

Ciclo de vida Analítico del dato

Práctica 1: Pig y Hive

Daniel Pérez Efremova

Índice

Aclaraciones	3
Creación de directorios en HDFS	4
Creación de la base de datos para Hive	4
Consultas	5
Primera consulta (C1)	5
Segunda consulta (C2)	6
Tercera consulta (C3)	6
Cuarta consulta (C4)	7
Resultados	9
Conslusiones	10
Anexo	11

Aclaraciones

La práctica se ha realizado con un pc personal, es decir, sin la máquina virtual del curso. No se ha utilizado dockers.

Aunque no se pedía para el ejercicio, se han desarrollado scripts para facilitar y agilizar el proceso de pruebas. Como demostración de autoría y por sencillez en la exposición, se adjuntan en los anexos varias capturas en función del programa:

- **Pig local:** Capturas de pantalla de la consola donde se ha ejecutado la consulta.
- **Pig mapreduce:** captura de pantalla de la interfaz de hadoop donde se ve la ruta donde se almacenan los correspondientes resultados en formato csv.
- **Hive:** captura de pantalla de la interfaz de hadoop donde se ven las bases de datos creadas y su contenido.

Para facilitar la corrección del ejercicio, se ha preparado el script bash *ejecuciones.sh*. Para usarlo, basta ejecutarlo pasando como argumento la ruta local al fichero *athleteEvents.csv* para los scripts de Pig local. Este script se encarga de subir a HDFS los datos, reiniciar el metastore, levantar el servicio Derby y ejecutar todas las consultas secuencialmente.

Creación de directorios en HDFS

En primer lugar se crean los directorios *hadoop/dataset* y *hadoop/out*.

Código 1: Creación de los directorios de trabajo

```
daniel@daniel:~$ hdfs dfs -ls /

drwxr-xr-x - daniel supergroup 0 2021-09-18 17:53 /test
drwx----- - daniel supergroup 0 2021-09-19 11:55 /tmp
drwxr-xr-x - daniel supergroup 0 2021-09-18 18:12 /user

daniel@daniel:~$ hdfs dfs -mkdir /hadoop/dataset
daniel@daniel:~$ hdfs dfs -mkdir /hadoop/out

daniel@daniel:~$ hdfs dfs -ls /hadoop/
Found 2 items
drwxr-xr-x - daniel supergroup 0 2021-09-19 17:52 /hadoop/dataset
drwxr-xr-x - daniel supergroup 0 2021-09-19 18:01 /hadoop/out
```

Después, se copia el archivo *athleteEvents.csv* en HDFS, se comprueba que efectivamente está en el directorio y se listan las primeras líneas. En el código 2 se ve que el archivo está cargado con un tamaño aproximado de 42 Mb.

Código 2: Ingesta de datos

```
daniel@daniel:~$ hdfs dfs -put ~/Downloads/athleteEvents.csv
/hadoop/dataset

daniel@daniel:~$ hdfs dfs -ls /hadoop/dataset
Found 1 items
-rw-r--r-- 1 daniel supergroup 41500688 2021-09-19 18:08
/hadoop/dataset/athleteEvents.csv

daniel@daniel:~$ hdfs dfs -cat /hadoop/dataset/athleteEvents.csv | head

"ID","Name","Sex","Age","Height","Weight","Team","NOC","Games","Year","
Season","City","Sport","Event","Medal"
"1","A Dijiang","M",24,180,80,"China","CHN","1992 Summer",1992,"Summer","
Barcelona","Basketball","Basketball Men's Basketball",NA
"2","A Lamusi","M",23,170,60,"China","CHN","2012 Summer",2012,"Summer","
London","Judo","Judo Men's Extra-Lightweight",NA
"3","Gunnar Nielsen Aaby","M",24,NA,NA,"Denmark","DEN","1920 Summer",1920,
"Summer","Antwerpen","Football","Football Men's Football",NA ...
```

Creación de la base de datos para Hive

Para crear la base de datos, se usa el script *BD_practicas.hql*. Este script se encarga de crear dos bases de datos:

- **practicas**. Almacena los datos del ejercicio.
- **results_P1_pighive**. Almacena los resultados de las consultas en formato .csv.

Código 3: Creación de bases de datos

```
CREATE DATABASE IF NOT EXISTS practicas;  
CREATE DATABASE IF NOT EXISTS results_P1_pighive;
```

Por otro lado, se crea la tabla atletas con los datos del ejercicio en *practicas* y se sobrescribe con los datos desde HDFS. Finalmente, se muestran por consola las bases de datos y las primeras líneas de los datos, para controlar que todo se ha ejecutado correctamente.

Código 4: Creación de bases de datos

```
SET input='hdfs:///hadoop/dataset/athleteEvents.csv';  
CREATE TABLE IF NOT EXISTS practicas.atletas (  
ID int,  
Name string,  
Sex string,  
Age int,  
Height int,  
Weight int,  
Team string,  
NOC string,  
Games string,  
Year int,  
Season string,  
City string,  
Sport string,  
Event string,  
Medal string  
)  
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'  
WITH SERDEPROPERTIES (  
"separatorChar" = ",",  
"quoteChar" = "\""  
)  
STORED AS TEXTFILE  
TBLPROPERTIES("skip.header.line.count"="1");  
  
LOAD DATA INPATH {hiveconf:input}  
OVERWRITE INTO TABLE practicas.atletas;  
  
-- Comprobacion creacion DB  
SHOW DATABASES;  
  
-- Comprobacion de la carga bde datos  
SELECT * FROM practicas.atletas LIMIT 5;
```

Consultas

Esta sección se divide en cuatro partes donde se presentan y explican las ideas detrás de cada script solución. En cada código se presentan la consultas, tanto para Pig como para Hive.

Primera consulta (C1)

La idea en los scripts consiste en crear un subconjunto del dataset que contenga sólo las medallas de Oro. Después, calcular el mínimo de la columna *Age* sobre este subconjunto. Finalmente, se filtra por la edad mínima.

Código 5: C1 para PIG

```
medallas_oro = FILTER atletas BY Medal == 'Gold';

edad_minima = FOREACH (GROUP medallas_oro ALL)
GENERATE MIN(medallas_oro.Age) AS valor:int;

medallistas = FOREACH (
FILTER medallas_oro BY Age == edad_minima.valor)
GENERATE Name, Games, Year, Sport, Event;
DUMP medallistas;

STORE medallistas INTO '$output' USING PigStorage (',');
```

Código 6: C1 para Hive

```
SELECT name, age, year, games, sport, event FROM practicas.atletas WHERE
    Age IN
(
SELECT MIN(Age) FROM practicas.atletas WHERE Medal = 'Gold'
) AND Medal = 'Gold';
```

Segunda consulta (C2)

La estrategia consiste en filtrar los datos de atletas por los que tengan a España por equipo. Después, basta ordenarlos descendientemente por la variable nombre.

Código 7: C2 para Pig

```
espanoles = FILTER atletas BY Team == 'Spain';
espanoles_ord = ORDER espanoles BY Name;
DUMP espanoles_ord;
```

Código 8: C2 para Hive

```
SELECT DISTINCT * FROM practicas.atletas WHERE Team='Spain' ORDER BY Name
DESC;
```

Tercera consulta (C3)

En esta consulta se distinguen dos partes, la primera para obtener el mejor atleta en la disciplina de vela y la segunda para obtener todas las olimpiadas por continente.

Para la primera parte, se considera el mejor atleta al que tenga más medallas de oro en toda la historia de los juegos olímpicos. Primero, se filtran los datos de atletas por la disciplina sailing y por medallas de oro¹. Después, para obtener el número de medallas de cada atleta se hace un recuento de las veces que aparece cada deportista en los datos filtrados. Lo siguiente es calcular el máximo de medallas sobre el recuento². Finalmente, se seleccionan los deportistas que tengan un número de medallas de oro igual al máximo calculado.

Código 9: C3 para el mejor atleta Sailing en Pig

```
-- filtro por medallas de oro y sailing
```

¹Este paso puede ser conflictivo, porque existen valores perdidos en la columna *Medal*.

²Si pudiésemos garantizar que el máximo es único, bastaría con ordenar los datos descendientemente y coger el primer valor, simplificando mucho la consulta. Este no es el caso y se recurre al cálculo explícito del máximo.

```

deportistas = FOREACH (
FILTER atletas BY Sport == 'Sailing' and Medal=='Gold')
GENERATE Name;
group_deportistas = GROUP deportistas BY Name;
deportistas_count = FOREACH group_deportistas
GENERATE group AS Name, COUNT(deportistas) AS Medallas;

-- calculo del maximo de medallas de oro
max_medallas = FOREACH (GROUP deportistas_count ALL)
GENERATE MAX(deportistas_count.Medallas) AS valor:int;

mejores = FOREACH (
FILTER deportistas_count BY Medallas == max_medallas.valor)
GENERATE Name, Medallas;
DUMP mejores;

```

Código 10: C3 para el mejor atleta en Hive

```

-- tabla auxiliar para contar medallas de oro en sailing

CREATE TEMPORARY TABLE recuento_sailing AS
SELECT Name, COUNT(*) AS recuento FROM practicas.atletas
WHERE Medal = 'Gold' AND Sport = 'Sailing'
GROUP BY Medal, Name;

SELECT Name, recuento FROM recuento_sailing
WHERE recuento IN (SELECT MAX(recuento) FROM recuento_sailing);

```

Para la segunda parte, se seleccionan los valores distintos de las variables que contienen el año, ciudad y estación. Para esto, basta agrupar por las tres variables y generar las ternas de valores de cada grupo.

Código 11: C3 para las olimpiadas por continente en Pig

```

olimp_agrup = GROUP atletas BY (Year, Season, City);
olimpiadas = FOREACH olimp_agrup GENERATE
group.Year,
group.City,
group.Season;

DUMP olimpiadas;

```

Código 12: C3 para las olimpiadas por continente en Hive

```

SELECT DISTINCT Year, City, Season FROM practicas.atletas;

```

Cuarta consulta (C4)

La estrategia en esta consulta consiste en contar cuántas medallas tiene cada deportista (de oro, plata y bronce). Para realizar el recuento, se crean 3 columnas dummy, una para cada tipo de medalla. Cada columna tomará el valor 1 si el campo *medalla* toma ese valor y 0 en caso contrario. Una vez se tienen estas nuevas columnas, basta agrupar los datos por el nombre del atleta y sumar las tres columnas. Por último, para obtener el total de medallas se suman las tres columnas.

Código 13: C4, recuento de medallas en Pig

```

dummies = FOREACH atletas GENERATE

```

```

Name,
(CASE WHEN Medal == 'Gold' THEN 1 ELSE 0 END) AS Gold:int,
(CASE WHEN Medal == 'Silver' THEN 1 ELSE 0 END) AS Silver:int,
(CASE WHEN Medal == 'Bronze' THEN 1 ELSE 0 END) AS Bronze:int;
--(CASE WHEN Medal == 'NA' THEN 1 ELSE 0 END) AS NA:int;

dummies_group = GROUP dummies BY Name;
dummies_count = FOREACH dummies_group GENERATE
group as Name,
SUM(dummies.Gold) AS Medallas_oro:int,
SUM(dummies.Silver) AS Medallas_plata:int,
SUM(dummies.Bronze) AS Medallas_bronce:int;

total_medallas = FOREACH dummies_count GENERATE
Name,
Medallas_oro,
Medallas_plata,
Medallas_bronce,
--Medallas_NA,
(Medallas_oro+Medallas_plata+Medallas_bronce) AS Total;

```

Código 14: C4, recuento de medallas en Hive

```

-- tabla auxiliar para contar las medallas de cada atleta
CREATE TEMPORARY TABLE recuento_medallas AS
SELECT Name,
SUM(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold,
SUM(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver,
SUM(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze
FROM practicas.atletas
GROUP BY Name;

-- tabla auxiliar para contar el total de medallas por atleta
CREATE TEMPORARY TABLE recuento_medallas_total AS
SELECT *, Gold+Silver+Bronze AS Total FROM recuento_medallas;

```

Finalmente, para obtener el atleta con más medallas, se realiza un proceso similar al de la consulta C3, se calcula el máximo de la columna que almacena el total de medallas y se filtra por ese valor³. Para calcular el atleta con más medallas de oro, se realiza el mismo proceso anterior pero sobre la columna recuento de medallas de oro.

Código 15: C4 filtro por el deportista con más medallas en total y de oro en Pig

```

max_medallas = FOREACH (GROUP total_medallas ALL)
GENERATE MAX(total_medallas.Total) AS valor:int;

max_medallas_oro = FOREACH (GROUP total_medallas ALL)
GENERATE MAX(total_medallas.Medallas_oro) AS valor:int;

mejores = FOREACH (
FILTER total_medallas BY Total == max_medallas.valor)
GENERATE Name, Medallas_oro, Medallas_plata, Medallas_bronce;

mejores_oro = FOREACH (

```

³Nuevamente, no podemos garantizar que el máximo sea único y se recurre al cálculo explícito del máximo


```

FILTER total_medallas BY Medallas_oro == max_medallas_oro.valor)
GENERATE Name, Medallas_oro, Medallas_plata, Medallas_bronce;

DUMP mejores;
DUMP mejores_oro;

```

Código 16: C4 filtro por el deportista con más medallas en total y de oro en Hive

```

SELECT Name, Gold FROM recuento_medallas WHERE Gold IN (SELECT MAX(Gold)
FROM recuento_medallas);

SELECT Name, Total FROM recuento_medallas_total WHERE Total IN (SELECT MAX
(Total) FROM recuento_medallas_total);

```

Resultados

En esta sección se presenta una comparativa de tiempos de ejecución para cada consulta en cada una de las tecnologías y, en las consultas que aplique, la respuesta.

Los tiempos de ejecución de cada script se midieron usando el comando Linux *time*. En la Figura 1 se observan los resultados. Los tiempos más bajos se obtienen para Pig local. Sin embargo, los tiempos de ejecución para Pig mapreduce⁴ y Hive no son significativamente distintos.

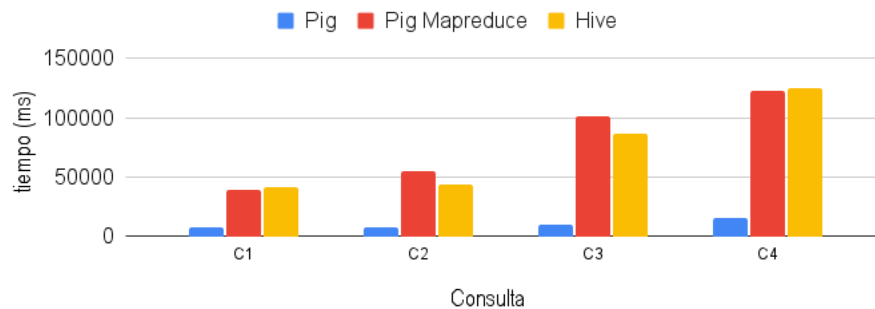


Figura 1: Gráfico de tiempos de ejecución

En la tabla 1 se muestra la cantidad de trabajos. En el primer valor los *map* y en el segundo los *reduce*. Se ve que Hive realiza la menor cantidad de tareas para cada consulta, seguido de pig Mapreduce y, finalmente, Pig local es el que más tareas realiza.

Tabla 1: Número de trabajos realizados.

Consulta	Pig local	Pig mapreduce	Hive
C1	(3, 1)	(2, 1)	(2, 1)
C2	(4, 2)	(3, 2)	(2, 2)
C3	(6, 3)	(5, 4)	(4, 3)
C4	(8, 4)	(6, 4)	(6, 3)

Las respuestas a las consultas son:

⁴El tiempo de ejecución en este caso se midió sin tener en cuenta el tiempo de almacenamiento de los resultados en HDFS

- **C3:** Charles Benedict Ainslie y Paul Bert Elvström. Ambos con cuatro medallas de oro.
- **C4:** Michael Phelps con 28 medallas, 23 de ellas de Oro. Es el mejor por ambos criterios, total de medallas y total de medallas de oro.

Conclusiones

A grandes rasgos, en esta práctica hay dos aspectos importantes a comparar entre Pig y Hive:

- La sintaxis para el desarrollo de consultas o creación de conjuntos de datos.
- Eficiencia en el cómputo de las consultas.

Respecto a la sintaxis, queda patente que Hive (HQL), por su gran similitud con SQL y legibilidad, es más accesible y cómodo para el público general. Sin embargo, Pig se demuestra más versátil que Hive, y podría integrarse en procesos *batch* para generar nuevos conjuntos de datos desde datos en bruto, al contrario que Hive, que está más orientado a *reporting* y datos ya estructurados. En este aspecto, se concluye que Pig es una herramienta para el entorno técnico mientras que Hive se centra en el entorno funcional de negocio.

Por otro lado, respecto a la eficiencia de cómputo, el tamaño de los datos no permite llegar a conclusiones claras. Pig local ha sido el más rápido, ya que sobre ficheros de pequeño tamaño, el acceso a disco es mucho más rápido a local que a HDFS. Respecto a Pig mapreduce y Hive, con este experimento no se puede afirmar que uno sea más eficiente que otro. Sin embargo, no se espera que los tiempos de ejecución aumenten significativamente al variar el tamaño del archivo, algo que en Pig local sí cabe esperar.

Pese a que se ha realizado en modo *standalone*, no cabe duda que ambas tecnologías constituyen una herramienta muy potente para el tratamiento de datos masivos y, en particular, el punto fuerte es que consiguen dar acceso a un gran público al paradigma mapreduce.

Anexo

```
daniel@daniel: ~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida
2021-10-15 23:08:29,899 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
3.3.0 0.17.0 daniel 2021-10-15 23:08:27 2021-10-15 23:08:29 GROUP_BY,FILTER

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Featu
re Outputs
job_local1966259457_0003 2 1 n/a n/a n/a n/a n/a n/a n/a 1-9,atletas,edad_minima,medallas_oro MULTI_QUERY,COMBINER
job_local1995447293_0004 1 0 n/a n/a n/a n/a n/a 0 0 1-10,medallistas MAP_ONLY file:///home/daniel/D
esktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C1_local,

Input(s):
Successfully read 271116 records from: "file:///home/daniel/Downloads/athleteEvents.csv"

Output(s):
Successfully stored 7 records in: "file:///home/daniel/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C1_local"

Counters:
Total records written : 7
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1966259457_0003 -> job_local1995447293_0004,
job_local1995447293_0004

2021-10-15 23:08:29,900 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:08:29,902 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:08:29,903 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:08:29,908 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:08:29,910 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:08:29,911 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:08:29,913 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAI
LED 132520 time(s).
2021-10-15 23:08:29,913 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-10-15 23:08:29,941 [main] INFO org.apache.pig.Main - Pig script completed in 7 seconds and 662 milliseconds (7662 ms)
(base) daniel@daniel:~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida$
```

Figura 2: Salida por consola C1 pig local

```
daniel@daniel: ~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida
Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Featu
re Outputs
job_local1575153481_0004 2 0 n/a n/a n/a n/a 0 0 0 0 atletas,espanoles MAP_ONLY
job_local665791577_0005 1 1 n/a n/a n/a n/a n/a n/a n/a n/a espanoles_ord SAMPLER
job_local932485370_0006 1 1 n/a n/a n/a n/a n/a n/a n/a n/a espanoles_ord ORDER_BY file:///home/daniel/Desktop/UAM_Msc_B
igData/ciclo_vida/practica1_ciclo_vida/resultados_local/C2_local,

Input(s):
Successfully read 271116 records from: "file:///home/daniel/Downloads/athleteEvents.csv"

Output(s):
Successfully stored 5224 records in: "file:///home/daniel/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C2_local"

Counters:
Total records written : 5224
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1575153481_0004 -> job_local665791577_0005,
job_local665791577_0005 -> job_local932485370_0006,
job_local932485370_0006

2021-10-15 23:10:24,027 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,028 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,029 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,031 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,032 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,033 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,036 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,039 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,040 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:24,042 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAI
LED 132520 time(s).
2021-10-15 23:10:24,042 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-10-15 23:10:24,062 [main] INFO org.apache.pig.Main - Pig script completed in 7 seconds and 991 milliseconds (7991 ms)
(base) daniel@daniel:~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida$
```

Figura 3: Salida por consola C2 pig local

```
daniel@daniel: ~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida
2021-10-15 23:10:57,656 [pool-25-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local1559498175_0008_r_000000_0
2021-10-15 23:10:57,656 [Thread-48] INFO org.apache.hadoop.mapred.LocalJobRunner - reduce task executor complete.
2021-10-15 23:10:57,694 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_local1559498175_0008]
2021-10-15 23:10:57,810 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:57,813 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:57,815 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:57,824 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2021-10-15 23:10:57,825 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
3.3.0 0.17.0 daniel 2021-10-15 23:10:56 2021-10-15 23:10:57 GROUP_BY

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Featu
re Outputs
job_local1559498175_0008 2 1 n/a n/a n/a n/a n/a n/a n/a atletas,olimp_agrup,olimpiadas GROUP_BY file:///home/
daniel/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C3_local_bis,

Input(s):
Successfully read 271116 records from: "file:///home/daniel/Downloads/athleteEvents.csv"

Output(s):
Successfully stored 52 records in: "file:///home/daniel/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C3_local_bis"

Counters:
Total records written : 52
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1559498175_0008

2021-10-15 23:10:57,828 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:57,831 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:57,833 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:10:57,839 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-10-15 23:10:57,866 [main] INFO org.apache.pig.Main - Pig script completed in 10 seconds and 526 milliseconds (10526 ms)
(base) daniel@daniel:~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida$
```

Figura 4: Salida por consola C3 pig local

```
daniel@daniel: ~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida
AM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C4_local,
job_local1267418590_0007 2 1 n/a n/a n/a n/a n/a n/a atletas,dummies,dummies_count,dummies_group,total_medallas G
ROUP_BY,COMBINER
job_local1473218330_0010 1 0 n/a n/a n/a n/a 0 0 0 1-28,mejores_oro,total_medallas MAP_ONLY file:///home/
daniel/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C4_local_oro,
job_local1534280372_0008 1 1 n/a n/a n/a n/a n/a n/a 1-25,1-26,max_medallas,max_medallas_oro MULTI_QUERY,COMBINER

Input(s):
Successfully read 271116 records from: "file:///home/daniel/Downloads/athleteEvents.csv"

Output(s):
Successfully stored 1 records in: "file:///home/daniel/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C4_local"
Successfully stored 1 records in: "file:///home/daniel/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida/resultados_local/C4_local_oro"

Counters:
Total records written : 2
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1267418590_0007 -> job_local1534280372_0008,job_local1083857694_0009,job_local1473218330_0010,
job_local1534280372_0008 -> job_local1083857694_0009,job_local1473218330_0010,
job_local1083857694_0009
job_local1473218330_0010

2021-10-15 23:11:48,875 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,877 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,878 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,882 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,884 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,886 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,890 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,892 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,893 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,895 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,897 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,898 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-10-15 23:11:48,899 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-10-15 23:11:48,924 [main] INFO org.apache.pig.Main - Pig script completed in 15 seconds and 758 milliseconds (15758 ms)
(base) daniel@daniel:~/Desktop/UAM_Msc_BigData/ciclo_vida/practica1_ciclo_vida$
```

Figura 5: Salida por consola C4 pig local

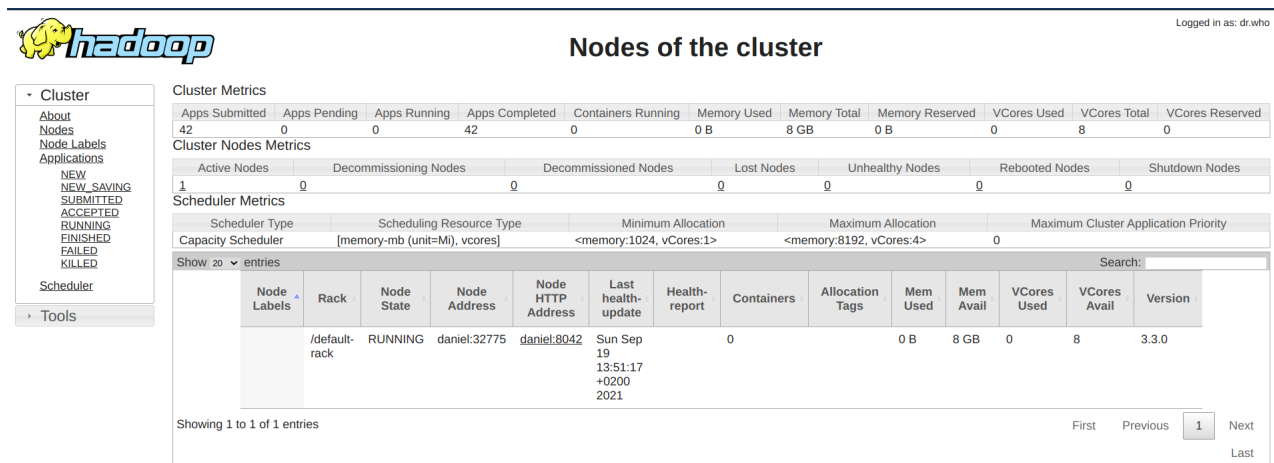


Figura 6: Resumen del servicio Hadoop funcionando correctamente

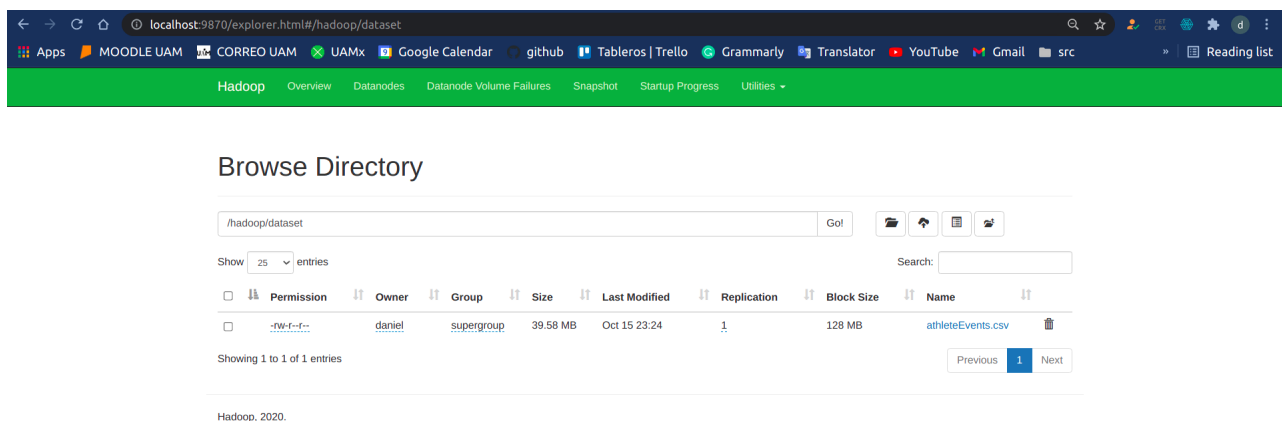


Figura 7: ruta a los datos en HDFS

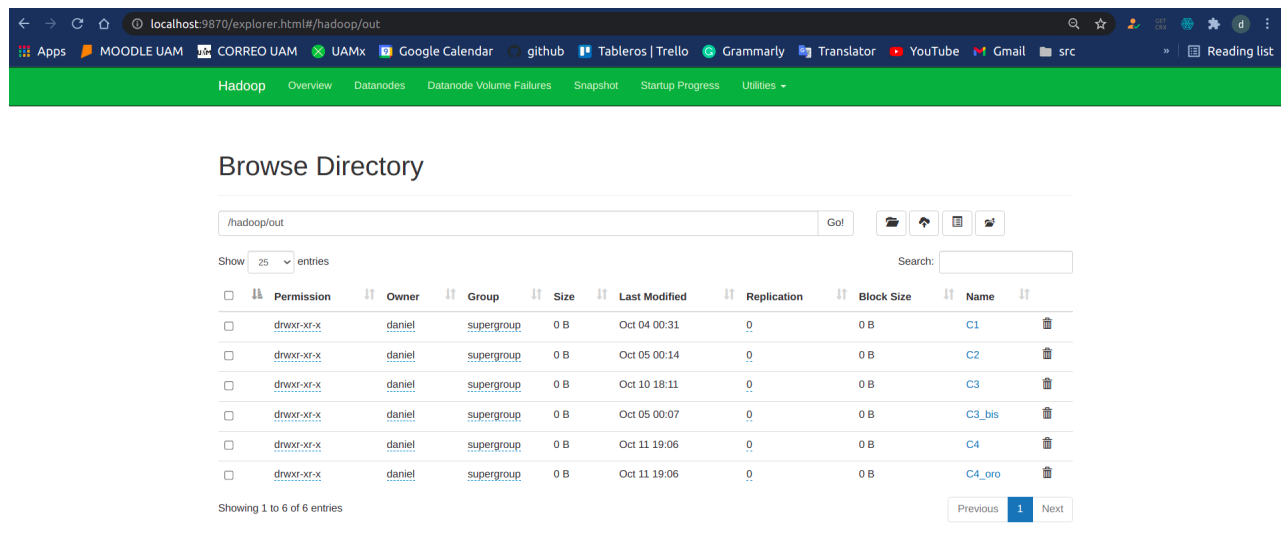


Figura 8: ruta a los resultados de PIG mapreduce en HDFS

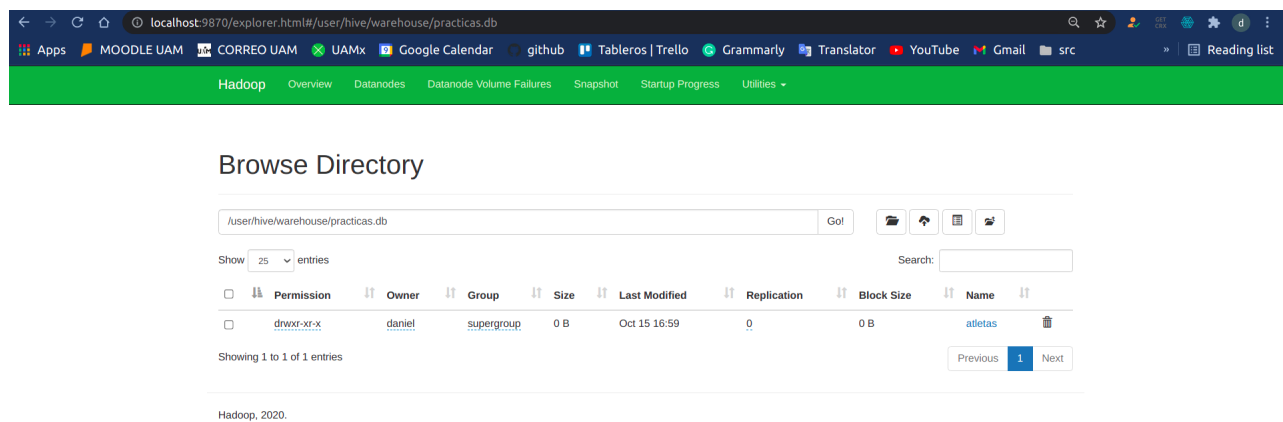


Figura 9: Contenido de la base de datos *practicass*

localhost:9870/explorer.html#/user/hive/warehouse/results_p1_pighive.db

Apps MOODLE UAM CORREO UAM UAMx Google Calendar github Tableros | Trello Grammarly Translator YouTube Gmail src Reading list

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

Browse Directory

/user/hive/warehouse/results_p1_pighive.db Go!

Show 25 entries Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	drwxr-xr-x	daniel	supergroup	0 B	Oct 15 17:16	0	0 B	c1	
<input type="checkbox"/>	drwxr-xr-x	daniel	supergroup	0 B	Oct 15 17:18	0	0 B	c2	
<input type="checkbox"/>	drwxr-xr-x	daniel	supergroup	0 B	Oct 15 17:23	0	0 B	c3_geograf	
<input type="checkbox"/>	drwxr-xr-x	daniel	supergroup	0 B	Oct 15 17:23	0	0 B	c3_sailing	
<input type="checkbox"/>	drwxr-xr-x	daniel	supergroup	0 B	Oct 15 17:26	0	0 B	c4_mejores	
<input type="checkbox"/>	drwxr-xr-x	daniel	supergroup	0 B	Oct 15 17:26	0	0 B	c4_mejores_oro	

Showing 1 to 6 of 6 entries

Previous 1 Next

Hadoop, 2020.

Figura 10: Contenido de la base de datos *results_p1_pighive*