# Data Modeling with Postgres

| REVIEW | CODE REVIEW 4 | HISTORY |
|---|---|---|

## Meets Specifications

Congratulations on passing the project. I hope you enjoyed and learned a lot with it. Keep the good work and happy learning.

## Table Creation

✓    The script, `create_tables.py`, runs in the terminal without errors. The script successfully connects to the Sparkify database, drops any tables if they exist, and creates the tables.

✓    CREATE statements in `sql_queries.py` specify all columns for each of the five tables with the right data types and conditions.

## ETL

✓ The script, `etl.py`, runs in the terminal without errors. The script connects to the Sparkify database, extracts and processes the `log_data` and `song_data`, and loads data into the five tables.

Since this is a subset of the much larger dataset, the solution dataset will only have 1 row with values for value containing ID for both `songid` and `artistid` in the fact table. Those are the only 2 values that the query in the `sql_queries.py` will return that are not-NONE. The rest of the rows will have NONE values for those two variables.

Excellent. Your ETL pipeline works perfectly!

✓ INSERT statements are correctly written for each table, and handle existing records where appropriate. `songs` and `artists` tables are used to retrieve the correct information for the `songplays` INSERT.

Very good using DO UPDATE to update the level column. This is important to know if the user changed their plan.

## Code Quality

✓ The README file includes a summary of the project, how to run the Python scripts, and an explanation of the files in the repository. Comments are used effectively and each function has a docstring.

Very nice README, congratulations!

✓ Scripts have an intuitive, easy-to-follow structure with code separated into logical functions. Naming for variables and functions follows the PEP8 style guidelines.

Perfect use of docstring and organization of the code!

⬇ DOWNLOAD PROJECT

4  CODE REVIEW COMMENTS ❯

Rate this review
★ ★ ★ ★ ★

DISCUSS ON STUDENT HUB ›

# Data Modeling with Postgres

| REVIEW | CODE REVIEW **4** | HISTORY |
|---|---|---|

▼ sql_queries.py  **3**

```
1  # DROP TABLES
```

**AWESOME**

Perfect use of IF EXISTS.

```
2
3  songplay_table_drop = "DROP TABLE IF EXISTS songplays;"
4  user_table_drop = "DROP TABLE IF EXISTS users;"
5  song_table_drop = "DROP TABLE IF EXISTS songs;"
6  artist_table_drop = "DROP TABLE IF EXISTS artists;"
7  time_table_drop = "DROP TABLE IF EXISTS time;"
8
9  # CREATE TABLES
10
11 user_table_create = ("""
12     CREATE TABLE IF NOT EXISTS users (
13         user_id INT PRIMARY KEY,
14         first_name VARCHAR,
15         last_name VARCHAR,
16         gender VARCHAR,
17         level VARCHAR
18 );
19 """)
20
21 song_table_create = ("""
22     CREATE TABLE IF NOT EXISTS songs (
23         song_id VARCHAR NOT NULL PRIMARY KEY,
24         title VARCHAR(255),
25         artist_id VARCHAR,
26         year INT CHECK (year >= 0),
```

Rate this review
☆☆☆☆☆

```
26        year INT CHECK (year >= 0),
27        duration FLOAT
28    );
29    """)
30
31 artist_table_create = ("""
32    CREATE TABLE IF NOT EXISTS artists (
33        artist_id VARCHAR CONSTRAINT artist_pk PRIMARY KEY,
34        name VARCHAR(255) NOT NULL,
35        location VARCHAR,
36        latitude FLOAT,
37        longitude FLOAT
38    );
39    """)
40
41 time_table_create = ("""
42    CREATE TABLE IF NOT EXISTS time (
43        start_time BIGINT NOT NULL PRIMARY KEY,
44        hour INT NOT NULL CHECK (hour >= 0),
45        day INT NOT NULL CHECK (day >= 0),
46        week INT NOT NULL CHECK (week >= 0),
47        month INT NOT NULL CHECK (month >= 0),
48        year INT NOT NULL CHECK (year >= 0),
49        weekday VARCHAR NOT NULL
50    );
51    """)
52
53 songplay_table_create = ("""
54    CREATE TABLE IF NOT EXISTS songplays (
55        songplay_id BIGSERIAL PRIMARY KEY,
```

**AWESOME**

Nice job defining the PRIMARY KEY.

```
56        start_time BIGINT NOT NULL,
57        user_id INT NOT NULL REFERENCES users(user_id),
58        level VARCHAR NOT NULL,
59        song_id VARCHAR REFERENCES songs(song_id),
60        artist_id VARCHAR REFERENCES artists(artist_id),
61        session_id INT NOT NULL,
62        location VARCHAR,
63        user_agent VARCHAR
64    );
65    """)
66
67
68 # INSERT RECORDS
69
70 # On conflict update user level
71 user_table_insert = ("""
72    INSERT INTO users(user_id, first_name, last_name, gender, level)
73    VALUES (%s, %s, %s, %s, %s)
74    ON CONFLICT (user_id) DO UPDATE SET level = EXCLUDED.level;
```

**AWESOME**

Excellent! 👏

```python
75  """)
76
77  song_table_insert = ("""
78      INSERT INTO songs(song_id, title, artist_id, year, duration)
79      VALUES (%s, %s, %s, %s, %s)
80      ON CONFLICT (song_id) DO NOTHING;
81  """)
82
83  # Updates might be needed to location, latitude and longitude
84  artist_table_insert = ("""
85      INSERT INTO artists(artist_id, name, location, latitude, longitude)
86      VALUES (%s, %s, %s, %s, %s)
87      ON CONFLICT (artist_id) DO UPDATE SET
88      location = EXCLUDED.location,
89      latitude = EXCLUDED.latitude,
90      longitude = EXCLUDED.longitude;
91  """)
92
93  time_table_insert = ("""
94      INSERT INTO time(start_time, hour, day, week, month, year, weekday) VALUES (%s, %s, %s, %s, %s
95      ON CONFLICT (start_time) DO NOTHING;
96  """)
97
98  songplay_table_insert = ("""
99      INSERT INTO songplays(songplay_id, start_time, user_id, level, song_id, artist_id, session_id
100     VALUES (DEFAULT, %s, %s, %s, %s, %s, %s, %s, %s);
101 """)
102
103 # FIND SONGS
104
105 song_select = ("""
106     SELECT s.song_id, s.artist_id
107     FROM songs s
108     JOIN artists a ON s.artist_id = a.artist_id
109     WHERE s.title = %s AND a.name = %s AND s.duration = %s
110 """)
111
112 # QUERY LISTS
113
114 create_table_queries = [user_table_create, artist_table_create, song_table_create, time_table_crea
115 drop_table_queries = [songplay_table_drop, user_table_drop, song_table_drop, artist_table_drop, ti
```

```python
1  import os
2  import glob
3  import psycopg2
4  import pandas as pd
5  from sql_queries import *
6
7
8  def process_song_file(cur, filepath):
9
10     """ Process_song_file takes the cursor and the list of song files as arguments and processes t
11
12     It reads the song json files and stores them in a dataframe. Then it extracts the Songs and Ar
13     """
```

**AWESOME**

Nice job providing docstring

```python
14
15     # open song file
16
17     df = pd.read_json(filepath, lines=True)
18
19     # insert song record
20
21     song_data = list(df[['song_id', 'title', 'artist_id', 'year', 'duration']].values[0])
22     cur.execute(song_table_insert, song_data)
23
24     # insert artist record
25
26     artist_data = list(df[['artist_id', 'artist_name', 'artist_location', 'artist_latitude', 'arti
27     cur.execute(artist_table_insert, artist_data)
28
29
30  def process_log_file(cur, filepath):
31
32     """  Process_log_file takes the cursor and the list of log files as arguments and processes th
33
34     It captures all json formatted log files and stores them in a dataframe. It formats the time i
35
36     It captures all user related columns from the main dataframe and stores them in a temporary da
37
38     It triggers an SQL statement to capture the song_id, Artist_id and length from the Songs and t
39     """
40
41     # open log file
42
43     df = pd.read_json(filepath, lines=True)
44
45     # filter by NextSong action
46
47     df = df.loc[df['page'] == 'NextSong']
48
49     # convert timestamp column to datetime
50
51     t = pd.to_datetime(df['ts'], unit='ms')
52
```

Rate this review
★★★★★

```python
        # insert time data records
        # isocalendar is now the new format in pandas

        time_data = [df.ts.values, t.dt.hour.values, t.dt.day.values, t.dt.isocalendar().week.values,

        # old version uses weekofyear. Deprecated in pandas since version 1.1.0
        #     time_data = [df.ts.values, t.dt.hour.values, t.dt.day.values, t.dt.weekofyear.values, t.d

        column_labels = ['start_time', 'hour', 'day', 'week', 'month', 'year', 'weekday']
        time_df = pd.DataFrame(dict(zip(column_labels, time_data)))

        for i, row in time_df.iterrows():
            cur.execute(time_table_insert, list(row))

        # load user table

        user_df = df[['userId', 'firstName', 'lastName', 'gender', 'level']]

        # insert user records

        for i, row in user_df.iterrows():
            cur.execute(user_table_insert, row)

        # insert songplay records

        for index, row in df.iterrows():

            # get songid and artistid from song and artist tables

            cur.execute(song_select, (row.song, row.artist, row.length))
            result = cur.fetchone()
            if result:
                songid, artistid = result
            else:
                songid, artistid = None, None

            # insert songplay record

            songplay_data = [row.ts, row.userId, row.level, songid, artistid, row.sessionId, row.locat
            cur.execute(songplay_table_insert, songplay_data)


def process_data(cur, conn, filepath, func):

    """ Process_data function uses the cursor and connection to collect all the files and call fun

    It collects all the files in the available directory and stores them in a list which it passes
    """

    # get all files matching extension from directory

    all_files = []
    for root, dirs, files in os.walk(filepath):
        files = glob.glob(os.path.join(root, '*.json'))
        for f in files:
            all_files.append(os.path.abspath(f))

    # get total number of files found

    num_files = len(all_files)
    print('{} files found in {}'.format(num_files, filepath))

    # iterate over files and process

    for i, datafile in enumerate(all_files, 1):
        func(cur, datafile)
        conn.commit()
        print('{}/{} files processed.'.format(i, num_files))


def main():
    """ Main function which calls the process_data function for processing the Songs and Logs file

    It passes the cursor, connection, filepaths and the function names as arguments.
    """

    conn = psycopg2.connect('host=127.0.0.1 dbname=sparkifydb user=student password=student')
    cur = conn.cursor()

    process_data(cur, conn, filepath='data/song_data', func=process_song_file)
    process_data(cur, conn, filepath='data/log_data', func=process_log_file)

    conn.close()


if __name__ == '__main__':
    main()
```