

## **ETL Challenge Report**

**(Francis Odo)**

### **Background**

Most business entities, government institutions and agencies maintain some type of database in their day-to-day operations. Depending on business type, industry category, as well as other need factors, the nature of data collection varies widely. There is the necessity to proactively keep up with the accumulation and maintenance brought about by business operations and other associated factors. One of the ways of achieving this task is the Extract Transformation and Load (ETL) process. This project takes us through the process from beginning to the end.

### **Objective**

The objective is to automate the process of Extract Transformation and Load (ETL) in fulfilling the need of an ideal business entity wanting to maintain an information database. The process takes in a set of files in raw data format, then:

- a) Extract - Take in the information in the best available format(file) and condition
- b) Transformation - Perform various types of cleaning, rearrangement and updates
- c) Load - Load the data into a functional database platform in a clean and easy to understand format.

### **Development Environment**

Code Development – Jupyter Notebook

Programming Language – Python

Libraries – Pandas

Database – PostgresDB

ORM - Sql Alchemy

Code File - Module\_8\_Challenge.ipynb

Config File – config.ny

### **Code Plan**

- (1) Import all the necessary dependencies and set up file paths to raw datafile location
- (2) Define engine and database connection string
- (3) Create a main function in Python using “def etl\_pipe\_func(file 1, file 2, file 3)”.  
The three raw files are:
  - (a) wikipedia.movies.json
  - (b) movies\_metadata.csv
  - (c) ratings.csv

- (4) Open and read the files with the appropriate and applicable Pandas method. Verify that data is being read correctly within the function created for automation purposes.
- (5) Perform transformation which includes cleaning, avoiding duplicates, merging columns, regular expression and other necessary steps.
- (6) Create and configure the “movies\_data” database in PostgreSQL.
- (7) Remove existing data from the table
- (8) Load the clean and transformed data into PostgreSQL database. Apply **try-except** to catch errors in the process.

### **Approach and Assumptions**

- a) Conditions of raw data can vary in so many ways and may require different steps/type of cleaning. The approach is to demonstrate the fundamental principle and methods of cleaning in limited fashion with emphasis on automating the process.
- b) One key assumption is that Opening and Reading files, Transforming/Cleaning and Loading to the Database steps are primary essential processes.
- c) Data cleaning varies depending on how dirty the data may be. There is no limit to cleaning as much as it adds improvement.

### **Risks**

1. Ratings file is large. Loading the file into the database is best done in chunks. Data integrity could be a challenge
2. The process is memory intensive, and tends to be slow sometimes. Loading process may halt and the process will have to be restarted.
3. Ensure that PostgresDB is up and running.

### **Conclusion**

There are other ways of creating data pipelines that involve other methods and applications. ETL is fairly straight-forward and can handle average sized data successfully and conveniently. However, these methods are outside the scope of this project.

This application can be tailored or modified for any data environment to fulfill the basic needs. The main advantage here is to eliminate repetitive processes in maintaining data.

Verification of tables being created and populated in PostgresDB.

Movies Table:

Berkeley/class\_folder/Movie X ETL\_Merging\_Split2 - Jupyter X ETL\_Parsing\_Split1 - Jupyter X Module\_8\_Challenge - Jupyter X Home Page - Select or create X pgAdmin 4

127.0.0.1:57431/browser/

pgAdmin File Object Tools Help

Browser

- FTS Templates
- Foreign Tables
- Functions
- Materialized Views
- Procedures
- Sequences
- Tables (2)
  - movies
    - Columns (41)
      - index
      - wikipedia\_url
      - year
      - imdb\_link
      - based\_on
      - starring
      - cinematography
      - Release date
      - Running time
      - country
      - Budget
      - director
      - distributor
      - editors
      - composers
      - producers
      - writers
      - imdb\_id
      - box\_office
      - belongs\_to\_collection
      - budget
      - genres

Dashboard Properties SQL Statistics Dependencies Dependents movie\_data/postgres@PostgreSQL 11

Query Editor Query History Scratch Pad

```
1 select * from movies;
```

Data Output Explain Messages Notifications

	index bigint	wikipedia_url text	year bigint	imdb_link text	based_on text	starring text	cinematography text	Release date text	Running time text	country text
1	0	https://en.wikipedi...	1990	https://www.L...	(Characters'by...	(Andrew Dic...	Oliver Wood	(July 11, 1990',1...	102 minutes	Unite
2	1	https://en.wikipedi...	1990	https://www.L...	(the novel'Alt...	(Jason Patri...	Mark Plummer	(May 17, 1990',1...	114 minutes	Unite
3	2	https://en.wikipedi...	1990	https://www.L...	(Air America'b...	(Mel Gibson...	Roger Deakins	(August 10, 1990',...	113 minutes	Unite
4	3	https://en.wikipedi...	1990	https://www.L...	[null]	(Alec Baldwi...	Carlo Di Palma	(December 25, 19...	106 minutes	Unite
5	4	https://en.wikipedi...	1990	https://www.L...	[null]	(Paul Hogan...	Russell Boyd	December 19, 1990	95 minutes	US
6	5	https://en.wikipedi...	1990	https://www.L...	[null]	(Eric Robert...	Jacques Haikin	(March 22, 1990',...	95 minutes	Unite
7	6	https://en.wikipedi...	1990	https://www.L...	[null]	[null]	(Tom Hurwitz',Mathie...	(October 6, 1990',...	100 minutes	(Unit

Type here to search

1:15 AM 8/18/2020

## Ratings Table:

Berkeley/class\_folder/Movie X ETL\_Merging\_Split2 - Jupyter X ETL\_Parsing\_Split1 - Jupyter X Module\_8\_Challenge - Jupyter X Home Page - Select or create X pgAdmin 4

127.0.0.1:57431/browser/

pgAdmin File Object Tools Help

Browser

- revenue
- runtime
- spoken\_languages
- status
- tagline
- video
- vote\_average
- vote\_count
- Constraints
- Indexes
- Rules
- Triggers
- ratings
  - Columns (5)
    - index
    - userid
    - movieid
    - rating
    - timestamp
  - Constraints
  - Indexes
  - Rules
  - Triggers
  - Trigger Functions
  - Types
  - Views
- postgres
- unit\_assessment\_db
- Login/Group Roles
- Tablespaces

Dashboard Properties SQL Statistics Dependencies Dependents movie\_data/postgres@PostgreSQL 11

Query Editor Query History Scratch Pad

```
1 select * from ratings;
```

Data Output Explain Messages Notifications

	index bigint	userid bigint	movieid bigint	rating double precision	timestamp bigint
1	0	1	110	1	1425941529
2	1	1	147	4.5	1425942435
3	2	1	858	5	1425941523
4	3	1	1221	5	1425941546
5	4	1	1246	5	1425941556
6	5	1	1968	4	1425942148
7	6	1	2762	4.5	1425941300
8	7	1	2918	5	1425941593

Type here to search

1:29 AM 8/18/2020