

CREACION DE CLUSTER EMR

Juan Pablo Gómez Triana¹

Resumen

Crear un Cluster AWS EMR en Amazon que contenga tecnologías como Hadoop, JupyterHub, Hive, Zeppelin, Hue y Spark; además de constar de un nodo maestro y 2 nodos core.

Palabras Clave: Cloud Computing, AWS, Cluster, EMR, Hadoop, JupyterHub, Hive, Zeppelin, Hue y Spark, Nodo Maestro, Nodo Core.

Introducción

En este documento se numeran los pasos requeridos para desplegar y configurar el Cluster planteado en los laboratorios 1 y 2 de la unidad 3.

¹ Estudiante de Ingeniería de Sistemas de la Universidad EAFIT, Medellín, Antioquia. E-mail: jpgomezt@eafit.edu.co

Creación del Bucket con S3

En esta sección mostraremos los pasos que fueron requeridos para la creación del Bucket con el sistema de AWS S3. Este bucket será utilizado para almacenar los datos de entrada o de salida de los programas que se ejecuten en el cluster. En este caso, crearemos el bucket para los datos del notebook.

Para iniciar, se instancio un Bucket en el servicio S3 de AWS:

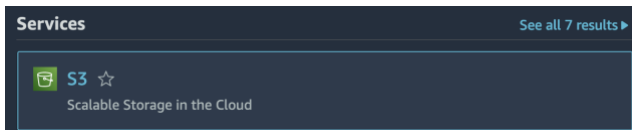


Figura 1. Servicio S3.

En este servicio, seleccionamos la opción de “Buckets” y seleccionamos “Create Bucket”.

Solo es necesario nombrar nuestro bucket, en este caso se llamará “notebooksjpgomezt”. Finalizamos seleccionando “Create Bucket”:

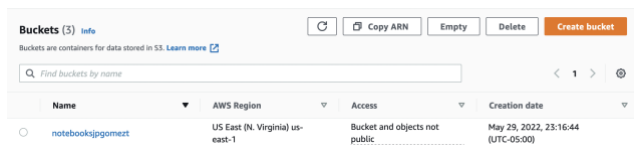


Figura 2. Bucket creado.

Creación del Cluster EMR

En esta sección mostraremos los pasos que fueron requeridos para la creación del Cluster con el sistema de AWS EMR. Para esto es importante primero poseer una clave para conectarse con el nodo mediante SSH. Para esto reutilizaremos la clave “AWS-KEY.pem” utilizada en laboratorios anteriores:

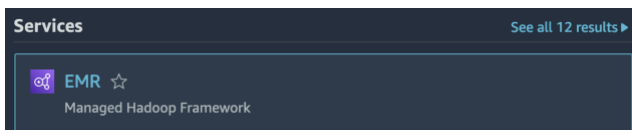


Figura 3. Servicio EMR.

En este servicio, seleccionamos la opción de “Clusters” y seleccionamos “Create cluster”.

Luego, es necesario seleccionar las opciones avanzadas, para poder configurar los softwares que utilizara el cluster.

Primero, seleccionamos la versión del EMR que utilizaremos. En este caso, utilizamos la versión “emr-6.3.1”. Luego seleccionamos los softwares que se instalaran en nuestro cluster:

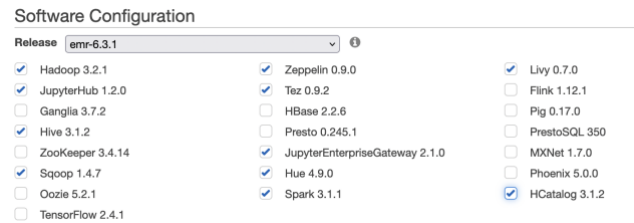


Figura 4. Software Configuration del Cluster.

Luego, hacemos integración entre el catálogo de Glue, lo que nos permite que las tablas que se creen sean visibles tanto desde Hive como desde Spark:

AWS Glue Data Catalog settings (optional)

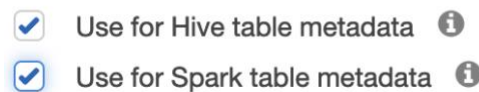


Figura 5. Integración de Glue con Hive y Spark.

Ahora configuramos la persistencia de los notebooks que se crearan con JupyterHub, utilizando el bucket que creamos con S3. Para esto, en “Edit software settings”, debemos entregar la siguiente configuración:

```
[
  {
    "Classification": "jupyter-s3-conf",
    "Properties": {
      "s3.persistance.enabled": "true",
      "s3.persistance.bucket": "notebooksjpgomezt"
    }
  }
]
```

Para el Hardware de nuestro cluster, en la opción “Cluster Nodes and Instances”, debemos cambiar el tipo de instancias que creara el servicio EMR de “m5.xlarge” a una “m4.xlarge”, y seleccionaremos la opción “Spot”:

Cluster Nodes and Instances

Choose the instance type, number of instances, and a purchasing option. [Learn more about instance purchasing options](#)

Console options for automatic scaling have changed. [Learn more](#)

Node type	Instance type	Instance count	Purchasing option
Master Master - 1	m4.xlarge 4 vCores, 16 GiB memory, EBS only storage EBS Storage: 64 GiB Add configuration settings	1 Instances	<input type="radio"/> On-demand <input checked="" type="radio"/> Spot Use on-demand as max price
Core Core - 2	m4.xlarge 4 vCores, 16 GiB memory, EBS only storage EBS Storage: 64 GiB Add configuration settings	2 Instances	<input type="radio"/> On-demand <input checked="" type="radio"/> Spot Use on-demand as max price

Figura 6. Hardware del Cluster.

Configuramos también la opción de “Auto-termination” para que el cluster sea destruido después de 1 hora de inactividad, y configuramos 20 GiB de almacenamiento:

Auto-termination

Select a time to have the cluster terminate after the cluster becomes idle. Choose a minimum of 1 minute or a max of 24 hours. [Learn more](#)

Auto-termination ☒ Enable auto-termination

Terminate cluster when it is idle after hours minutes

EBS Root Volume

Specify the root device volume size up to 100 GiB. This sizing applies to all instances in the cluster. [Learn more](#)

Root device EBS volume size GiB

Figura 7. Auto-termination y Almacenamiento.

Nombramos nuestro cluster en “Cluster name”. En este caso el cluster se llamará “My cluster jpgomez”. Finalmente, en “Security Options”, le asignamos la clave anteriormente mencionada “AWS-KEY.pem”. Una vez configurada estas opciones, podemos seleccionar “Create Cluster”:

Filter	All clusters	Filter clusters	2 clusters (all loaded)
Name	ID	Status	Creation time (UTC-8)
My cluster jpgomez	j-3TLNLSJWRLN	Starting	2022-05-30 12:23 (UTC-8)
			Elapsed time
			Normalized instance hours

Figura 8. Creación del Cluster.

Debemos esperar entre 20-30 minutos a que el cluster se instancie para poder interactuar con el cluster. Ahora configuraremos los puertos que el cluster tendrá abiertos.

Puertos del Cluster

Existen dos configuraciones para abrir los puertos requeridos para interactuar con las aplicaciones del cluster. Primero, vamos a la opción “Block public Access” dentro del servicio EMR. En esta podemos editar los puertos por los que queremos permitir el acceso:

Block public access settings

Block public access

On [Change](#)

Exceptions

A cluster can launch with security group rules that allow i

[Edit](#)

Port range ▲

22

8888-8888

8890-8890

9443-9443

Figura 9. Puertos permitidos en Block Public Access.

Ahora, si seleccionamos nuestro cluster, podemos ir a la opción “Security and Access”, y desde esta podemos configurar el security group del nodo master. Debemos añadir el puerto 8888, 9443, 8890, y el 22:

Port range ▼	Source ▼
8443	54.240.217.8/29
8443	54.240.217.64/28
8443	207.171.167.26/32
8443	72.21.217.0/24
8443	207.171.167.101/32
8443	207.171.172.6/32
All	sg-062f9b05756f7a71...
8888	0.0.0.0/0
All	sg-03b7e20072519dcf...
0 - 65535	sg-062f9b05756f7a71...
8443	72.21.198.64/29
8443	54.239.98.0/24
9443	0.0.0.0/0
22	0.0.0.0/0
8443	207.171.167.25/32
8443	54.240.217.16/29
0 - 65535	sg-03b7e20072519dcf...
8443	54.240.217.80/29

Figura 10. Puertos permitidos en el Security Group.

Configuración del Cluster EMR

Una vez nuestro cluster esté en funcionamiento, nos conectaremos por SSH para poder configurar Sqoop, de forma que lo podamos utilizar con la

interfaz web de Hue. Para conectarnos utilizamos el DNS público que ofrece el cluster para el nodo maestro:

```
ssh -i AWS-KEY.pem hadoop@ec2-54-162-195-99.compute-1.amazonaws.com
```

Ahora, necesitamos conocer el nombre del directorio lib donde se encuentran ubicados los componentes de oozie. Para esto podemos correr el siguiente comando:

```
hdfs dfs -ls /user/oozie/share/lib/
```

Esta operación nos lanzara de resultado el nombre del directorio:

[illegible]

Figura 11. Directorio Lib del EMR.

Podemos observar que el directorio tiene el nombre de “lib_20220530174018”. Ahora solo debemos correr los siguientes comandos para habilitar Sqoop en la interfaz web Hue:

```
hdfs dfs -put /usr/share/java/mysql-connector-  
java.jar  
/user/oozie/share/lib/lib_20220530174018/sqoop/
```

```
hdfs      dfs      -chown      oozie
/user/oozie/share/lib/lib_20220530174018/sqoop/
mysql-connector-java.jar
```

```
hdfs      dfs      -chgrp      oozie
/user/oozie/share/lib/lib_20220530174018/sqoop/
mysql-connector-java.jar
```

<i>hdfs</i>	<i>dfs</i>	<i>-cp</i>
<code>/user/oozie/share/lib/lib_20220530174018/hive/*</code>		
<code>/user/oozie/share/lib/lib_20220530174018/sqoop/</code>		

```
hdfs      dfs      -chown      oozie
/user/oozie/share/lib/lib_20220530174018/sqoop/
*
```

```
hdfs          dfs          -chgrp          oozie
/user/oozie/share/lib/lib_20220530174018/sqoop/
*
```

Finalmente, verificamos que el proceso haya sido exitoso y no ocurriera ningún error:

```
oozie admin -sharelibupdate
```

```
[ShareLib update status]
sharelibIdOld = hdfsf://ip-172-31-36-120.ec2.internal:8020/user/oozie/share/lib/lib_20220530174018
host = http://ip-172-31-36-120.ec2.internal:11000/oozie
sharelibIdNew = hdfsf://ip-172-31-36-120.ec2.internal:8020/user/oozie/share/lib/lib_20220530174018
status = Successful
```

Figura 12. Configuración de Sqoop.

Configuración de las Aplicaciones Web

Una vez nuestro cluster esté en funcionamiento, debemos acceder a la consola web de HUE. Para esto, utilizamos la dirección que nos ofrece el panel de “Application user interfaces”:

On-cluster application user interfaces

On-cluster UIs are available only while clusters are running. Because they are hosted on the master node, on-cluster UI require a connection via SSH tunneling. Set up SSH tunneling before accessing these application UI. [Learn more](#)

Application	User interface URL	Status
HCFS Name Node	http://ec2-54-162-195-99.compute-1.amazonaws.com:9870/	Available
Hue	http://ec2-54-162-195-99.compute-1.amazonaws.com:8888/	Available
JupyterLab	https://ec2-54-162-195-99.compute-1.amazonaws.com:9413/	Available
Zeppelin	https://ec2-54-162-195-99.compute-1.amazonaws.com:8880/	Available
Tel UI	http://ec2-54-162-195-99.compute-1.amazonaws.com:8080/hazelcast/	Available
Spark History Server	http://ec2-54-162-195-99.compute-1.amazonaws.com:18080/	Available
Livy	http://ec2-54-162-195-99.compute-1.amazonaws.com:8998/	Available
Resource Manager	http://ec2-54-162-195-99.compute-1.amazonaws.com:8086/	Available

Figura 12. Dirección de la interfaz web de HUE.

Una vez ingresamos, debemos crear un usuario y contraseña, y una vez hecho esto, podemos ingresar a la interfaz web de HUE:

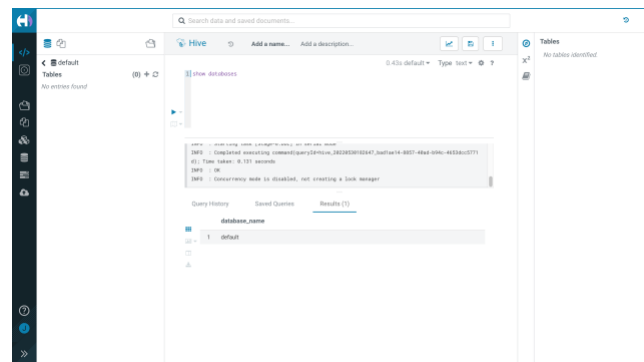


Figura 13. Interfaz web de HUE.

También podemos ingresar a la consola web de JupyterHub, buscando la dirección en el mismo

panel de “Application user interfaces”. Una vez ingresamos, debemos proveer un usuario y clave para ingresar. Por defecto, estos son “jovyan” y “jupyter” respectivamente. Una vez adentro, podemos empezar a crear notebooks:

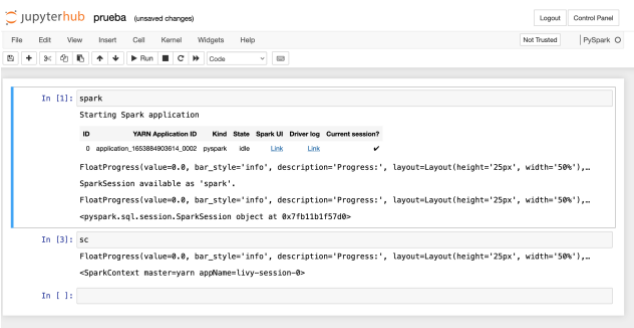


Figura 14. Notebook con PySpark en JupyterHub.

Finalmente, podemos ingresar a la consola web de Zeppelin, buscando la dirección en el mismo panel de “Application user interfaces”.

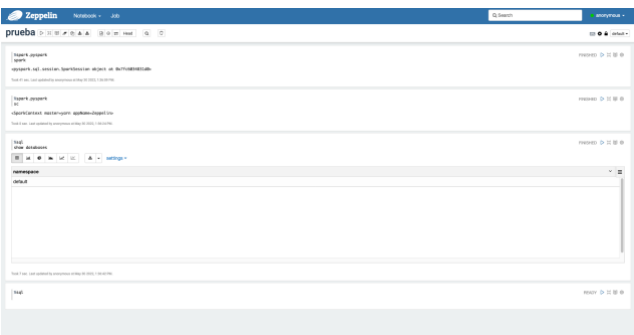


Figura 15. Notebook en Zeppelin.

Con esto ya tendríamos nuestro Cluster EMR configurado:

Name	ID	Status	Creation time (UTC-5)	Elapsed time	Normalized instance hours
My cluster jopagmet	j-STLDNUGJWPRLN	Waiting Cluster ready	2022-05-30 12:23 (UTC-5)	1 hour, 14 minutes	0

Figura 16. Cluster EMR Corriendo.

Referencias

[1] AWS. (s. f.-a). Adding Jupyter Notebook users and administrators - Amazon EMR. <https://docs.aws.amazon.com/emr/latest/ReleaseGuide/emr-jupyterhub-user-access.html>

[2] AWS. (s. f.-b). Configuring persistence for notebooks in Amazon S3 - Amazon EMR. <https://docs.aws.amazon.com/emr/latest/ReleaseGuide/emr-jupyterhub-s3.html>

[3] Montoya, E. N. [Edwin Nelson Montoya]. (2021a, noviembre 3). AWS EMR 6 3 1 parte1 20211103 [Vídeo]. YouTube. <https://www.youtube.com/watch?v=MyXSwxN5Zdk>

[4] Montoya, E. N. [Edwin Nelson Montoya]. (2021b, noviembre 3). AWS EMR 6 3 1 Parte2 20211103 [Vídeo]. YouTube. <https://www.youtube.com/watch?v=3sao-qJG34Y>

[5] Montoya, J. C. (2022, 25 mayo). ST0263/st0263-2022-1. GitHub. <https://github.com/ST0263/st0263-2022-1>