

```
In [1]: 1 import spotipy
        2 from spotipy.oauth2 import SpotifyClientCredentials
        3 import pandas as pd
        4 import numpy as np
        5 import matplotlib.pyplot as plt
        6 import scipy.stats as stats
        7 import statsmodels.api as sm
```

```
In [2]: 1 #connecting to spotify API using Spotipy
        2 cid = "12b051049eb14a05ac3dd252770b6bd3"
        3 secret = "8844a8e5f63f4fa3998ba69811b2e005"
        4
        5 client_credentials_manager = SpotifyClientCredentials(client_id=cid, client_secret=secret)
        6 sp = spotipy.Spotify(client_credentials_manager=client_credentials_manager)
```

```
In [3]: 1 artist_name = []
        2 track_name = []
        3 popularity = []
        4 track_id = []
        5
        6 #loading in 1000 different tracks from 2020 with Spotipy
        7 for i in range(0,1000,10):
        8     track_results = sp.search(q='year:2020', type='track', limit=10,offset=i)
        9     for i, t in enumerate(track_results['tracks']['items']):
       10         #populating lists we will use for DF later
       11         artist_name.append(t['artists'][0]['name'])
       12         track_name.append(t['name'])
       13         track_id.append(t['id'])
       14         popularity.append(t['popularity'])
```

```
In [4]: 1 #storing data in dataframe
        2 df_tracks = pd.DataFrame({'artist_name':artist_name,'track_name':track_name,'track_id':track_id
```

```
In [5]: 1 rows = []
        2 batchsize = 100
        3
        4 # Create list of audio features for each song
        5 for i in range(0,len(df_tracks['track_id']),batchsize):
        6     batch = df_tracks['track_id'][i:i+batchsize]
        7     feature_results = sp.audio_features(batch)
        8     for i, t in enumerate(feature_results):
        9         if t != None:
        10             rows.append(t)
        11
```

In [6]:

```
1 # Storing audio feature data of each song in a dataframe
2 df_audio_features = pd.DataFrame.from_dict(rows,orient='columns')
3 print("Shape of the dataset:", df_audio_features.shape)
4 df_audio_features.head()
```

Shape of the dataset: (1000, 18)

	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness	liveness	valence	tempo	
0	0.653	0.524	11	-9.016	0	0.0502	0.1120	0.000000	0.2030	0.553	83.970	audio_f
1	0.699	0.558	9	-14.713	1	0.0541	0.1650	0.883000	0.0932	0.398	99.996	audio_f
2	0.680	0.826	0	-5.487	1	0.0309	0.0212	0.000012	0.5430	0.644	118.051	audio_f
3	0.642	0.749	6	-7.060	0	0.2890	0.0299	0.000000	0.2540	0.561	81.862	audio_f
4	0.436	0.655	1	-8.370	0	0.0583	0.4990	0.000008	0.6880	0.412	121.002	audio_f

In [7]:

```
1  # EM algo with 2 dimensions
2  # Needs to be modified to fit more features (~7) in the future
3
4  def EM(dat, k):
5
6      p_class=np.zeros(k)
7      means=np.zeros((k,2))
8      covars=np.zeros((k,2,2))
9      mean_dist=np.array(0)
10     p_data_given_class=np.zeros((len(dat),k))
11     #initializations
12     init_idx=np.random.choice(range(len(dat)), size=k, replace=False)
13
14     for dim in range(k):
15         covars[dim,:,:]=np.cov(np.transpose(dat))
16         means[dim,:]=dat.iloc[init_idx[dim]]
17         p_class[dim]=1/k
18
19     for step in range(50):
20         #Bayes stuff: pdfs then pdf*mixtures, then normalize
21         for dim in range(k):
22             p_data_given_class[:,dim]= np.array([stats.multivariate_normal.pdf(x=dat, mean=mea
23             p_class_given_data=p_data_given_class*p_class
24
25             sums=np.sum(p_class_given_data, axis=1)
26             for dim in range(k):
27                 p_class_given_data[:,dim]=p_class_given_data[:,dim]*(1/sums)
28             n_class = np.sum(p_class_given_data, axis=0)
29             p_class=n_class/len(dat)
30
31         # mean and covar updates
32         for dim in range(k):
33             means[dim,0]=np.sum(p_class_given_data[:,dim]*dat.iloc[:,0])*(1/n_class[dim])
```

```

34         means[dim,1]=np.sum(p_class_given_data[:,dim]*dat.iloc[:,1])*(1/n_class[dim])
35         covars[dim,0,0]=np.sum(p_class_given_data[:,dim]*((dat.iloc[:,0]-means[dim,0])**2))
36         covars[dim,1,1]=np.sum(p_class_given_data[:,dim]*((dat.iloc[:,1]-means[dim,1])**2))
37         covars[dim,0,1]=np.sum(p_class_given_data[:,dim]*(dat.iloc[:,1]-means[dim,1])*(dat.
38             covars[dim,1,0]=np.sum(p_class_given_data[:,dim]*(dat.iloc[:,1]-means[dim,1])*(dat.
39     mean_dist=0
40
41     for pt in range(len(dat)):
42         for dim in range(k):
43             #for each datum-mean pair, compute their prob-weighted distance apart
44             mean_dist+=np.sqrt(np.sum((means[dim,:]-np.array(dat.iloc[pt]))**2)*p_class_given_d
45     mean_dist=mean_dist/(len(dat)*k)
46     return p_class, means, covars, mean_dist, p_class_given_data

```

```

In [8]: 1 # Method to compute the probabiltiy a new song belongs to each genre (cluster)
2 def p_class_given_data(song, k, means, covars, p_class):
3     p_song_given_cluster = []
4     p_cluster_given_song = []
5     for cluster in range(k):
6         p_song_given_cluster.append( stats.multivariate_normal.pdf(x=song, mean=means[cluster],
7             p_cluster_given_song.append( p_song_given_cluster[cluster] * p_class[cluster] )
8
9     summ = sum(p_cluster_given_song)
10
11     for cluster in range(k):
12         p_cluster_given_song[cluster] = p_cluster_given_song[cluster] / summ
13
14     return p_cluster_given_song

```

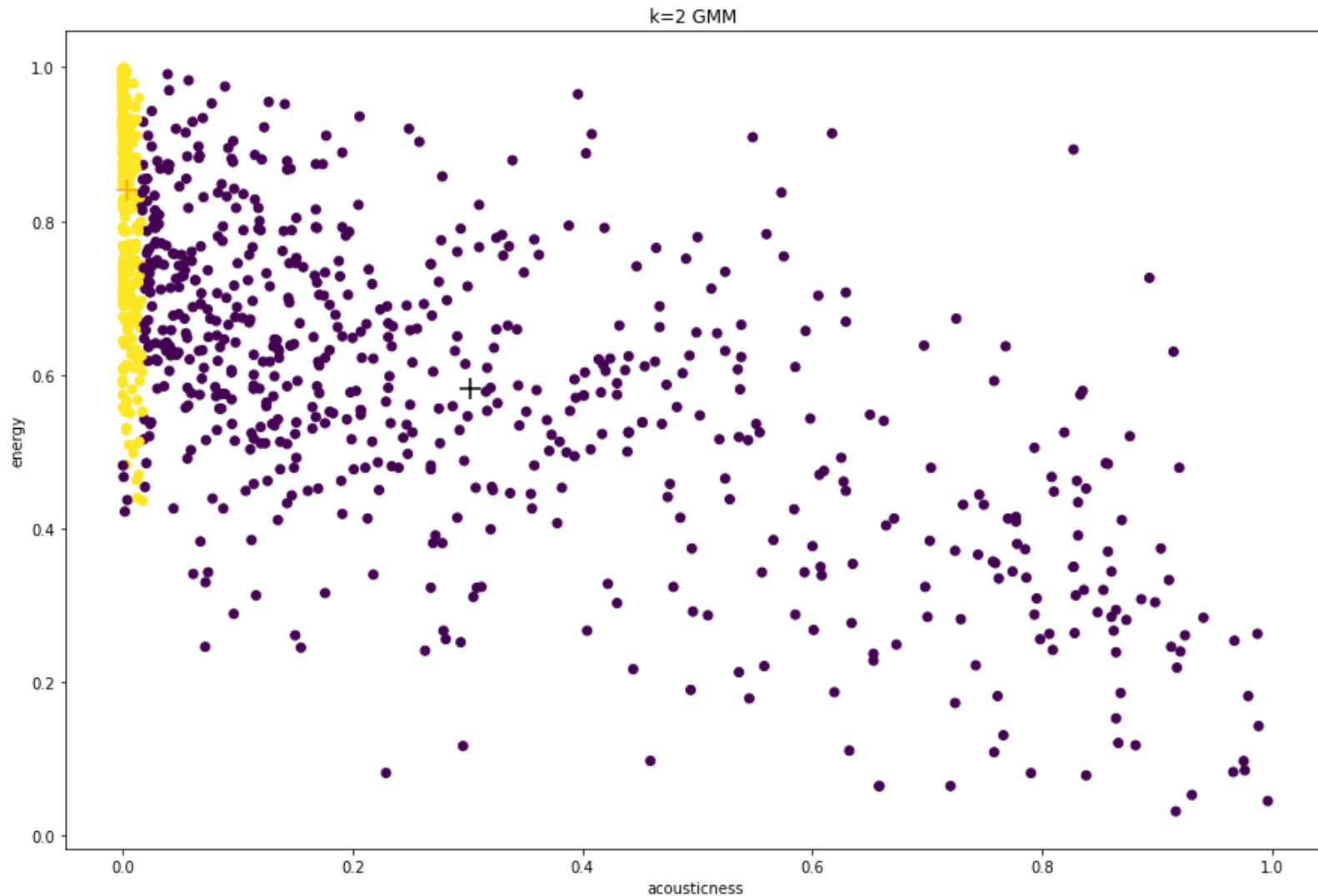
```
In [9]: 1 # Preliminary model using only 2 features
2 test_df = df_audio_features[['acousticness', 'energy']]
3 test_df
4 p2, m2, c2, d2, pc2 = EM(test_df, 2)
5
6 clusters = []
7 for pt in pc2:
8     if(pt[0]>pt[1]):
9         clusters.append(0)
10    else:
11        clusters.append(1)
```

```
In [10]: 1 test_song = [0.000112, 0.506]
2 probs = p_class_given_data(test_song, 2, m2, c2, p2)
3 print("Probability new song belongs to each cluster: ", probs)
```

Probability new song belongs to each cluster: [0.6534666004710983, 0.3465333995289018]

```
In [11]:
```

```
1 # Preliminary GMM model using only 2 features
2 fig, ax = plt.subplots(figsize=(15,10))
3 ax.set_title('k=2 GMM ')
4 ax.scatter(test_df['acousticness'],test_df['energy'], c=clusters)
5 ax.scatter(m2[0][0], m2[0][1], c='black', marker='+', s = 200)
6 ax.scatter(m2[1][0], m2[1][1], c='orange', marker='+', s = 200)
7 ax.set_xlabel('acousticness')
8 ax.set_ylabel('energy');
```



Project Check In

Our data so far is two dataframes containing song data from popular songs in 2020 from the Spotify API. `df_tracks` has the basic song data including the artist name, song name, track id and a popularity rating. `df_audio_features` contains the audio features for each song in the `df_tracks` data frame. This includes acousticness, danceability, energy, etc.

We have a preliminary GMM with $k=2$ using acousticness and energy as our features. The crosses on the plot represent the mean of each cluster. With this preliminary data, we have a proof of concept of being able to run the EM algorithm on our data set. We are also able to take a new song, and predict the probability it belongs to each cluster. For example a song with 0.000012 acousticness and 0.503 energy has a 65% chance of belonging to "cluster 1" with our preliminary GMM.

The remaining work we have is to modify the EM method above to account for all 9 audio features instead of just two. Additionally, we need to extract audio features of each song using Aubio, and run the EM algorithm on that data as well as the spotify data. Finally, we need to analyze and compare the Spotify API GMM and the Aubio GMM to see how closely we mimicked Spotify's song classification.