**Tutorial uso de mapreduce y hdfs**

Pablo Moreno Quintero

Escuela de Ciencias e Ingeniería, Universidad EAFIT

Pregrado en Ingeniería de Sistemas

Edwin Nelson Montoya Munera
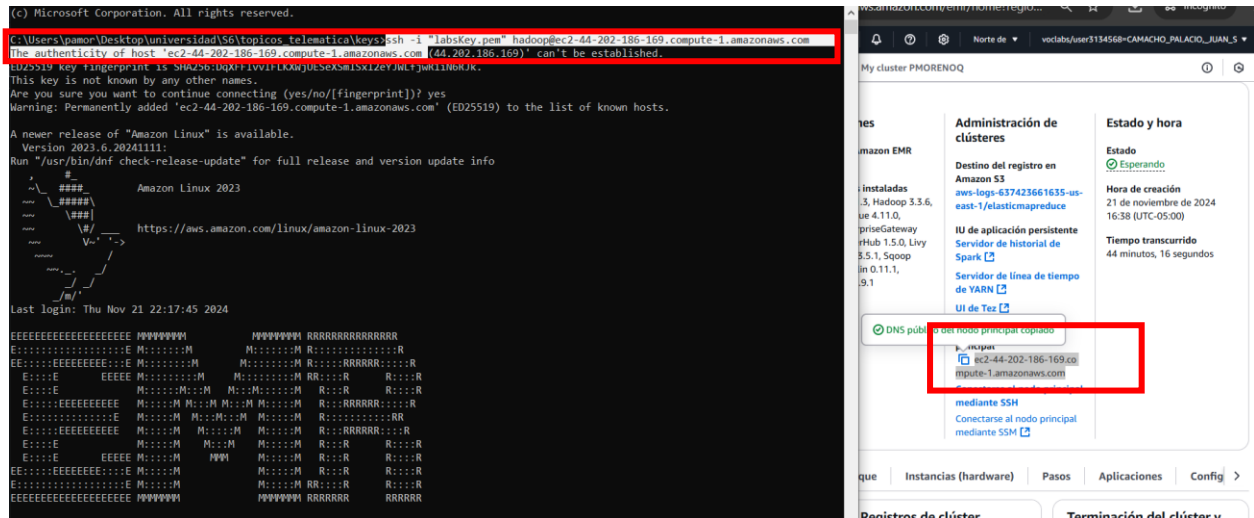
23 de noviembre de 2024

1. Acceder a nuestra instancia EC2 del nodo máster.

En consola

ssh -i "<ubicacion de la clave .pem>" hadoop@<DNS clúster principal>

Posteriormente nos debe mostrar la siguiente pantalla en nuestra consola



Luego de acceder a nuestra instancia master debemos descargar git

Damos "yes" cuando lo pida.

```
[hadoop@ip-172-31-15-33 ~]$ sudo yum install git
Last metadata expiration check: 0:50:47 ago on Thu Nov 21 21:39:25 2024.
Dependencies resolved.
================================================================================
 Package              Architecture    Version                  Repository    Size
================================================================================
Installing:
 git                  x86_64          2.40.1-1.amzn2023.0.3    amazonlinux   54 k
Installing dependencies:
 git-core             x86_64          2.40.1-1.amzn2023.0.3    amazonlinux   4.3 M
 git-core-doc         noarch          2.40.1-1.amzn2023.0.3    amazonlinux   2.6 M
 perl-Error           noarch          1:0.17029-5.amzn2023.0.2 amazonlinux   41 k
 perl-Git             noarch          2.40.1-1.amzn2023.0.3    amazonlinux   42 k
 perl-TermReadKey     x86_64          2.38-9.amzn2023.0.2      amazonlinux   36 k
 perl-lib             x86_64          0.65-477.amzn2023.0.6    amazonlinux   15 k

Transaction Summary
================================================================================
Install  7 Packages

Total download size: 7.1 M
Installed size: 34 M
Is this ok [y/N]: y
Downloading Packages:
(1/7): git-2.40.1-1.amzn2023.0.3.x86_64.rpm          941 kB/s |  54 kB   00:00
(2/7): git-core-doc-2.40.1-1.amzn2023.0.3.noarch.rpm  29 MB/s | 2.6 MB   00:00
(3/7): perl-Error-0.17029-5.amzn2023.0.2.noarch.rpm  1.2 MB/s |  41 kB   00:00
(4/7): git-core-2.40.1-1.amzn2023.0.3.x86_64.rpm      34 MB/s | 4.3 MB   00:00
(5/7): perl-Git-2.40.1-1.amzn2023.0.3.noarch.rpm     1.1 MB/s |  42 kB   00:00
(6/7): perl-TermReadKey-2.38-9.amzn2023.0.2.x86_64.rpm 956 kB/s | 36 kB  00:00
(7/7): perl-lib-0.65-477.amzn2023.0.6.x86_64.rpm     737 kB/s |  15 kB   00:00
--------------------------------------------------------------------------------
Total                                                 35 MB/s | 7.1 MB   00:00
Running transaction check
Transaction check succeeded.
Running transaction test
Transaction test succeeded.
```

Clonamos el repositorio

Git clone https://github.com/st0263eafit/st0263-242

```
[hadoop@ip-172-31-15-33 ~]$ git clone https://github.com/st0263eafit/st0263-242.git
```

Utilizamos los siguientes comandos para crear directorios de hdfs

Hdfs dfs -mkdir <ruta>

Y el siguiente comando para verificar su creación

Hdfs dfs -ls <ruta>

```
[hadoop@ip-172-31-12-162 ~]$ hdfs dfs -mkdir /user/hadoop/datasets
[hadoop@ip-172-31-12-162 ~]$ hdfs dfs -ls /user/hadoop/datasets
[hadoop@ip-172-31-12-162 ~]$ hdfs dfs -mkdir /user/hadoop/datasets/gutenberg-small
[hadoop@ip-172-31-12-162 ~]$ hdfs dfs -ls  /user/hadoop/datasets/gutenberg-small
```

Movemos los archivos descargados a hdfs

hdfs dfs -put /home/ec2-home/st0263-242/bigdata/datasets/gutenberg-small/*.txt /user/hadoop/datasets/gutenberg-small/

```
[hadoop@ip-172-31-15-33 ~]$ hdfs dfs -ls /
Found 4 items
drwxr-xr-x   - hdfs hdfsadmingroup          0 2024-11-21 21:49 /apps
drwxrwxrwt   - hdfs hdfsadmingroup          0 2024-11-21 21:51 /tmp
drwxr-xr-x   - hdfs hdfsadmingroup          0 2024-11-21 21:49 /user
drwxr-xr-x   - hdfs hdfsadmingroup          0 2024-11-21 21:49 /var
[hadoop@ip-172-31-15-33 ~]$ hdfs dfs -ls /user
Found 9 items
drwxrwxrwx   - hadoop    hdfsadmingroup     0 2024-11-21 21:49 /user/hadoop
drwxr-xr-x   - mapred    mapred             0 2024-11-21 21:49 /user/history
drwxrwxrwx   - hdfs      hdfsadmingroup     0 2024-11-21 21:49 /user/hive
drwxrwxrwx   - hue       hue                0 2024-11-21 21:49 /user/hue
drwxrwxrwx   - livy      livy               0 2024-11-21 21:49 /user/livy
drwxrwxrwx   - oozie     oozie              0 2024-11-21 21:50 /user/oozie
drwxrwxrwx   - root      hdfsadmingroup     0 2024-11-21 21:49 /user/root
drwxrwxrwx   - spark     spark              0 2024-11-21 21:49 /user/spark
drwxrwxrwx   - zeppelin  hdfsadmingroup     0 2024-11-21 21:49 /user/zeppelin
```

hdfs dfs -get /user/hadoop/datasets/gutenberg-small/* user/datasets/

verificamos

```
[hadoop@ip-172-31-12-162 ~]$ sudo ls /home/hadoop/st0263-242/bigdata/datasets/gutenberg-small/
AbrahamLincoln___LincolnLetters.txt
AbrahamLincoln___LincolnsFirstInauguralAddress.txt
AbrahamLincoln___LincolnsGettysburgAddressGivenNovember-19-1863.txt
AbrahamLincoln___LincolnsInauguralsAddressesandLettersSelections.txt
AbrahamLincoln___LincolnsSecondInauguralAddress.txt
AbrahamLincoln___SpeechesandLettersofAbrahamLincoln1832-1865.txt
AbrahamLincoln___StateoftheUnionAddresses.txt
AbrahamLincoln___TheEmancipationProclamation.txt
AbrahamLincoln___TheLifeandPublicServiceofGeneralZacharyTaylorAnAddress.txt
AbrahamLincoln___TheWritingsofAbrahamLincolnVolume1.txt
AbrahamLincoln___TheWritingsofAbrahamLincolnVolume2.txt
AbrahamLincoln___TheWritingsofAbrahamLincolnVolume3.txt
AbrahamLincoln___TheWritingsofAbrahamLincolnVolume4.txt
AbrahamLincoln___TheWritingsofAbrahamLincolnVolume5.txt
AbrahamLincoln___TheWritingsofAbrahamLincolnVolume6.txt
AbrahamLincoln___TheWritingsofAbrahamLincolnVolume7.txt
```

```
[hadoop@ip-172-31-12-162 ~]$ hdfs dfs -ls /user/hadoop/datasets
Found 1 items
drwxr-xr-x   - hadoop hdfsadmingroup          0 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small
[hadoop@ip-172-31-12-162 ~]$ hdfs dfs -ls /user/hadoop/datasets/gutenberg-small/
Found 16 items
-rw-r--r--   1 hadoop hdfsadmingroup       5878 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___LincolnLetters.txt
-rw-r--r--   1 hadoop hdfsadmingroup      21586 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___LincolnsFirstInauguralAddress.txt
-rw-r--r--   1 hadoop hdfsadmingroup       1653 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___LincolnsGettysburgAddressGivenNovember-19-1863.txt
-rw-r--r--   1 hadoop hdfsadmingroup     262083 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___LincolnsInauguralsAddressesandLettersSelections.txt
-rw-r--r--   1 hadoop hdfsadmingroup       4093 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___LincolnsSecondInauguralAddress.txt
-rw-r--r--   1 hadoop hdfsadmingroup     516298 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___SpeechesandLettersofAbrahamLincoln1832-1865.txt
-rw-r--r--   1 hadoop hdfsadmingroup     167895 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___StateoftheUnionAddresses.txt
-rw-r--r--   1 hadoop hdfsadmingroup       3928 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheEmancipationProclamation.txt
-rw-r--r--   1 hadoop hdfsadmingroup      45664 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheLifeandPublicServiceofGeneralZacharyTaylorAnAddress.txt
-rw-r--r--   1 hadoop hdfsadmingroup     459006 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheWritingsofAbrahamLincolnVolume1.txt
-rw-r--r--   1 hadoop hdfsadmingroup     505150 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheWritingsofAbrahamLincolnVolume2.txt
-rw-r--r--   1 hadoop hdfsadmingroup     254941 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheWritingsofAbrahamLincolnVolume3.txt
-rw-r--r--   1 hadoop hdfsadmingroup     209643 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheWritingsofAbrahamLincolnVolume4.txt
-rw-r--r--   1 hadoop hdfsadmingroup     692051 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheWritingsofAbrahamLincolnVolume5.txt
-rw-r--r--   1 hadoop hdfsadmingroup     601102 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheWritingsofAbrahamLincolnVolume6.txt
-rw-r--r--   1 hadoop hdfsadmingroup     478689 2024-11-22 23:35 /user/hadoop/datasets/gutenberg-small/AbrahamLincoln___TheWritingsofAbrahamLincolnVolume7.txt
```

2. Gestion de archivos mediante HUE en Amazon EMR

Accedemos al servicio de HUE en las aplicaciones de nuestro clúster



Luego de iniciar sesión veremos la siguiente interfaz

Nos dirigimos en la barra de navegación a "files"



Navegamos hasta los archivos que compartimos

Vemos un archivo cargado en el sistema.