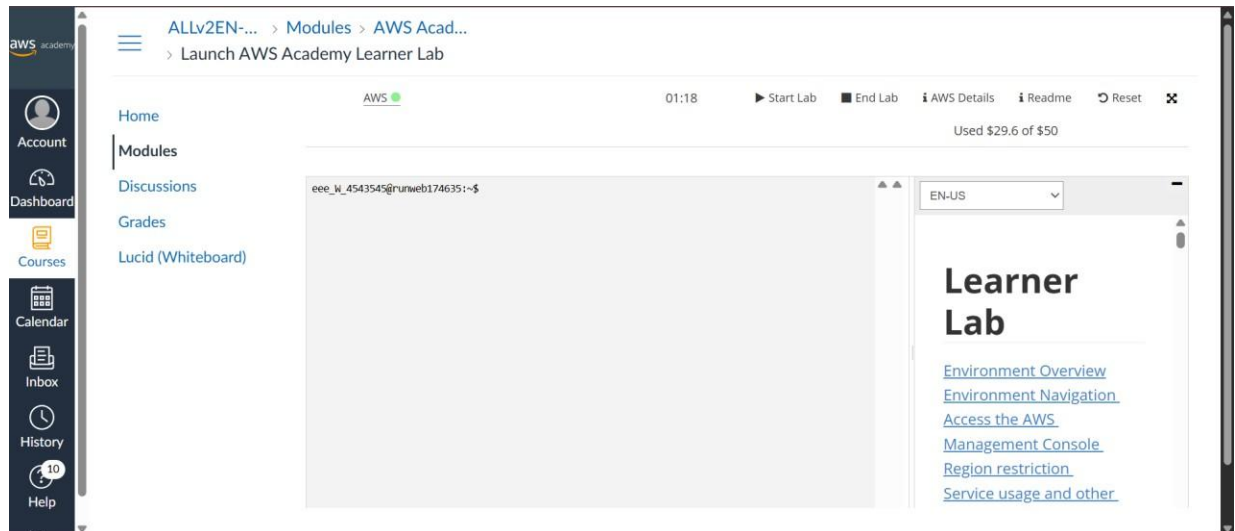
Línea de énfasis en Ciencias de los Datos

ST1800 – Almacenamiento y recuperación de la información

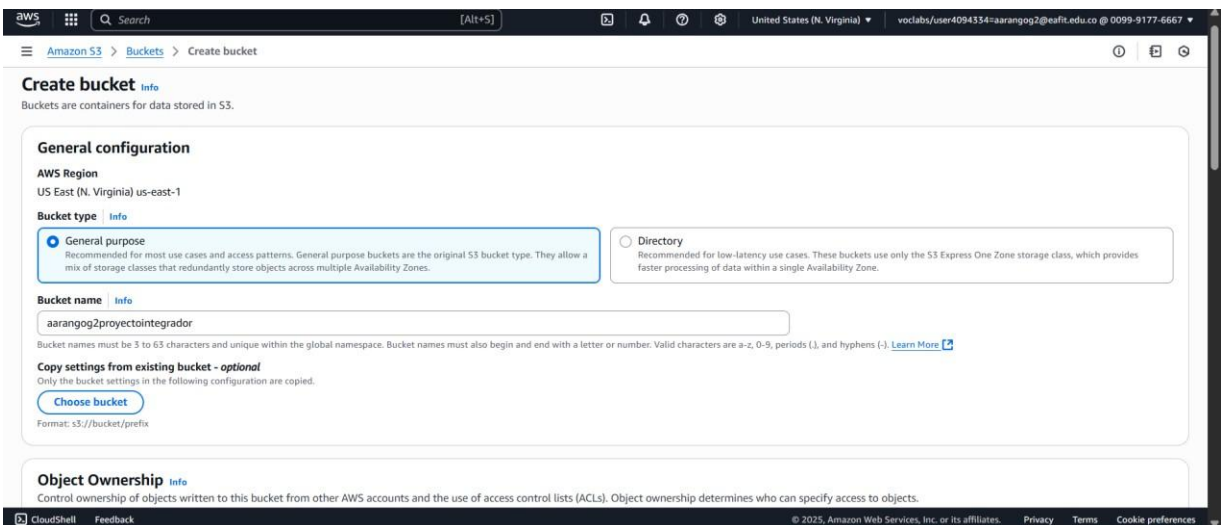
UNIVERSIDAD EAFIT
ESCUELA DE CIENCIAS APLICADAS E INGENIERÍA
CARRERA DE INGENIERÍA MECÁNICA
MEDELLÍN
2025 - 1

Para reproducir este proyecto, se deben seguir los siguientes pasos de manera secuencial:

1. Activar una sesión de AWS Academy Learner Lab, o, si se cuenta con otras credenciales, activar la sesión.



2. Entrar a la interfaz web de AWS y buscar el servicio de AWS S3. Oprimir en la opción “Create bucket” y mantener la siguiente configuración:



Si se cambia el nombre del bucket, se deben cambiar las rutas definidas en los códigos para leer y guardar los archivos allí almacenados.

aws [Search] [Alt+S] United States (N. Virginia) voclabs/user4094334=aarangog2@eafit.edu.co @ 0099-9177-6667

Amazon S3 Buckets Create bucket

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. Access to this bucket and its objects is specified using only policies.

☐ **ACLs enabled**
Objects in this bucket can be owned by other AWS accounts. Access to this bucket and its objects can be specified using ACLs.

Object Ownership
 Bucket owner enforced

Block Public Access settings for this bucket

Public access is granted to buckets and objects through access control lists (ACLs), bucket policies, access point policies, or all. In order to ensure that public access to this bucket and its objects is blocked, turn on Block all public access. These settings apply only to this bucket and its access points. AWS recommends that you turn on Block all public access, but before applying any of these settings, ensure that your applications will work correctly without public access. If you require some level of public access to this bucket or objects within, you can customize the individual settings below to suit your specific storage use cases. [Learn more](#)

☐ **Block all public access**
Turning this setting on is the same as turning on all four settings below. Each of the following settings are independent of one another.

☐ **Block public access to buckets and objects granted through new access control lists (ACLs)**
S3 will block public access permissions applied to newly added buckets or objects, and prevent the creation of new public access ACLs for existing buckets and objects. This setting doesn't change any existing permissions that allow public access to S3 resources using ACLs.

☐ **Block public access to buckets and objects granted through any access control lists (ACLs)**
S3 will ignore all ACLs that grant public access to buckets and objects.

☐ **Block public access to buckets and objects granted through new public bucket or access point policies**
S3 will block new bucket and access point policies that grant public access to buckets and objects. This setting doesn't change any existing policies that allow public access to S3 resources.

☐ **Block public and cross-account access to buckets and objects through any public bucket or access point policies**
S3 will ignore public and cross-account access for buckets or access points with policies that grant public access to buckets and objects.

Turning off block all public access might result in this bucket and the objects within becoming public
 AWS recommends that you turn on block all public access, unless public access is required for specific and verified use cases such as static website hosting.

☒ I acknowledge that the current settings might result in this bucket and the objects within becoming public.

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

aws [Search] [Alt+S] United States (N. Virginia) voclabs/user4094334=aarangog2@eafit.edu.co @ 0099-9177-6667

Amazon S3 Buckets Create bucket

Bucket Versioning

Versioning is a means of keeping multiple variants of an object in the same bucket. You can use versioning to preserve, retrieve, and restore every version of every object stored in your Amazon S3 bucket. With versioning, you can easily recover from both unintended user actions and application failures. [Learn more](#)

Bucket Versioning

☒ **Disable**

☐ **Enable**

Tags - optional (0)

You can use bucket tags to track storage costs and organize buckets. [Learn more](#)

No tags associated with this bucket.

[Add tag](#)

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**

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

3. Una vez creado el bucket “aarangog2proyectointegrador”, se crean las carpetas: zona raw, zona trusted y zona refined.

aws [Search] [Alt+S] United States (N. Virginia) voclabs/user4094334=aarangog2@eafit.edu.co @ 0099-9177-66

Amazon S3

Amazon S3

- General purpose buckets
- Directory buckets
- Table buckets
- Access Grants
- Access Points
- Object Lambda Access Points
- Multi-Region Access Points
- Batch Operations
- IAM Access Analyzer for S3

Block Public Access settings for this account

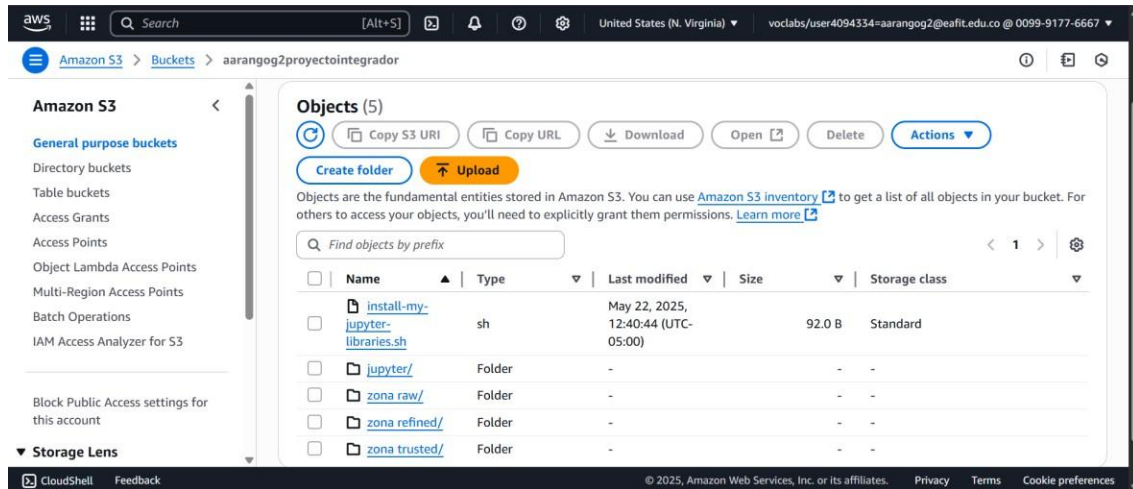
General purpose buckets (3) info **All AWS Regions**

[Refresh](#) [Copy ARN](#) [Empty](#) [Delete](#) [Create bucket](#)

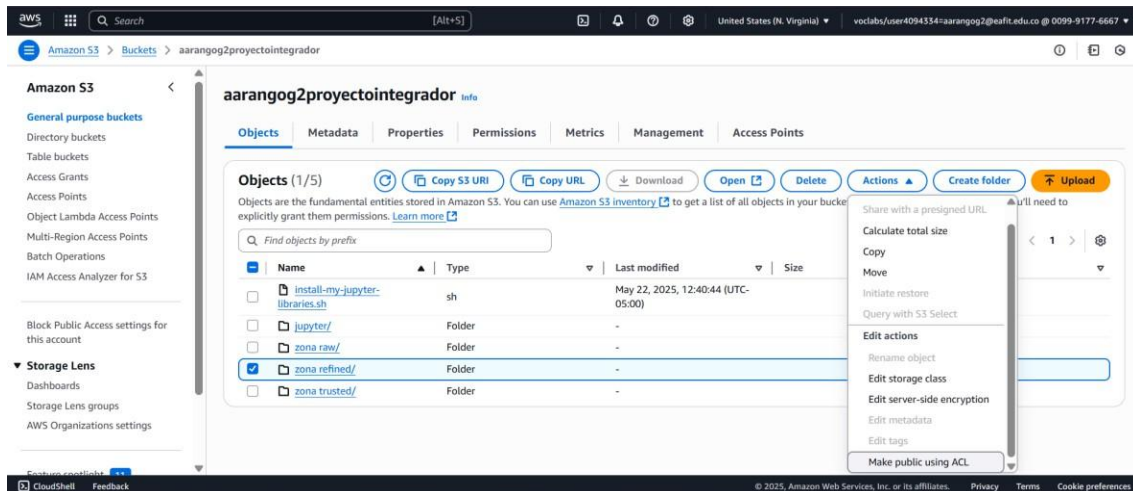
Buckets are containers for data stored in S3.

	Name	AWS Region	IAM Access Analyzer	Creation date
<input type="radio"/>	aarangog2lab1	US East (N. Virginia) us-east-1	View analyzer for us-east-1	May 17, 2025, 16:04:58 (UTC-05:00)
<input type="radio"/>	aarangog2proyectointegrador	US East (N. Virginia) us-east-1	View analyzer for us-east-1	May 20, 2025, 18:08:42 (UTC-05:00)
<input type="radio"/>	aws-logs-009991776667-us-east-1	US East (N. Virginia) us-east-1	View analyzer for us-east-1	May 20, 2025, 18:21:36 (UTC-05:00)

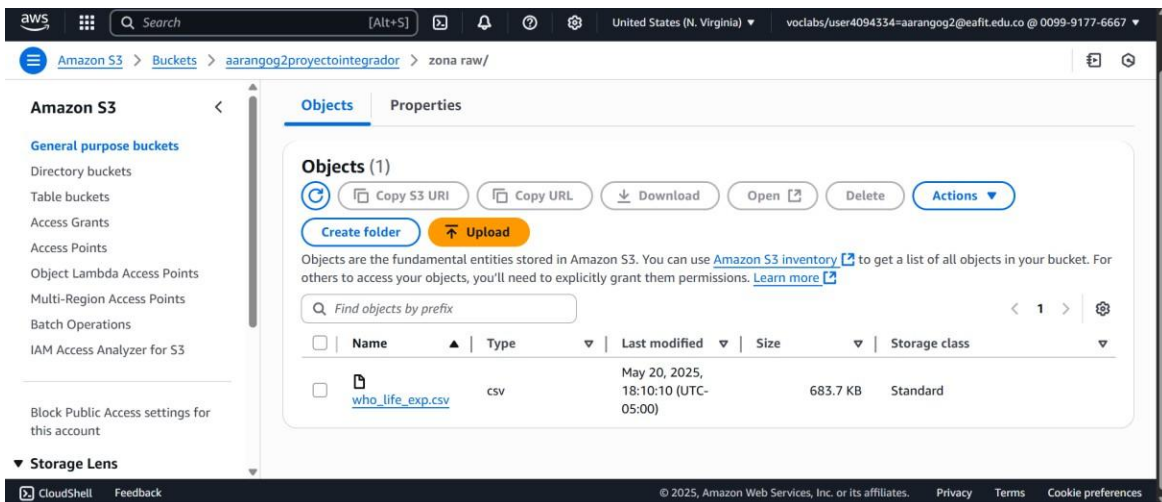
CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences



Las zonas “trusted” y “refined” deben volverse públicas para su futura consulta. Esto se realiza por medio de la acción “Make public using ACL”.



4. En la carpeta “zona raw”, se deben cargar manualmente los datos crudos. En nuestro caso, sería el archivo “who_life_exp.csv”.



5. Buscar en la interfaz web el servicio EMR y crear un cluster con la siguiente configuración:

This screenshot shows the first step of the AWS EMR console 'Create cluster' wizard. The page is titled 'Clone "aarangog2 Proyecto Integrador" info'. The 'Name and applications - required' section includes a text input for the cluster name 'aarangog2 Proyecto Integrador', a dropdown for the Amazon EMR release label 'emr-7.8.0', and a grid of application bundles. The 'Custom' bundle is selected. Below the bundles, there are checkboxes for various applications like AmazonCloudWatchAgent, Hive, and Tez. The 'AWS Glue Data Catalog settings' section is also visible. On the right, a 'Summary' panel shows the configured details. At the bottom right, there are 'Cancel' and 'Clone cluster' buttons.

Name and applications - required

Name: aarangog2 Proyecto Integrador

Amazon EMR release: emr-7.8.0

Application bundle: Custom (HCatalog 3.1.3, Hadoop 3.4.1, Hive 3.1.3, Hue 4.11.0, JupyterEnterpriseGateway 2...)

Cluster configuration - required

Uniform instance groups: Primary (m5.xlarge), Core (m5.xlarge), Task (m5.xlarge)

Cluster scaling and provisioning - required

Provisioning configuration: Core size: 1 instance, Task size: 1 instance

This screenshot shows the second step of the AWS EMR console 'Create cluster' wizard. The 'Cluster configuration - required' section is active. It shows the 'Uniform instance groups' option selected. Under 'Uniform instance groups', the 'Primary' group is configured with an 'm5.xlarge' EC2 instance type. The 'Node configuration - optional' section is also visible. On the right, the 'Summary' panel is updated with the new configuration. At the bottom right, there are 'Cancel' and 'Clone cluster' buttons.

Cluster configuration - required

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

Uniform instance groups

Primary: Choose EC2 instance type: m5.xlarge

Node configuration - optional

Core: Choose EC2 instance type: m5.xlarge

Summary

Name: aarangog2 Proyecto Integrador

Amazon EMR release: emr-7.8.0

Application bundle: Custom (HCatalog 3.1.3, Hadoop 3.4.1, Hive 3.1.3, Hue 4.11.0, JupyterEnterpriseGateway 2...)

Cluster configuration - required

Uniform instance groups: Primary (m5.xlarge), Core (m5.xlarge), Task (m5.xlarge)

Cluster scaling and provisioning - required

Provisioning configuration: Core size: 1 instance, Task size: 1 instance

This screenshot shows the third step of the AWS EMR console 'Create cluster' wizard. The 'Core' and 'Task' node configurations are being set. Both are configured with an 'm5.xlarge' EC2 instance type. The 'EBS root volume' section is also visible, showing a size of 15 GB. On the right, the 'Summary' panel is updated with the new configuration. At the bottom right, there are 'Cancel' and 'Clone cluster' buttons.

Core

Choose EC2 instance type: m5.xlarge

Task 1 of 1

Choose EC2 instance type: m5.xlarge

EBS root volume

Size (GiB): 15, IOPS: 3000, Throughput (MiB/s): 125

Summary

Name: aarangog2 Proyecto Integrador

Amazon EMR release: emr-7.8.0

Application bundle: Custom (HCatalog 3.1.3, Hadoop 3.4.1, Hive 3.1.3, Hue 4.11.0, JupyterEnterpriseGateway 2...)

Cluster configuration - required

Uniform instance groups: Primary (m5.xlarge), Core (m5.xlarge), Task (m5.xlarge)

Cluster scaling and provisioning - required

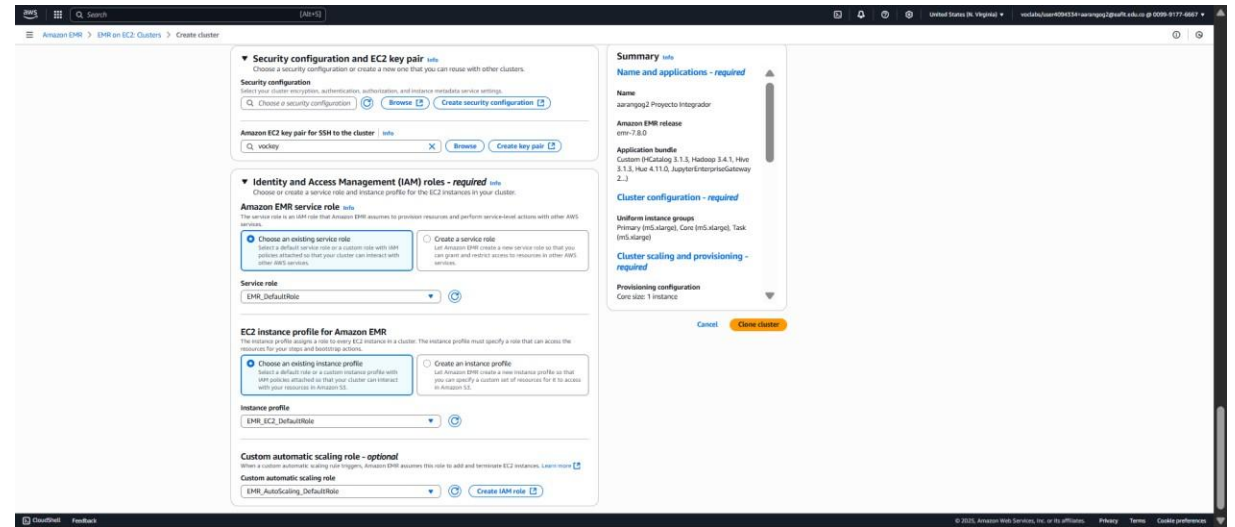
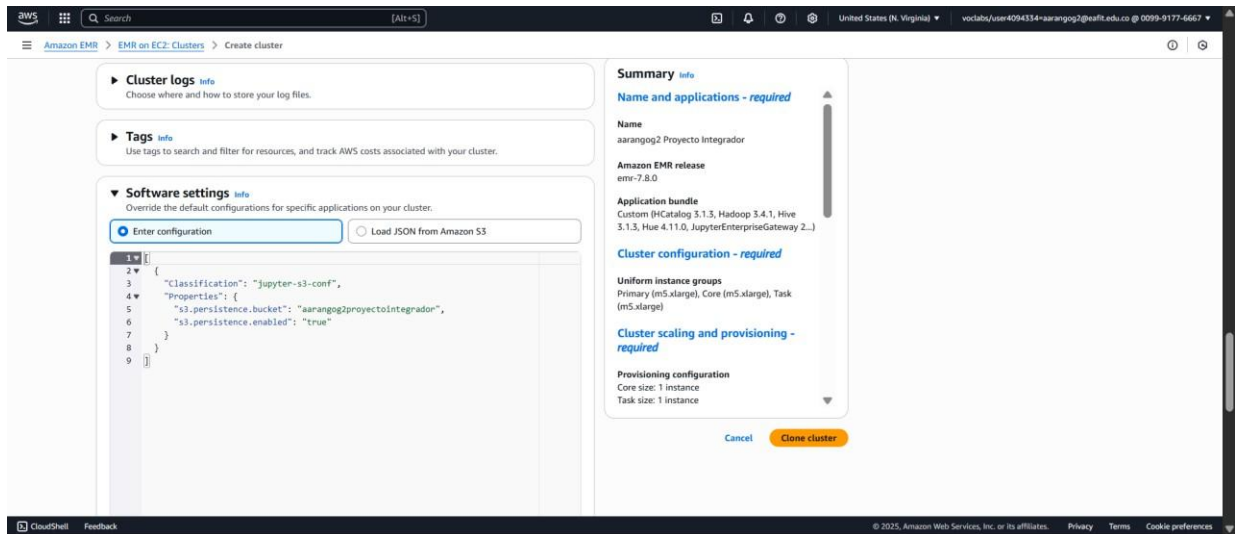
Provisioning configuration: Core size: 1 instance, Task size: 1 instance

Adicionar “Bootstrap actions” es opcional en este caso de uso, puesto que las librerías utilizadas en las diferentes fases son SparkML y SparkSQL, las cuales forman parte de Apache Spark. Sin embargo, de necesitarse otras librería, se puede crear un código como el siguiente, guardarlo con la extensión “sh”, y añadirlo a la sección “Bootstrap actions” en la configuración del cluster.

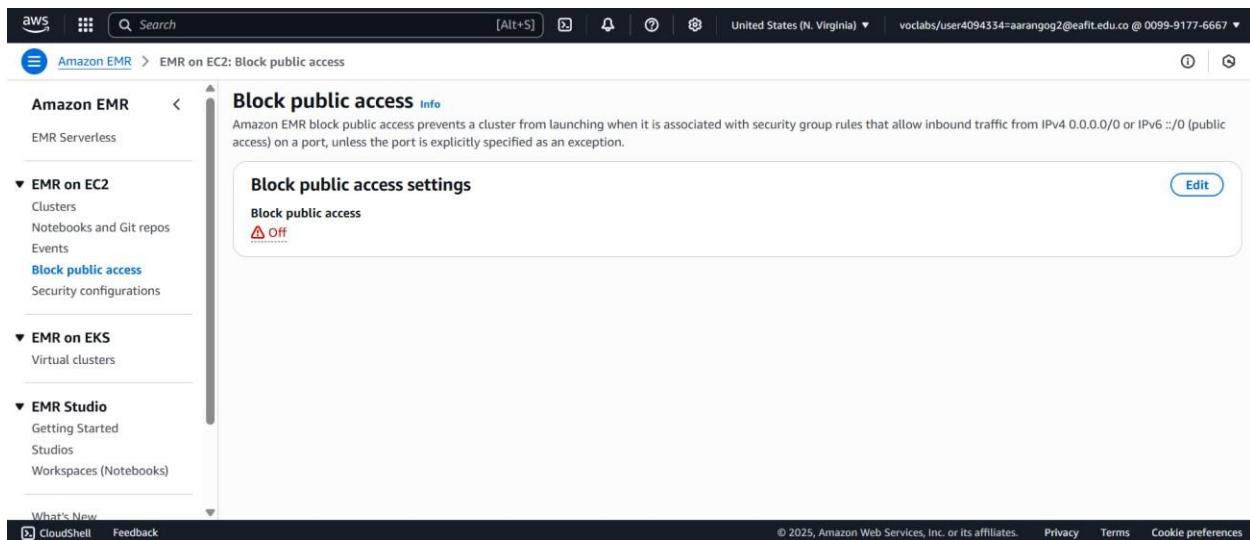
```

$ install-my-jupyter-libraries.sh
C:\Users\aleje\OneDrive - Universidad EAFIT\Escritorio\EAFIT\2025-1\Proyecto Integrador\Entregas> $ install-my-jupyter-libraries.sh
1 #!/bin/bash
2
3 sudo python3 -m pip install shap seaborn matplotlib scipy scikit-learn pandas

```

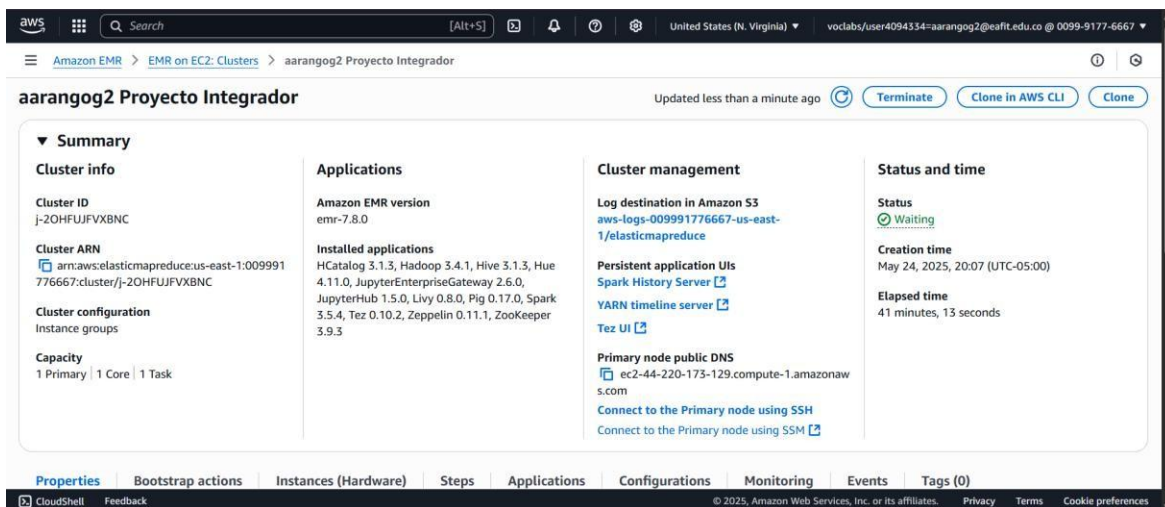



6. Entrar a la opción “Block public access” del menú de la izquierda y abrir todos los puertos TCP para acceso al clúster de la siguiente manera:

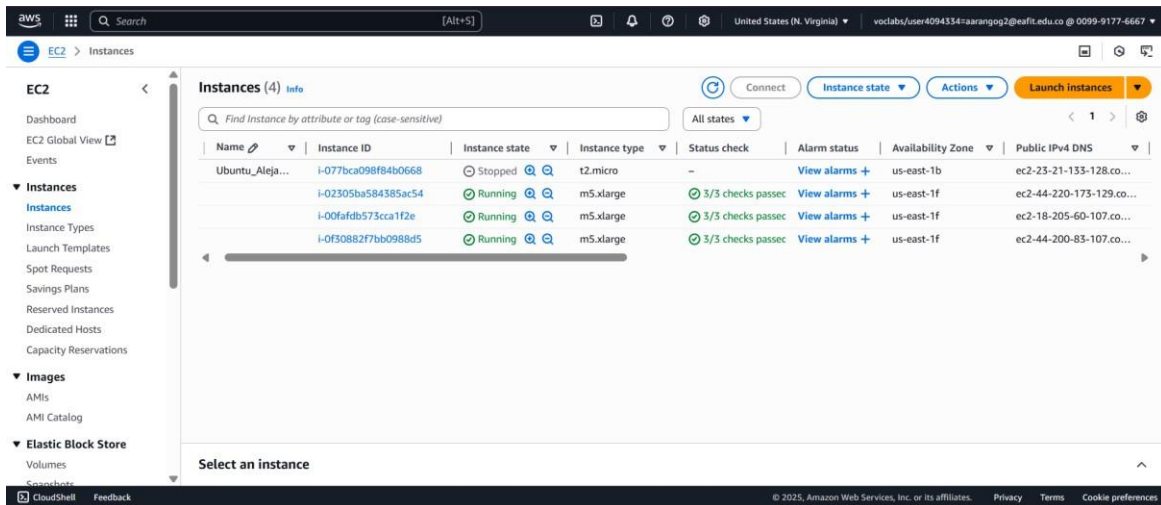


7. Abrir los puertos de las aplicaciones de hadoop/Spark en el Security Group del nodo MASTER del clúster como se muestra a continuación:

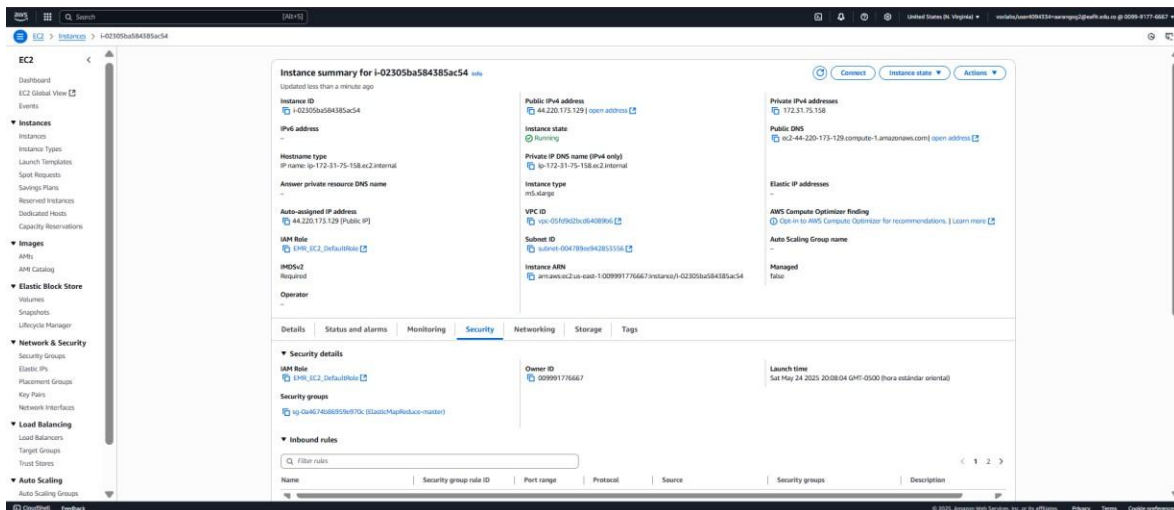
7.1 Identificar el nodo primario del cluster recién creado, el cual se muestra en “Primary node public DNS”.



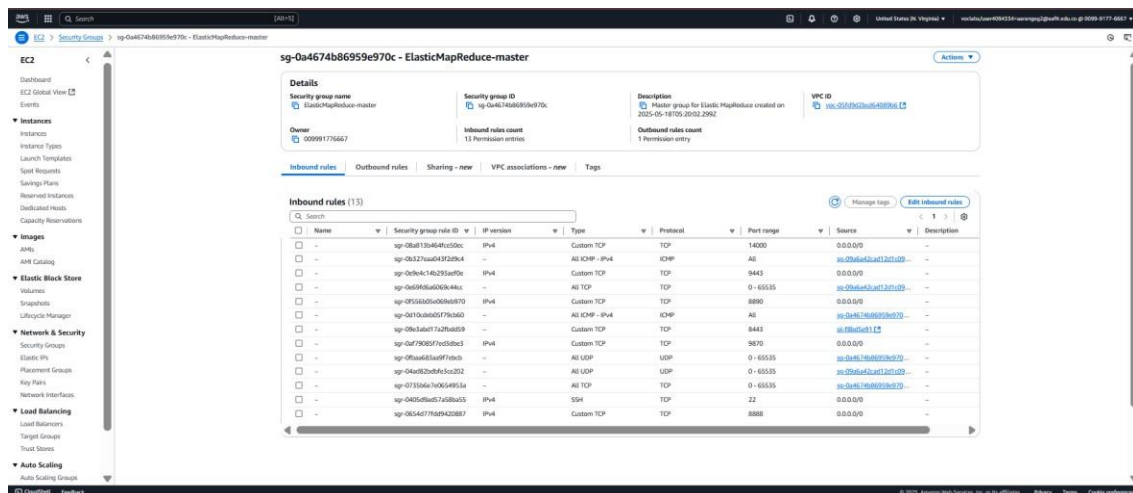
7.2 Buscar en la interfaz web el servicio E2C, donde se encontrarán tres máquinas. Abrir la que tenga el valor de la columna “Public IPv4 DNS” igual al valor del “Primary node public DNS” del cluster.



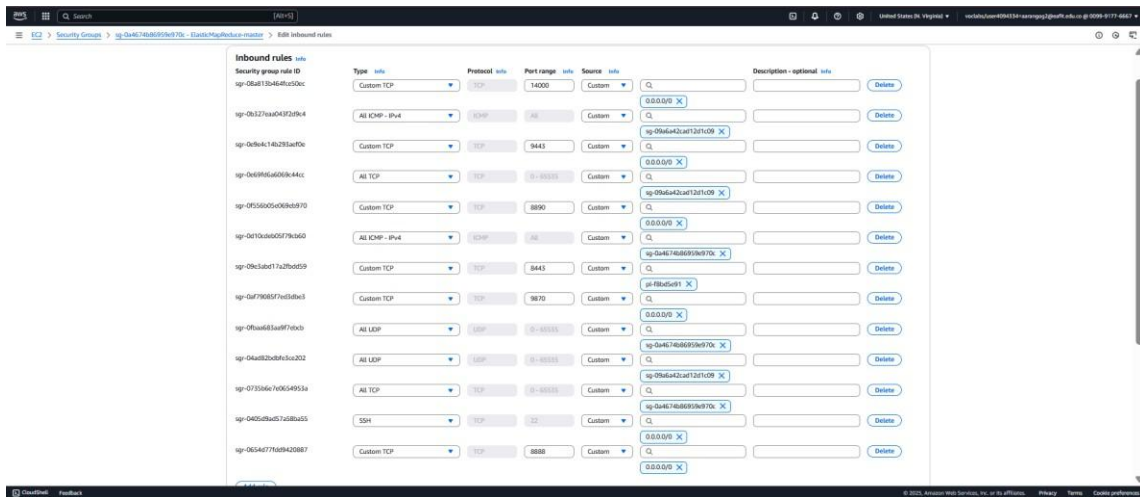
7.3 Entrar a la pestaña de seguridad de la Instancia EC2 del nodo master y abrir la opción “Security groups”.



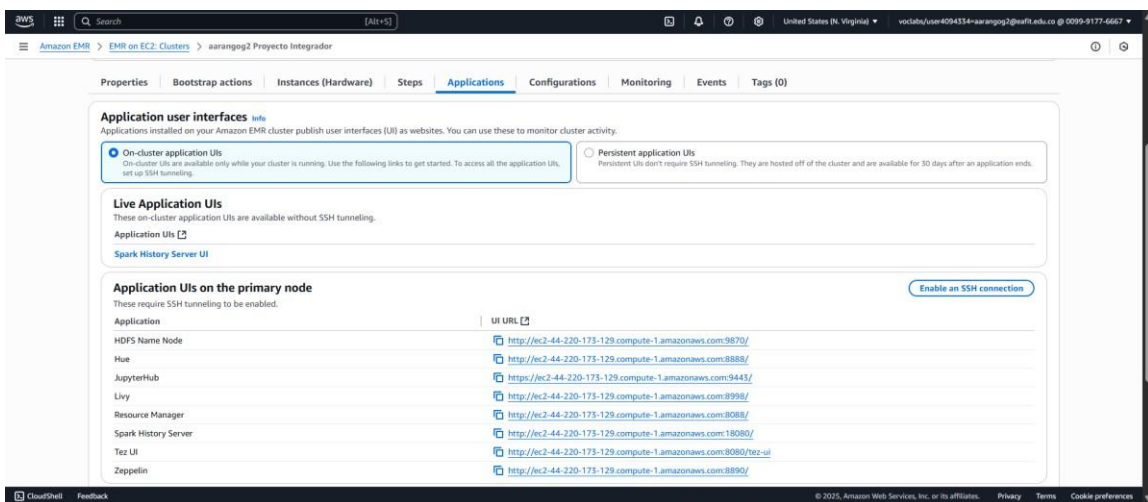
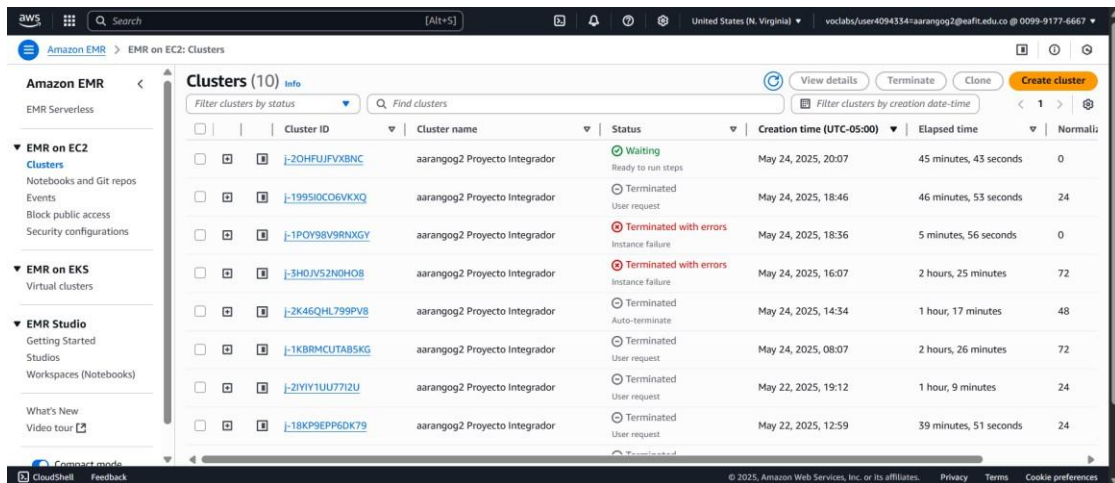
7.4 Entrar a la opción “Edit inbound rules”.



7.5 Habilitar los nodos: 22, 14000, 9870, 8888, 9443, y 8890.



8. Con el cluster en estado “Waiting”, seleccionarlo y en el menú “Applications”, oprimir el UI que lleva al servicio de Jupyterhub.



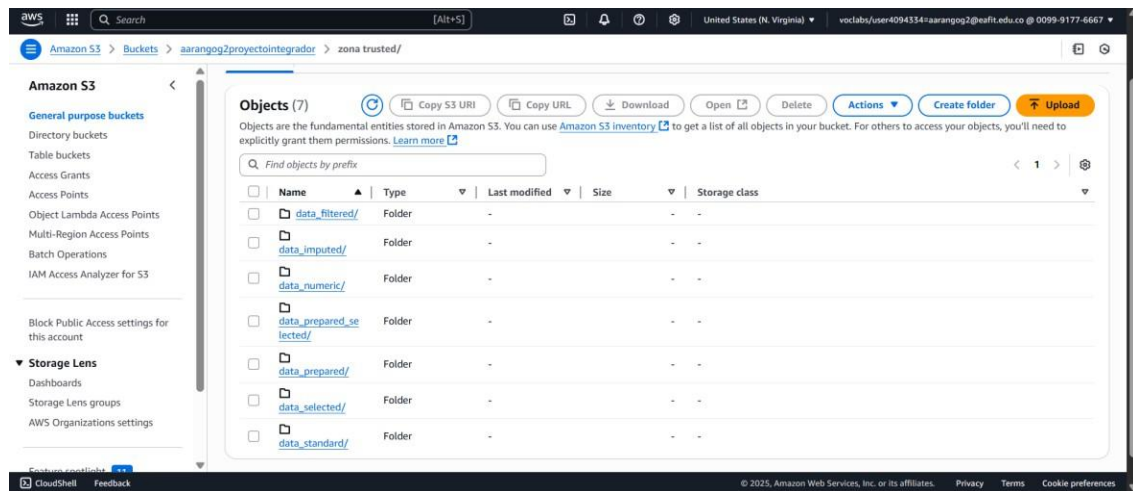
9. Ingresar con las credenciales:

- **Username:** jovyan
- **Password:** jupyter

Cargar y ejecutar el notebook: “Data prep.ipynb”.

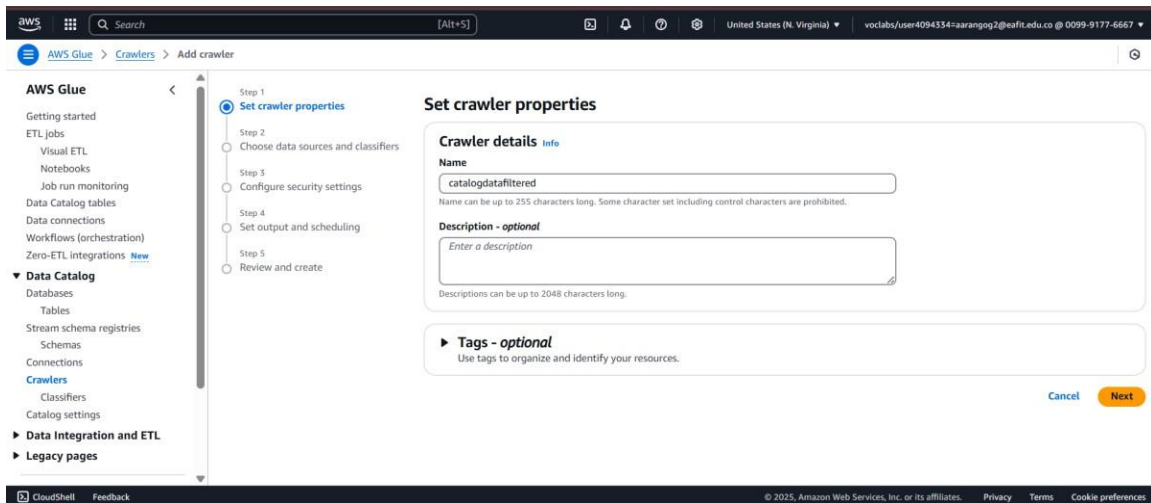


Tras ejecutar este notebook, se generarán los siguientes archivos con formato parquet en la zona trusted del S3:

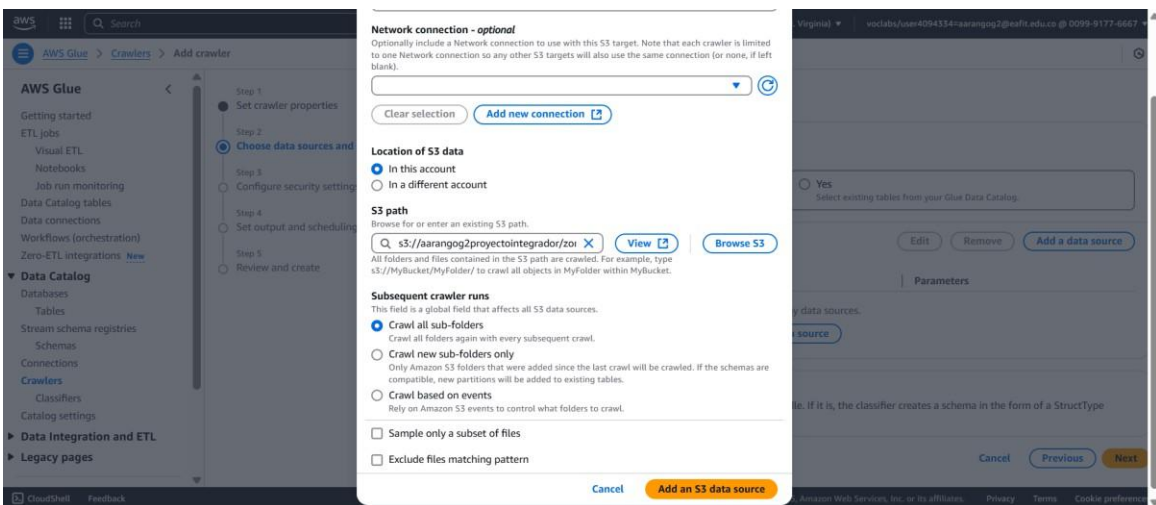
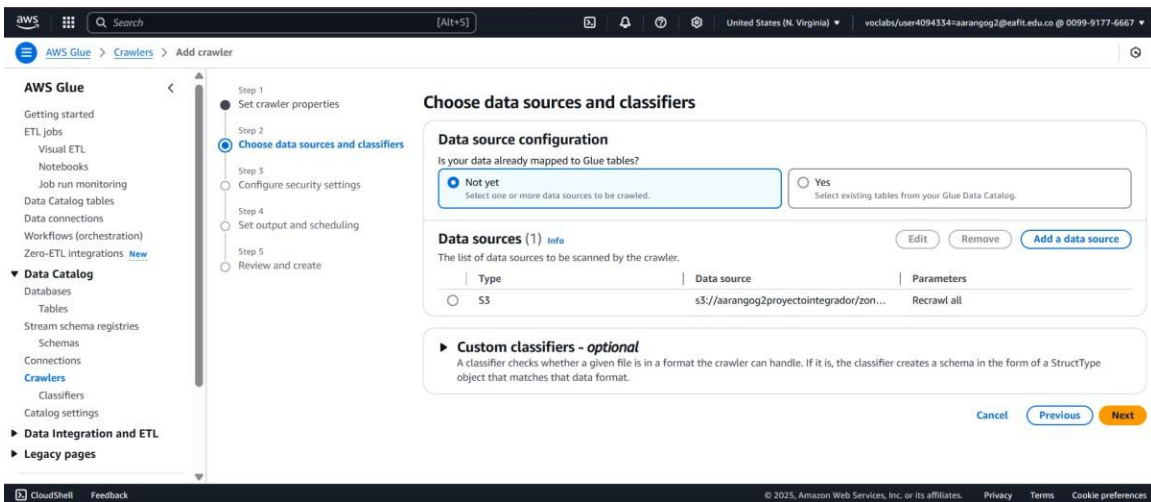


10. Catalogar los resultados del notebook “Data prep.ipynb” con el servicio de AWS Glue, creando un crawler por cada uno de los resultados. Para la creación de un crawler, se debe utilizar la siguiente configuración, y replicarla para cada uno de ellos:

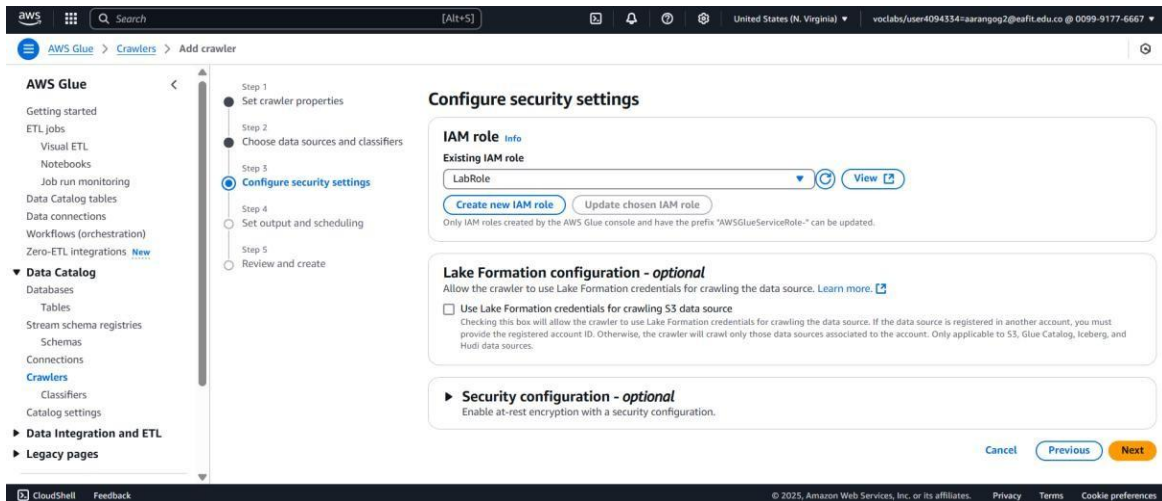
10.1 Seleccionar el nombre del crawler.



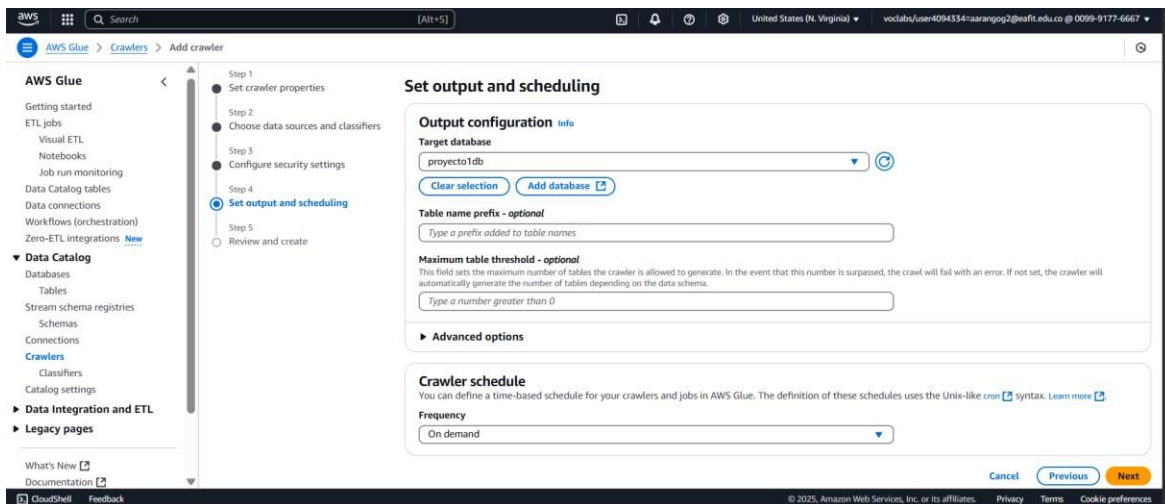
10.2 Seleccionar la ruta al archivo almacenado en S3 que se quiere catalogar.



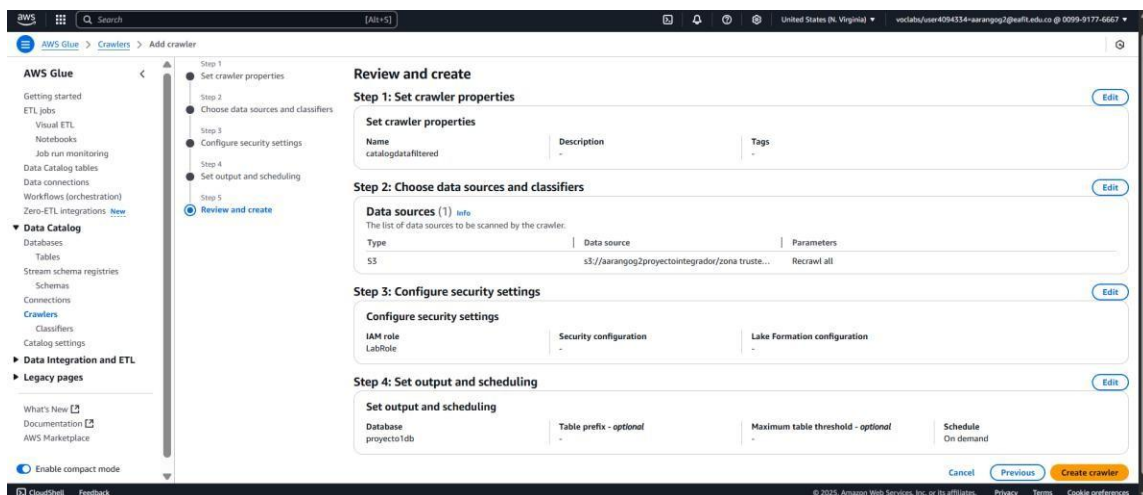
10.3 Seleccionar el rol “LabRole” para el “IAM role”.



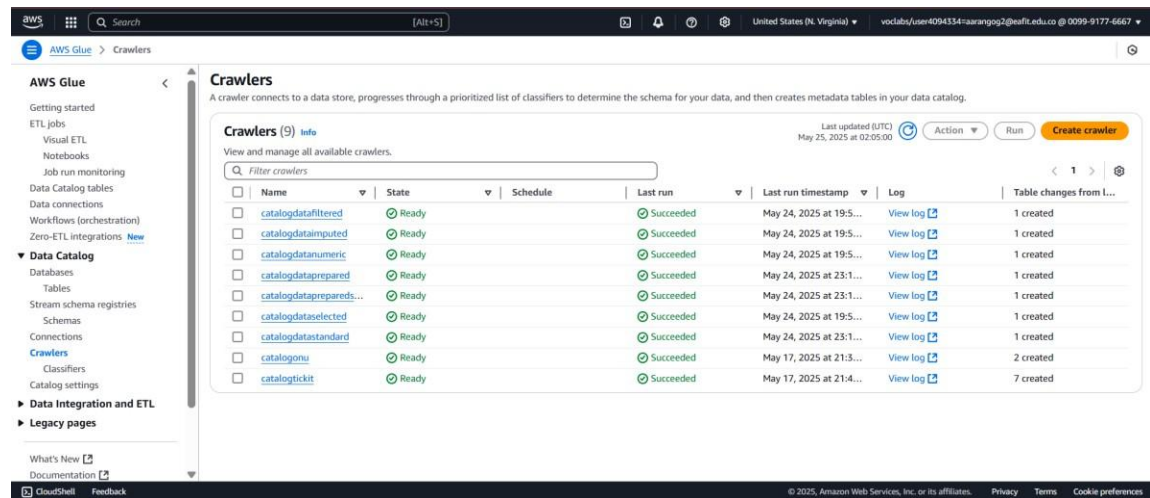
10.4 Crear una nueva base de datos llamada “proyecto1db” y seleccionarla para el almacenamiento de los esquemas.



10.5 Crear y correr el crawler.

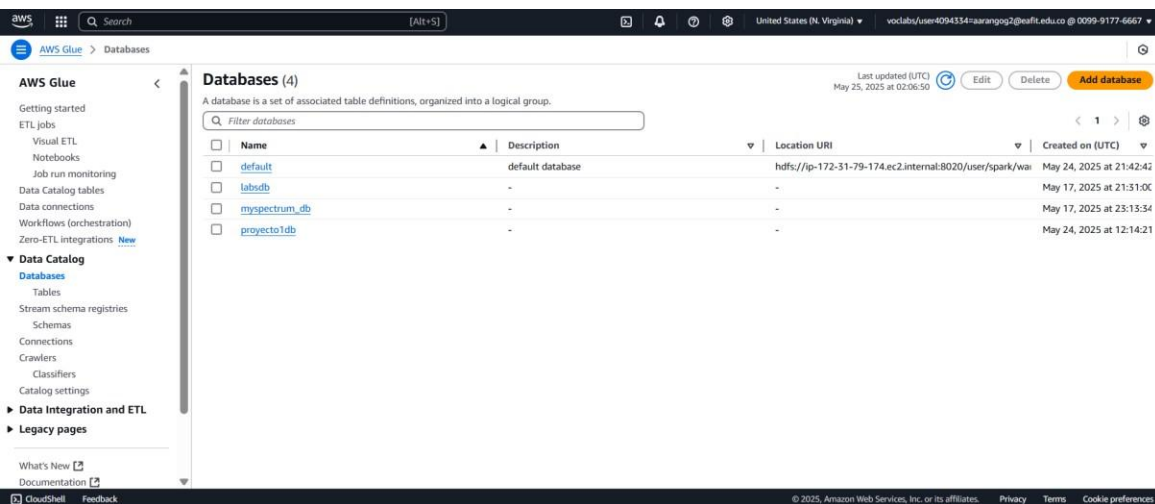


11. Correr todos los crawlers para obtener las siguientes tablas en la base de datos “proyecto1db”:



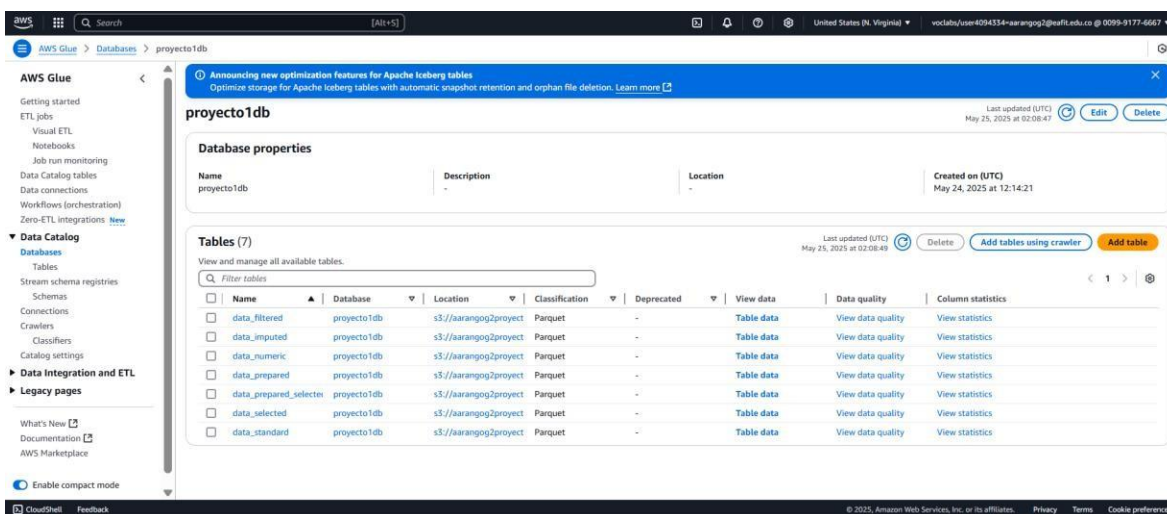
The screenshot shows the AWS Glue console with the 'Crawlers' tab selected. A list of 9 crawlers is displayed, all in a 'Ready' state and having completed their last run successfully. The crawlers are designed to extract data from various sources into the 'proyecto1db' database.

Name	State	Schedule	Last run	Last run timestamp	Log	Table changes from last run
catalogdatafiltered	Ready		Succeeded	May 24, 2025 at 19:5...	View log	1 created
catalogdataimputed	Ready		Succeeded	May 24, 2025 at 19:5...	View log	1 created
catalogdatanumeric	Ready		Succeeded	May 24, 2025 at 19:5...	View log	1 created
catalogdataprepared	Ready		Succeeded	May 24, 2025 at 23:1...	View log	1 created
catalogdataprepareds...	Ready		Succeeded	May 24, 2025 at 23:1...	View log	1 created
catalogdataselected	Ready		Succeeded	May 24, 2025 at 19:5...	View log	1 created
catalogdatastandard	Ready		Succeeded	May 24, 2025 at 23:1...	View log	1 created
cataloggonu	Ready		Succeeded	May 17, 2025 at 21:3...	View log	2 created
catalogtickit	Ready		Succeeded	May 17, 2025 at 21:4...	View log	7 created



The screenshot shows the AWS Glue console with the 'Databases' tab selected. A list of 4 databases is displayed. The 'proyecto1db' database is highlighted, showing its location URI and creation timestamp.

Name	Description	Location URI	Created on (UTC)
default	default database	hdfs://ip-172-31-79-174.ec2.internal:8020/user/spark/wai	May 24, 2025 at 21:42:42
labsdb	-	-	May 17, 2025 at 21:31:06
myspectrum_db	-	-	May 17, 2025 at 23:13:34
proyecto1db	-	-	May 24, 2025 at 12:14:21



The screenshot shows the AWS Glue console with the 'projecto1db' database selected. The 'Database properties' section shows the name, description, location, and creation timestamp. Below, the 'Tables (7)' section lists the tables created by the crawlers, including their names, databases, locations, classifications, and data quality.

Name	Database	Location	Classification	Deprecated	View data	Data quality	Column statistics
data_filtered	projecto1db	s3://aarangog2project	Parquet	-	Table data	View data quality	View statistics
data_imputed	projecto1db	s3://aarangog2project	Parquet	-	Table data	View data quality	View statistics
data_numeric	projecto1db	s3://aarangog2project	Parquet	-	Table data	View data quality	View statistics
data_prepared	projecto1db	s3://aarangog2project	Parquet	-	Table data	View data quality	View statistics
data_prepared_selects	projecto1db	s3://aarangog2project	Parquet	-	Table data	View data quality	View statistics
data_selected	projecto1db	s3://aarangog2project	Parquet	-	Table data	View data quality	View statistics
data_standard	projecto1db	s3://aarangog2project	Parquet	-	Table data	View data quality	View statistics

12. Ejecutar los notebooks “EDA.ipynb”, y “Train Model.ipynb” en el EMR para obtener los siguientes resultados en la zona refined. Todos, a excepción del archivo “scaler” que es creado en “Data prep.ipynb”, son generados en el EDA y el entrenamiento del modelo.

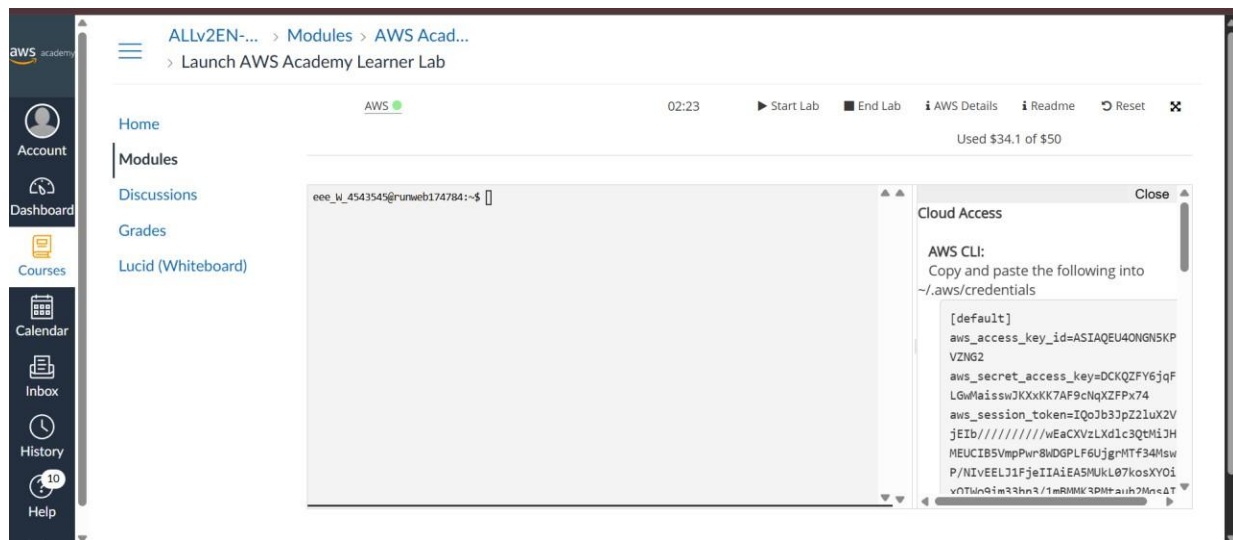
Name	Type	Last modified	Size	Storage class
et_eda	Folder	-	-	-
et_eda_model	Folder	-	-	-
et_eda_model_2	Folder	-	-	-
et_eda_model_3	Folder	-	-	-
et_eda_model_4	Folder	-	-	-
et_eda_model_5	Folder	-	-	-
et_eda_model_6	Folder	-	-	-
et_eda_model_7	Folder	-	-	-
et_eda_model_8	Folder	-	-	-
et_eda_model_9	Folder	-	-	-
et_eda_model_10	Folder	-	-	-
et_eda_model_11	Folder	-	-	-
et_eda_model_12	Folder	-	-	-
et_eda_model_13	Folder	-	-	-
et_eda_model_14	Folder	-	-	-
et_eda_model_15	Folder	-	-	-
et_eda_model_16	Folder	-	-	-
et_eda_model_17	Folder	-	-	-

13. Ejecutar en Google Colab el notebook “Visualizaciones.ipynb” para visualizar los resultados del EDA y el desempeño del modelo.

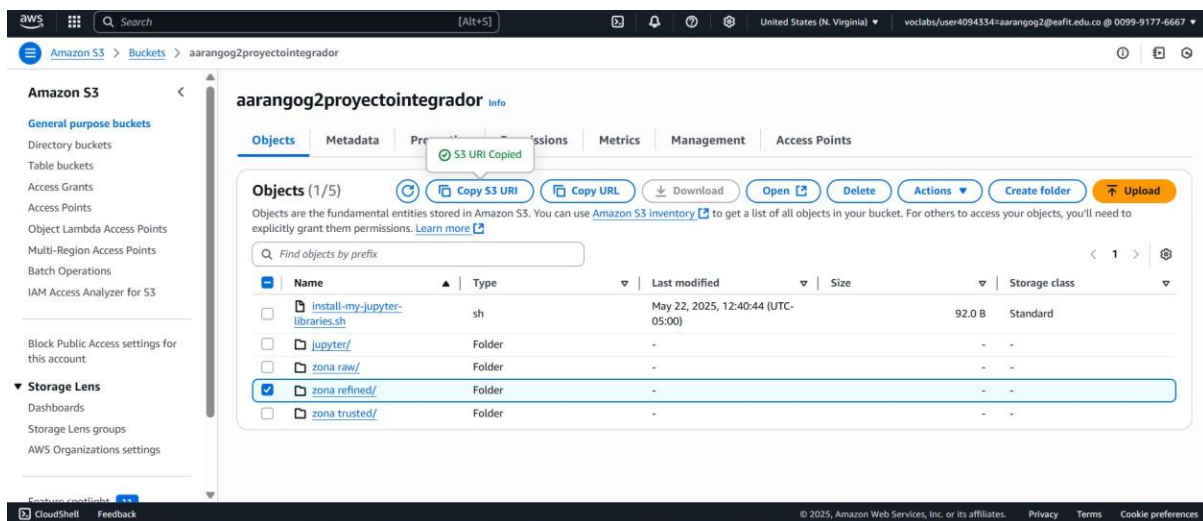
Para la ejecución exitosa del notebook, se deben modificar los siguientes parámetros:

- aws_access_key_id
- aws_secret_access_key
- aws_session_token

Se encuentran en la sección “AWS Details” de la terminal de la sesión de AWS creada en el “AWS Academy Learner Lab”.



El parámetro “s3_path” se puede obtener copiando la URI de la zona refined tras seleccionar la carpeta “zona refined/” y oprimiendo la opción “Copy S3 URI”.



Asumiendo los archivos se guardaron con los mismos nombres definidos en los notebooks “EDA.ipynb” y “Train Model.ipynb”, solo se debe correr el notebook para obtener las visualizaciones para “data_filtered”, “data_selected” y el desempeño del modelo.

