

# L3 Exercise 3 - Parallel ETL - Solution

July 10, 2021

## 1 Exercise 3: Parallel ETL

```
In [ ]: %load_ext sql
In [ ]: from time import time
import configparser
import matplotlib.pyplot as plt
import pandas as pd
```

## 2 STEP 1: Get the params of the created redshift cluster

- We need:
  - The redshift cluster endpoint
  - The IAM role ARN that give access to Redshift to read from S3

```
In [ ]: config = configparser.ConfigParser()
config.read_file(open('dwh.cfg'))
KEY=config.get('AWS','key')
SECRET= config.get('AWS','secret')

DWH_DB= config.get("DWH","DWH_DB")
DWH_DB_USER= config.get("DWH","DWH_DB_USER")
DWH_DB_PASSWORD= config.get("DWH","DWH_DB_PASSWORD")
DWH_PORT = config.get("DWH","DWH_PORT")

In [ ]: # FILL IN THE REDSHIFT ENPOINT HERE
# e.g. DWH_ENDPOINT="redshift-cluster-1.csmamz5zxmle.us-west-2.redshift.amazonaws.com"
DWH_ENDPOINT=""

#FILL IN THE IAM ROLE ARN you got in step 2.2 of the previous exercise
#e.g DWH_ROLE_ARN="arn:aws:iam::988332130976:role/dwhRole"
DWH_ROLE_ARN=""
```

## 3 STEP 2: Connect to the Redshift Cluster

```
In [ ]: conn_string="postgresql://{user}:{password}@{host}/{db}".format(DWH_DB_USER, DWH_DB_PASSWORD, DWH_ENDPOINT, DWH_DB)
print(conn_string)
%sql $conn_string
```

```
In [ ]: import boto3

s3 = boto3.resource('s3',
                    region_name="us-west-2",
                    aws_access_key_id=KEY,
                    aws_secret_access_key=SECRET
                    )

sampleDbBucket = s3.Bucket("udacity-labs")

for obj in sampleDbBucket.objects.filter(Prefix="tickets"):
    print(obj)
```

## 4 STEP 3: Create Tables

```
In [ ]: %%sql
DROP TABLE IF EXISTS "sporting_event_ticket";
CREATE TABLE "sporting_event_ticket" (
    "id" double precision DEFAULT nextval('sporting_event_ticket_seq') NOT NULL,
    "sporting_event_id" double precision NOT NULL,
    "sport_location_id" double precision NOT NULL,
    "seat_level" numeric(1,0) NOT NULL,
    "seat_section" character varying(15) NOT NULL,
    "seat_row" character varying(10) NOT NULL,
    "seat" character varying(10) NOT NULL,
    "ticketholder_id" double precision,
    "ticket_price" numeric(8,2) NOT NULL
);
```

## 5 STEP 4: Load Partitioned data into the cluster

```
In [ ]: %%time
qry = """
    copy sporting_event_ticket from 's3://udacity-labs/tickets/split/part'
    credentials 'aws_iam_role={}'
    gzip delimiter ';' compupdate off region 'us-west-2';
    """.format(DWH_ROLE_ARN)

%sql $qry
```

## 6 STEP 4: Create Tables for the non-partitioned data

```
In [ ]: %%sql
DROP TABLE IF EXISTS "sporting_event_ticket_full";
CREATE TABLE "sporting_event_ticket_full" (
    "id" double precision DEFAULT nextval('sporting_event_ticket_seq') NOT NULL,
```

```

        "sporting_event_id" double precision NOT NULL,
        "sport_location_id" double precision NOT NULL,
        "seat_level" numeric(1,0) NOT NULL,
        "seat_section" character varying(15) NOT NULL,
        "seat_row" character varying(10) NOT NULL,
        "seat" character varying(10) NOT NULL,
        "ticketholder_id" double precision,
        "ticket_price" numeric(8,2) NOT NULL
    );

```

## 7 STEP 5: Load non-partitioned data into the cluster

- Note how it's slower than loading partitioned data

```
In [ ]: %%time
```

```

qry = """
    copy sporting_event_ticket_full from 's3://udacity-labs/tickets/full/full.csv.gz'
    credentials 'aws_iam_role={}'
    gzip delimiter ';' compupdate off region 'us-west-2';
    """.format(DWH_ROLE_ARN)

%sql $qry

```

```
In [ ]:
```