

## airflow\dags\spotify\_dag.py

```
1 # 685.652, Spring 2025 - Group 6 Final Project
2 # spotify_dag.py
3
4 # This DAG gets spotify data
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, BranchPythonOperator
12 from airflow.models import Variable
13 from airflow.providers.postgres.hooks.postgres import PostgresHook
14 from datetime import datetime, timedelta
15 import requests
16 import base64
17 from urllib.parse import urlencode
18 import time
19 import pandas as pd
20 from datetime import datetime
21 import os
22 import unicodedata
23 import uuid
24
25
26 default_args = {
27     'owner': 'group6',
28     'depends_on_past': False,
29     'start_date': datetime(2025, 1, 1),
30     'email_on_failure': False,
31     'email_on_retry': False,
32     'retries': 1,
33     'retry_delay': timedelta(minutes=5),
34 }
35
36
37
38 # Retrieve tracks from specific playlists from Spotify API
39 def get_spot_tracks():
40
41     # These need to be set in variables.json
42     client_id = Variable.get("SPOTIFY_CLIENT_ID")
43     client_secret = Variable.get("SPOTIFY_CLIENT_SECRET")
44
45     # Get a Spotify access token
46     auth_header = base64.b64encode(f"{client_id}:{client_secret}".encode()).decode()
47     headers = {
48         "Authorization": f"Basic {auth_header}",
```

```
49     "Content-Type": "application/x-www-form-urlencoded",
50 }
51 data = {"grant_type": "client_credentials"}
52
53 print("Requesting an access token from Spotify")
54 response = requests.post("https://accounts.spotify.com/api/token", headers=headers,
data=data)
55 if response.status_code == 200:
56     token = response.json().get("access_token")
57     token_type = response.json().get("token_type")
58     expires_in = response.json().get("expires_in")
59     print(f"Received a Spotify {token_type} access token that expires in {expires_in}
seconds\n")
60 else:
61     print("Error:", response.status_code, response.text)
62
63
64 headers = {"Authorization": f"Bearer {token}"}
65
66 # Harcoding good playlists for our purposes
67 # Need to do query to get the playlist ID, Spotify seems to rotate playlist IDs
68 playlist_queries = [
69     "Billboard Hot 100: All #1 hit songs 1958-2024",
70     "Most well known songs ever",
71     "100 Greatest Rock Songs",
72     "Top 100 hip hop hits of all time",
73     "100 Greatest Pop Songs",
74     "Top 100 Most Popular Electronic Songs Of All Time",
75     "Top 100 Alternative Rock Songs",
76     "Top 100 Jazz Songs",
77     "The 100 Greatest Heavy Metal Songs of All Time",
78     "100 Best Folk Songs",
79     "Top hits of the 2000s",
80     "90s hits top 100 songs",
81     "Top hits of the 2010s",
82     "80s hits top 100 songs",
83     "70s hits top 100 songs"
84 ]
85
86 # Dictionary to store playlist information
87 playlist_info = {}
88
89 # Search for each playlist and retrieve its ID
90 for query in playlist_queries:
91     time.sleep(1)
92     print(f"Searching for playlist: \"{query}\"...")
93
94     # Set up search parameters
95     search_url = "https://api.spotify.com/v1/search"
96     params = {
```

```
97         "q": query,
98         "type": "playlist",
99         "limit": 1
100     }
101
102     # Make the search request
103     response = requests.get(f"{search_url}?{urlencode(params)}", headers=headers)
104
105     if response.status_code == 200:
106         search_data = response.json()
107
108         # Check if any playlists were found
109         if search_data['playlists']['items']:
110             # Get the first matching playlist
111             first_result = search_data['playlists']['items'][0]
112             playlist_name = first_result['name']
113             playlist_id = first_result['id']
114
115             print(f"Found playlist: \"{playlist_name}\" with ID: {playlist_id}")
116
117             # Store the playlist info
118             playlist_info[playlist_name] = playlist_id
119         else:
120             print(f"No playlists found for query: \"{query}\"")
121     else:
122         print(f"Error searching for playlist: {response.status_code}")
123         print(response.text)
124
125     print(f"Found {len(playlist_info)} playlists\n")
126
127
128     all_tracks = []
129
130     # Get tracks from the playlist
131     for playlist_name, playlist_id in playlist_info.items():
132         playlist_tracks_url = f"https://api.spotify.com/v1/playlists/{playlist_id}/tracks"
133         page = 0
134
135         print(f"Retrieving tracks from playlist \"{playlist_name}\"...")
136
137         # Get all pages
138         while playlist_tracks_url:
139             page += 1
140             if page > 50: # Just a failsafe to prevent excessive API calls
141                 break
142
143             # To comply with Spotify API rate limits
144             time.sleep(1)
145
146             response = requests.get(playlist_tracks_url, headers=headers)
```

```
147         if response.status_code == 200:
148             data = response.json()
149             track_items = data.get("items", [])
150
151             for track_info in track_items:
152                 track = track_info.get("track", {})
153                 all_tracks.append(track)
154                 print(f"Getting page {page} of tracks from playlist \"{playlist_name}\"...")
155                 playlist_tracks_url = data.get("next")
156
157         else:
158             print("Error fetching tracks:", response.status_code)
159             break
160
161     print(f"Retrieved all tracks from Spotify playlist(s): ")
162     print(f"Total Tracks Retrieved: {len(all_tracks)}\n")
163
164     return all_tracks
165
166
167 # Takes columns we want to keep, and returns a dataframe
168 def parse_spotify_tracks(spot_tracks):
169     print(f"Parsing Spotify data to keep relevant columns...\n")
170     all_track_details = []
171
172     # Only keep columns with data we find interesting for tables
173     for track in spot_tracks:
174         track_details = {}
175
176         # Keep top artist and list of all artists separately
177         artists = track.get("artists", None)
178         if artists:
179             top_artist = artists[0]["name"]
180             track_details["top_artist"] = top_artist
181             artist_names = [artist.get("name", "") for artist in artists] # Empty is OK here
182             track_details["artists"] = ", ".join(artist_names)
183         else:
184             track_details["top_artist"] = None
185             track_details["artists"] = None
186
187         track_details["song_name"] = track.get("name", None)
188         track_details["duration"] = track.get("duration_ms", None)
189         track_details["popularity"] = track.get("popularity", None)
190         track_details["spotify_id"] = track.get("id", None)
191
192
193         album = track.get("album", None)
194         if album:
195             album_name = album.get("name", None)
196             track_details["album_name"] = album_name
```

```

197         track_details["album_id"] = album.get("id", None)
198         track_details["album_release_date"] = album.get("release_date", None)
199         track_details['album_release_date_precision'] = album.get("release_date_precision", None)
200         images = album.get("images", [])
201         if len(images) > 1:
202             track_details["album_image"] = images[1].get("url", None)
203         else:
204             track_details["album_image"] = None
205     else:
206         track_details["album_name"] = None
207         track_details["album_id"] = None
208         track_details["album_release_date"] = None
209         track_details['album_release_date_precision'] = None
210         track_details["album_image"] = None
211
212     track_details["explicit_lyrics"] = track.get("explicit", None)
213     track_details["isrc"] = track.get("external_ids", {}).get("isrc", None)
214     track_details["spotify_url"] = track.get("external_urls", {}).get("spotify", None)
215     track_details["available_markets"] = ", ".join(track.get("available_markets", None))
216
217     all_track_details.append(track_details)
218
219 df = pd.DataFrame(all_track_details)
220
221 # Make a timestamped csv for record
222 # Saves to data folder
223 data_dir = os.path.join(os.getcwd(), 'data')
224 os.makedirs(data_dir, exist_ok=True)
225 timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
226 spot_tracks_file_name = f"spotify_tracks_{timestamp}.csv"
227 spot_tracks_file_path = os.path.join(data_dir, spot_tracks_file_name)
228 df.to_csv(spot_tracks_file_path, index=False, encoding='utf-8-sig')
229 print(f"Successfully wrote to {spot_tracks_file_name}\n")
230
231 return df
232
233 # A few cases where only year is listed
234 # Need to prepend day and month
235 # Since to_datetime doesn't work with just year
236 def convert_date(date_str):
237     if pd.isna(date_str):
238         return None
239     date_str = str(date_str).strip()
240     if len(date_str) == 4 and date_str.isdigit():
241         return pd.to_datetime(f"01/01/{date_str}", errors='coerce')
242     return pd.to_datetime(date_str, errors='coerce')
243
244
245 # Function to replace accented characters

```

```

246 def replace_accented_characters(s):
247     if isinstance(s, str): # Check if the value is a string
248         normalized_string = unicodedata.normalize('NFD', s)
249         return ''.join(c for c in normalized_string if unicodedata.category(c) != 'Mn')
250     return s
251
252
253 # Further clean the spotify tracks dataframe
254 def clean_spotify_tracks(df):
255     print(f"Cleaning Spotify data...\n")
256
257     spot_df = df.copy()
258     spot_df.replace('', None, inplace=True)
259
260     print(f"Current number of tracks: {len(spot_df)}")
261     print(f"Dropping duplicates where playlists overlap...")
262
263     # isrc is great duplicate key for this spotify data
264     duplicate_count = spot_df.duplicated(subset='isrc', keep=False).sum()
265     print(f"Number of 'isrc' duplicates: {duplicate_count}")
266     spot_df = spot_df.drop_duplicates(subset='isrc', keep='first')
267     print(f"Removed {duplicate_count} duplicates")
268     print(f"Remaining number of tracks: {len(spot_df)}\n")
269
270     # Clean a few specific columns
271     cols_to_clean = ['top_artist', 'artists', 'song_name', 'album_name']
272     for col in cols_to_clean:
273         if col in spot_df.columns:
274             spot_df[col] = spot_df[col].str.lower().str.strip()
275             spot_df[col] = spot_df[col].apply(replace_accented_characters)
276
277     # Make sure integer columns are ints
278     int_columns = ['duration', 'popularity']
279     for col in int_columns:
280         spot_df[col] = pd.to_numeric(spot_df[col], errors='coerce')
281     spot_df[int_columns] = spot_df[int_columns].where(pd.notnull(spot_df[int_columns]), None)
282
283     # Force int columns to Int64 (needed for one column that was being read as float)
284     for col in int_columns:
285         spot_df[col] = spot_df[col].astype('Int64')
286
287     # Replace & with and - helps with matching
288     spot_df['top_artist'] = spot_df['top_artist'].str.replace('&', 'and')
289     spot_df['artists'] = spot_df['artists'].str.replace('&', 'and')
290
291     # Helps with matching - removes " (feat. *)"
292     spot_df['song_name'] = spot_df['song_name'].str.replace(r' \((feat\..*?)\)', '',
293     regex=True)
294
295     # Delete if no artist

```

```
295     spot_df = spot_df[spot_df['artists'].notna() & (spot_df['artists'] != '')] # Delete if no
artist
296
297     spot_df['album_release_date'] = spot_df['album_release_date'].apply(convert_date)
298
299     # Checking for any artists, song_name duplicates
300     print("Checking for any duplicates by top_artist and song_name...")
301     duplicate_count = spot_df.duplicated(subset=['top_artist', 'song_name']).sum()
302     print(f"Number of duplicates in 'top_artist' and 'song_name': {duplicate_count}")
303     spot_df = spot_df.drop_duplicates(subset=['top_artist', 'song_name'])
304     print(f"Removed {duplicate_count} duplicates")
305     print(f"Remaining number of tracks: {len(spot_df)}\n")
306
307     # Sort by top artist and song name
308     # Not needed for storage, but helps with debugging
309     spot_df.sort_values(by=['top_artist', 'song_name'], inplace=True)
310
311     # Check for similar entries (adjacent rows after sorting)
312     print("Checking for similar songs based on character matching...")
313     indices_to_drop = []
314     similar_count = 0
315
316     # Efficient way to check adjacent rows for duplicates
317     # Keep the more popular one
318     for i in range(len(spot_df) - 1): # Iterate through all rows except the last one
319         current_row = spot_df.iloc[i]
320         next_row = spot_df.iloc[i + 1]
321
322         # Get the values to compare
323         current_artist = str(current_row['top_artist'])
324         next_artist = str(next_row['top_artist'])
325         current_song = str(current_row['song_name'])
326         next_song = str(next_row['song_name'])
327
328         # Compare full length or first 10 characters - 10 is arbitrary
329         artist_compare_len = min(len(current_artist), len(next_artist), 10)
330         song_compare_len = min(len(current_song), len(next_song), 10)
331         if song_compare_len == 0: # Skip if either song is empty
332             continue
333
334         # Check if both beginning parts match
335         if (current_artist[:artist_compare_len] == next_artist[:artist_compare_len] and
336             current_song[:song_compare_len] == next_song[:song_compare_len]):
337
338             # Found a potential duplicate, decide which one to keep based on popularity
339             current_popularity = current_row.get('popularity', 0)
340             next_popularity = next_row.get('popularity', 0)
341
342             # If next row has higher popularity, drop current row
343             if next_popularity > current_popularity:
```

```

344         indices_to_drop.append(spot_df.index[i])
345         similar_count += 1
346         print(f"Similar songs found: '{current_song}' and '{next_song}' by
'{current_artist}'")
347         print(f" Keeping '{next_song}' (popularity: {next_popularity})")
348         print(f" Dropping '{current_song}' (popularity: {current_popularity})")
349         # If current row has higher or equal popularity, drop next row
350         else:
351             indices_to_drop.append(spot_df.index[i + 1])
352             similar_count += 1
353             print(f"Similar songs found: '{current_song}' and '{next_song}' by
'{current_artist}'")
354             print(f" Keeping '{current_song}' (popularity: {current_popularity})")
355             print(f" Dropping '{next_song}' (popularity: {next_popularity})")
356
357         # Drop the identified duplicates
358         if similar_count > 0:
359             spot_df = spot_df.drop(indices_to_drop)
360             print(f"Removed {similar_count} similar songs based on character matching")
361         else:
362             print("No similar songs found")
363
364         print(f"Remaining number of Spotify tracks: {len(spot_df)}\n")
365
366
367         data_dir = os.path.join(os.getcwd(), 'data')
368         os.makedirs(data_dir, exist_ok=True)
369         timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
370         spot_clean_file_name = f"spotify_tracks_cleaned_{timestamp}.csv"
371         spot_clean_file_path = os.path.join(data_dir, spot_clean_file_name)
372         spot_df.to_csv(spot_clean_file_path, index=False, encoding='utf-8-sig')
373         print(f"Successfully processed Spotify tracks...")
374         print(f"Successfully wrote record to {spot_clean_file_name}\n")
375
376         return spot_df
377
378
379
380 # Calls other functions to get and clean spotify data
381 # Then loads it into postgres
382 def load_spotify_data():
383
384     raw_spot_tracks = get_spot_tracks()
385     df = parse_spotify_tracks(raw_spot_tracks)
386     clean_spot_df = clean_spotify_tracks(df)
387     num_tracks = len(clean_spot_df)
388
389     # Get postgres connection
390     pg_hook = PostgresHook(postgres_conn_id='pg_group6')
391     with pg_hook.get_conn() as conn:

```



```

392         with conn.cursor() as cur:
393             for index, row in clean_spot_df.iterrows():
394                 group6_id = uuid.uuid5(uuid.NAMESPACE_DNS, str(row['song_name'].strip() +
row['top_artist'].strip()))
395
396                 cur.execute("""INSERT INTO spotify_tracks (group6_id, top_artist, artists,
song_name, duration,
397                             popularity, spotify_id, album_name, album_id, album_release_date,
album_release_date_precision,
398                             album_image, explicit_lyrics, isrc, spotify_url,
available_markets)
399                             VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s,
%s, %s)""",
400
401                             (group6_id,
402                             row['top_artist'],
403                             row['artists'].split(','),
404                             row['song_name'],
405                             row['duration'],
406                             row['popularity'],
407                             row['spotify_id'],
408                             row['album_name'],
409                             row['album_id'],
410                             row['album_release_date'],
411                             row['album_release_date_precision'],
412                             row['album_image'],
413                             row['explicit_lyrics'],
414                             row['isrc'],
415                             row['spotify_url'],
416                             row['available_markets']))
417
418                 conn.commit()
419
420             print(f"Saved {num_tracks} Spotify tracks to database")
421
422 with DAG(
423     'spotify_dag',
424     default_args=default_args,
425     description='Dag for spotify data',
426     schedule_interval=None,
427     catchup=False
428 ) as spotify_dag:
429
430     start_task = EmptyOperator(task_id='start')
431
432     create_tables = PostgresOperator(
433         task_id='create_tables',
434         postgres_conn_id='pg_group6',
435         sql='sql/spotify_create.sql'
436     )
437
438     load_spotify_data = PythonOperator(

```

```
438         task_id='load_spotify_data',
439         python_callable=load_spotify_data
440     )
441
442     end_task = EmptyOperator(task_id='end')
443
444     start_task >> create_tables >> load_spotify_data >> end_task
445
446
```