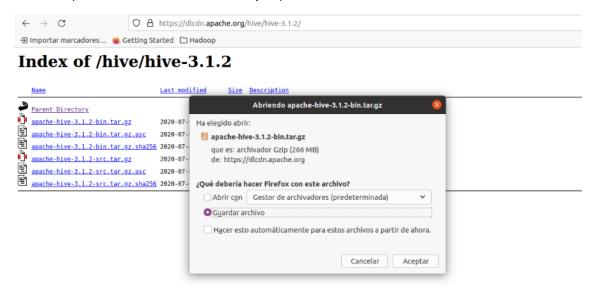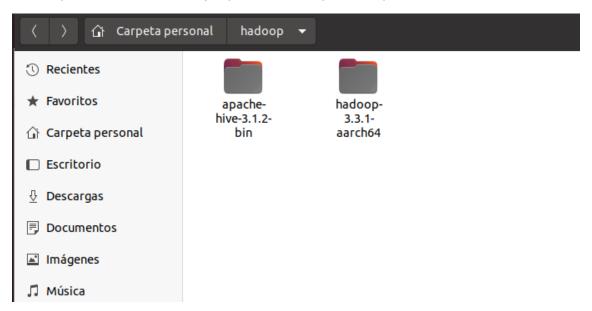# 6. HIVE

## 6.1. Instalación

Descargar a partir del sitio https://hive.apache.org/downloads.html. Verificar la versión de HIVE compatible con HADOOP, en este ejemplo, a versión de Hive 3.1.2



Descomprimir el archivo (en el ejemplo utilizo la carpeta hdoop)



Crear la variable de entorno HIVE_HOME y agregarla al PATH. Abrir el archivo de configuración "~/.profile" y agregar:

```
PATH="/home/hadoop/hadoop/apache-hive-3.1.2-bin/bin:$PATH"

export HIVE_HOME="/home/hadoop/hadoop/apache-hive-3.1.2-bin"
```

Quedando el fichero ~/.profile de la siguiente manera.

```
28
29 #agrego las variables para hadoop y modifico el path
30 PATH="/home/hadoop/hadoop/hadoop-3.3.1-aarch64/hadoop-3.3.1/bin:$PATH"
31 PATH="/home/hadoop/hadoop/hadoop-3.3.1-aarch64/hadoop-3.3.1/sbin:$PATH"
32 PATH="/home/hadoop/hadoop/apache-hive-3.1.2-bin/bin:$PATH"
33
34
35 export HADOOP_HOME="/home/hadoop/hadoop/hadoop-3.3.1-aarch64/hadoop-3.3.1"
36 export HADOOP_MAPRED_HOME=$HADOOP_HOME
37 export HADOOP_COMMON_HOME=$HADOOP_HOME
38 export HADOOP_HDFS_HOME=$HADOOP_HOME
39 export YARN_HOME=$HADOOP_HOME
40 export HIVE_HOME="/home/hadoop/hadoop/apache-hive-3.1.2-bin"
41
42
```

Reiniciar session para que se carguen las nuevas variables de entorno y se actualice el path

Iniciar el shell de Hive con

```
$ hive
```

```
hadoop@hadoop2:~$ hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/hadoop/hadoop/apache-hive-3.1.2-bin/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf
4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/hadoop/hadoop/hadoop-3.3.1-aarch64/hadoop-3.3.1/share/hadoop/common/lib/slf4
j-log4j12-1.7.30.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Hive Session ID = ac808d14-75b5-4f02-8c07-e78c5c33ca16

Logging initialized using configuration in jar:file:/home/hadoop/hadoop/apache-hive-3.1.2-bin/lib/hive-common-3.1.2.
jar!/hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execu
tion engine (i.e. spark, tez) or using Hive 1.X releases.
hive>
```

Verificar funcionamiento

```
hive> show tables hiv;
```

```
hive> show tables;
OK
Time taken: 0.644 seconds
hive>
```

En el caso de error, reiniciar la base de datos desde el fuera del Shell de hive y volver a ejecutar

```
$ rm -rf metastore_db/

$ $HIVE_HOME/bin/schematool -initSchema -dbType derby
```

## 6.2. Funcionamiento básico

Además del Shell de Hive CLI tenemos otro Shell BeeLine, que se basa en el anterior pero basado en conexión JDBCexit

### 6.2.1. Creación de una tabla de ejemplo

Crear una tabla basada en el fichero /examples/test.dat

```
CREATE TABLE src(valor INT);
```



Cargar los datos desde el fichero test.dats

```
LOAD DATA LOCAL INPATH '/home/hadoop/hadoop/apache-hive-3.1.2-
bin/examples/files/test.dat' OVERWRITE INTO TABLE src;
```



### 6.2.2. Ejecutar instrucciones

- -f permite ejecutar los comandos desde un archivo



- -e , permite ejecutar el comando sin iniciar el Shell

```
hive -e "select * from src"
```

Notas sintaxis HiveQL   :

- Las instrucciones no son casesensitive
- Las sentencias tienen que finalizar con ";"
- La tecla TAB completa la sintaxis

## 6.3. Utilización de tablas externas

Para este ejemplo tenemos varios archivos con los campos separados por # y con la siguiente estructura

movies:

Estructura: id#name#genre

Ejemplo :
1#Toy Story (1995)#Animation|Children's|Comedy
2#Jumanji (1995)#Adventure|Children's|Fantasy
3#Grumpier Old Men (1995)#Comedy|Romance

users:

Estructura:id#gender#age#occupationid#zipcode

Ejemplo:
1#F#1#10#48067
2#M#56#16#70072
3#M#25#15#55117

ratings:

estructura:userid#movieid#rating#tmstmp

Ejemplo:
1#1193#5#978300760
1#661#3#978302109
1#914#3#978301968

occupations:

Estructura:id#category#place

Ejemplo:
0#other/not specified
1#academic/educator
2#artist

### 6.3.1.  Creo la estructura y copio los ficheros al sistema HDFS

Creo la estructura en HDFS

```
$ hdfs dfs -mkdir /hive

$ hdfs dfs -mkdir /hive/data

$ hdfs dfs -mkdir /hive/data/occupation
```

```
$ hdfs dfs -mkdir /hive/data/movie

$ hdfs dfs -mkdir /hive/data/rating

$ hdfs dfs -mkdir /hive/data/user
```

Copio los archivos (ejemplo suponiendo que estoy en el mismo directorio)

```
$ hdfs dfs -put occupations.dat /hive/data/occupation

$ hdfs dfs -put movies.dat /hive/data/movie

$ hdfs dfs -put ratings.dat /hive/data/rating

$ hdfs dfs -put users.dat /hive/data/user
```

Teniendo un resultado similar a:

### 6.3.2.  Creación de la base de datos y tablas

A partir de este conjunto de datos, estructuramos la base de datos y tablas en Hive para poder hacer las consultas.

Creamos la base de datos

```
hive>   create database movielens;

hive>        use movielens;
```

```
hive>   create database movielens;
OK
Time taken: 1.761 seconds
hive>          use movielens;
OK
Time taken: 0.072 seconds
hive>
```

Creamos las tablas enlazadas con los datos externos al metastore de hive

```
        CREATE EXTERNAL TABLE ratings (
                userid INT,
                movieid INT,
                rating INT,
                tstamp STRING
                 ) ROW FORMAT DELIMITED
                 FIELDS TERMINATED BY '#'
                 STORED AS TEXTFILE
                 LOCATION '/hive/data/rating';
```

```
hive>          CREATE EXTERNAL TABLE ratings (
    >                    userid INT,
    >                    movieid INT,
    >                    rating INT,
    >                    tstamp STRING
    >                     ) ROW FORMAT DELIMITED
    >                    FIELDS TERMINATED BY '#'
    >                    STORED AS TEXTFILE
    >                    LOCATION '/hive/data/rating';
OK
Time taken: 1.2 seconds
hive>
```

```
CREATE EXTERNAL TABLE movies (
        movieid INT,
        title STRING,
        genres ARRAY<STRING>
) ROW FORMAT DELIMITED
FIELDS TERMINATED BY '#'
COLLECTION ITEMS TERMINATED BY "|"
STORED AS TEXTFILE
LOCATION '/hive/data/movie';
```

```
hive>          CREATE EXTERNAL TABLE movies (
    >                    movieid INT,
    >                    title STRING,
    >                    genres ARRAY<STRING>
    >                     ) ROW FORMAT DELIMITED
    >                    FIELDS TERMINATED BY '#'
    >                    COLLECTION ITEMS TERMINATED BY "|"
    >                    STORED AS TEXTFILE
    >                    LOCATION '/hive/data/movie';
OK
Time taken: 0.213 seconds
```

```
CREATE EXTERNAL TABLE users (
        userid INT,
        gender STRING,
        age INT,
        occupation_id INT,
        zipcode STRING
) ROW FORMAT DELIMITED
FIELDS TERMINATED BY '#'
STORED AS TEXTFILE
LOCATION '/hive/data/user';
```

```
hive>
    >              CREATE EXTERNAL TABLE users (
    >                      userid INT,
    >                      gender STRING,
    >                      age INT,
    >                      occupation_id INT,
    >                      zipcode STRING
    >                      ) ROW FORMAT DELIMITED
    >                      FIELDS TERMINATED BY '#'
    >                      STORED AS TEXTFILE
    >                      LOCATION '/hive/data/user';
OK
Time taken: 0.19 seconds
```

```
CREATE EXTERNAL TABLE occupations (
     id INT,
     occupation STRING
) ROW FORMAT DELIMITED
FIELDS TERMINATED BY '#'
STORED AS TEXTFILE
LOCATION '/hive/data/occupation';
```

```
hive>
    >              CREATE EXTERNAL TABLE occupations (
    >                      id INT,
    >                      occupation STRING
    >                      ) ROW FORMAT DELIMITED
    >                      FIELDS TERMINATED BY '#'
    >                      STORED AS TEXTFILE
    >                      LOCATION '/hive/data/occupation';
OK
Time taken: 0.168 seconds
```

### 6.3.3. Verificar datos de las tablas

Mediante Hive, realizar una consulta simple obteniendo los tres primeros registros de cada tabla

```
hive> select * from <tabla> limit 3;
```

## 6.4. Consultas con HiveSQ

A partir de las tablas anteriores, algunos ejemplos de consulta

### 6.4.1. Encontrar ocupación de todos los usuarios

```
hive>select u.*, o.occupation from users u, occupations o where
u.occupation_id= o.id limit 10;
```

### 6.4.2. Buscar el número de menores que han valorado películas

```
hive> select count(*) from users where age < 18;
```



### 6.4.3. Buscar el número de usuarios con la misma ocupación y con más de 25 años. Mostrar los detalles de la ocupación

```
hive> select o.occupation, count(1)  from users u join occupations o where
u.occupation_id= o.id AND u.age > 24 group by o.occupation;
```

### 6.4.4. Encontrar la edad el usuario con más valoraciones

```
hive> select u.userid, u.age, x.count from users u join ( select r.userid,
count(rating) count from ratings r group by (r.userid) order by count DESC
limit 1) x where u.userid = x.userid;
```