### Clase 7

Consigna: Por cada ejercicio, escribir el código y agregar una captura de pantalla del resultado obtenido.

## Diccionario de datos:

https://www.nyc.gov/assets/tlc/downloads/pdf/data\_dictionary\_trip\_records\_vellow.pdf

- 1. En Hive, crear la siguiente tabla (externa) en la base de datos tripdata:
  - a. airport\_trips(tpep\_pickup\_datetetime, airport\_fee, payment\_type, tolls\_amount, total\_amount)

Dentro de Hive, escribimos (la base de datos ha sido creada previamente)

2. En Hive, mostrar el esquema de airport trips

```
hive > describe formatted airport_trips;
OK
# col_name
                          data_type
                                                   comment
tpep_pickup_datetetime date
airport_fee
                          float
payment_type
                          int
tolls_amount
                          float
total_amount
                          float
# Detailed Table Information
Database:
                         tripdata
                         hadoop
Owner:
                         Sat Oct 18 15:22:56 ART 2025
CreateTime:
LastAccessTime:
                         UNKNOWN
Retention:
                         hdfs://172.17.0.2:9000/tables/external/tripdata
Location:
Table Type:
                         EXTERNAL_TABLE
Table Parameters:
        EXTERNAL
                                  TRUE
        transient_lastDdlTime
                                  1760811776
# Storage Information
SerDe Library:
                         org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
                         org.apache.hadoop.mapred.TextInputFormat org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
InputFormat:
OutputFormat:
                         No
-1
[]
Compressed:
Num Buckets:
Bucket Columns:
Sort Columns:
Storage Desc Params:
        field.delim
        serialization.format
Time taken: 0.087 seconds, Fetched: 31 row(s)
```

3. Crear un archivo .bash que permita descargar los archivos mencionados abajo e ingestarlos en HDFS:

https://data-engineer-edvai-public.s3.amazonaws.com/yellow\_tripdata\_2021-01.parquet https://data-engineer-edvai-public.s3.amazonaws.com/yellow\_tripdata\_2021-02.parquet

# Entramos al directorio /home/hadoop/scripts

```
hadoop@615cf53bef6c:~$ cd /home/hadoop/scripts
hadoop@615cf53bef6c:~/scripts$ ls
Notebook derby.log ingest.sh landing.sh pyspark_jupyter.sh spark-warehouse start-services.sh transformation.py
hadoop@615cf53bef6c:~/scripts$ |
```

# Creamos el script *landing\_2.sh*, que va a descargar los archivos e ingestarlos en HDFS

```
hadoop@615cf53bef6c:~/scripts$ cat > landing_2.sh

# Punto 3 (script: landing_2.sh)

# Descarga datos desde el repositorio al directorio /home/hadoop/landing
wget -P /home/hadoop/landing https://data-engineer-edvai-public.s3.amazonaws.com/yellow_tripdata_2021-01.parquet
wget -P /home/hadoop/landing https://data-engineer-edvai-public.s3.amazonaws.com/yellow_tripdata_2021-02.parquet

# Lleva el archivo a HDFS al directorio /ingest
hdfs dfs -put /home/hadoop/landing/yellow_tripdata_2021-01.parquet /ingest
hdfs dfs -put /home/hadoop/landing/yellow_tripdata_2021-02.parquet /ingest

# Borra el archivo starwars.csv del directorio /home/hadoop/landing/
rm /home/hadoop/landing/yellow_tripdata_2021-01.parquet
rm /home/hadoop/landing/yellow_tripdata_2021-02.parquet
hadoop@615cf53bef6c:~/scripts$ |
```

## Consultamos permisos

```
hadoop@615cf53bef6c:~/scripts$ ls -l
total 36
drwxrwxr-x 3 hadoop hadoop 4096 Sep 29 19:46 Notebook
-rw-rw-r-- 1 hadoop hadoop 670 Feb 28 2022 derby.log
-rwxrwxr-x 1 hadoop hadoop 272 Sep 22 01:14 ingest.sh
-r-xr-xr-x 1 hadoop hadoop 398 Sep 17 16:57 landing.sh
-rw-rw-r-- 1 hadoop hadoop 716 Oct 18 16:24 landing_2.sh
-rwxrwxrwx 1 hadoop hadoop 273 Sep 28 12:14 pyspark_jupyter.sh
drwxr-xr-x 2 hadoop hadoop 4096 Feb 9 2022 spark-warehouse
-rwxrwxrwx 1 hadoop hadoop 1089 May 9 2022 start-services.sh
-rwxrwxrwx 1 hadoop hadoop 1058 May 9 2022 transformation.py
hadoop@615cf53bef6c:~/scripts$
```

Asignamos permisos de ejecución a *landing 2.sh* y verificamos que se pueda ejecutar

```
hadoop@615cf53bef6c:~/scripts$ chmod 777 landing_2.sh
hadoop@615cf53bef6c:~/scripts$ ls -l
total 36
drwxrwxr-x 3 hadoop hadoop 4096 Sep 29 19:46 Notebook
-rw-rw-r-- 1 hadoop hadoop 670 Feb 28 2022 derby.log
-rwxrwxr-x 1 hadoop hadoop 272 Sep 22 01:14 ingest.sh
-r-xr-xr-x 1 hadoop hadoop 398 Sep 17 16:57 landing.sh
-rwxrwxrwx 1 hadoop hadoop 716 Oct 18 16:24 landing_2.sh
-rwxrwxrwx 1 hadoop hadoop 273 Sep 28 12:14 pyspark_jupyter.sh
drwxr-xr-x 2 hadoop hadoop 4096 Feb 9 2022 spark-warehouse
-rwxrwxrwx 1 hadoop hadoop 1089 May 9 2022 start-services.sh
-rwxrwxrwx 1 hadoop hadoop 1058 May 9 2022 transformation.py
hadoop@615cf53bef6c:~/scripts$
```

# Ejecutamos el archivo landing 2.sh

hadoop@615cf53bef6c:~\$

Verificamos que realmente se borraron los archivos de /home/hadoop/landing y que fueron correctamente ingresados en HDFS

```
hadoop@615cf53bef6c:~$ ls /home/hadoop/landing
starwars.csv 'yellow-taxi-ny?resource=download' yellow_tripdata_2021-01.csv

hadoop@615cf53bef6c:~$ hdfs dfs -ls /ingest
Found 4 items
-rw-r---- 1 hadoop supergroup 5462 2025-09-26 18:19 /ingest/starwars.csv
-rw-r---- 1 hadoop supergroup 125981363 2025-09-28 12:55 /ingest/yellow_tripdata_2021-01.csv
-rw-r---- 1 hadoop supergroup 21686067 2025-10-18 16:34 /ingest/yellow_tripdata_2021-01.parquet
-rw-r---- 1 hadoop supergroup 21777258 2025-10-18 16:34 /ingest/yellow_tripdata_2021-02.parquet
```

4. Crear un archivo .py que permita, mediante Spark, crear un data frame uniendo los viajes del mes 01 y mes 02 del año 2021 y luego Insertar en la tabla airport\_trips los viajes que tuvieron como inicio o destino aeropuertos, que hayan pagado con dinero.

####### # Archivo: transform\_load\_2.py # Tarea: Transformar Yellow Taxi (2021-01/02) y cargar en Hive (tripdata.airport\_trips) # Requisitos: # - La base de datos y tabla Hive existen: tripdata.airport\_trips # - Esquema de la tabla: # tpep\_pickup\_datetime DATE **FLOAT** # airport fee # payment\_type INT # tolls amount FLOAT total\_amount **FLOAT** ####### from pyspark.sql import SparkSession from pyspark.sql.functions import col, coalesce, lit def main(): # 1) Crear la sesión de Spark con soporte Hive spark = ( SparkSession.builder .appName("transform\_load") .enableHiveSupport()

```
#.config("hive.metastore.uris", "thrift://metastore:9083") # si usas metastore remoto
    #.config("spark.sql.warehouse.dir", "/user/hive/warehouse") # ajusta si aplica
    .getOrCreate()
)
# 2) Rutas de entrada (Parquet en HDFS)
p1 = "hdfs://172.17.0.2:9000/ingest/yellow_tripdata_2021-01.parquet"
p2 = "hdfs://172.17.0.2:9000/ingest/yellow_tripdata_2021-02.parquet"
df1 = spark.read.parquet(p1)
df2 = spark.read.parquet(p2)
# 3) Filtro de viajes de aeropuerto:
# - airport fee > 0 (con COALESCE para evitar null)
# - PULocationID o DOLocationID en {132 (JFK), 138 (LGA)}
# - RatecodeID en {2, 3} (tarifas típicas asociadas a aeropuertos)
def airport filter(df):
  return df.where(
    (coalesce(col("airport_fee"), lit(0.0)) > 0.0)
    (col("PULocationID").isin(132, 138)) |
    (col("DOLocationID").isin(132, 138)) |
   (col("RatecodeID").isin(2, 3))
 )
df1_f = airport_filter(df1)
df2_f = airport_filter(df2)
```

```
# 4) Unir meses asegurando el match por nombre de columnas
 df_union = df1_f.unionByName(df2_f)
 # 5) Filtrar sólo pagos en efectivo (payment_type = 2)
 df_cash = df_union.where(col("payment_type") == 2)
 # 6) Seleccionar y castear columnas EXACTAMENTE como la tabla Hive destino
 df_load = df_cash.selectExpr(
   "CAST(tpep_pickup_datetime AS DATE) AS tpep_pickup_datetime",
   "CAST(airport_fee
                        AS FLOAT) AS airport_fee",
   "CAST(payment_type
                           AS INT) AS payment_type",
   "CAST(tolls_amount
                          AS FLOAT) AS tolls_amount",
   "CAST(total amount
                         AS FLOAT) AS total_amount"
 )
 #7) Insertar en tabla externa (append). Requiere que exista tripdata.airport trips
 df_load.write.mode("append").insertInto("tripdata.airport_trips")
 # (Opcional) Log rápido
 inserted = df load.count()
 print(f"[OK] Filas insertadas en tripdata.airport_trips: {inserted}")
 spark.stop()
if __name__ == "__main__":
 main()
```

5. Realizar un proceso automático en Airflow que orqueste los archivos creados en los puntos 3 y 4. Correrlo y mostrar una captura de pantalla (del DAG y del resultado en la base de datos)

Creamos el DAG dentro del directorio donde se alojan en Airflow: /home/hadoop/airflow/dags/

```
hadoop@615cf53bef6c:~/airflow/dags$ cat > DAG_ejercicio_5.py
from datetime import timedelta
from airflow import DAG
from airflow.operators.bash import BashOperator
from airflow.operators.dummy import DummyOperator # en 2.8+ preferir EmptyOperator
from airflow.utils.dates import days_ago
args = {
 "owner": "airflow",
 "retries": 0,
with DAG(
 dag_id="ingest-transform-2",
 default_args=args,
 schedule_interval="0 0 * * *", # diario a medianoche
 start_date=days_ago(2),
 catchup=False,
 dagrun_timeout=timedelta(minutes=60),
 max_active_runs=1,
 tags=["ingest", "transform"],
) as dag:
 inicio_proceso = DummyOperator(task_id="inicio_proceso")
 finaliza_proceso = DummyOperator(task_id="finaliza_proceso")
 ingest = BashOperator(
```

task\_id="ingest",

bash\_command="""

```
set-e
   /bin/bash /home/hadoop/scripts/landing_2.sh
transform = BashOperator(
 task_id="transform",
 bash_command="""
   set-e
   /home/hadoop/spark/bin/spark-submit \
    --master local[*] \
    --deploy-mode client \
    --files /home/hadoop/hive/conf/hive-site.xml \
    /home/hadoop/scripts/transform_load_2.py
 # Si necesitas pasar conf/vars explícitas, descomenta:
 # env={
 # "HADOOP_CONF_DIR": "/home/hadoop/hadoop/etc/hadoop",
 # "JAVA_HOME": "/usr/lib/jvm/java-8-openjdk-amd64",
 # "SPARK_HOME": "/home/hadoop/spark",
 # "PATH": "/home/hadoop/spark/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin",
 # }
inicio_proceso >> ingest >> transform >> finaliza_proceso
```

Verificamos que el DAG se creo correctamente en el directorio especificado:

```
hadoop@615cf53bef6c:~/airflow/dags$ ls
DAG_ejercicio_5.py __pycache__ example-DAG.py ingest-transform.py
```

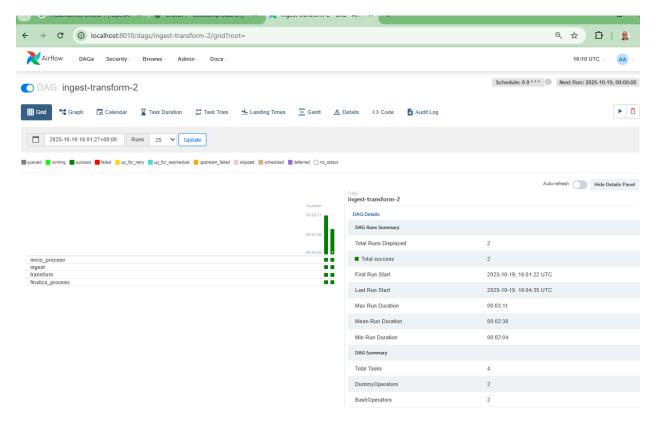
Comprobamos que tenga permisos

```
hadoop@615cf53bef6c:~/airflow/dags$ ls -l
total 16
-rw-rw-r-- 1 hadoop hadoop 1142 Oct 18 21:21 DAG_ejercicio_5.py
drwxrwxr-x 1 hadoop hadoop 4096 May 5 2022 __pycache__
-rw-rw-r-- 1 hadoop hadoop 1079 May 1 2022 example-DAG.py
-rw-rw-r-- 1 hadoop hadoop 1024 May 5 2022 ingest-transform.py
```

### Reiniciamos Airflow

## El DAG se llama ingest-transform-2:





Mostramos que la tabla se escribió correctamente en destino.

```
hive> select * from airport_trips limit 10;
OK
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
2020-12-31
                                              33.92
                  NULL
                           2
                                    6.12
2020-12-31
                  NULL
                           2
                                    6.12
                                              59.42
                                              40.3
2020-12-31
                  NULL
                                    0.0
2020-12-31
                  NULL
                           2
                                    0.0
                                              13.3
2020-12-31
                           2
                  NULL
                                    0.0
                                              20.3
2020-12-31
                  NULL
                           2
                                    0.0
                                              47.3
2020-12-31
                           2
                                              37.42
                  NULL
                                    6.12
2020-12-31
                           2
                  NULL
                                    0.0
                                              27.3
                           2
2020-12-31
                                              19.8
                  NULL
                                    0.0
2020-12-31
                  NULL
                           2
                                    0.0
                                              42.3
Time taken: 1.232 seconds, Fetched: 10 row(s)
hive>
```

### Desde DBeaver se puede ver esto también

