**airflow\dags\kaggle_dag.py**

```
1    # 685.652, Spring 2025 - Group 6 Final Project
2    # kaggle_dag.py
3
4    # This DAG gets billboard hot 100 chart data from kaggle
5    # Cleans and transforms it
6    # Then loads it into postgres
7
8    from airflow import DAG
9    from airflow.operators.empty import EmptyOperator
10   from airflow.providers.postgres.operators.postgres import PostgresOperator
11   from airflow.operators.python_operator import PythonOperator
12   from airflow.models import Variable
13   from airflow.providers.postgres.hooks.postgres import PostgresHook
14   from datetime import datetime, timedelta
15   import os
16   import zipfile
17   import subprocess
18   import pandas as pd
19   import uuid
20
21   default_args = {
22       'owner': 'group6',
23       'depends_on_past': False,
24       'start_date': datetime(2024, 1, 1),
25       'email_on_failure': False,
26       'email_on_retry': False,
27       'retries': 1,
28       'retry_delay': timedelta(minutes=5),
29   }
30
31
32   # Retrieves a specific file from kaggle via API
33   def get_kag_file():
34
35       # User needs to set up account on kaggle and generate API key
36       # Set the KAG_USERNAME and KAG_KEY in variables.json
37       os.environ['KAGGLE_USERNAME'] = Variable.get("KAG_USERNAME")
38       os.environ['KAGGLE_KEY'] = Variable.get("KAG_KEY")
39
40       # Retrieving a specific dataset
41       dataset_name = 'elizabethearhart/billboard-hot-1001958-2024'
42       zip_file_name = f"{dataset_name.split('/')[-1]}.zip"
43
44       # Data is to go into data subdirectory
45       data_dir = os.path.join(os.getcwd(), 'data')
46       os.makedirs(data_dir, exist_ok=True)
47
48       # Download the dataset by running the kaggle command
```

```python
49          # The kaggle package gets installed on build
50          subprocess.run([
51              'kaggle', 'datasets', 'download',
52              '-d', dataset_name,
53              '--force', '-p', data_dir])
54
55          # Unzip to data subdirectory
56          zip_path = os.path.join(data_dir, zip_file_name)
57          try:
58
59              with zipfile.ZipFile(zip_path, 'r') as zip_ref:
60                  zip_ref.extractall(data_dir)
61                  extracted_files = zip_ref.namelist()
62
63              if extracted_files:
64                  extracted_file_name = extracted_files[0]
65                  print(f"Extracted Kaggle file: {extracted_file_name} to data directory")
66
67          except PermissionError:
68              print("Error extracting ZIP file: Permission denied")
69              print("Close Kaggle data csv if you already have it open in another program")
70              return None
71
72          # Extract first file
73          extracted_file_name = extracted_files[0] if extracted_files else None
74          print(f"Extracted Kaggle file: {extracted_file_name} to data directory")
75
76          # Keep extracted file, delete zip
77          os.remove(zip_path)
78          print(f"The ZIP file {zip_file_name} has been deleted\n")
79
80          return extracted_file_name
81
82
83  # Processes kaggle file in data directory and returns a dataframe
84  def process_kag_file(file_name):
85      print(f"\nProcessing Kaggle file: {file_name} ...\n"
86            "File contains Bilboard Hot 100 chart data from 1958 to 2024\n"
87            "It contains - chart_week, title, performer, current_week, last_week, peak_pos,
    wks_on_chart\n"
88            "Duplicate song title, performer pairs will be removed\n")
89
90      data_dir = os.path.join(os.getcwd(), 'data')
91      os.makedirs(data_dir, exist_ok=True)
92      kaggle_file_path = os.path.join(data_dir, file_name)
93
94      # Read kaggle file from data directory
95      df = pd.read_csv(kaggle_file_path, encoding='utf-8')
96
97      # Strip and lowercase string cols
```

```python
 98        df['title'] = df['title'].str.strip().str.lower()
 99        df['performer'] = df['performer'].str.strip().str.lower()
100
101        # Force chart_wk to date for sorting
102        df['chart_week'] = pd.to_datetime(df['chart_week'], errors='coerce')
103
104        # Replace empty strings with None
105        df.replace('', None, inplace=True)
106
107        df_kaggle = df.copy()
108
109        # Replace some strings to help with matching to Spotify data
110        df_kaggle['performer'] = df_kaggle['performer'].str.replace(r' \(featuring ([^)]*)\)',
    r',\1', regex=True)
111        df_kaggle['performer'] = df_kaggle['performer'].str.replace(r' \(feat. ([^)]*)\)',
    r',\1', regex=True)
112        df_kaggle['performer'] = df_kaggle['performer'].str.replace('featuring', ',')
113        df_kaggle['performer'] = df_kaggle['performer'].str.replace('feat.', ',')
114        df_kaggle['performer'] = df_kaggle['performer'].str.replace('&', 'and')
115        df_kaggle['performer'] = df_kaggle['performer'].str.replace(r' ,', r',', regex=True)
116        df_kaggle['title'] = df_kaggle['title'].str.replace(r' \(featuring.*$', '', regex=True)
117
118        df_kaggle['top_artist'] = df_kaggle['performer'].str.split(',').str[0]
119        df_kaggle['performer'] = df_kaggle['performer'].str.split(',')
120
121        # Force int columns to numeric, force any NaNs to None
122        int_columns = ['current_week', 'last_week', 'peak_pos', 'wks_on_chart']
123        for col in int_columns:
124            df_kaggle[col] = pd.to_numeric(df_kaggle[col], errors='coerce')
125        df_kaggle[int_columns] = df_kaggle[int_columns].where(pd.notnull(df_kaggle[int_columns]),
    None)
126
127        # Force int columns to Int64 (needed for one column that was being read as float)
128        # Replace 0 with None
129        for col in int_columns:
130            df_kaggle[col] = df_kaggle[col].astype('Int64')
131        for col in int_columns:
132            df_kaggle[col] = df_kaggle[col].replace(0, None)
133
134
135        print(f"Number of Kaggle dataset rows before removing duplicates: {len(df_kaggle)}")
136
137        # Group by top_artist and title, then aggregate
138        date_groups = df_kaggle.groupby(['top_artist', 'title']).agg(
139            performer=('performer', 'first'),          # Keep the first performer string
140            chart_week_min=('chart_week', 'min'),      # Rename agg results directly
141            chart_week_max=('chart_week', 'max'),
142            current_week=('current_week', 'last'),
143            peak_pos=('peak_pos', 'min'),
144            wks_on_chart=('wks_on_chart', 'max')
145        ).reset_index() # Makes 'top_artist' and 'title' columns again
```

```
146
147        # Rename the columns after aggregation
148        date_groups.rename(columns={
149            'performer': 'artists',          # Now rename the carried-over performer
150            'title': 'song_name',            # Rename title
151            'chart_week_min': 'wk_first_charted',
152            'chart_week_max': 'wk_last_charted',
153            'current_week': 'last_chart_pos',
154            'peak_pos': 'peak_chart_pos',
155            'wks_on_chart': 'total_wks_on_chart'
156        }, inplace=True)
157
158        # Assign the correctly aggregated and renamed DataFrame back
159        df_kaggle = date_groups
160
161        print(f"Number of Kaggle dataset rows after removing duplicates: {len(df_kaggle)}\n")
162
163        # This column reordering should now work
164        df_kaggle = df_kaggle[['top_artist', 'artists', 'song_name', 'peak_chart_pos',
165                               'last_chart_pos', 'total_wks_on_chart',
166                               'wk_first_charted', 'wk_last_charted']]
167
168
169        # Write to CSV for record
170        timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
171        unique_hot100 = f"unique_hot100_{timestamp}.csv"
172        kag_file_name = os.path.join(data_dir, unique_hot100)
173
174        df_kaggle.to_csv(kag_file_name, index=False,
175                         encoding='utf-8-sig')
176        print(f"Kaggle dataset processed and saved to {unique_hot100}\n")
177
178        count = (df_kaggle['peak_chart_pos'] == 1).sum()
179        print(f"Processed Kaggle dataset contains {len(df_kaggle)} songs in total and "
180              f"{count} songs that peaked at #1")
181        print(f"Dataframe has columns: {', '.join(df_kaggle.columns)}\n")
182
183        return df_kaggle
184
185
186
187 # Loads kaggle data into postgres
188 def load_billboard_data():
189        kag_name = get_kag_file()
190        df_kaggle = process_kag_file(kag_name)
191
192        # get postgres connection
193        pg_hook = PostgresHook(postgres_conn_id='pg_group6')
194        with pg_hook.get_conn() as conn:
195            with conn.cursor() as cur:
```

```python
196                for index, row in df_kaggle.iterrows():
197                    group6_id = uuid.uuid5(uuid.NAMESPACE_DNS, str(row['song_name'].strip() +
       row['top_artist'].strip()))
198
199                    cur.execute("""INSERT INTO billboard_chart_data (group6_id, top_artist,
       artists, song_name, peak_chart_pos,
200                                    last_chart_pos, total_wks_on_chart, wk_first_charted,
       wk_last_charted)
201                                    VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s)""",
202                                    (group6_id,
203                                    row['top_artist'],
204                                    row['artists'],
205                                    row['song_name'],
206                                    row['peak_chart_pos'],
207                                    row['last_chart_pos'],
208                                    row['total_wks_on_chart'],
209                                    row['wk_first_charted'],
210                                    row['wk_last_charted']))
211
212
213   with DAG(
214           'kaggle_dag',
215           default_args=default_args,
216           description='Dag for kaggle billboard data',
217           schedule_interval=None,
218           catchup=False
219       ) as kaggle_dag:
220
221       start_task = EmptyOperator(task_id='start')
222
223       create_tables = PostgresOperator(
224           task_id='create_tables',
225           postgres_conn_id='pg_group6',
226           sql='sql/billboard_create.sql'
227       )
228
229       load_billboard_data = PythonOperator(
230           task_id='load_billboard_data',
231           python_callable=load_billboard_data
232       )
233
234       end_task = EmptyOperator(task_id='end')
235
236       start_task >> create_tables >> load_billboard_data >> end_task
237
```