

Final Project Submission

- Student name: Justin Grisanti
- Student pace: self-paced
- Scheduled project review date/time: 6/20/2022
- Instructor name: Claude Fried
- Blog post URL: https://justingrisanti.github.io/spotify_recommendation_system
(https://justingrisanti.github.io/spotify_recommendation_system)

Section 1: Business Understanding

Context:

Spotify is an audio streaming and media services platform, created in 2006. It is one of the largest music streaming service providers with over 406 million monthly active users, including 180 million paying subscribers, as of December 2021.

Spotify offers digital copyright restricted recorded music and podcasts, including more than 82 million songs, from record labels and media companies. As a freemium service, basic features are free with advertisements and limited control, while additional features, such as offline listening and commercial-free listening, are offered via paid subscriptions. Spotify is currently available in 180+ countries as of October 2021. Users can search for music based on artist, album, or genre, and can create, edit, and share playlists.

Two of the most important aspects of Spotify that has led to its popularity are its music discovery functionalities, and playlist creation fostering a new social aspect to music listening. In a 2021 How-To Geek article called "6 Awesome Spotify Features You Should Be Using," 3 of the 6 features are related to playlists, and one speaks about music discovery. One of these features related to music discovery is called "Enhance." Enhance allows you to discover new tracks that might best fit one of your existing playlists. For example, if you have a playlist of a collection of 80s rock songs, Enhance might suggest that you add "Eye of the Tiger" by Survivor.

Business Problem:

Music is now one of the easiest forms of media to consume, due to smartphone technology allowing users to access millions of songs at their fingertips. However, with all of these choices, it is hard to know where to begin when it comes to finding new music that matches someone's preferences, and users may become overwhelmed when trying to find music.

The solution to this problem is to create a recommendation system from scratch that can reperform the functionality of Enhance, which is to use a selection of songs from a playlist and use content-based filtering to suggest a list of songs that are similar.

The stakeholders of this project are Spotify, music-listeners, DJs, and other music-related occupations.

The main purpose of this recommendation system is inferential, meaning that this model should be able infer information about songs from a given playlist and then to predict songs that a user will likely add to that same playlist.

Section 2: Data Understanding

Now that I have developed a overall business understanding, I will take a deeper dive into the data that I will be using for this project.

```
In [1]: 1 # Import relevant libraries
2
3 import pandas as pd
4 from numpy.random import seed
5 seed(123)
6 import numpy as np
7 import random
8 import shutil
9 import math
10 import statistics as stat
11 import os
12 from os.path import exists
13 import datetime
14 import seaborn as sns
15 from sklearn.metrics.pairwise import cosine_similarity as cos
16 from sklearn.impute import SimpleImputer
17 import matplotlib.pyplot as plt
18 import sklearn as sk
19 from sklearn.preprocessing import OneHotEncoder, MinMaxScaler
20 from sklearn.feature_extraction.text import TfidfVectorizer
21 import json
22 from pandas.io.json import json_normalize
23 import spotipy
24 from spotipy.oauth2 import SpotifyClientCredentials
25 import time
26 import config
27 sp = spotipy.Spotify(auth_manager=SpotifyClientCredentials(client_
28                                                         client_
29
30 import sqlite3
31 conn = sqlite3.connect('Data/music_recs.db')
32 cur = conn.cursor()
33 from pyspark import SparkContext
34 from pyspark.sql import SparkSession
35 import dataframe_image as dfi
36
37 import sys
38 import re
39 import collections
```

```
In [2]: 1 # Set input path for os functions
        2
        3 input_path = 'Data/spotify_million_playlist_dataset/data/'
```

```
In [3]: 1 # Check how many items are in the data folder. There are 1000 files
        2
        3 list_files = os.listdir(input_path)
        4 number_files = len(list_files)
        5 print(number_files)

1019
```

```
In [4]: 1 # View a handful of files to see their types
        2
        3 os.listdir(input_path)[0:5]
```

```
Out[4]: ['mpd.slice.549000-549999.json',
         'mpd.slice.613000-613999.json',
         'mpd.slice.115000-115999.json',
         'mpd.slice.778000-778999.json',
         'mpd.slice.290000-290999.json']
```

As we can see above, it appears that we have 1000 json files in our data folder. Each file appears to have 1000 records each, which means there are 1 million songs in our dataset. After trying to load the first file, it was too large for my computer to handle. Instead, I will load the first playlist to see its contents.

```
In [5]: 1 # Open json file to see formatting
        2
        3 test = open(input_path+'mpd.slice.0-999.json')
        4 data_0_999 = json.load(test)
        5 data_0_999['playlists'][0]
```

```
Out[5]: {'name': 'Throwbacks',
         'collaborative': 'false',
         'pid': 0,
         'modified_at': 1493424000,
         'num_tracks': 52,
         'num_albums': 47,
         'num_followers': 1,
         'tracks': [{ 'pos': 0,
                       'artist_name': 'Missy Elliott',
                       'track_uri': 'spotify:track:0UaMYEvWZi0ZqiD0oHU3YI',
                       'artist_uri': 'spotify:artist:2wIVse2owClT7go1WT98tk',
                       'track_name': 'Lose Control (feat. Ciara & Fat Man Scoop)',
                       'album_uri': 'spotify:album:6vV5UrXcfyQD1wu4Qo2I9K',
                       'duration_ms': 226863,
                       'album_name': 'The Cookbook'},
                    { 'pos': 1,
                       'artist_name': 'Britney Spears',
                       'track_uri': 'spotify:track:6I9VzXrHx09rA9A5euc8Ak',
                       'artist_uri': 'spotify:artist:26dSoYclwsYLMAKD3tp0r4',
                       'track_name': 'Torn'}
```

Looking at this playlist above, we can see that we have the following features:

Playlist Attributes

- Playlist Name
- Playlist Type
- Number of Tracks
- Number of Unique Albums
- Number of Followers
- Number of Edits
- Duration in Milliseconds
- Number of Artists

Song Attributes

- Artist Name
- Track URI
- Artist URI
- Track Name
- Album URI
- Duration in Milliseconds
- Album Name

Some features that will be important to our model will be track name, artist name, album name, and their respective URIs. These will help us get more detail about a song from Spotipy. Another piece that could be helpful is the number of followers. This shows us interest in a given playlist. If a playlist has a jumble of random songs that don't form a cohesive playlist, it will probably have less followers than a well-crafted playlist.

Below is a (reformatted) summary from the readme file provided with the [data \(https://www.aicrowd.com/challenges/spotify-million-playlist-dataset-challenge/dataset_files\)](https://www.aicrowd.com/challenges/spotify-million-playlist-dataset-challenge/dataset_files).

General Info

- Number of Playlists: 1,000,000
- Number of Tracks: 66,346,428
- Number of Unique Tracks: 2,262,292
- Number of Unique Albums: 734,684
- Number of Unique Artists: 295,860
- Number of Unique Titles: 92,944
- Number of Playlists with Descriptions: 18,760
- Number of Unique Normalized Titles: 17,381
- Average Playlist Length: 66.346428

Top Playlist Titles

- 10,000 country
- 10,000 chill
- 8,493 rap
- 8,481 workout
- 8,146 oldies

Top Tracks

- 46,574 HUMBLE. by Kendrick Lamar
- 43,447 One Dance by Drake
- 41,309 Broccoli (feat. Lil Yachty) by DRAM
- 41,079 Closer by The Chainsmokers
- 39,987 Congratulations by Post Malone

Top Artists

- 847,160 Drake
- 413,297 Kanye West
- 353,624 Kendrick Lamar
- 339,570 Rihanna
- 316,603 The Weeknd

As we see in the summary above, there are 1 million playlists with over 66 million songs. Within these playlists, there are 2.2 million unique songs, 734k unique albums, and 300k unique artists. There is a lot of data to work with here. The playlists are sorted into categories, with country and chill being the top, followed by rap and workout. Drake, Kanye West, and Kendrick Lamar are the 3 top artists. The next step in this process is to convert all of this json data to be compatible with python.

Section 3: Data Preparation

The first steps are to convert our JSON data to python and DataFrames. Once we have that, we need to unwrap our track data so it can be converted to a DataFrame, as well.

```
In [6]: 1 """
2 This code was used to open 100 of the files in the dataset and app
3 We then select the top 100 playlists according to number of follow
4 cohesive playlists
5 """
6
7 if os.path.exists('Data/Spotipy Custom DataFrames/spotify_playlist
8     pass
9 else:
10     ## Generate DataFrame for playlist data
11     spotify_playlist_data = pd.DataFrame()
12
13     ## For Loop to open and append 100 json files to a DataFrame
14     for item in range(0,100):
15         open_file = open(input_path+sorted(os.listdir(input_path))
16         load_file = json.load(open_file)
17         spotify_playlist_data = spotify_playlist_data.append(load_
18
19     ## Send our DataFrame to a pickle file so we can call it inste
20     spotify_playlist_data.to_pickle('Data/Spotipy Custom DataFrame
21
22     ## Creating a top 100 playlist DataFrame to select the top 100
23     spotify_playlist_top_100 = spotify_playlist_data.sort_values(b
24     spotify_playlist_top_100 = spotify_playlist_top_100.set_index(
25     spotify_playlist_top_100.to_pickle('Data/Spotipy Custom DataFr
```

```
In [7]: 1 # Open pickle file with our top 100 playlist DataFrame
2 spd_top100 = pd.read_pickle('Data/Spotipy Custom DataFrames/spotif
3 spd_top100
```

```
Out [7]:      name collaborative  modified_at  num_tracks  num_albums  num_followers
pid
```


180831	My Little Pony	false	1478908800	85	9	31539	[[{'artist_': 'Applet', 'track_': 'My Little Pony'}]]
159077	Rock Hits	false	1509408000	56	54	22102	[[{'artist_': 'Fighter', 'track_': 'Rock Hits'}]]
101121	Workout Playlist	false	1456531200	26	23	11745	[[{'artist_': 'Be', 'track_': 'tra', 'album_': 'Workout Playlist'}]]
147486	J cole	false	1491523200	139	55	7912	[[{'artist_': 'J.', 'track_': 'tra', 'album_': 'J cole'}]]
17675	raggaeton	false	1501718400	107	71	2994	[[{'artist_': 'R', 'track_': 'Ken-', 'album_': 'raggaeton'}]]
...
100819	Dance!	false	1440115200	10	6	74	[[{'artist_': 'Márqu', 'track_': 'Dance!', 'album_': 'Dance!' }]]
127212	Demi Lovato	false	1470009600	84	23	72	[[{'artist_': 'Lovato', 'track_': 'Demi Lovato', 'album_': 'Demi Lovato'}]]
154640	Thankful.	false	1479168000	63	62	72	[[{'artist_': 'tob', 'track_': 'tra', 'album_': 'Thankful.' }]]
143072	Dutch!	false	1490227200	216	197	68	[[{'artist_': 'Discl', 'track_': 'Dutch!', 'album_': 'Dutch!' }]]
147046	punk goes pop	false	1354147200	71	12	68	[[{'artist_': 'Mε', 'track_': 'May', 'album_': 'punk goes pop'}]]

100 rows × 11 columns

Now that we have aggregated our data, we can see that our track data is nested, which makes sense given we are using json data. In order to get this track data into a pandas dataframe, we will be using similar code to normalize our track data into its own dataframe. We will also pull in the playlist id so we know which track relates to which playlist for when we analyze or use SQL.

```
In [8]: 1 if os.path.exists('Data/Spotipy Custom DataFrames/unwrapped_track_
2         pass
3     else:
4         # Creating DataFrame for the tracks column in our playlist dat
5         unwrapped_track_data = pd.DataFrame()
6
7         # Adding playlist IDs to a list to be appended to the unwrapp
8         index_list = spd_top100.index.tolist()
9         pid = []
10
11        # Using json_normalize to reformat the json data the DataFrame
12        for item in range(0, len(spd_top100)):
13            unwrapped_track_data = unwrapped_track_data.append(pd.jso
14
15        # Appending the playlist IDs to the pid list
16        for row in range(0, len(spd_top100.index)):
17            track=0
18            while track<spd_top100['num_tracks'].iloc[row]:
19                track+=1
20                pid.append(index_list[row])
21
22        # Adding the playlist IDs to the unwrapped track data
23        unwrapped_track_data['pid'] = pid
24
25        # Save the data as a pickle file
26        unwrapped_track_data.to_pickle('Data/Spotipy Custom DataFrames
```

```
In [9]: 1 # Call the pickle file for data use
        2 unwrapped_track_data = pd.read_pickle('Data/Spotipy Custom DataFra
        3 unwrapped_track_data
```

```
Out[9]:
```

	pos	artist_name	track_uri	
0	0	Applebloom	spotify:track:527IbJxFcjjTix0ONdxDdS	spotify:artist:7ggsXdK95oJBki
1	1	Apple Jack	spotify:track:1mOnMHXxt2vGI0b804eQsy	spotify:artist:1r0v3fdCiqrr9m
2	2	Twilight Sparkle	spotify:track:4GciJR91Tj8a7dLJ12WFvr	spotify:artist:53CQUfjaBNRwVz
3	3	Twilight	spotify:track:68dCOK4xvIMF5Nhr44TevS	spotify:artist:53CQUfjaBNRwVz

```
In [10]: 1 # Ensure the sum of the number of tracks for each playlist equals
        2 spd_top100['num_tracks'].sum() == unwrapped_track_data.shape[0]
```

```
Out[10]: True
```

```
In [11]: 1 # Stripping the URI fields to leave only the URI itself
        2 for item in ['track', 'artist', 'album']:
        3     placeholder = item
        4     unwrapped_track_data[f'{placeholder}_uri'] = unwrapped_track_c
```

In [12]:

```
1 unwrapped_track_data = unwrapped_track_data[['pos', 'track_uri', 'tr
2 unwrapped_track_data
```

Out[12]:

	pos	track_uri	track_name	artist_name	artist_uri	al
0	0	527IbJxFcjjTix0ONdxDdS	Hearts as Strong as Horses	Applebloom	7ggsXdK95oJBkuZu1txVjC	
1	1	1mOnMHXxt2vGI0b804eQsy	Apples to the Core	Apple Jack	1r0v3fdCiqrr9mYtvbCccT	
2	2	4GciJR91Tj8a7dLJ12WFvr	Ballad of the Crystal Ponies	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	
3	3	68dC0K4xglMF5Nhr44TevS	Find a Way	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	
4	4	0vhKFkvwgdPBrrNr9gUbVa	A True, True Friend	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	
...
9318	66	3G0PWfSGIsUrl4EI8u46EX	All Or Nothing	Fake ID	4GwCpkFWajqx3KSm0MVh2a	
9319	67	5qpQhP4hyjCKC95oRtTqni	Borderline	Showoff	6IKhTkyp4EJ0ocidcwafs6	
9320	68	7wwXG5FOebfhVBotX4vTXo	Send Me An Angel	Thrice	3NChzMpu9exTINPiqUQ2DE	
9321	69	1lbfP4HOrAS0fB8eloEYko	Baby One More Time	Nicotine	0p3U0uLx2oSf0yn8i5XZki	
9322	70	6ktA174mwmCqcb9hfdL178	Heaven Is A Place On Earth	Student Rick	6AluDNoFgmeUTnOc7DXYIN	

9323 rows × 7 columns

Now that our track data is unwrapped, we will remove duplicates to get a DataFrame of unique songs.

```
In [13]: 1 # Drop duplicates, reset index, and isolate important columns.
2 unique_track_data = unwrapped_track_data.drop_duplicates(subset='track_name')
3 unique_track_data = unique_track_data.reset_index()
4 unique_track_data = unique_track_data[['track_uri', 'track_name', 'artist_name', 'album_name']]
5 unique_track_data.to_sql('unique_track_data', con = conn, if_exists='replace')
6 unique_track_data
```

```
Out[13]:
```

	track_uri	track_name	artist_name	artist_uri	album_name
0	527IbJxFcjjTix0ONdxDdS	Hearts as Strong as Horses	Applebloom	7ggsXdK95oJBkuZu1txVjC	Soi Poi (Music the O
1	1mOnMHXxt2vGI0b804eQsy	Apples to the Core	Apple Jack	1r0v3fdCiqrr9mYtvbCccT	Soi Poi (Music the O
2	4GciJR91Tj8a7dLJ12WFvr	Ballad of the Crystal Ponies	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	Soi Poi (Music the O
3	68dC0K4xgIMF5Nhr44TevS	Find a Way	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	Soi Poi (Music the O
4	0vhKFkwgdPBrrNr9gUbVa	A True, True Friend	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	Soi Poi (Music the O
...
8076	3G0PWfSGIsUrI4EI8u46EX	All Or Nothing	Fake ID	4GwCpkFWajqx3KSm0MVh2a	Punk
8077	5qpQhP4hyjCKC95oRtTqni	Borderline	Showoff	6lKhTkyp4EJ0ocidcwafs6	Punk
8078	7wwXG5FOebfhVBotX4vTXo	Send Me An Angel	Thrice	3NChzMpu9exTINPiqUQ2DE	Punk
8079	1lbfP4HOrAS0fB8eloEYko	Baby One More Time	Nicotine	0p3U0uLx2oSf0yn8i5XZki	Punk
8080	6ktA174mwmCqcb9hfdL178	Heaven Is A Place On	Student	6AluDNoFgmeUTnOc7DYXIN	Punk

Earth

Rick

8081 rows × 6 columns

Now that I have my data prepared into tables, we will need to import some relevant features that describe the type of each song. While we have already have many features that describe the music, such as song name, length and artist, none of these describe the *characteristics* of each song. To do this, I have imported the Spotipy library. Per Spotipy's website, this library is described as "a lightweight Python library for the Spotify Web API. With Spotipy you get full access to all of the music data provided by the Spotify platform."

We can use this library to get more information about each song, and therefore classify the type of playlist that we are listening to. The reason we have to do this is because the playlist names might not always have a good description of the type of playlist we have. For example, playlist 159077 above is labelled "Rock Hits!", but playlist 154640 is just labelled an arbitrary title, "Thankful.". We can name Rock songs, but "Thankful" songs can be more subjective, and it is harder to classify a sound or vibe from the title alone. Once we have the type of playlist that we are looking at, we can begin suggesting songs from a similar genre or "vibe".

We will be using the API to import thhe following:

1. Audio features that analyze the sound/vibe of the song
2. Genre for each song
3. Popularity of each song

3.1 Audio Features DataFrames

```
In [14]: 1 # Using the audio_features method to get characteristics of our song
         2 sp.audio_features('0UaMYEvWZi0ZqiD0oHU3YI')
```

```
Out[14]: [{'danceability': 0.904,
            'energy': 0.813,
            'key': 4,
            'loudness': -7.105,
            'mode': 0,
            'speechiness': 0.121,
            'acousticness': 0.0311,
            'instrumentalness': 0.00697,
            'liveness': 0.0471,
            'valence': 0.81,
            'tempo': 125.461,
            'type': 'audio_features',
            'id': '0UaMYEvWZi0ZqiD0oHU3YI',
            'uri': 'spotify:track:0UaMYEvWZi0ZqiD0oHU3YI',
            'track_href': 'https://api.spotify.com/v1/tracks/0UaMYEvWZi0ZqiD0oHU3YI',
            'analysis_url': 'https://api.spotify.com/v1/audio-analysis/0UaMYEvWZi0ZqiD0oHU3YI',
            'duration_ms': 226864,
            'time_signature': 4}]
```

Above, Spotipy gives us relevant audio features for each song:

- **Danceability:** describes how suitable a track is for dancing based on a combination of musical elements including tempo, rhythm stability, beat strength, and overall regularity. A value of 0.0 is least danceable and 1.0 is most danceable.
- **Acousticness:** A measure from 0.0 to 1.0 of whether the track is acoustic.
- **Energy:** a measure from 0.0 to 1.0 and represents a perceptual measure of intensity and activity. Typically, energetic tracks feel fast, loud, and noisy.
- **Instrumentalness:** Predicts whether a track contains no vocals. The closer the instrumentalness value is to 1.0, the greater likelihood the track contains no vocal content.
- **Liveness:** Detects the presence of an audience in the recording. Higher liveness values represent an increased probability that the track was performed live.
- **Loudness:** The overall loudness of a track in decibels (dB). Loudness values are averaged across the entire track. Values typical range between -60 and 0 db.
- **Speechiness:** detects the presence of spoken words in a track. The more exclusively speech-like the recording (e.g. talk show, audio book, poetry), the closer to 1.0 the attribute value.
- **Tempo:** The overall estimated tempo of a track in beats per minute (BPM). In musical terminology, tempo is the speed or pace of a given piece and derives directly from the average beat duration.
- **Valence:** A measure from 0.0 to 1.0 describing the musical positiveness conveyed by a track. Tracks with high valence sound more positive (e.g. happy, cheerful, euphoric), while tracks with low valence sound more negative (e.g. sad, depressed, angry).

In [15]:

```

1  '''
2  Creating a DataFrame for all unique tracks and their audio features
3  file to save time. I also found I was getting rate limited from the
4  '''
5
6  if os.path.exists('Data/Spotipy Custom DataFrames/audio_features_df.pkl'):
7      pass
8  else:
9      audio_features_df = pd.DataFrame()
10
11     # For loop to append audio features for each track
12     for track in range(0, len(unique_track_data)):
13         audio_features_df = audio_features_df.append(sp.audio_features(track_id))
14
15     # Save to pickle to save loading time
16     audio_features_df.to_pickle('Data/Spotipy Custom DataFrames/audio_features_df.pkl')

```

In [16]:

```

1  ## Call file from pickle
2  audio_features_df = pd.read_pickle('Data/Spotipy Custom DataFrames/audio_features_df.pkl')
3  audio_features_df

```

Out[16]:

	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness
0	0.794	0.826	1	-4.384	1	0.0294	0.242000	0.0000
1	0.802	0.814	0	-2.489	1	0.0286	0.232000	0.0000
2	0.632	0.482	0	-7.480	1	0.0289	0.419000	0.0000
3	0.530	0.392	10	-8.648	1	0.0341	0.881000	0.0000
4	0.678	0.735	10	-6.135	1	0.0355	0.305000	0.0000
...
8076	0.360	0.909	0	-4.105	1	0.0587	0.017700	0.0000
8077	0.540	0.901	1	-5.300	1	0.0460	0.014000	0.0006

3.2 Genre DataFrames

Spotipy does not have a method to call genre for an individual song, so we will need to try our best in order to extract genre for an artist's music. I will create a function that gets the overall artist genre, which we can then apply to their songs. One possible downfall is if an artist goes across multiple genres, the songs might be mapped incorrectly.

```
In [17]: 1 unique_artist_data = unique_track_data.drop_duplicates(subset='artist_name')
         2 unique_artist_data
```

Out[17]:

	track_uri	track_name	artist_name	artist_uri	album_name
0	527IbJxFcjjTix0ONdxDdS	Hearts as Strong as Horses	Applebloom	7ggsXdK95oJBkuZu1txVjC	Soi Poi (Music the O
1	1mOnMHXxt2vGI0b804eQsy	Apples to the Core	Apple Jack	1r0v3fdCiqr9mYtvbCccT	Soi Poi (Music the O
2	4GciJR91Tj8a7dLJ12WFvr	Ballad of the Crystal Ponies	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	Soi Poi (Music the O
8	1rBQeB5zwAWf9mXL2HfMjt	Make a Wish - Extended Version	Pinkie Pie	7ExZeMNpyKhYSokWo9riU5	Soi Poi (Music the O
11	3I4r9SZcAhEydPIS5eS5Sz	Becoming Popular	Rarity	6PqIHmHPCKrZoyLMf98era	Soi Frier and I (Music
...
8075	4krBGFZoDTDtjtGZC8rG9g	Sometimes	Reach The Sky	7masWBjicrrYn9G2iZezdv	Punk
8076	3G0PWfSGIsUrI4EI8u46EX	All Or Nothing	Fake ID	4GwCpkFWajqx3KSm0MVh2a	Punk
8077	5qpQhP4hyjCKC95oRtTqni	Borderline	Showoff	6IKhTkyp4EJ0ocidcwafs6	Punk
8079	1lbfP4HOrAS0fB8eloEYko	Baby One More Time	Nicotine	0p3U0uLx2oSf0yn8i5XZki	Punk
8080	6ktA174mwmCqcb9hfdL178	Heaven Is A Place On Earth	Student Rick	6AluDNoFgmeUTnOc7DYXIN	Punk

3148 rows × 6 columns

```

In [18]: 1 # Extract artist data from spotipy using the search method and loc
          2
          3 def genre_extract(artist):
          4     extract = pd.DataFrame()
          5     result = sp.search(artist)
          6     if len(result['tracks']['items']) != 0:
          7         track = result['tracks']['items'][0]
          8         artist = sp.artist(track["artists"][0]["external_urls"]["s
          9         extract['genres'] = artist["genres"]
         10         extract['artist_id'] = artist["id"]
         11         return extract
         12     else:
         13         pass

```

```

In [19]: 1 genre_extract('Taylor Swift')

```

```

Out[19]:
   genres  artist_id
0    pop  06HL4z0CvFAxyc27GXpf02

```

We will now create a DataFrame with the genres for each artist.

```

In [20]: 1 # Creating DataFrame for genre data
          2
          3 if os.path.exists('Data/Spotipy Custom DataFrames/unique_genre_dat
          4     pass
          5 else:
          6     unique_genre_data = pd.DataFrame()
          7
          8     # Using function created above to append the genre data to the
          9     for artist in unique_artist_data['artist_name']:
         10         unique_genre_data = unique_genre_data.append(genre_extract
         11
         12     # Save to pickle to save runtime
         13     unique_genre_data.to_pickle('Data/Spotipy Custom DataFrames/ur

```

In [21]:

```

1 # Call data from pickle
2 unique_genre_data = pd.read_pickle('Data/Spotipy Custom DataFrames')
3 unique_genre_data

```

Out [21]:

	genres	artist_id
0	pony	7ggsXdK95oJBkuZu1txVjC
0	trap queen	1ziRj7e5Tm72Qf2ag6jHed
0	pony	53CQUfjaBNRwV2nFro1nac
0	pony	53CQUfjaBNRwV2nFro1nac
0	alternative emo	2ElhbnEc2cvYIAsXXbo9tg
...
4	tropical house	7i9j813KFoSBMldGqlh2Z1
5	uk dance	7i9j813KFoSBMldGqlh2Z1
0	bow pop	4zeHJ3kiJyjYXIIOcG4MA7
1	pop violin	4zeHJ3kiJyjYXIIOcG4MA7
0	modern rock	20JZFwl6HVI6yg8a4H3ZqK

11796 rows × 2 columns

In [22]:

```

1 # Set index and group data by artist ID
2 unique_genre_data = unique_genre_data.set_index('artist_id')
3 unique_genre_data = unique_genre_data.groupby('artist_id').agg({'g
4 unique_genre_data

```

Out [22]:

genres	
artist_id	
001aJOc7CSQVo3XzoLG4DK	[classic soul, disco, electro, funk, post-disc...
00FQb4jTyendYWaN8pK0wa	[art pop, pop]
00RJAKLnjGx4kVWVJbOJx1	[opm]
00TKPo9MxwZ0j4oovelxWZ	[alt z, dance pop, electropop, indie pooptimism...
00Z3UDoAQwzvGu13HoAM7J	[indie pop rap, pop rap, underground hip hop, ...
...	...
7z55f4aJkaPR4EF2BXqsq7	[gaming edm]
7z5WFjZAIYejWy0NI5lv4T	[contemporary country, country, pop]
7zICaxnDB9ZprDSiFpvbbW	[dirty south rap, gangster rap, hip hop, houst...
7zmk5lkmCMVvfvwF3H8FWC	[hip pop, neo soul, pop r&b, r&b, urban contem...
7zsin6lgVsR1rqSRCNYDwq	[stomp and holler]

2499 rows × 1 columns

In [23]:

```

1 # Create new column reformatting data from list to string
2 unique_genre_data['genres_string'] = ['',''].join(map(str, l)) for l
3 unique_genre_data

```

Out [23]:

artist_id	genres	genres_string
001aJOc7CSQVo3XzoLG4DK	[classic soul, disco, electro, funk, post-disc...	classic soul,disco,electro,funk,post-disco,qui...
00FQb4jTyendYWaN8pK0wa	[art pop, pop]	art pop,pop
00RJAKLnjGx4kVWVJbOJx1	[opm]	opm
00TKPo9MxwZ0j4oovelxWZ	[alt z, dance pop, electropop, indie poptimism...	alt z,dance pop,electropop,indie poptimism,nyc...
00Z3UDoAQwzvGu13HoAM7J	[indie pop rap, pop rap, underground hip hop, ...	indie pop rap,pop rap,underground hip hop,indi...
...
7z55f4aJkaPR4EF2BXqsq7	[gaming edm]	gaming edm
7z5WFjZAIYejWy0NI5lv4T	[contemporary country, country, pop]	contemporary country,country,pop
7zICaxnDB9ZprDSiFpvbbW	[dirty south rap, gangster rap, hip hop, houst...	dirty south rap,gangster rap,hip hop,houston r...
7zmk5lkmCMVvfvwF3H8FWC	[hip pop, neo soul, pop r&b, r&b, urban contem...	hip pop,neo soul,pop r&b,r&b,urban contemporary
7zsin6lgVsR1rqSRCNYDwq	[stomp and holler]	stomp and holler

2499 rows × 2 columns

In [24]:

```

1 # Split string column to explode it across all columns
2 unique_genre_data = unique_genre_data['genres_string'].str.split(' ')
3 unique_genre_data

```

Out [24]:

	0	1	2	3	4	
artist_id						
001aJOc7CSQVo3XzoLG4DK	classic soul	disco	electro	funk	post-disco	qi
00FQb4jTyendYWaN8pK0wa	art pop	pop	None	None	None	
00RJAKLnjGx4kVWVJbOJx1	opm	None	None	None	None	
00TKPo9MxwZ0j4oovelxWZ	alt z	dance pop	electropop	indie poptimism	nyc pop	
00Z3UDoAQwzvGu13HoAM7J	indie pop rap	pop rap	underground hip hop	indie pop rap	pop rap	unc
...	
7z55f4aJkaPR4EF2BXqsq7	gaming edm	None	None	None	None	
7z5WFjZAIYejWy0NI5lv4T	contemporary country	country	pop	None	None	
7zICaxnDB9ZprDSiFpvbbW	dirty south rap	gangster rap	hip hop	houston rap	new orleans rap	
7zmk5lkmCMVvfvwF3H8FWC	hip pop	neo soul	pop r&b	r&b	urban contemporary	
7zsin6lgVsR1rqSRCNYDwq	stomp and holler	None	None	None	None	

2499 rows × 42 columns

In [25]:

```

1 # Strip columns to remove any unwanted characters
2
3 for column in range(0,41):
4     unique_genre_data[column] = unique_genre_data[column].apply(lambda x:
5     unique_genre_data
6

```

Out [25]:

	0	1	2	3	4	
	artist_id					
001aJOc7CSQVo3XzoLG4DK	classic soul	disco	electro	funk	post-disco	qi
00FQb4jTyendYWaN8pK0wa	art pop	pop	None	None	None	
00RJAKLnjGx4kVWVJbOJx1	opm	None	None	None	None	
00TKPo9MxwZ0j4oovelxWZ	alt z	dance pop	electropop	indie pop optimism	nyc pop	
00Z3UDoAQwzvGu13HoAM7J	indie pop rap	pop rap	underground hip hop	indie pop rap	pop rap	unc
...	
7z55f4aJkaPR4EF2BXqsq7	gaming edm	None	None	None	None	
7z5WFjZAIYejWy0NI5lv4T	contemporary country	country	pop	None	None	
7zICaxnDB9ZprDSiFpvbbW	dirty south rap	gangster rap	hip hop	houston rap	new orleans rap	
7zmk5lkmCMVvfvwF3H8FWC	hip pop	neo soul	pop r&b	r&b	urban contemporary	
7zsin6lgVsR1rqSRCNYDwq	stomp and holler	None	None	None	None	

2499 rows × 42 columns

3.3 Popularity DataFrames

In [26]:

```

1 # We can get popularity of a track from the track method
2 sp.track('0KKkJNfGyhkQ5aFogxQAPU')

```

Out[26]:

```

{'album': {'album_type': 'album',
  'artists': [{'external_urls': {'spotify': 'https://open.spotify.c
om/artist/0du5cEVh5yTK9QJze8zA0C'}},
  'href': 'https://api.spotify.com/v1/artists/0du5cEVh5yTK9QJze8z
A0C',
  'id': '0du5cEVh5yTK9QJze8zA0C',
  'name': 'Bruno Mars',
  'type': 'artist',
  'uri': 'spotify:artist:0du5cEVh5yTK9QJze8zA0C'}],
  'available_markets': ['AD',
  'AE',
  'AG',
  'AL',
  'AM',
  'AO',
  'AR',
  'AT',
  'AU',
  'AZ',
  'BA',
  'BB',
  'BD',
  'BE',
  'BF',
  'BG',
  'BH',
  'BI',
  'BJ',
  'BM',
  'BN',
  'BO',
  'BR',
  'BS',
  'BT',
  'BV',
  'BW',
  'BY',
  'BZ',
  'CA',
  'CC',
  'CD',
  'CF',
  'CG',
  'CH',
  'CI',
  'CK',
  'CL',
  'CM',
  'CN',
  'CO',
  'CR',
  'CU',
  'CV',
  'CW',
  'CX',
  'CY',
  'CZ',
  'DE',
  'DG',
  'DK',
  'DM',
  'DO',
  'DZ',
  'EC',
  'EE',
  'EG',
  'EH',
  'ER',
  'ES',
  'ET',
  'FI',
  'FJ',
  'FK',
  'FM',
  'FO',
  'FR',
  'GA',
  'GB',
  'GD',
  'GE',
  'GF',
  'GG',
  'GH',
  'GI',
  'GL',
  'GM',
  'GN',
  'GP',
  'GQ',
  'GR',
  'GS',
  'GT',
  'GU',
  'GW',
  'GY',
  'HK',
  'HM',
  'HN',
  'HO',
  'HR',
  'HT',
  'HU',
  'ID',
  'IE',
  'IL',
  'IM',
  'IN',
  'IO',
  'IQ',
  'IR',
  'IS',
  'IT',
  'JE',
  'JM',
  'JO',
  'JP',
  'KE',
  'KG',
  'KH',
  'KI',
  'KM',
  'KN',
  'KR',
  'KW',
  'KY',
  'KZ',
  'LA',
  'LB',
  'LC',
  'LI',
  'LK',
  'LR',
  'LS',
  'LT',
  'LU',
  'LV',
  'LY',
  'MA',
  'MC',
  'MD',
  'ME',
  'MG',
  'MH',
  'MK',
  'ML',
  'MM',
  'MN',
  'MO',
  'MP',
  'MQ',
  'MR',
  'MS',
  'MT',
  'MU',
  'MV',
  'MW',
  'MX',
  'MY',
  'MZ',
  'NA',
  'NC',
  'NE',
  'NF',
  'NG',
  'NI',
  'NL',
  'NO',
  'NP',
  'NR',
  'NU',
  'NZ',
  'OM',
  'PA',
  'PE',
  'PF',
  'PG',
  'PH',
  'PK',
  'PL',
  'PM',
  'PN',
  'PR',
  'PS',
  'PT',
  'PW',
  'PY',
  'QA',
  'RE',
  'RO',
  'RS',
  'RU',
  'RW',
  'SA',
  'SB',
  'SC',
  'SD',
  'SE',
  'SG',
  'SH',
  'SI',
  'SJ',
  'SK',
  'SL',
  'SM',
  'SN',
  'SO',
  'SR',
  'SS',
  'ST',
  'SV',
  'SX',
  'SY',
  'SZ',
  'TC',
  'TD',
  'TF',
  'TG',
  'TH',
  'TJ',
  'TK',
  'TL',
  'TM',
  'TN',
  'TO',
  'TR',
  'TT',
  'TV',
  'TW',
  'TZ',
  'UA',
  'UG',
  'UM',
  'US',
  'UY',
  'UZ',
  'VA',
  'VC',
  'VE',
  'VG',
  'VI',
  'VN',
  'VU',
  'WF',
  'WS',
  'YE',
  'YT',
  'ZA',
  'ZM',
  'ZW']}]

```

In [27]:

```

1 sp.track('0KKkJNfGyhkQ5aFogxQAPU')['popularity']

```

Out[27]: 83

In [28]:

```

1 # Define function to extract the popularity from a given track. Pa
2 # don't break the function
3
4 def popularity_extract(track):
5     extract = pd.DataFrame()
6     pop_list = []
7     try:
8         result = sp.track(track)
9         track = result['id']
10        popularity = result['popularity']
11        pop_list.append(int(popularity))
12        extract['popularity'] = pop_list
13        extract['track_id'] = result['id']
14        return extract
15    except Exception:
16        pass

```

In [29]:

```

1 # Test function
2 popularity_extract('0KKkJNfGyhkQ5aFogxQAPU')

```

Out[29]:

	popularity	track_id
0	83	0KKkJNfGyhkQ5aFogxQAPU

In [30]:

```

1 # Test function if track_id doesn't exist
2 popularity_extract('example_error')

```

HTTP Error for GET to https://api.spotify.com/v1/tracks/example_error (https://api.spotify.com/v1/tracks/example_error) with Params: {'market': None} returned 400 due to invalid id

In [31]:

```

1 # Creating DataFrame for track popularity
2
3 if os.path.exists('Data/Spotipy Custom DataFrames/unique_track_popularity.pkl'):
4     pass
5 else:
6     unique_track_popularity = pd.DataFrame()
7
8     # Iterating over track data to apply function and append it to DataFrame
9     for track in range(0, len(unique_track_data)):
10         track_id = unique_track_data['track_uri'][track]
11         unique_track_popularity = unique_track_popularity.append(popularity_extract(track_id))
12
13     # Save time by exporting DataFrame to pickle
14     unique_track_popularity.to_pickle('Data/Spotipy Custom DataFrames/unique_track_popularity.pkl')

```

In [32]:

```

1 # Call pickle file and remove 0 popularity songs
2 unique_track_popularity = pd.read_pickle('Data/Spotipy Custom Data')
3 unique_track_popularity.drop(unique_track_popularity[unique_track_
4 unique_track_popularity

```

Out [32]:

	popularity	track_id
82	64	3kdMzXOcrDIdSWLdONHNK5
83	69	7aOor99o8NNLZYEIOXIBG1
86	60	45HAjqRWiNv6mMPw4NvZrU
87	56	5y1jgbDNgTfxoWXv3FhH2Q
90	67	2UZtl2HUyLRzqBjodvcUmY
...
8007	30	76ictxnZf8a4MAmaeNqvbU
8009	39	0MHJ3ObkdI3EN29A8nv6uz
8042	43	4POJUFV0qevJyeAX0j2mxR
8061	46	34ccBqL3xNaCzPxr0UqoEw
8069	33	6HcSRCF0R0DYRNY6vG0448

Next I am going to create an overall track_features file that combines all of our tables together.

In [33]:

```

1 # Send DataFrames to SQL
2 audio_features_df.to_sql('audio_features_df', con = conn, if_exists='replace')
3 unique_track_data.to_sql('unique_track_data', con = conn, if_exists='replace')
4 unwrapped_track_data.to_sql('unwrapped_track_data', con = conn, if_exists='replace')
5 spd_top100 = spd_top100.applymap(str)
6 spd_top100.to_sql('spd_top100', con = conn, if_exists='replace')
7 unique_genre_data.to_sql('unique_genre_data', con = conn, if_exists='replace')
8 unique_track_popularity.to_sql('unique_track_popularity', con = conn, if_exists='replace')

```

In [34]:

```
1 # Query and join all of our relevant DataFrames
2
3 track_features = """SELECT *
4                     FROM unwrapped_track_data as utd
5                     LEFT JOIN audio_features_df as af
6                         ON utd.track_uri = af.id
7                     LEFT JOIN unique_genre_data as ugd
8                         ON utd.artist_uri = ugd.artist_id
9                     LEFT JOIN unique_track_popularity as utp
10                        ON utd.track_uri = utp.track_id"""
```

In [35]:

```
1 # Convert it to a DataFrame
2
3 track_features = pd.read_sql(track_features, con = conn)
```

In [36]:

```

1  # Reformat the DataFrame, dropping columns, removing duplicates, k
2
3  track_features = track_features.drop(columns=['pos','index','type']
4  track_features = track_features.loc[:,~track_features.columns.duplic
5  track_features = track_features.drop(track_features.iloc[:, 21:62]
6  track_features = track_features.rename(columns={'pid':'playlist_id')
7  track_features = track_features.set_index('playlist_id')
8
9
10 scaler = MinMaxScaler()
11 minmax_columns = track_features.columns.tolist()[5:18] + ['popular
12 track_features[minmax_columns]=scaler.fit_transform(track_features
13 track_features.info()

```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Int64Index: 9323 entries, 180831 to 147046
```

```
Data columns (total 21 columns):
```

#	Column	Non-Null Count	Dtype
0	track_uri	9323 non-null	object
1	track_name	9323 non-null	object
2	artist_name	9323 non-null	object
3	artist_uri	9323 non-null	object
4	album_name	9323 non-null	object
5	danceability	9323 non-null	float64
6	energy	9323 non-null	float64
7	key	9323 non-null	float64
8	loudness	9323 non-null	float64
9	mode	9323 non-null	float64
10	speechiness	9323 non-null	float64
11	acousticness	9323 non-null	float64
12	instrumentalness	9323 non-null	float64
13	liveness	9323 non-null	float64
14	valence	9323 non-null	float64

In [37]: 1 track_features

Out[37]:

		track_uri	track_name	artist_name	artist_uri	all
playlist_id						
180831	527lbJxFcjjTix0ONdxDdS		Hearts as Strong as Horses	Applebloom	7ggsXdK95oJBkuZu1txVjC	(
180831	1mOnMHXxt2vGI0b804eQsy		Apples to the Core	Apple Jack	1r0v3fdCiqr9mYtvbCccT	(
180831	4GciJR91Tj8a7dLJ12WFvr		Ballad of the Crystal Ponies	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	(
180831	68dC0K4xglMF5Nhr44TevS		Find a Way	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	(
180831	0vhKFkvwgdPBrrNr9gUbVa		A True, True Friend	Twilight Sparkle	53CQUfjaBNRwV2nFro1nac	(
...
147046	3G0PWfSGIsUrI4EI8u46EX		All Or Nothing	Fake ID	4GwCpkFWajqx3KSm0MVh2a	
147046	5qpQhP4hyjCKC95oRtTqni		Borderline	Showoff	6lKhTkyp4EJ0ocidcwafs6	
147046	7wwXG5FOebfhVBotX4vTXo		Send Me An Angel	Thrice	3NChzMpu9exTINPiQ2DE	
147046	1lbfp4HOrAS0fB8eloEYko		Baby One More Time	Nicotine	0p3U0uLx2oSf0yn8i5XZki	
147046	6ktA174mwmCqcb9hfdL178		Heaven Is A Place On Earth	Student Rick	6AluDNoFgmeUTnOc7DYXIN	

9323 rows × 21 columns

```
In [38]: 1 # Send it to SQL
          2
          3 track_features.to_sql('track_features', con = conn, if_exists='rep
```

3.4 Create Playlist/Track Vectors for Modeling

Now that we have our general information tables, it is time to create our vectors for our recommendation system. We will make a Features, popularity, and genre vector for both our track data, and then we will use spark to aggregate the data and group it by playlist. All columns will be numeric.

3.4.1.1 Track Features Vector

```
In [39]: 1 # Combine track features with base list to get track_uri
          2
          3 playlist_features_by_track = """SELECT *
          4             FROM track_features as tf
          5             LEFT JOIN spd_top100 as spd
          6             ON tf.playlist_id = spd.pid"""
```

```
In [40]: 1 # Convert to DataFrame
          2 playlist_features_by_track = pd.read_sql(playlist_features_by_track
```

```
In [41]: 1 # View DataFrame to determine columns to drop
          2 playlist_features_by_track
```

```
Out[41]:
```

	playlist_id	track_uri	track_name	artist_name	artist_
0	180831	527lbJxFcjjTix0ONdxDdS	Hearts as Strong as Horses	Applebloom	7ggsXdK95oJBkuZu1tx\
1	180831	1mOnMHXxt2vGI0b804eQsy	Apples to the Core	Apple Jack	1r0v3fdCiqrr9mYtvbCc
2	180831	4GciJR91Tj8a7dLJ12WFvr	Ballad of the Crystal Ponies	Twilight Sparkle	53CQUfjaBNRwV2nFro1r

3	180831	68dC0K4xgIMF5Nhr44TevS	Find a Way	Twilight Sparkle	53CQUfjaBNRwV2nFro1r
4	180831	0vhKFkvwdPBrrNr9gUbVa	A True, True Friend	Twilight Sparkle	53CQUfjaBNRwV2nFro1r
...
9318	147046	3G0PWfSGIsUrl4EI8u46EX	All Or Nothing	Fake ID	4GwCpkFWajqx3KSm0MVr
9319	147046	5qpQhP4hyjCKC95oRtTqni	Borderline	Showoff	6lKhTkyp4EJ0ocidcwa
9320	147046	7wwXG5FOebfhVBotX4vTXo	Send Me An Angel	Thrice	3NChzMpu9exTINPiQUQ2
9321	147046	1lbfp4HOrAS0fB8eloEYko	Baby One More Time	Nicotine	0p3U0uLx2oSf0yn8i5X
9322	147046	6ktA174mwmCqcb9hfdL178	Heaven Is A Place On Earth	Student Rick	6AluDNoFgmeUTnOc7DY

9323 rows × 34 columns

In [42]:

```

1 # Drop irrelevant columns that we will not use
2 playlist_features_by_track = playlist_features_by_track.drop(column

```

In [43]:

```

1 # Create Master Vector for Audio Features (1/6)
2
3 master_track_audio_features = playlist_features_by_track
4 master_track_audio_features = master_track_audio_features.set_index('track_uri')
5 master_track_audio_features = master_track_audio_features.iloc[:,5:]
6 master_track_audio_features

```

Out [43]:

	danceability	energy	key	loudness	mode	speechiness
track_uri						
5271bJxFcjjTix0ONdxDdS	0.811861	0.827313	0.090909	0.905208	1.0	0.033832
1mOnMHXxt2vGI0b804eQsy	0.820041	0.815265	0.000000	0.951044	1.0	0.032911
4GciJR91Tj8a7dLJ12WFvr	0.646217	0.481938	0.000000	0.830322	1.0	0.033257
68dC0K4xgIMF5Nhr44TevS	0.541922	0.391578	0.909091	0.802070	1.0	0.039241
0vhKFkvwgdPBrrNr9gUbVa	0.693252	0.735949	0.909091	0.862855	1.0	0.040852
...
3G0PWfSGIsUri4EI8u46EX	0.368098	0.910644	0.000000	0.911956	1.0	0.067549
5qpQhP4hyjCKC95oRtTqni	0.561350	0.902612	0.363636	0.882834	1.0	0.053970
7wwXG5FOebfhVBotX4vTXo	0.267894	0.943776	0.545455	0.920809	1.0	0.059839
1lbfu4HOrAS0fB8eloEYko	0.601227	0.886548	0.818182	0.884164	0.0	0.038780

3.4.1.2 Playlist Features Vector

In [44]:

```

1 # Create Spark Session
2 spark = SparkSession.builder.master('local').getOrCreate()

```

In [45]:

```

1 # Create first spark DataFrame
2 spark_df = spark.createDataFrame(playlist_features_by_track)

```

In [46]:

```
1  # Assign aggregate type to the columns. We will be using mean
2
3  aggregate_features = {'danceability': 'mean',
4                        'energy' : 'mean',
5                        'key' : 'mean',
6                        'loudness' : 'mean',
7                        'mode' : 'mean',
8                        'speechiness' : 'mean',
9                        'acousticness' : 'mean',
10                       'instrumentalness' : 'mean',
11                       'liveness' : 'mean',
12                       'valence' : 'mean',
13                       'tempo' : 'mean',
14                       'time_signature' : 'mean'}
15
16 # Create aggregate DataFrame for playlist
17 playlist_features_aggregate_spark = spark_df.groupBy('playlist_id')
```

In [47]:

```

1  # Convert Spark DF to DataFrame in Pandas
2  playlist_features_aggregate = playlist_features_aggregate_spark.to
3  playlist_features_aggregate = playlist_features_aggregate.set_inde
4
5  # Scale our data
6  scaler.fit(playlist_features_aggregate)
7  playlist_features_aggregate[['avg(tempo)', 'avg(valence)', 'avg(er
8      'avg(speechiness)', 'avg(acousticness)', 'avg(key)',
9      'avg(time_signature)', 'avg(danceability)', 'avg(mode)',
10     'avg(loudness)', 'avg(instrumentalness)']] = scaler.fit_tra
11     'avg(speechiness)', 'avg(acousticness)', 'avg(key)',
12     'avg(time_signature)', 'avg(danceability)', 'avg(mode)',
13     'avg(loudness)', 'avg(instrumentalness)']]
14 # Final Master playlist audio features (2/6)
15 playlist_features_aggregate

```

Out[47]:

	avg(tempo)	avg(valence)	avg(energy)	avg(liveness)	avg(speechiness)	avg(acousticn
playlist_id						
161637	0.644151	0.498094	0.708264	0.178966	0.111970	0.155
102585	0.680578	0.569183	0.823455	0.117857	0.196751	0.145
108597	0.715631	0.151800	0.579258	0.219829	0.014027	0.195
108807	0.220464	0.778609	0.655196	0.069106	0.127409	0.415
154640	0.427408	0.487836	0.562870	0.229843	0.092631	0.335
...
184707	0.581173	0.480149	0.827352	0.141990	0.052378	0.085
101224	0.736112	0.771093	0.871079	0.199587	0.147407	0.085
182403	0.922561	0.714557	0.631823	0.054815	0.031743	0.275
17675	0.423771	0.962023	0.840775	0.170312	0.213944	0.195
147261	0.562398	0.484761	0.521355	0.089924	0.050583	0.255

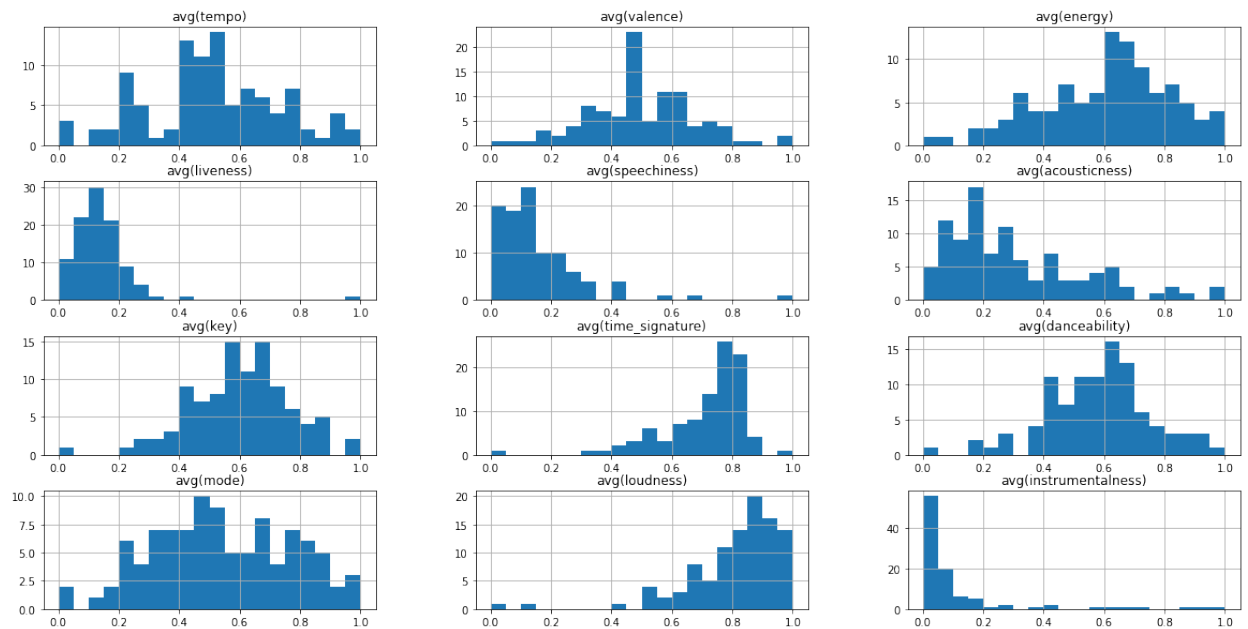
100 rows × 12 columns

In [48]:

```

1 # Create Visualization for each feature type
2
3 fig, axes = plt.subplots(len(playlist_features_aggregate.columns)/
4
5 i = 0
6 for triaxis in axes:
7     for axis in triaxis:
8         playlist_features_aggregate.hist(column = playlist_features_aggregate.columns[i],
9         i = i+1

```



3.4.2.1 Popularity Track Vector

In [49]:

```

1 # Create df for all popularity data that have no nulls, set index
2 playlist_popularity_no_nulls = playlist_features_by_track.dropna(a
3 master_track_popularity_data = playlist_popularity_no_nulls[['trac
4 master_track_popularity_data = master_track_popularity_data.set_in

```

```
In [50]: 1 # Master Popularity data for Tracks (3/6)
          2 master_track_popularity_data
```

```
Out[50]:
```

	popularity
track_uri	
3kdMzXOcrDIdSWLdONHNK5	0.724138
7aOor99o8NNLZYEIOXIBG1	0.781609
45HAjqRWiNv6mMPw4NvZrU	0.678161
5y1jgbDNgtfxoWXv3FhH2Q	0.632184
2UZtl2HUyLRzqBjodvcUmY	0.758621
...	...
76ictxnZf8a4MAmaeNqvbU	0.333333
0MHJ3Obkdl3EN29A8nv6uz	0.436782
4POJUFV0qevJyeAX0j2mxR	0.482759
34ccBqL3xNaCzPxr0UqoEw	0.517241

3.4.2.2 Popularity Playlist Vector

```
In [51]: 1 # Create second Spark df for popularity dataset
          2 spark_df2 = spark.createDataFrame(playlist_popularity_no_nulls)
```

```
In [52]: 1 # Group our popularity data by playlist
          2 playlist_popularity_aggregate_spark = spark_df2.groupBy('playlist_
```

In [53]:

```
1 # Send our spark df to Pandas and set index to playlist id (4/6)
2 playlist_popularity_aggregate = playlist_popularity_aggregate_spar
3 playlist_popularity_aggregate = playlist_popularity_aggregate.set_
4 playlist_popularity_aggregate
```

Out [53]:

avg(popularity)	
playlist_id	
161637	0.326230
102585	0.665113
108597	0.471436
108807	0.445680
154640	0.218391
...	...
184707	0.451393
101224	0.556975
182403	0.678161
17675	0.616240
147261	0.527203

99 rows × 1 columns

3.4.3.1 Genre Track Vector

In [54]:

```

1 # Prepare dataframe for genre for tracks
2
3 playlist_features_by_track[playlist_features_by_track['genre_1'].r

```

Out [54]:

	playlist_id	track_uri	track_name	artist_name	artist
0	180831	527IbJxFcjjTix0ONdxDdS	Hearts as Strong as Horses	Applebloom	7ggsXdK95oJBkuZu1tx\
2	180831	4GciJR91Tj8a7dLJ12WFvr	Ballad of the Crystal Ponies	Twilight Sparkle	53CQUfjaBNRwV2nFro1r
3	180831	68dC0K4xglMF5Nhr44TevS	Find a Way	Twilight Sparkle	53CQUfjaBNRwV2nFro1r
4	180831	0vhKFkvwgdPBrrNr9gUbVa	A True, True Friend	Twilight Sparkle	53CQUfjaBNRwV2nFro1r
6	180831	1XBQVELMITpvMal5Pn1UpO	Babs Seed	Applebloom	7ggsXdK95oJBkuZu1tx\
...
9310	147046	2jbupzHNwz0VgORg1uJb3D	Like A Prayer	Rufio	0HjoylTAvSVktTCjXUa4
9311	147046	6HcSRCF0R0DYRNY6vG0448	Bye, Bye, Bye	Further Seems Forever	1Enp9WKfk0aI9CFi2YGE
9313	147046	57uO0ogaSWb7t20CY85CfD	I'm Like A Bird	Element 101	6pndtpE63q5pHaGhDko
9315	147046	5Al3VJyLLmB8hEloCAVCIm	I'm Real	The Starting Line	3E3xrZtBU5ORqcmX78v5
9320	147046	7wwXG5FOebfhVBotX4vTXo	Send Me An Angel	Thrice	3NChzMpu9exTINPiqUQ2

7207 rows × 21 columns


```
In [55]: 1 # Select genre data and set index to track_uri
2 genre_data_no_nulls = playlist_features_by_track[playlist_features
3
4 # OneHotEncode our data so it can be read by our model
5 ohe = OneHotEncoder()
6 X = ohe.fit_transform(genre_data_no_nulls['genre_1'].values.reshape
7 y = ohe.get_feature_names(['genre_1'])
8 master_track_genre_data = pd.DataFrame(X, columns = y)
9 master_track_genre_data.index = genre_data_no_nulls.index
10
11 # Final master genre data by track (5/6)
12 master_track_genre_data
```

Out[55]:

	genre_1_21st century classical	genre_1_a cappella	genre_1_abstract beats	genre_1_abstract hip hop	genre_1_rock
track_uri					
527lbJxFcjjTix0ONdxDdS	0.0	0.0	0.0	0.0	0.0
4GciJR91Tj8a7dLJ12WFvr	0.0	0.0	0.0	0.0	0.0
68dC0K4xglMF5Nhr44TevS	0.0	0.0	0.0	0.0	0.0
0vhKFkvwgdPBrrNr9gUbVa	0.0	0.0	0.0	0.0	0.0
1XBQVELMITpvMal5Pn1UpO	0.0	0.0	0.0	0.0	0.0
...
2jbupzHNwz0VgORg1uJb3D	0.0	0.0	0.0	0.0	0.0
6HcSRCF0R0DYRNY6vG0448	0.0	0.0	0.0	0.0	0.0
57uO0ogaSWb7t20CY85CfD	0.0	0.0	0.0	0.0	0.0
5AI3VJyLLmB8hElOCAVCIm	0.0	0.0	0.0	0.0	0.0
7wwXG5FOebfhVBotX4vTXo	0.0	0.0	0.0	0.0	0.0

7207 rows × 500 columns

3.4.3.2 Genre Playlist Vector

```
In [56]: 1 # Prep our track genre data
2 playlist_genre_prep = master_track_genre_data.join(genre_data_no_r
```

```
In [57]: 1 # Create our third spark dataframe
          2 spark_df3 = spark.createDataFrame(playlist_genre_prep)
```

```
In [58]: 1 # Aggregate our genre data by sum
          2
          3 genre_keys = playlist_genre_prep.columns.to_list()
          4 genre_values = ['sum'] * 500
          5 aggregate = dict(zip(genre_keys, genre_values))
          6
          7 playlist_genre_aggregate_spark = spark_df3.groupBy('playlist_id').
```

```
In [59]: 1 # MinMax scale our sums to weight our most common genres on a scale
          2 playlist_genre_aggregate = playlist_genre_aggregate_spark.toPandas()
          3 playlist_genre_aggregate = playlist_genre_aggregate.set_index('playlist_id')
          4 playlist_pie_chart = playlist_genre_aggregate.copy()
          5 minmax_columns2 = playlist_genre_aggregate.columns.tolist()
          6 playlist_genre_aggregate[minmax_columns2] = scaler.fit_transform(playlist_pie_chart[minmax_columns2])
          7
          8 # Final Master playlist data for genre (6/6)
          9 playlist_genre_aggregate.index.get_loc(102585)
```

Out[59]: 1

```
In [60]: 1 playlist_genre_aggregate.iloc[1:2]
```

Out[60]:

	sum(genre_1_la pop)	sum(genre_1_lovers rock)	sum(genre_1_lo- fi beats)	sum(genre_1_french shoegaze)	sum(genre_1_classic)
playlist_id					
102585	0.0	0.0	0.0	0.0	0.0

1 rows × 500 columns

```
In [61]: 1 # End Spark Session
          2 SparkSession.stop(spark)
```

3.5 Final Vectors

Track Vector

In [62]:

```

1 # Send master track vectors to SQL
2
3 master_track_audio_features.to_sql('master_track_audio_features',
4 master_track_popularity_data.to_sql('master_track_popularity_data',
5 master_track_genre_data.to_sql('master_track_genre_data', con = conn)

```

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.
method=method,

In [63]:

```

1 # Join track vectors using SQL
2
3 final_track_join = """SELECT *
4 FROM master_track_audio_features as taf
5 LEFT JOIN master_track_popularity_data as mtp
6 ON taf.track_uri = mtp.track_uri
7 LEFT JOIN master_track_genre_data as mtg
8 ON taf.track_uri = mtg.track_uri"""
9 # Read into Pandas and set index, impute values for most frequent
10 final_track_vector = pd.read_sql(final_track_join, con = conn)
11 final_track_vector = final_track_vector.set_index('track_uri')
12 impute = SimpleImputer(missing_values=np.nan, strategy='most_frequent')
13 final_track_vector[final_track_vector.columns.to_list()] = impute.fit_transform(
14 final_track_vector[final_track_vector.columns.to_list()])
15 # Fix index (formatting was strange)
16 new_index = []
17
18 for item in range(0, len(final_track_vector.index)):
19     new_index.append(final_track_vector.index[item][0])
20
21 final_track_vector['track_uri_fix'] = new_index
22 final_track_vector = final_track_vector.set_index('track_uri_fix')
23 final_track_vector.index = final_track_vector.index.rename(name = 'track_uri')
24
25 # Drop categorical features that shouldn't be used in modeling
26 final_track_vector = final_track_vector.drop(['key', 'mode', 'time_signature', 'tempo', 'duration', 'explicit_lyrics', 'explicit_content_lyrics', 'explicit_content_lyrics_text', 'explicit_content_lyrics_text_text', 'explicit_content_lyrics_text_text_text'])
27
28 # Final track vector ready for modelling shape is (8081 x 513)
29 final_track_vector = final_track_vector.drop_duplicates()
30 final_track_vector

```

Out[63]:

	danceability	energy	loudness	speechiness	acousticness	instrumentalness
5271bJxFcjjTix0ONdxDdS	0.811861	0.827313	0.905208	0.033832	0.242971	

1mOnMHXxt2vGI0b804eQsy	0.820041	0.815265	0.951044	0.032911	0.232930
4GciJR91Tj8a7dLJ12WFvr	0.646217	0.481938	0.830322	0.033257	0.420682
68dC0K4xglMF5Nhr44TevS	0.541922	0.391578	0.802070	0.039241	0.884538
0vhKFkvwgdPBrrNr9gUbVa	0.693252	0.735949	0.862855	0.040852	0.306224
...
3G0PWfSGIsUri4EI8u46EX	0.368098	0.910644	0.911956	0.067549	0.017769
5qpQhP4hyjCKC95oRtTqni	0.561350	0.902612	0.882834	0.053970	0.014055
7wwXG5FOebfhVBotX4vTXo	0.267894	0.943776	0.920809	0.059839	0.000006
1lbfp4HOrAS0fB8eloEYko	0.601227	0.886548	0.884164	0.038780	0.001866
6ktA174mwmCqcb9hfdL178	0.627812	0.794181	0.895920	0.035328	0.000684

8081 rows × 510 columns

Playlist Vector

```
In [64]: 1 # Send playlist vectors to SQL
2 playlist_features_aggregate.to_sql('playlist_features_aggregate',
3 playlist_popularity_aggregate.to_sql('playlist_popularity_aggregate',
4 playlist_genre_aggregate.to_sql('playlist_genre_aggregate', con =
```

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.
method=method,

```
In [65]: 1 # Join playlist vectors using SQL
2 final_playlist_join = """SELECT *
3
4 FROM playlist_features_aggregate as pfa
5 LEFT JOIN playlist_popularity_aggregate as ppa
6 ON pfa.playlist_id = ppa.playlist_id
7 LEFT JOIN playlist_genre_aggregate as pga
8 ON pfa.playlist_id = pga.playlist_id"""
9
10 # Bring vector into pandas and set index/reformat
11 final_playlist_vector = pd.read_sql(final_playlist_join, con = con)
12 final_playlist_vector = final_playlist_vector.T.groupby(level=0).first()
13 final_playlist_vector = final_playlist_vector.set_index('playlist_id')
14 final_playlist_vector.index = final_playlist_vector.index.astype('int')
15
16 # Impute missing values by mean
17 impute = SimpleImputer(missing_values=np.nan, strategy='mean')
```

```

17 final_playlist_vector[final_playlist_vector.columns.to_list()] = i
18 final_playlist_vector.columns = final_playlist_vector.columns.strl
19
20 # Drop categorical features that shouldn't be used in modeling
21 final_playlist_vector = final_playlist_vector.drop(['key', 'mode', '
22
23 # Final playlist vector ready for modelling shape is (100 x 513)
24 final_playlist_vector

```

Out [65]:

	acousticness	danceability	energy	instrumentalness	liveness	loudness	popularity
playlist_id							
161637	0.152660	0.447993	0.708264	0.018751	0.178966	0.878009	0.3262
102585	0.143031	0.588997	0.823455	0.024649	0.117857	0.941022	0.6651
108597	0.199158	0.272991	0.579258	0.000178	0.219829	0.859859	0.4714
108807	0.413473	0.676067	0.655196	0.018121	0.069106	0.793770	0.4456
154640	0.335231	0.492703	0.562870	0.021785	0.229843	0.785069	0.2183
...
184707	0.084968	0.517632	0.827352	0.076839	0.141990	0.925287	0.4513
101224	0.085839	0.633265	0.871079	0.021512	0.199587	0.971021	0.5569
182403	0.271278	0.675073	0.631823	0.044869	0.054815	0.755928	0.6781
17675	0.190159	0.900861	0.840775	0.002419	0.170312	0.946618	0.6162
147261	0.250436	0.504570	0.521355	0.083943	0.089924	0.768521	0.5272

100 rows × 510 columns

Section 4: Modeling

Content-Based Filtering using Cosine Similarity

Now that we have our vectors with our audio features, popularity, and genre data, I will use cosine similarity to build our recommendation system. Here is how it will work:

- Step 1: Our model will take a look at a selected playlist id and its vector to analyze the overall features of the songs within.
- Step 2: Our model will then compare the playlist's vector to our unique track data vector (audio features, popularity, and genre) to generate a list of x number of songs (user-inputted) that are closely comparable to the playlist vectors, using cosine similarity. The greater the cosine similarity, the more similar the vectors are related.
- Step 3: We will then analyze the features of the top song generated for reasonableness.

In [66]:

```

1  # Recommenders
2
3  def song_recommender(playlist_id,num_recs):
4      '''
5          This function serves as the recommendation system behind our
6          Inputs pulled from spotify_recommendation function below:
7
8          playlist_id: spotify playlist id from top 100 popular playlis
9          num_recs: number of recomendations desired
10         '''
11
12         # Remove songs from track vector that appear in the selected
13         playlist_tracks_uri_key = pd.DataFrame(playlist_info(playlist_id))
14         master_track_key = pd.DataFrame(final_track_vector.reset_index())
15         unique_track_key = pd.DataFrame(final_track_vector.reset_index())
16         unique_track_key = unique_track_key.append(playlist_tracks_uri_key)
17         new_final_track_vector_no_playlist_tracks = pd.merge(unique_track_key,
18         new_final_track_vector_no_playlist_tracks = new_final_track_vector_no_playlist_tracks
19
20         # Create consine similarity model between unique tracks and p
21         row_number = final_playlist_vector.index.get_loc(playlist_id)
22         playlist_selection_vector = final_playlist_vector.iloc[row_number]
23         playlist_feature_cos = pd.DataFrame(cos(new_final_track_vector, playlist_selection_vector))
24         playlist_feature_cos.index = new_final_track_vector_no_playlist_tracks.index
25
26         # Reformat and sort based on strongest cosine similarity value
27         playlist_feature_cos = playlist_feature_cos.rename(columns = {'cosine_similarity': 'similarity'})
28         playlist_feature_cos = playlist_feature_cos.sort_values('similarity', ascending = False)
29
30         # Use SQL to bring in track name, artist name and album name
31         playlist_feature_cos.to_sql('playlist_feature_cos', con = con, if_exists = 'replace')
32         clean_recommendations = """SELECT *
33                                 FROM playlist_feature_cos as pfc
34                                 LEFT JOIN unique_track_data as utd

```

```

34         LEFT JOIN unique_track_data as utd
35             ON pfc.track_uri = utd.track_uri"""
36
37     # Send to Pandas and reformat
38     clean_recommendations_df = pd.read_sql(clean_recommendations,
39     clean_recommendations_df = clean_recommendations_df.drop(colu
40     clean_recommendations_df = clean_recommendations_df.iloc[:,2:
41     return clean_recommendations_df
42
43 def playlist_info(playlist_id):
44     playlist_info = unwrapped_track_data[unwrapped_track_data['pi
45     playlist_info = playlist_info[['track_uri','track_name','arti
46     dfi.export(playlist_info.head(10), 'Visualizations/playlistsn
47     return playlist_info
48
49 def spotify_recommendation(playlist_id,num_recs):
50     '''
51     This function provides us with our general interface for our
52
53     Inputs:
54     playlist_id: spotify playlist id from top 100 popular playlis
55     num_recs: number of recomendations desired
56     '''
57
58     dfi.export(song_recommender(playlist_id,num_recs), 'Visualiza
59     return print('\nThank you for inputting your playlist. Please
60
61 # Validation
62 def check_duplicates(playlist_id,num_recs):
63     '''
64     This function is used to check for any duplicate recommendati
65     recommendations for tracks that aren't already in the playlis
66     playlist and our recommendations, the returned set will not b
67     the set will be empty.
68     '''
69     rec_check = song_recommender(playlist_id,num_recs)['track_uri
70     playlist_check = playlist_info(playlist_id)['track_uri'].valu
71     if len(set(rec_check) & set(playlist_check))==0:
72         return print("No duplicates being recommended!")
73     else:
74         return print("Duplicate recommendations found!")
75
76 # Comparisons
77 def song_vs_playlist_comparison(playlist_id, num_recs):
78
79     # This function gives us the sample variance between our sele
80     rnum = final_playlist_vector.index.get_loc(playlist_id)
81     playlist_snapshot = final_playlist_vector.iloc[rnum:rnum+1].r
82     rec_tracks = pd.DataFrame()
83     for item in range(0,num_recs):
84         rec_tracks = rec_tracks.append(song_recommender(playlist

```

```

85     rec_tracks = rec_tracks.drop(columns=['track_name', 'artist_name'])
86     final_track_vector.to_sql('final_track_vector', con = conn, if_exists='replace')
87     rec_tracks.to_sql('rec_tracks', con = conn, if_exists='replace')
88     rec_tracks_features = """SELECT *
89                             FROM rec_tracks as rt
90                             LEFT JOIN final_track_vector as ftv
91                             ON rt.track_uri = ftv.track_uri"""
92     rec_tracks_features_df = pd.read_sql(rec_tracks_features, con = conn)
93     rec_tracks_features_df = rec_tracks_features_df.set_index('track_uri')
94     rec_tracks_features_df['mean'] = rec_tracks_features_df['mean'].mean()
95
96
97     # Manually calculated variance
98     abs_variance_df = playlist_snapshot - rec_tracks_features_df
99     abs_variance_df['sum difference'] = abs_variance_df.sum(axis=1)
100    abs_variance_df = abs_variance_df.drop(columns='index', axis=1)
101    return abs_variance_df
102
103    def song_vs_playlist_visual_comparison(playlist_id, num_recs):
104
105        '''This function is for creating a cluster barchart that will
106        and playlist '''
107
108        # Get track data formatted for bar graph
109        variance_bar = song_vs_playlist_comparison(playlist_id, num_recs)
110        variance_bar = variance_bar.loc[:, (variance_bar != 0).any(axis=1)]
111        variance_bar = variance_bar.drop(columns='sum difference')
112        variance_bar = variance_bar.drop(columns=variance_bar.iloc[:, 0])
113        size_track = len(variance_bar.columns)
114        variance_bar = variance_bar.to_dict('records')[0]
115        variance_bar = dict(sorted(variance_bar.items(), key=lambda i: i[1]))
116        variance_bar_values = list(variance_bar.values())
117        variance_bar_columns = {k.title() for k in variance_bar.keys()}
118        variance_bar_columns = list(variance_bar_columns)
119
120        x = np.arange(len(variance_bar_columns)) # the label locations
121        width = 0.35 # the width of the bars
122
123        # Create plot for clustered bar graph
124        fig, ax = plt.subplots(figsize=(30, 15), facecolor='white')
125        ax.bar(range(len(variance_bar_columns)), height=variance_bar_values)
126        ax.set_xlabel('Feature', fontsize=30)
127        ax.set_ylabel('Difference (min: -1; max: 1)', fontsize=20)
128        ax.set_title('Features Comparison: Playlist vs. Recommended Tracks',
129                    fontsize=20, rotation=45)
130        ax.legend(loc="lower right")
131        ax.legend(fontsize=20)
132        plt.ylim(-1, 1)
133        plt.axhline(y=0, color='black')
134        ax.set_xticks(x)

```



```

135     sns.set_context('poster')
136     plt.savefig('Visualizations/songvsplaylist'+str(playlist_id)+
137     return plt.show()
138
139 # Genre population breakdown by playlist
140 def playlist_genre_breakdown(playlist_id):
141     playlist_pie_chart_df = playlist_pie_chart
142     row_number = playlist_pie_chart_df.index.get_loc(playlist_id)
143     playlist_pie_chart_df = playlist_pie_chart_df.iloc[row_number]
144     playlist_pie_chart_df = playlist_pie_chart_df.loc[:, (playlist_pie_chart_df.columns != 'genre')]
145     playlist_pie_chart_df = playlist_pie_chart_df.loc[:, ~playlist_pie_chart_df.columns.str.startswith('genre')]
146     rows = playlist_pie_chart_df.shape[1]
147     playlist_pie_chart_values = np.reshape(playlist_pie_chart_df.values, (rows,))
148     playlist_pie_chart_labels = playlist_pie_chart_df.columns.tolist()
149     playlist_pie_chart_labels = [label.title() for label in playlist_pie_chart_labels]
150     playlist_pie_chart_labels = [label[12:-1] for label in playlist_pie_chart_labels]
151     fig, ax = plt.subplots(figsize=(20,10), facecolor='white')
152     sns.set_context('poster')
153     ax.pie(playlist_pie_chart_values, labels=playlist_pie_chart_labels)
154     ax.set_title('Song Genre in Playlist (Most Common)', fontsize=16)
155     plt.savefig('Visualizations/PlaylistGenreBreakdown'+str(playlist_id)+'.png')
156     return fig, ax

```

How can I review results when there aren't many metrics available for content-based filtering? First, I will review the overall features for our selected playlist and compare them to the top track selected by our recommendation system for reasonableness. Second, I will plot the feature scores against each other to visualize the magnitude of difference between the average playlist scores and the top recommended song. Lastly, I will listen to a selection of 3 songs on each playlist, and then I will listen to and analyze the tracks provided by our recommendation system to subjectively determine whether I think they could belong on the playlist.

Playlist #102585 Reasonableness Test

In [67]:

```

1 # Call function for recommender
2
3 spotify_recommendation(102585,5)

```

Thank you for inputting your playlist. Please see playlist tracks below:

	track_uri	track_name	artist_name	album_name
pos				

0	67WTwafOMgegV6ABnBQxcE	Some Nights	fun.	Some Nights
1	6Ep6BzlOB9tz3P4sWqiiAB	Radioactive	Imagine Dragons	Night Visions
2	4wCmqSrbyCgxEXROQE6vtV	Somebody That I Used To Know	Gotye	Making Mirrors
3	5j9iuo3tMmQlfnEEQOOjxh	Best Day Of My Life	American Authors	Oh, What A Life
4	6cpk00i5TxCqSeqNi2Hule	One More Night	Maroon 5	Overexposed Track By Track
...
137	1YGvv0iH1TEjMrp0oRPB5a	Louder Than Your Love	Andy Black	The Shadow Side
138	1DtLI0k0zMrFqcW0pquxpf	Broken Pieces	Andy Black	The Shadow Side
139	4WTesnTwFArtC6fhXuX31	The Void	Andy Black	The Shadow Side
140	0NWQTyapmz4GuDTSN9xTB7	Candyman	Zedd	Candyman
141	2qPUnoasNe4Ep43emVXEig	Billionaire (feat. Bruno Mars)	Travie McCoy	Lazarus

142 rows × 4 columns

Here are 5 tracks that might fit this playlist:

	track_uri	track_name	artist_name	album_name
0	3s5ogvexUgA6XjNj37zpnP	Piledriver waltz	Alex Turner	Submarine - Original Songs From The Film By Al...
1	57LTVY8oPSGy2a7c8SzJpD	What Goes Around...	Alesana	Punk Goes Pop, Vol. 2
2	0fiB5cRA2IngpF06X1Ou4u	Bang Bang You're Dead	Dirty Pretty Things	Waterloo To Anywhere
3	5lDriBxJd22lhOH9zTcFrV	Dirty Little Secret	The All-American Rejects	Move Along
4	0X0Lz7LwpilWcdGqVWaxXD	Mess Around	Cage The Elephant	Tell Me I'm Pretty

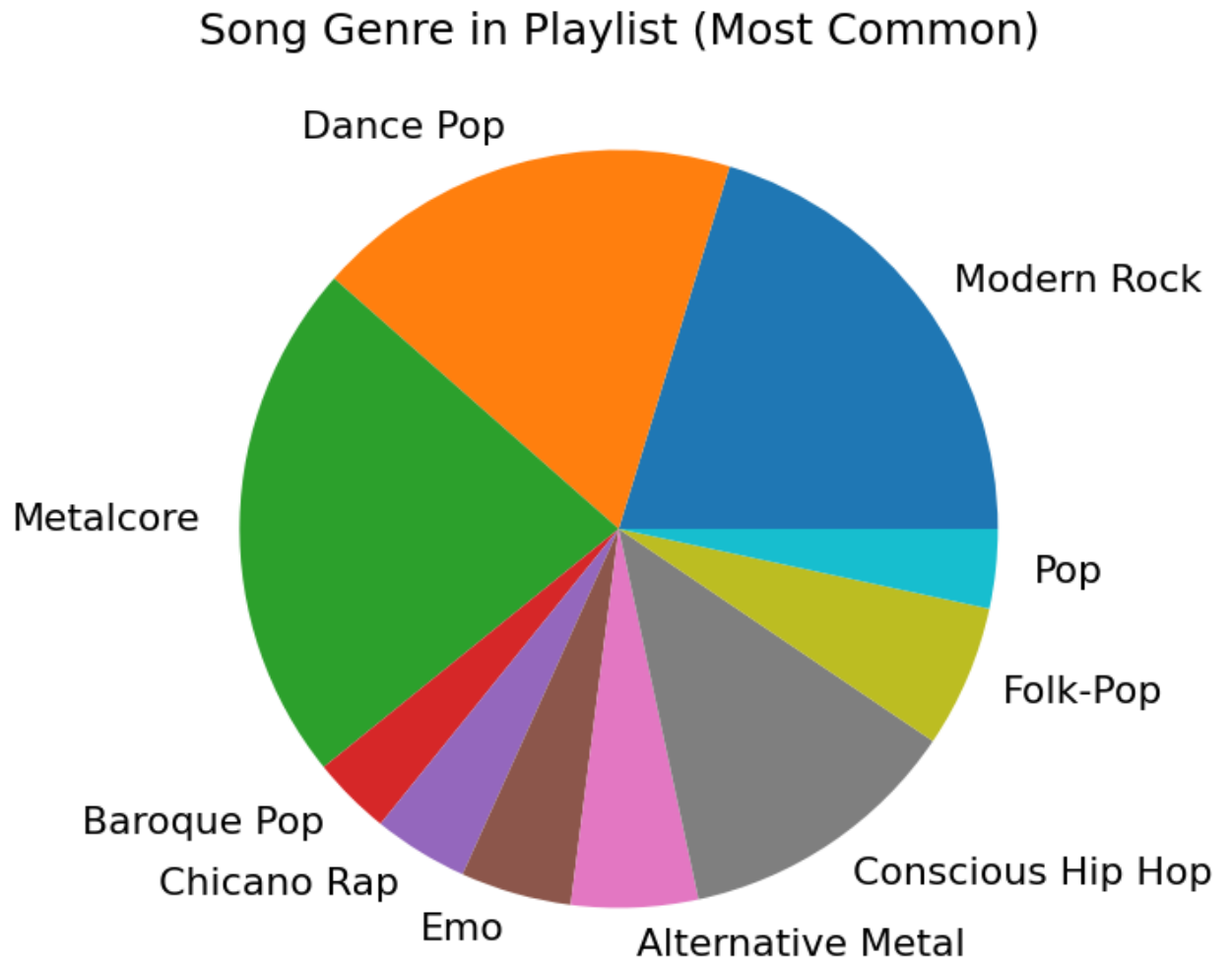
Out[67]: (None, None, None, None)

In [68]: `1 check_duplicates(102585,5)`

No duplicates being recommended!

```
In [69]: 1 playlist_genre_breakdown(102585)
```

```
Out[69]: (<Figure size 1440x720 with 1 Axes>,  
  <AxesSubplot:title={'center':'Song Genre in Playlist (Most Common)'}>)
```



In [70]:

```

1 # Call function for variance
2
3 song_vs_playlist_comparison(102585,5)

```

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.
method=method,

Out[70]:

	acousticness	danceability	energy	genre_1_21st century classical	genre_1_a cappella	genre_1_abstract beats	genre_1_a
0	0.115059	0.162003	-0.050844	0.0	0.125	0.0	

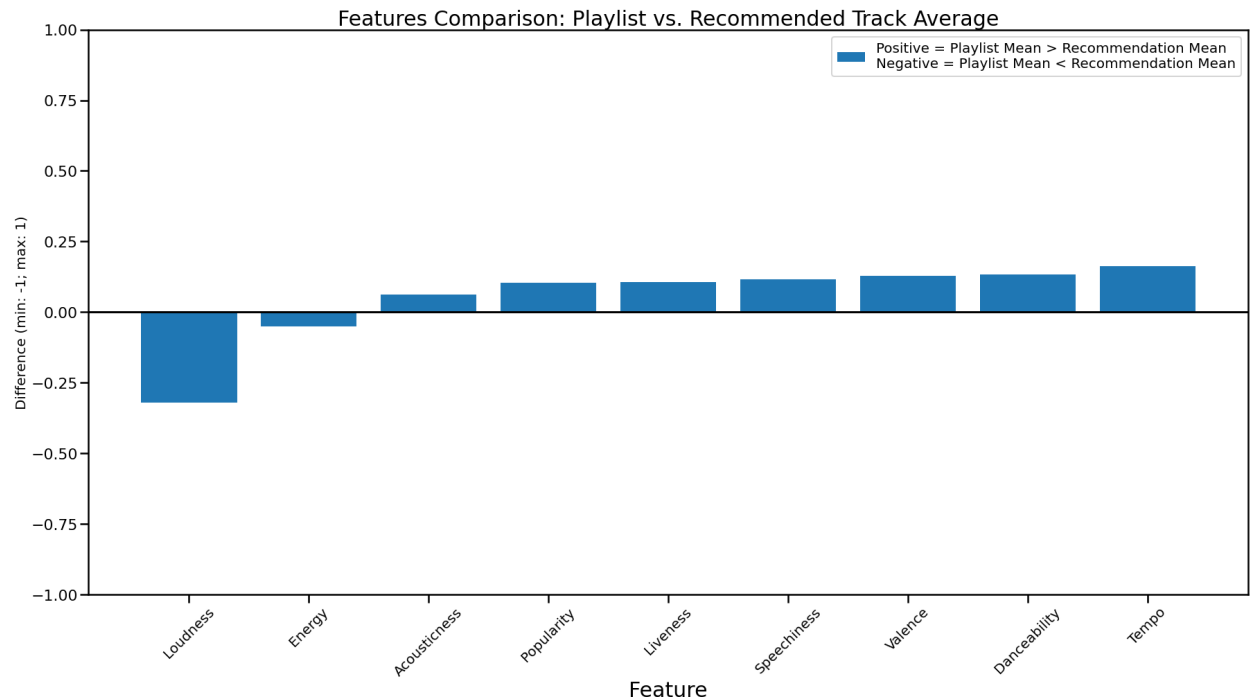
1 rows × 511 columns

```
In [71]: 1 # Call function for graphical comparison
        2 song_vs_playlist_visual_comparison(102585,5)
```

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.

method=method,

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/ipykernel_launcher.py:129: UserWarning: FixedFormatter should only be used together with FixedLocator



As we can see above, the aggregate difference between our recommended songs and our playlist average is 14, which is driven mainly by the genre data. This metric is only useful when compared to other aggregate differences, so I will analyze when I review another playlist. While it can be argued that our genre data is so specific that it is dominating our model, I personally think the specificity of the genre helps us narrow our recommendations down to suggest more accurate songs.

Looking at our visualization, for the most part it appears that our features for our recommendations are comparable to the playlist, with Tempo driving the main difference.

Now for the sound test. I will document my subjective findings below:

Some Nights, Radioactive, and Somebody That I Used to Know were all hits around 2013, and I personally have listened to these songs already. My description of Radioactive is an arena rock, bass thumping anthem that brings a high level of intensity. Some Nights and Somebody That I Used to Know are more chill and have a very nostalgic sound to them.

For our recommendations, What Goes Around..., Bang Bang You're Dead, Dirty Little Secret and Mess Around are all alternative/rock sounding tracks, whereas Piledriver Waltz is definitely more on the chill side.

Overall, based on our genre analysis and feature comparison, I would deem that these 5 recommended songs pass the reasonableness test.

Playlist #765 Reasonableness Test

In [72]: `spotify_recommendation(765,10)`

Thank you for inputting your playlist. Please see playlist tracks below:

	track_uri	track_name	artist_name	album_name
pos				
0	03xWMkKEbeO4SnylA53ipj	When Will My Life Begin - From "Tangled"/Sound...	Mandy Moore	Tangled
1	1IOSxJNCLvWm2bYaTcTSmK	Mother Knows Best - From "Tangled"/Soundtrack ...	Donna Murphy	Tangled
2	0TCt7OFRdD8PQ6vTRQxNgQ	I've Got a Dream - From "Tangled"/Soundtrack V...	Mandy Moore	Tangled
3	6klpXs2uAjagnZMFkt4qkl	I See the Light - From "Tangled"/Soundtrack Ve...	Mandy Moore	Tangled

4	75VVIB2x1h6BfxD2PqOO57	Healing Incantation - From "Tangled"/Soundtrac...	Mandy Moore	Tangled
...
76	6FHUBs8P5qcj7C2QHdEq	Tulou Tagaloe	Olivia Foa'i	Moana
77	3ZJnc1eGicPxRitBoC7eWZ	An Innocent Warrior	Vai Mahina	Moana
78	2bwSCluNtVrQPvddCi8sOW	Where You Are	Christopher Jackson	Moana
79	3C4WmF4klgfmb6GzW8DEdX	We Know The Way - From "Moana"	Opetiaia Foa'i	We Know The Way
80	2wCRJwiL1WSrW0Dwfc07Nj	Know Who You Are	Auli'i Cravalho	Moana

81 rows × 4 columns

Here are 10 tracks that might fit this playlist:

	track_uri	track_name	artist_name	album_name
0	7wjiMdiSgsL9Vkrwb10Num	On My Way	Phil Collins	Brother Bear
1	2stkLJ0JNcXklRDNF3ld6c	You've Got A Friend In Me - From "Toy Story"/ ...	Randy Newman	Toy Story
2	0odIT9B9BvOCnXfS0e4IB5	Bette Davis Eyes	Kim Carnes	Mistaken Identity
3	2s6wCS3vDZFPY9NOTIPXJZ	Take Me Home - 2016 Remastered	Phil Collins	No Jacket Required (Remastered)
4	4LCywJRtsZfr5qEJ2ckrTw	Tornado	Lea Michele	Places
5	27mSYFYtli1K3eXMWMmZVZ	The Rendezvous (feat. Madi Diaz)	Rob Cantor	Not a Trampoline
6	0qxtQ8rf3W1nld3D2r0xH4	I Just Can't Wait to Be King - From "The Lion ...	Jason Weaver	The Lion King
7	7G061Oqw7NXFr1NDTpXoI4	Happy Working Song - From "Enchanted"/Soundtra...	Amy Adams	Enchanted
8	6mDxu0xwhv5tn1oMTNUypu	Something There	Robby Benson	Beauty and the Beast
9	0QKHM0vKGZcgRpryeqtYkG	The Chipmunk Song	Alvin & The Chipmunks	Alvin & The Chipmunks / OST

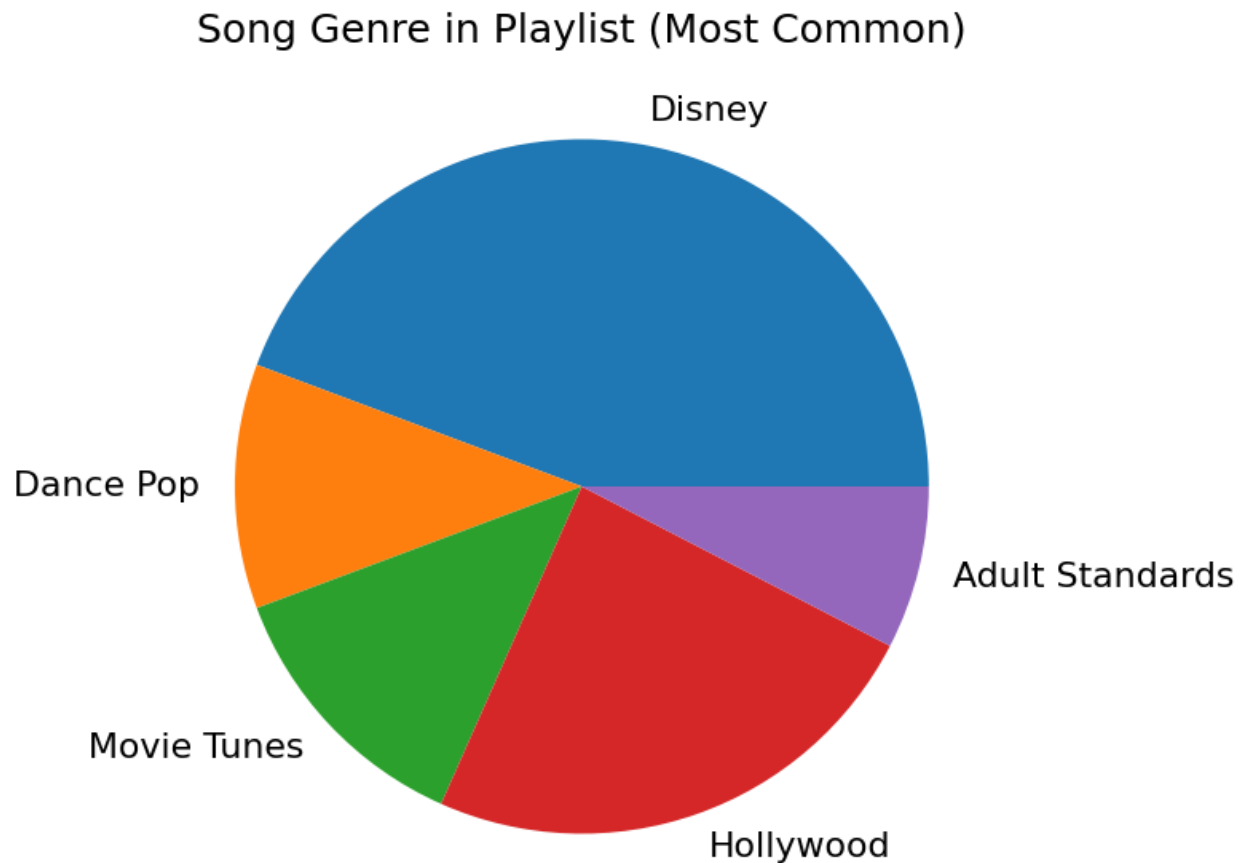
Out[72]: (None, None, None, None)

```
In [73]: 1 check_duplicates(765,10)
```

No duplicates being recommended!

```
In [74]: 1 playlist_genre_breakdown(765)
```

```
Out[74]: (<Figure size 1440x720 with 1 Axes>,  
<AxesSubplot:title={'center':'Song Genre in Playlist (Most Common)'}>)
```



In [75]: `1 song_vs_playlist_comparison(765,10)`

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.

method=method,

Out[75]:

	acousticness	danceability	energy	genre_1_21st century classical	genre_1_a cappella	genre_1_abstract beats	genre_1_ab hi
0	0.326452	-0.284049	-0.29488	0.0	0.0	0.0	

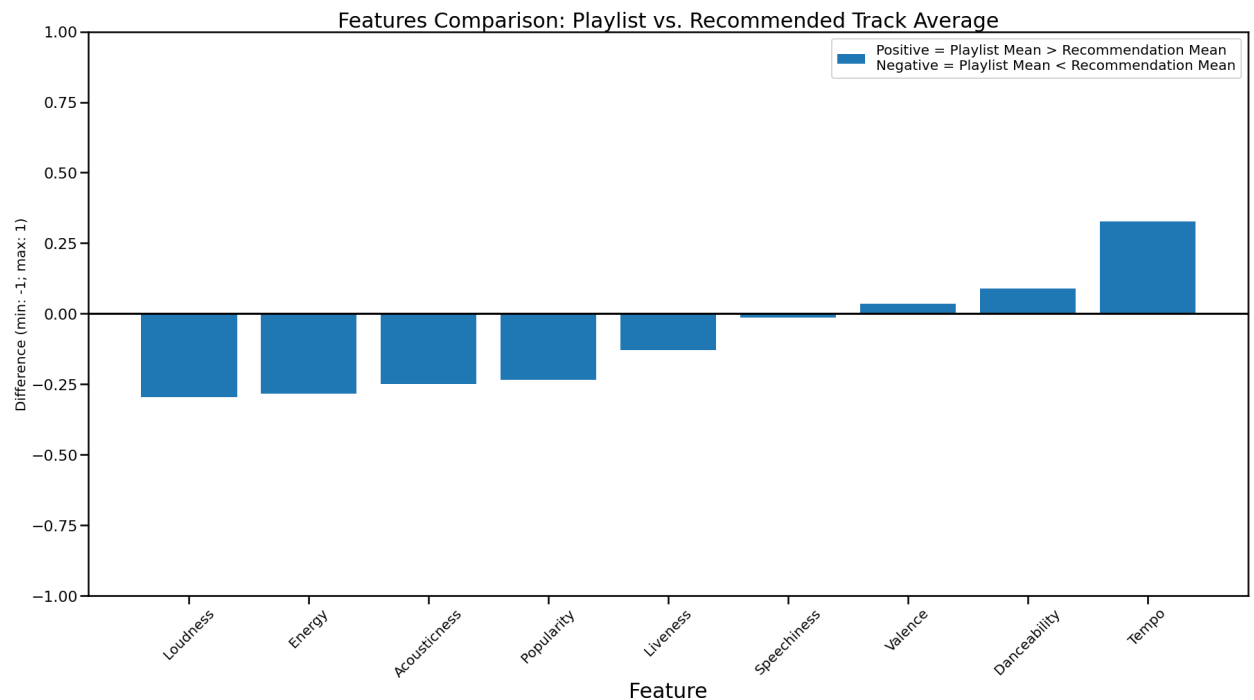
1 rows x 511 columns

In [76]: `1 song_vs_playlist_visual_comparison(765,10)`

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.

method=method,

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/ipykernel_launcher.py:129: UserWarning: FixedFormatter should only be used together with FixedLocator



In [77]: `1 unwrapped_track_data[unwrapped_track_data['pid']==765].head(50)`

Out[77]:

	pos	track_uri	track_name	artist_name	
1791	0	03xWMkKEbeO4SnyIA53ipj	When Will My Life Begin - From "Tangled"/Sound...	Mandy Moore	2LJxr7Pt3JnP6i
1792	1	1IOSxJNCLvWm2bYaTcTSmK	Mother Knows Best - From "Tangled"/Soundtrack ...	Donna Murphy	5BuTOT6mPoNZ5
1793	2	0TCt7OFRdD8PQ6vTRQxNgQ	I've Got a Dream - From "Tangled"/Soundtrack V...	Mandy Moore	2LJxr7Pt3JnP6i
1794	3	6klpXs2uAjagnZMFkt4qkl	I See the Light - From "Tangled"/Soundtrack Ve...	Mandy Moore	2LJxr7Pt3JnP6i
1795	4	75VVIB2x1h6BfxD2PqOO57	Healing Incantation - From "Tangled"/Soundtrac	Mandy Moore	2LJxr7Pt3JnP6i

As we can see above, the aggregate difference between our recommended songs and our playlist average is 6, which is driven mainly by the genre data. Most of our songs here are Disney/Soundtrack, so there seems to be little variation here. This aggregate difference is better than our first playlist, although it seems our features may be offsetting some of that difference.

Looking at our visualization, for the most part it appears all of the features are comparable, with tempo, energy, popularity, loudness and danceability having the largest difference. Seeing that the general theme of this playlist seems to be soundtrack for childrens movies, I believe that genre is the most important driver of the recommender here, and it seems like our genre hit the mark.

Now for the sound test. I will document my subjective findings below:

As we can see from our selection of songs, they all seem to be Disney songs from Moana and Tangled. I would expect that our recommendation system would suggest other kids songs accross childrens movies. After pulling more songs from the unwrapped track data, I can confirm that there are some songs that are not Disney.

For our recommendations, all 10 songs appear to be soundtrack songs, which is great, and they are all mainly from beloved kids movies. Because it seems like kids movies is the main category for this playlist, I will deem this recommendation adequate.

Overall, I would deem that these 10 recommended songs pass the reasonableness test.

Playlist #160101 Reasonableness Test

In [78]: `1 spotify_recommendation(160101,5)`

Thank you for inputting your playlist. Please see playlist tracks below:

	track_uri	track_name	artist_name	album_name
pos				
0	4w3dm0pGQ9otu7cG5uWy88	Not Just a Girl	She Wants Revenge	Valleyheart
1	37r6i0GTqgR05rGe5wNhmp	When They Fight, They Fight	Generational	Con Law
2	0grFc6klR3hxoHLcgCYsF4	Howlin' For You	The Black Keys	Brothers
3	6M23RkYPbVR91c4iWVNkcl	Changing	The Airborne Toxic Event	All At Once
4	5nHRIKsXDwUpse9gzaXLR	Oxford Comma	Vampire Weekend	Vampire Weekend
...
202	6T9ZJ2cJbtF4eDapnRHCux	Death Valley	My Jerusalem	Preachers
203	2QVmiA93GVhWNTWQctyY1K	Outro	M83	Hurry Up, We're Dreaming
204	7raMTVKDjLfTAyfDXKlfrz	Beneath The Surface	Demons Of Ruby Mae	Beneath The Surface
205	6y88SnrCoqRDZs1WTjKIZc	You're Loved & I'm Hated	Christopher Tyng	Suits (Original Television Soundtrack)
206	3uvsVUrAaGQJCTEUR1S3Sx	Bare	WILDES	Bare

207 rows × 4 columns

Here are 5 tracks that might fit this playlist:

	track_uri	track_name	artist_name	album_name
0	6mHw0o12JUfDodSxMwp8TI	Caves	Haux	All We've Known
1	629Cjw0fUyZUMkBjnjtDR	Places	Shlohmo	Bad Vibes
2	2BvDTv6akKbiKLT6ol2NXp	Coast to Coast	Waxahatchee	Cerulean Salt

3	2ByM6ejiN2EERm2NmXlq5t	Emerge From Smoke	Shlohmo	Dark Red
4	2VvfTbJWcguP88dvnPH7hl	Adult Diversion	Alvvays	Alvvays

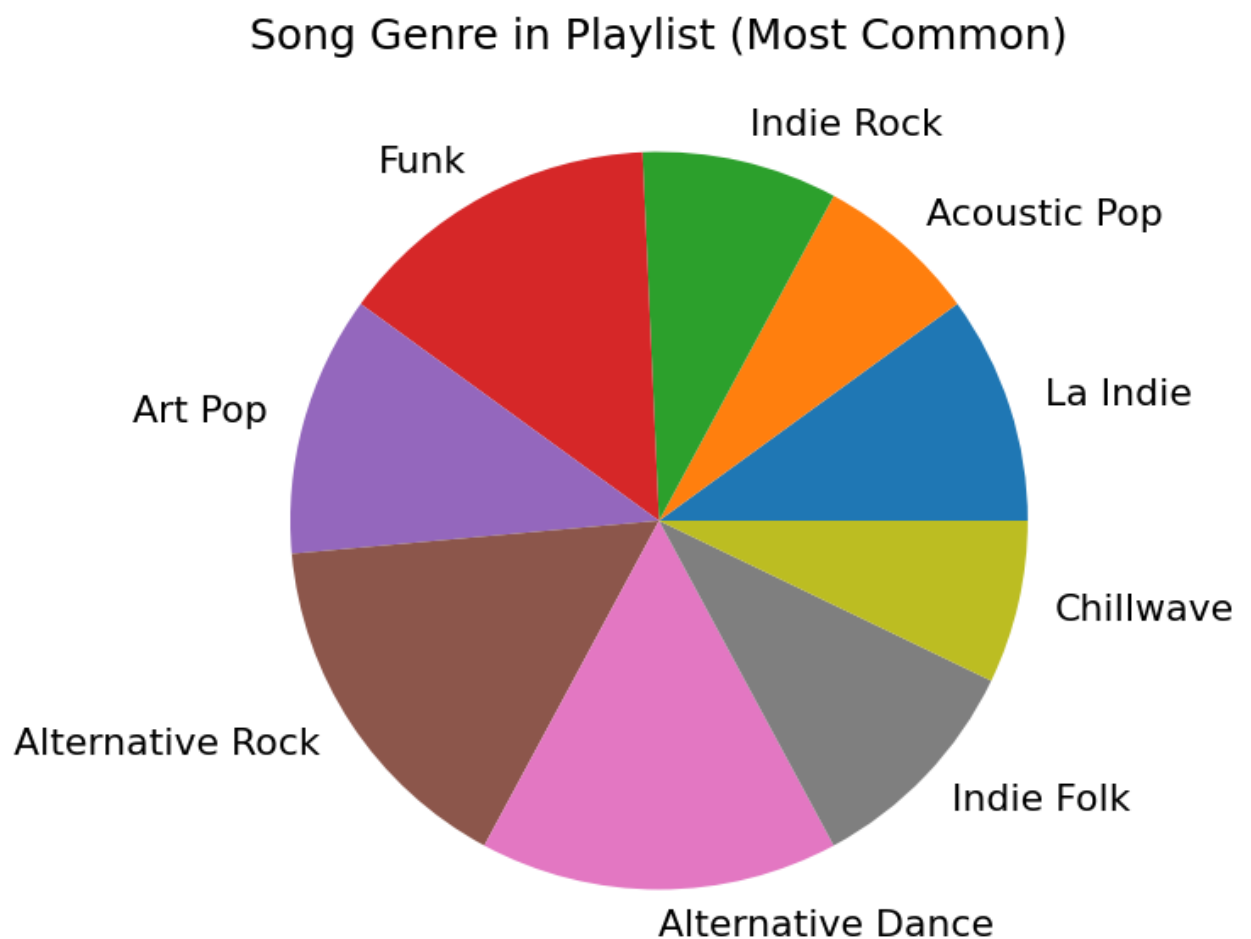
Out[78]: (None, None, None, None)

In [79]: 1 check_duplicates(160101,5)

No duplicates being recommended!

In [80]: 1 playlist_genre_breakdown(160101)

Out[80]: (<Figure size 1440x720 with 1 Axes>,
<AxesSubplot:title={'center':'Song Genre in Playlist (Most Common)'}>)



In [81]: `1 song_vs_playlist_comparison(160101,5)`

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.

method=method,

Out[81]:

	acousticness	danceability	energy	genre_1_21st century classical	genre_1_a cappella	genre_1_abstract beats	genre_1_a
0	0.234817	0.033102	-0.027902	0.0	0.0	0.0	

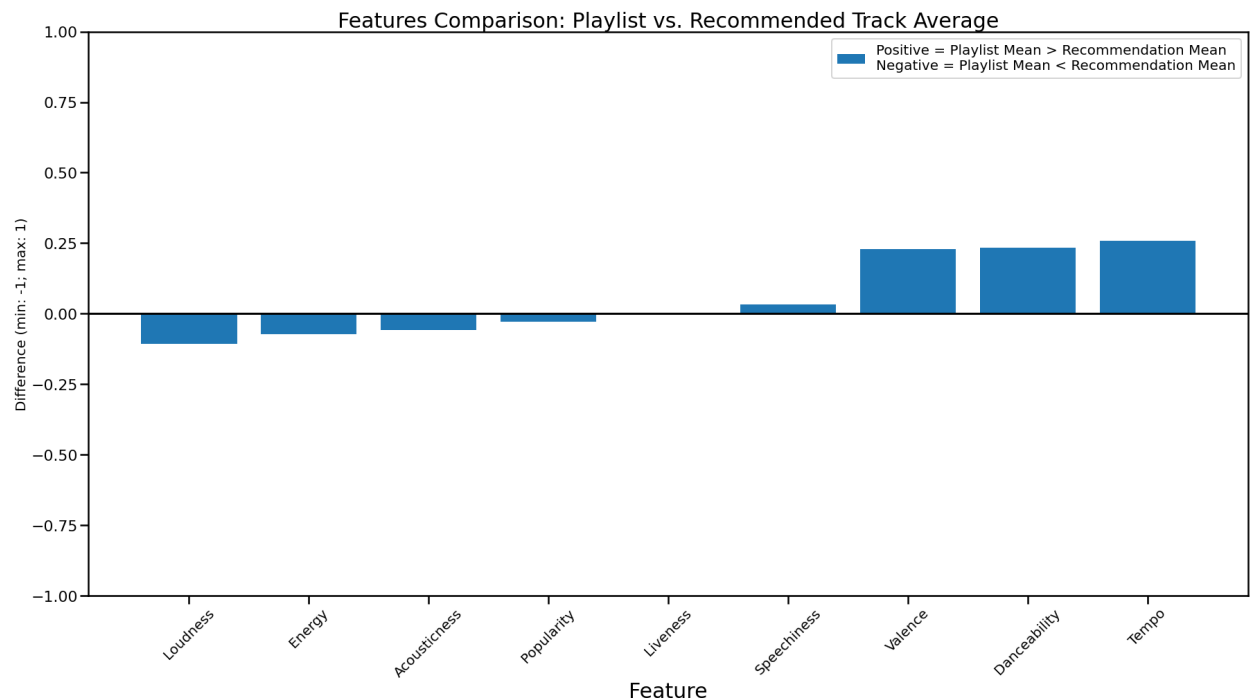
1 rows x 511 columns

In [82]: `1 song_vs_playlist_visual_comparison(160101,5)`

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/generic.py:2615: UserWarning: The spaces in these column names will not be changed. In pandas versions < 0.14, spaces were converted to underscores.

method=method,

/opt/anaconda3/envs/learn-env/lib/python3.6/site-packages/ipykernel_launcher.py:129: UserWarning: FixedFormatter should only be used together with FixedLocator



As we can see above, the aggregate difference between our recommendations and our playlist data is pretty large, at 44. I believe this is driven by our genre data, as it is very evenly spread out across multiple genres

According to our feature visualization, our differences are very minimal compared to the other playlists. The largest differences are Speechiness, Acousticness and Energy.

Now for the sound test. I will document my subjective findings below:

From our playlist, I already recognize M83, the Black Keys and Vampire Weekend. Outro by M83 is very spacy and slow, whereas Howlin' For You is more uptempo and rock themed. Oxford Comma is also uptempo and kind of spaced out, but overall sounds chill. Not Just a Girl is slower paced, but also sounds like alt-rock overall. For our recommendations, I would expect chiller alt-rock that is slower in pace and has mostly fleshed out instrumentals.

For our recommendations, I recognize none of the songs. Caves, Coast to Coast and Diversion are very Indie sounding, somewhat Alternative. The two songs by Shlohmo are dance-like, and sound more electronic. They are also instrumentals. Outro by M83 fits a similar vibe to these songs.

This playlist seems like there is more variety between the genres, and as such, my recommendations also seem varied. However, due to the feature difference being relatively small, and due to the genres fitting the broad scope of the playlist, I deem these recommendations sufficient.

Section 5: Results

To summarize: using cosine similarity, I was able to successfully create a recommendation system using a selection of data from the Spotify Million Playlist Dataset. Overall, I have identified 3 potential drawbacks that might affect this model: 1. Suggesting songs from the same basket that the machine is learning from may create bias, 2. There are few metrics to assess our system, which makes it harder to determine the broader success of the model, 3. I analyzed the model using subjectivity. Tastes can be very different and subjectivity begets bias, as well. With these in mind, however, I believe that this recommendation system is adequate at suggesting similar songs based on our tests in the results section.

As shown in the results section, our music recommendations seem to be driven mainly by the genre. Extracting Genre data from Spotify provided us with very specific genre data that might've been able to isolate our list to a very select handful of songs. It then used the audio features and popularity vectors to differentiate the songs when it came to using cosine similarity.

I believe some good improvements to this project would be to source more general genre data if it is eventually made available as an export from Spotify. As I had to manually export the data from Spotify using different means, it is hard to tell how accurate the genre data actually is. I would also like to include information on the year released, in case a playlist is related to a certain year, for example, 2013s hits. I would also like user-review data that shows how much a user liked a song, so I could try a collaborative-based filtering approach. Another thing I would like to do is research other forms of content-based recommendation systems, perhaps some form of clustering. Lastly, if I had a faster computer I would include all of the playlists in my analysis to capture more songs for my recommendation system.