**airflow\dags\lastfm_brainz_dag.py**

```python
 1  from airflow import DAG
 2  from airflow.operators.empty import EmptyOperator
 3  from airflow.providers.postgres.operators.postgres import PostgresOperator
 4  from airflow.operators.python_operator import PythonOperator, BranchPythonOperator
 5  from airflow.models import Variable
 6  from airflow.providers.postgres.hooks.postgres import PostgresHook
 7  import pandas as pd
 8  import requests
 9  import time
10  import json
11  from datetime import datetime, timedelta
12  import os
13  import numpy as np
14  import unicodedata
15  import uuid
16
17
18  default_args = {
19      'owner': 'group6',
20      'depends_on_past': False,
21      'start_date': datetime(2025, 1, 1),
22      'email_on_failure': False,
23      'email_on_retry': False,
24      'retries': None
25  }
26
27
28  user_agent = Variable.get("LASTFM_USER_AGENT")
29  api_key = Variable.get("LASTFM_API_KEY")
30  base_url = Variable.get("LASTFM_BASE_URL")
31
32  MAX_TAGS = 50
33  TRACKS_PER_TAG = 100
34  HEADERS = {"User-Agent": user_agent}
35
36  def get_top_tags():
37      print(f"Getting top {MAX_TAGS} tags from LastFM")
38
39      # Get postgres connection
40      pg_hook = PostgresHook(postgres_conn_id='pg_group6')
41
42      # Set up params for lastfm api
43      params = {
44          "method": "tag.getTopTags",
45          "api_key": api_key,
46          "format": "json",
47      }
48
```

```python
49        # Call lastfm api
50        top_tags_r = requests.get(base_url, params=params, headers=HEADERS)
51        top_tags_data = top_tags_r.json()
52        df_top_tags = pd.DataFrame(top_tags_data['toptags']['tag'][:MAX_TAGS])
53
54        # Write data to postgres
55        with pg_hook.get_conn() as conn:
56            with conn.cursor() as cur:
57                for index, row in df_top_tags.iterrows():
58                    cur.execute("""INSERT INTO lastfm_top_tags (tag_name, tag_count, tag_reach)
59                                VALUES (%s, %s, %s)""", (row['name'], row['count'],
   row['reach']))
60            conn.commit()
61        print("Saved tags to database")
62        return
63
64 # Gets basic track info from tag.getTopTracks endpoint
65 def get_track_basics():
66        print(f"Getting basics for {TRACKS_PER_TAG} tracks for each tag")
67
68        # Get api key and base url from airflow variables
69        api_key = Variable.get("LASTFM_API_KEY")
70        base_url = Variable.get("LASTFM_BASE_URL")
71
72        # Get postgres connection
73        pg_hook = PostgresHook(postgres_conn_id='pg_group6')
74
75
76
77        # Get top tags from db
78        with pg_hook.get_conn() as conn:
79            with conn.cursor() as cur:
80                cur.execute("SELECT tag_name FROM lastfm_top_tags")
81                top_tags = [tag[0] for tag in cur.fetchall()]
82
83        # List to store all tracks
84        all_tracks_initial = []
85        track_data = []
86
87        # Get tracks for each tag
88        for tag in top_tags:
89            print(f"Getting tracks for tag: {tag}")
90            tag_tracks = []
91
92            # Need to paginate
93            page = 1
94            while len(tag_tracks) < TRACKS_PER_TAG:
95                time.sleep(0.2)
96                params = {
97                    "method": "tag.getTopTracks",
```

```python
 98                    "tag": tag,
 99                    "api_key": api_key,
100                    "format": "json",
101                    "limit": TRACKS_PER_TAG,
102                    "page": page
103                }
104
105            response = requests.get(base_url, params=params, headers=HEADERS).json()
106            tracks = response.get("tracks", {}).get("track", [])
107            if not tracks:
108                break
109
110            # Add tag to each track
111            for track in tracks:
112                track["tag"] = tag
113
114            tag_tracks.extend(tracks)
115            page += 1
116        tag_tracks = tag_tracks[:TRACKS_PER_TAG]
117        all_tracks_initial.extend(tag_tracks)
118
119    print("Total number of tracks retrieved from LastFM:", len(all_tracks_initial))
120
121    # If duplicate mbid, consolidate tag rank
122    consolidated_tracks = {}
123    for track in all_tracks_initial:
124        mbid = track['mbid']
125        if not mbid:
126            continue
127        tag = track.get('tag', None)
128        rank = track.get('@attr', {}).get('rank', None)
129        rank = int(rank) if rank is not None else None
130        if mbid in consolidated_tracks:
131            consolidated_tracks[mbid]['tag_ranks'][tag] = rank
132        else:
133            track_basics = {
134                'mbid': mbid,
135                'tag_ranks': {tag: rank}
136            }
137            consolidated_tracks[mbid] = track_basics
138
139    # Convert to dataframe
140    track_data = list(consolidated_tracks.values())
141    df_lastfm_initial = pd.DataFrame(track_data)
142    initial_count_raw = len(all_tracks_initial)
143    final_count = len(df_lastfm_initial)
144    print(f"Retrieved basic info for {initial_count_raw} raw tracks")
145    print(f"After consolidating tracks by 'mbid', we have {final_count} unique tracks")
146    print(f"Removed {initial_count_raw - final_count} duplicate tracks")
147
```

```python
148         return df_lastfm_initial

149

150 # Gets detailed track info from track.getInfo endpoint
151 def get_track_details(df_initial):
152     request_count = 0
153     all_tracks_details = []

154

155     for index, row in df_initial.iterrows():
156         mbid = row['mbid']

157

158         params = {
159             "method": "track.getInfo",
160             "mbid": mbid,
161             "api_key": api_key,
162             "format": "json"
163         }

164

165         request_count += 1
166         try:
167             response = requests.get(base_url, params=params, headers=HEADERS)
168             if response.status_code == 200:
169                 response_data = response.json()
170                 if 'track' in response_data:
171                     track_data = response_data['track']
172                     track_data['tag_ranks'] = row['tag_ranks']
173                     all_tracks_details.append(track_data)
174                 if request_count % 10 == 0:
175                     print(f"Retrieved last.fm track info for {request_count} tracks...")
176             else:
177                 print(f"Failed to retrieve track info for {mbid}. Status code:
    {response.status_code}")
178         except Exception as e:
179             print(f"Error retrieving track info for {mbid}: {e}")

180

181         # Rate limit
182         # Last.fm doesn't publish their limits, but 0.2s was mentioned in some forums
183         time.sleep(0.2)

184

185     return all_tracks_details

186

187

188 def process_lastfm_tracks(all_tracks_details):
189     print("Processing last.fm tracks...")
190     lastfm_tracks = []
191     for track in all_tracks_details:
192         track_details = {}
193         track_details['artist'] = track.get('artist', {}).get('name', None)
194         track_details['song_name'] = track.get('name', None)
195         track_details['duration'] = track.get('duration', None)
196         track_details['listeners'] = track.get('listeners', None)
```

```python
197            track_details['playcount'] = track.get('playcount', None)
198            track_details['mbid'] = track.get('mbid', None)
199            track_details['album_name'] = track.get('album', {}).get('title', None)
200            track_details['url'] = track.get('url', None)
201            track_details['tag_ranks'] = json.dumps(track.get('tag_ranks', None))
202            toptags = track.get('toptags', {}).get('tag', [])
203            if toptags:
204                tag_names = [tag.get('name') for tag in toptags if tag.get('name')]
205                track_details['toptags'] = tag_names if tag_names else None
206            else:
207                track_details['toptags'] = None
208            track_details['wiki_summary'] = track.get('wiki', {}).get('summary', None)
209
210            lastfm_tracks.append(track_details)
211
212        df_lastfm_tracks = pd.DataFrame(lastfm_tracks)
213
214        df_lastfm_tracks.replace('', None, inplace=True)
215
216        # Clean a few columns
217        cols_to_clean = ['artist', 'song_name', 'album_name']
218        for col in cols_to_clean:
219            if col in df_lastfm_tracks.columns:
220                df_lastfm_tracks[col] = df_lastfm_tracks[col].str.lower().str.strip()
221
222        # Make sure integer columns are ints
223        int_columns = ['duration', 'listeners', 'playcount']
224        for col in int_columns:
225            df_lastfm_tracks[col] = pd.to_numeric(df_lastfm_tracks[col], errors='coerce')
226        df_lastfm_tracks[int_columns] =
    df_lastfm_tracks[int_columns].where(pd.notnull(df_lastfm_tracks[int_columns]), None)
227        for col in int_columns:
228            df_lastfm_tracks[col] = df_lastfm_tracks[col].astype('Int64')
229
230        print(f"Processed and cleaned {len(df_lastfm_tracks)} tracks")
231
232        # Make a timestamped csv for record
233        timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
234        data_dir = os.path.join(os.getcwd(), 'data')
235        os.makedirs(data_dir, exist_ok=True)
236        lastfm_clean_file_name = f"lastfm_clean_{timestamp}.csv"
237        lastfm_clean_file_path = os.path.join(data_dir, lastfm_clean_file_name)
238        df_lastfm_tracks.to_csv(lastfm_clean_file_path, index=False, encoding='utf-8-sig')
239        print(f"Successfully processed Last.fm tracks...")
240        print(f"Successfully wrote record to {lastfm_clean_file_name}\n")
241
242        return df_lastfm_tracks
243
244  # Calls the other functions to get last.fm tracks
245  # Then loads into postgres
```

```python
246  def get_and_load_lastfm_tracks():
247      print("Getting and loading last.fm tracks...")
248      df_initial = get_track_basics()
249      all_tracks_details = get_track_details(df_initial)
250      df_lastfm_tracks = process_lastfm_tracks(all_tracks_details)
251
252      # Get postgres connection
253      pg_hook = PostgresHook(postgres_conn_id='pg_group6')
254
255
256      # Write data to postgres
257      with pg_hook.get_conn() as conn:
258          with conn.cursor() as cur:
259              for index, row in df_lastfm_tracks.iterrows():
260
261                  # Generate group6_id
262                  group6_id = str(uuid.uuid5(uuid.NAMESPACE_DNS, str(row['song_name'].strip() +
     row['artist'].strip())))
263
264                  cur.execute("""INSERT INTO lastfm_tracks (mbid, group6_id, artist, song_name,
     duration, listeners,
265                                  playcount, album_name, url, tag_ranks, toptags, wiki_summary)
266                                  VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s)
267                                  ON CONFLICT DO NOTHING""",
268                                  (row['mbid'], group6_id,
269                                   row['artist'], row['song_name'], row['duration'],
     row['listeners'],
270                                   row['playcount'], row['album_name'], row['url'],
271                                   row['tag_ranks'], row['toptags'], row['wiki_summary']))
272          conn.commit()
273          print(f"Saved {len(df_lastfm_tracks)} tracks to postgres")
274
275      return
276
277
278  # Gets acousticbrainz features
279  # For all tracks with mbids currently in the lastfm_tracks table
280  def get_ab_features():
281      print("Getting acousticbrainz features")
282
283      # Get base url from airflow variables
284      base_url = Variable.get("AB_BASE_URL")
285      endpoint = base_url + "/api/v1/high-level"
286
287      # Get postgres connection
288      pg_hook = PostgresHook(postgres_conn_id='pg_group6')
289
290      # Get all mbid from lastfm_tracks
291      with pg_hook.get_conn() as conn:
292          with conn.cursor() as cur:
293              # Select both mbid and group6_id
```

```python
294             cur.execute("SELECT mbid, group6_id FROM lastfm_tracks WHERE mbid IS NOT NULL and
    mbid != ''")
295             mbid_results = cur.fetchall()
296
297             # Create dict mapping from mbid to group6_id
298             mbid_to_group6id = {row[0]: row[1] for row in mbid_results}
299
300             # List of just mbids for the API calls
301             mbid_list = list(mbid_to_group6id.keys())
302
303     batch_size = 25
304     print('Starting to get features for', len(mbid_list), f'tracks using batch_size =
    {batch_size}')
305
306     results = []
307
308     # Get features for each batch of mbid
309     for i in range(0, len(mbid_list), batch_size):
310         batch = mbid_list[i:i + batch_size]
311         params = {'recording_ids': ';'.join(batch)}
312
313         response = requests.get(endpoint, params=params, headers=HEADERS)
314         response.raise_for_status()  # Raise exception for HTTP errors
315
316         response_data = response.json()
317
318         # Process each mbid in the response
319         for mbid, data in response_data.items():
320             if mbid not in batch or '0' not in data:
321                 continue
322
323             group6_id = mbid_to_group6id.get(mbid)
324             if not group6_id:
325                 print(f"Warning: No group6_id found for mbid: {mbid}")
326                 continue
327
328             metadata = data['0']['metadata']
329             tags = metadata.get('tags', {})
330             highlevel = data['0']['highlevel']
331
332             # Simple, direct field extraction
333             features = {
334                 'mbid': mbid,
335                 'group6_id': group6_id,
336                 'artist': json.dumps(tags.get('artist', None)),
337                 'song_name': json.dumps(tags.get('title', None)),
338                 'album': json.dumps(tags.get('album', None)),
339                 'date': json.dumps(tags.get('date', None)),
340                 'isrcs': json.dumps(tags.get('isrc', None)),
341                 'bpm': json.dumps(tags.get('bpm', None)),
```

```
342                        'initialkey': tags.get('initialkey', None),
343                        'musicbrainz_albumid': tags.get('musicbrainz_albumid', None),
344                        'musicbrainz_artistid': tags.get('musicbrainz_artistid', None),
345                        'mood': json.dumps(tags.get('mood', None)),
346
347                        'danceability_danceable': highlevel.get('danceability', {}).get('all',
      {}).get('danceable', None),
348                        'danceability_not_danceable': highlevel.get('danceability', {}).get('all',
      {}).get('not_danceable', None),
349                        'danceability_max_class': highlevel.get('danceability', {}).get('value',
      None),
350
351                        'gender_female': highlevel.get('gender', {}).get('all', {}).get('female',
      None),
352                        'gender_male': highlevel.get('gender', {}).get('all', {}).get('male', None),
353                        'gender_max_class': highlevel.get('gender', {}).get('value', None),
354
355                        'genre_alternative': highlevel.get('genre_dortmund', {}).get('all',
      {}).get('alternative', None),
356                        'genre_blues': highlevel.get('genre_dortmund', {}).get('all',
      {}).get('blues', None),
357                        'genre_electronic': highlevel.get('genre_dortmund', {}).get('all',
      {}).get('electronic', None),
358                        'genre_folkcountry': highlevel.get('genre_dortmund', {}).get('all',
      {}).get('folkcountry', None),
359                        'genre_funksoulrnb': highlevel.get('genre_dortmund', {}).get('all',
      {}).get('funksoulrnb', None),
360                        'genre_jazz': highlevel.get('genre_dortmund', {}).get('all', {}).get('jazz',
      None),
361                        'genre_pop': highlevel.get('genre_dortmund', {}).get('all', {}).get('pop',
      None),
362                        'genre_raphiphop': highlevel.get('genre_dortmund', {}).get('all',
      {}).get('raphiphop', None),
363                        'genre_rock': highlevel.get('genre_dortmund', {}).get('all', {}).get('rock',
      None),
364                        'genre_dortmund_max_class': highlevel.get('genre_dortmund', {}).get('value',
      None),
365
366                        'voice_instrumental_instrumental': highlevel.get('voice_instrumental',
      {}).get('all', {}).get('instrumental', None),
367                        'voice_instrumental_voice': highlevel.get('voice_instrumental',
      {}).get('all', {}).get('voice', None),
368                        'voice_instrumental_max_class': highlevel.get('voice_instrumental',
      {}).get('value', None),
369                        'genre': json.dumps(tags.get('genre', None))
370                    }
371
372
373                results.append(features)
374
375
```

```python
376             print('Batch', i // batch_size + 1, 'of', (len(mbid_list) - 1) // batch_size + 1,
    'complete')
377             time.sleep(0.2)
378
379     if results:
380         results_df = pd.DataFrame(results)
381
382         results_df.replace('', None, inplace=True)
383
384         # Convert empty strings to None for all numeric fields
385         numeric_fields = ['danceability_danceable', 'danceability_not_danceable',
386                           'gender_female', 'gender_male',
387                           'genre_alternative', 'genre_blues', 'genre_electronic',
388                           'genre_folkcountry', 'genre_funksoulrnb', 'genre_jazz',
389                           'genre_pop', 'genre_raphiphop', 'genre_rock',
390                           'voice_instrumental_instrumental', 'voice_instrumental_voice']
391
392         # Convert numeric fields that need to be floats
393         for field in numeric_fields:
394             if field in results_df.columns:
395                 results_df[field] = pd.to_numeric(results_df[field], errors='coerce')
396
397         # Replace NaN with None
398         results_df = results_df.where(pd.notnull(results_df), None)
399
400         timestamp = datetime.now().strftime('%Y%m%d_%H%M%S')
401         data_dir = os.path.join(os.getcwd(), 'data')
402         filename = f"acousticbrainz_features_{timestamp}.csv"
403         file_path = os.path.join(data_dir, filename)
404         results_df.to_csv(file_path, index=False, encoding='utf-8-sig')
405         print(f"Saved {len(results_df)} records to {filename}")
406     else:
407         print("No features were retrieved")
408
409
410     # Now open a new database connection to insert the data
411     with pg_hook.get_conn() as conn:
412         with conn.cursor() as cur:
413             for track in results:
414
415                 if 'bpm' in track and (track['bpm'] == '' or track['bpm'] == []):
416                     track['bpm'] = None
417
418                 # Directly replace all empty strings in all features
419                 for key in track:
420                     if track[key] == '':
421                         track[key] = None
422
423                 cur.execute("""INSERT INTO acousticbrainz_features
424                             (mbid, group6_id, artist, song_name, album, date, isrcs,
```

```
425                              bpm, initialkey, musicbrainz_albumid, musicbrainz_artistid, mood,
426                              danceability_danceable, danceability_not_danceable,
    danceability_max_class,
427                              gender_female, gender_male, gender_max_class,
428                              genre_alternative, genre_blues, genre_electronic,
    genre_folkcountry, genre_funksoulrnb,
429                              genre_jazz, genre_pop, genre_raphiphop, genre_rock,
    genre_dortmund_maxclass,
430                              voice_instrumental_instrumental, voice_instrumental_voice,
    voice_instrumental_max_class, genre)
431                              VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s,
    %s, %s, %s, %s,
432                                  %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s)""",
433                          (track['mbid'],
434                          track['group6_id'],
435                          track['artist'],
436                          track['song_name'],
437                          track['album'],
438                          track['date'],
439                          track['isrcs'],
440                          track['bpm'],
441                          track['initialkey'],
442                          track['musicbrainz_albumid'],
443                          track['musicbrainz_artistid'],
444                          track['mood'],
445                          track['danceability_danceable'],
446                          track['danceability_not_danceable'],
447                          track['danceability_max_class'],
448                          track['gender_female'],
449                          track['gender_male'],
450                          track['gender_max_class'],
451                          track['genre_alternative'],
452                          track['genre_blues'],
453                          track['genre_electronic'],
454                          track['genre_folkcountry'],
455                          track['genre_funksoulrnb'],
456                          track['genre_jazz'],
457                          track['genre_pop'],
458                          track['genre_raphiphop'],
459                          track['genre_rock'],
460                          track['genre_dortmund_max_class'],
461                          track['voice_instrumental_instrumental'],
462                          track['voice_instrumental_voice'],
463                          track['voice_instrumental_max_class'],
464                          track['genre']))
465              conn.commit()
466          print(f"Inserted {len(results)} tracks into database")
467
468
469  with DAG(
470      'lastfm_brainz_dag',
```

```
471          default_args=default_args,
472          description='Dag for lastfm and AcousticBrainz data',
473          schedule_interval=None,
474          catchup=False
475      ) as lastfm_dag:
476
477      start_task = EmptyOperator(task_id='start')
478
479      create_tables = PostgresOperator(
480          task_id='create_tables',
481          postgres_conn_id='pg_group6',
482          sql='sql/lastfm_brainz_create_z.sql'
483      )
484
485      load_top_tags = PythonOperator(
486          task_id='load_top_tags',
487          python_callable=get_top_tags
488      )
489
490      load_lastfm_tracks = PythonOperator(
491          task_id='load_lastfm_tracks',
492          python_callable=get_and_load_lastfm_tracks
493      )
494
495      load_ab_features = PythonOperator(
496          task_id='load_ab_features',
497          python_callable=get_ab_features
498      )
499
500      end_task = EmptyOperator(task_id='end')
501
502      start_task >> create_tables >> load_top_tags >> load_lastfm_tracks >> load_ab_features >>
   end_task
503
504
```