

In [1]:

```
import os
import numpy as np
import pandas as pd

import seaborn as sns
import plotly.express as px
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.manifold import TSNE
from sklearn.decomposition import PCA
from sklearn.metrics import euclidean_distances
from scipy.spatial.distance import cdist

import warnings
warnings.filterwarnings("ignore")
```

In [2]:

```
data = pd.read_csv("data.csv.zip")
genre_data = pd.read_csv('data_by_genres.csv')
year_data = pd.read_csv('data_by_year.csv')
```

In [3]:

```
print(data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 170653 entries, 0 to 170652
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  -
0   valence                170653 non-null float64
1   year                  170653 non-null int64
2   acousticness          170653 non-null float64
3   artists               170653 non-null object
4   danceability          170653 non-null float64
5   duration_ms           170653 non-null int64
6   energy                170653 non-null float64
7   explicit              170653 non-null int64
8   id                    170653 non-null object
9   instrumentalness       170653 non-null float64
10  key                    170653 non-null int64
11  liveness              170653 non-null float64
12  loudness              170653 non-null float64
13  mode                  170653 non-null int64
14  name                  170653 non-null object
15  popularity             170653 non-null int64
16  release_date          170653 non-null object
17  speechiness           170653 non-null float64
18  tempo                 170653 non-null float64
dtypes: float64(9), int64(6), object(4)
memory usage: 24.7+ MB
None
```

In [4]:

```
print(genre_data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2973 entries, 0 to 2972
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   mode                  2973 non-null  int64
1   genres                2973 non-null  object
2   acousticness          2973 non-null  float64
3   danceability          2973 non-null  float64
4   duration_ms           2973 non-null  float64
5   energy                2973 non-null  float64
6   instrumentalness       2973 non-null  float64
7   liveness              2973 non-null  float64
8   loudness              2973 non-null  float64
9   speechiness           2973 non-null  float64
10  tempo                 2973 non-null  float64
11  valence               2973 non-null  float64
12  popularity             2973 non-null  float64
13  key                   2973 non-null  int64
dtypes: float64(11), int64(2), object(1)
memory usage: 325.3+ KB
None
```

In [5]:

```
print(year_data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 14 columns):
 #   Column                Non-Null Count  Dtype  
---  -
 0   mode                  100 non-null   int64  
 1   year                  100 non-null   int64  
 2   acousticness          100 non-null   float64 
 3   danceability          100 non-null   float64 
 4   duration_ms           100 non-null   float64 
 5   energy                100 non-null   float64 
 6   instrumentalness       100 non-null   float64 
 7   liveness              100 non-null   float64 
 8   loudness               100 non-null   float64 
 9   speechiness           100 non-null   float64 
10   tempo                 100 non-null   float64 
11   valence                100 non-null   float64 
12   popularity             100 non-null   float64 
13   key                   100 non-null   int64  
dtypes: float64(11), int64(3)
memory usage: 11.1 KB
None
```

In [6]:

```
!pip install yellowbrick
```

```
Requirement already satisfied: yellowbrick in c:\users\chari\anaconda3\lib\site-packages (1.5)
Requirement already satisfied: scikit-learn>=1.0.0 in c:\users\chari\anaconda3\lib\site-packages (from yellowbrick) (1.0.2)
Requirement already satisfied: matplotlib!=3.0.0,>=2.0.2 in c:\users\chari\anaconda3\lib\site-packages (from yellowbrick) (3.5.1)
Requirement already satisfied: scipy>=1.0.0 in c:\users\chari\anaconda3\lib\site-packages (from yellowbrick) (1.7.3)
Requirement already satisfied: cycler>=0.10.0 in c:\users\chari\anaconda3\lib\site-packages (from yellowbrick) (0.11.0)
Requirement already satisfied: numpy>=1.16.0 in c:\users\chari\anaconda3\lib\site-packages (from yellowbrick) (1.21.5)
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\chari\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (3.0.4)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\chari\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (2.8.2)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\chari\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (1.3.2)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\chari\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (4.25.0)
Requirement already satisfied: packaging>=20.0 in c:\users\chari\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (21.3)
Requirement already satisfied: pillow>=6.2.0 in c:\users\chari\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (9.0.1)
Requirement already satisfied: six>=1.5 in c:\users\chari\anaconda3\lib\site-packages (from python-dateutil>=2.7->matplotlib!=3.0.0,>=2.0.2->yellowbrick) (1.16.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\chari\anaconda3\lib\site-packages (from scikit-learn>=1.0.0->yellowbrick) (2.2.0)
Requirement already satisfied: joblib>=0.11 in c:\users\chari\anaconda3\lib\site-packages (from scikit-learn>=1.0.0->yellowbrick) (1.1.0)
```

In [7]:

```
from yellowbrick.target import FeatureCorrelation

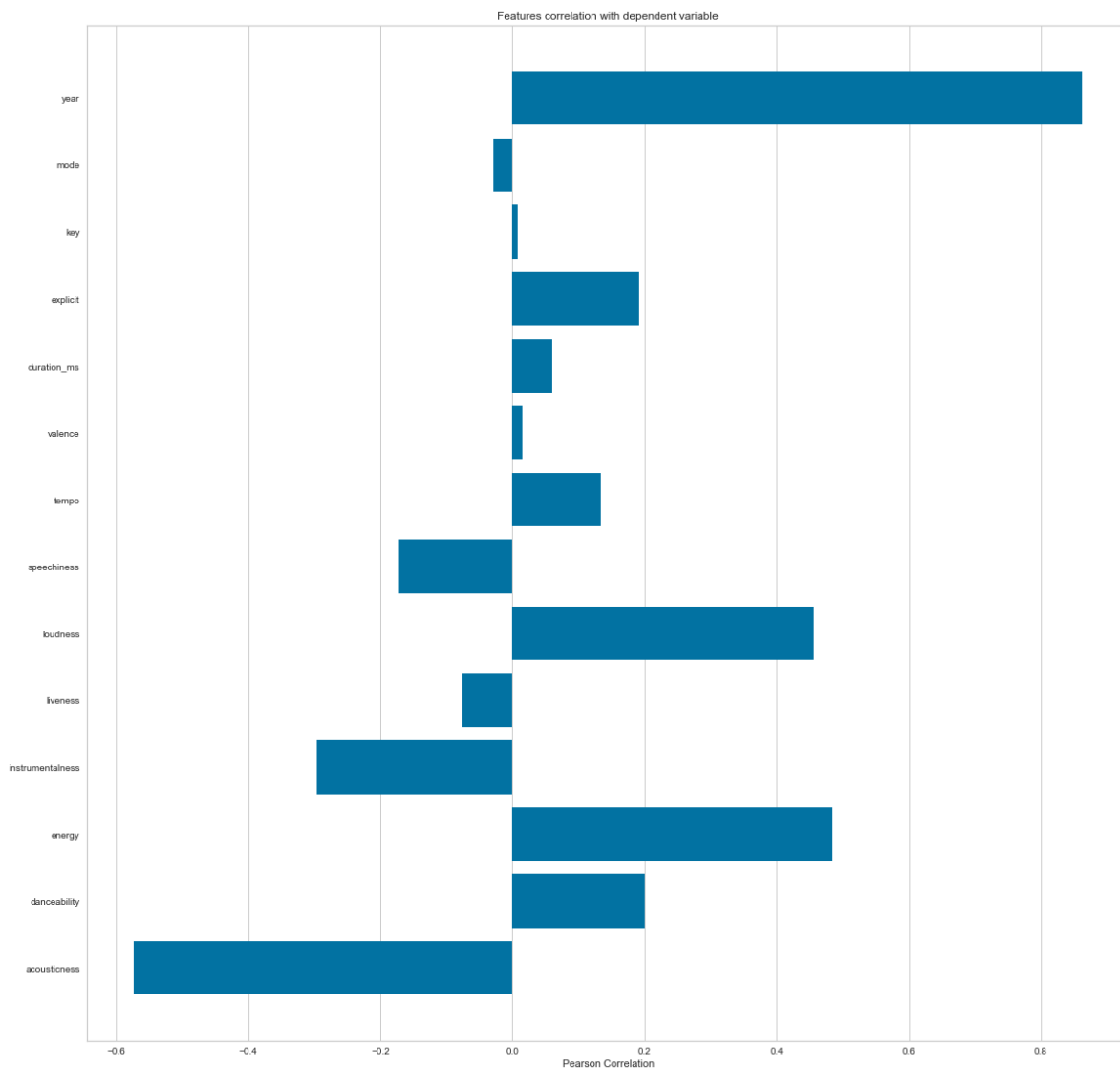
feature_names = ['acousticness', 'danceability', 'energy', 'instrumentalness',
                 'liveness', 'loudness', 'speechiness', 'tempo', 'valence', 'duration_ms', 'explicit',
                 ]

X, y = data[feature_names], data['popularity']

# Create a list of the feature names
features = np.array(feature_names)

# Instantiate the visualizer
visualizer = FeatureCorrelation(labels=features)

plt.rcParams['figure.figsize']=(20,20)
visualizer.fit(X, y)      # Fit the data to the visualizer
visualizer.show()
```



Out[7]:

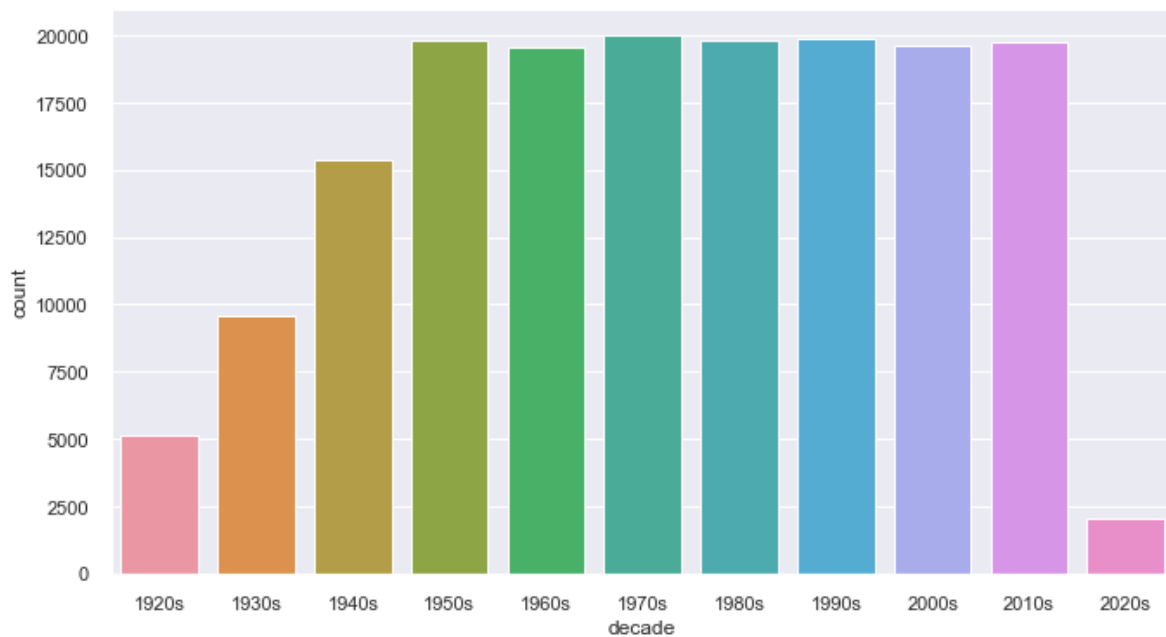
<AxesSubplot:title={'center':'Features correlation with dependent variable'}, xlabel='Pearson Correlation'>

In [8]:

```
def get_decade(year):  
    period_start = int(year/10) * 10  
    decade = '{}s'.format(period_start)  
    return decade  
  
data['decade'] = data['year'].apply(get_decade)  
  
sns.set(rc={'figure.figsize':(11 ,6)})  
sns.countplot(data['decade'])
```

Out[8]:

<AxesSubplot:xlabel='decade', ylabel='count'>



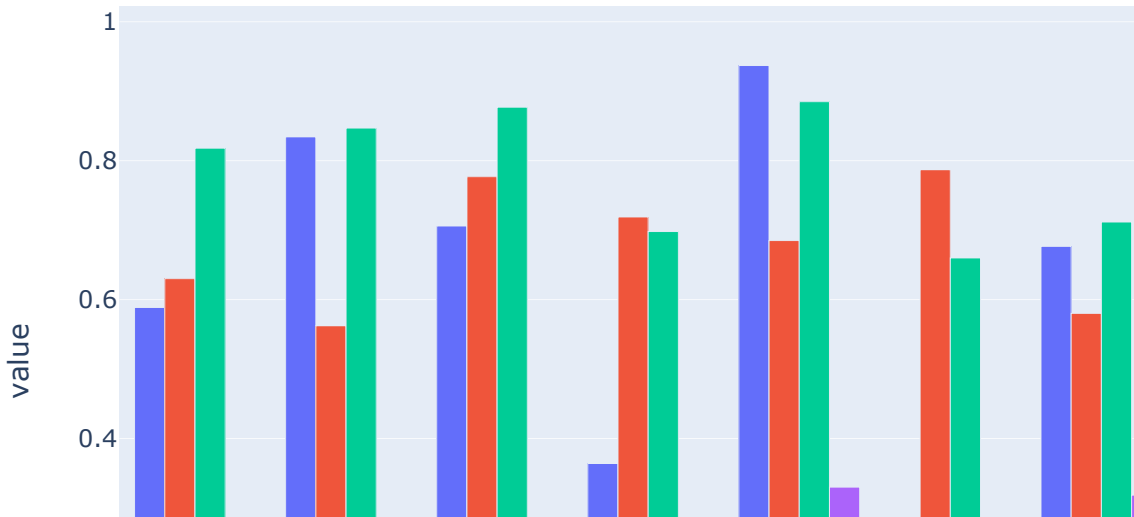
In [9]:

```
sound_features = ['acousticness', 'danceability', 'energy', 'instrumentalness', 'liveness',  
fig = px.line(year_data, x='year', y=sound_features)  
fig.show()
```



In [10]:

```
top10_genres = genre_data.nlargest(10, 'popularity')  
  
fig = px.bar(top10_genres, x='genres', y=['valence', 'energy', 'danceability', 'acousticness'])  
fig.show()
```



In [11]:

```
from sklearn.cluster import KMeans  
from sklearn.preprocessing import StandardScaler  
from sklearn.pipeline import Pipeline  
  
cluster_pipeline = Pipeline([('scaler', StandardScaler()), ('kmeans', KMeans(n_clusters=10))])  
X = genre_data.select_dtypes(np.number)  
cluster_pipeline.fit(X)  
genre_data['cluster'] = cluster_pipeline.predict(X)
```



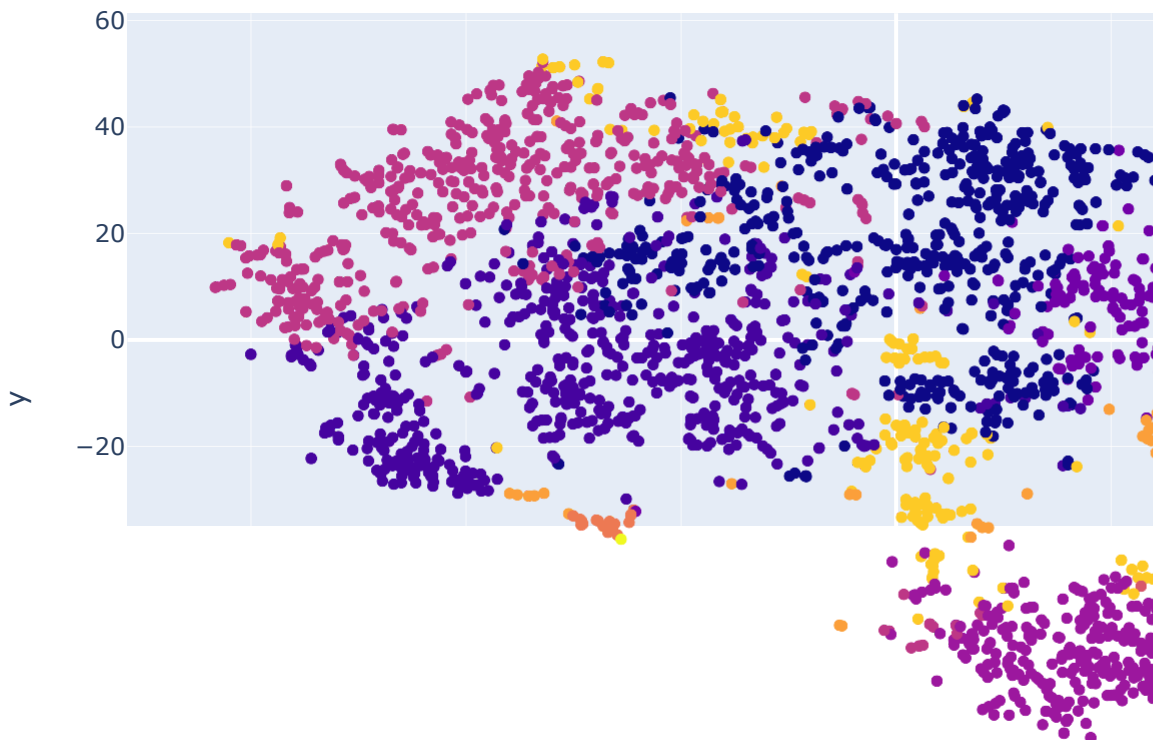
In [12]:

```
from sklearn.manifold import TSNE

tsne_pipeline = Pipeline([('scaler', StandardScaler()), ('tsne', TSNE(n_components=2, verbose=0))])
genre_embedding = tsne_pipeline.fit_transform(X)
projection = pd.DataFrame(columns=['x', 'y'], data=genre_embedding)
projection['genres'] = genre_data['genres']
projection['cluster'] = genre_data['cluster']

fig = px.scatter(
    projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'genres'])
fig.show()
```

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 2973 samples in 0.111s...
[t-SNE] Computed neighbors for 2973 samples in 0.404s...
[t-SNE] Computed conditional probabilities for sample 1000 / 2973
[t-SNE] Computed conditional probabilities for sample 2000 / 2973
[t-SNE] Computed conditional probabilities for sample 2973 / 2973
[t-SNE] Mean sigma: 0.777516
[t-SNE] KL divergence after 250 iterations with early exaggeration: 76.34357
5
[t-SNE] KL divergence after 1000 iterations: 1.407372
```



In [13]:

```
song_cluster_pipeline = Pipeline([('scaler', StandardScaler()),
                                   ('kmeans', KMeans(n_clusters=20,
                                                       verbose=False))
                                   ], verbose=False)

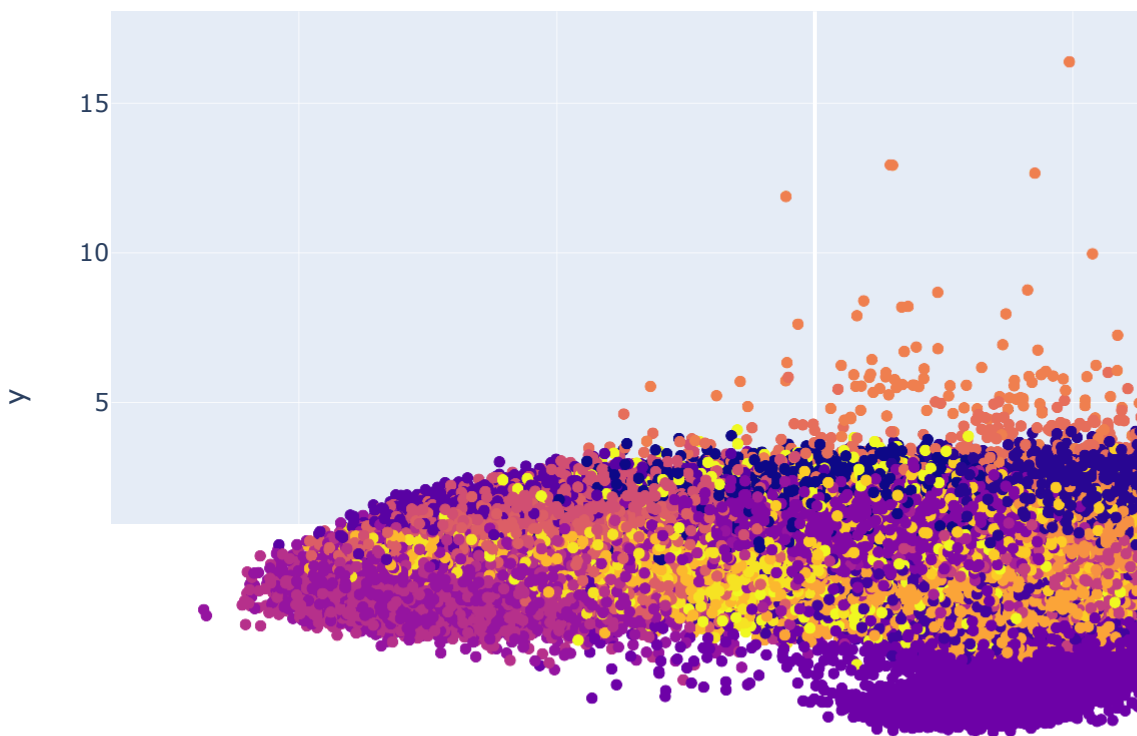
X = data.select_dtypes(np.number)
number_cols = list(X.columns)
song_cluster_pipeline.fit(X)
song_cluster_labels = song_cluster_pipeline.predict(X)
data['cluster_label'] = song_cluster_labels
```

In [14]:

```
from sklearn.decomposition import PCA

pca_pipeline = Pipeline([('scaler', StandardScaler()), ('PCA', PCA(n_components=2))])
song_embedding = pca_pipeline.fit_transform(X)
projection = pd.DataFrame(columns=['x', 'y'], data=song_embedding)
projection['title'] = data['name']
projection['cluster'] = data['cluster_label']

fig = px.scatter(
    projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'title'])
fig.show()
```



In [15]:

```
!pip install spotipy
```

```
Requirement already satisfied: spotipy in c:\users\chari\anaconda3\lib\site-packages (2.21.0)
Requirement already satisfied: requests>=2.25.0 in c:\users\chari\anaconda3\lib\site-packages (from spotipy) (2.27.1)
Requirement already satisfied: six>=1.15.0 in c:\users\chari\anaconda3\lib\site-packages (from spotipy) (1.16.0)
Requirement already satisfied: urllib3>=1.26.0 in c:\users\chari\anaconda3\lib\site-packages (from spotipy) (1.26.9)
Requirement already satisfied: redis>=3.5.3 in c:\users\chari\anaconda3\lib\site-packages (from spotipy) (4.3.4)
Requirement already satisfied: deprecated>=1.2.3 in c:\users\chari\anaconda3\lib\site-packages (from redis>=3.5.3->spotipy) (1.2.13)
Requirement already satisfied: packaging>=20.4 in c:\users\chari\anaconda3\lib\site-packages (from redis>=3.5.3->spotipy) (21.3)
Requirement already satisfied: async-timeout>=4.0.2 in c:\users\chari\anaconda3\lib\site-packages (from redis>=3.5.3->spotipy) (4.0.2)
Requirement already satisfied: wrapt<2,>=1.10 in c:\users\chari\anaconda3\lib\site-packages (from deprecated>=1.2.3->redis>=3.5.3->spotipy) (1.12.1)
Requirement already satisfied: pyparsing!=3.0.5,>=2.0.2 in c:\users\chari\anaconda3\lib\site-packages (from packaging>=20.4->redis>=3.5.3->spotipy) (3.0.4)
Requirement already satisfied: charset-normalizer~=2.0.0 in c:\users\chari\anaconda3\lib\site-packages (from requests>=2.25.0->spotipy) (2.0.4)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\chari\anaconda3\lib\site-packages (from requests>=2.25.0->spotipy) (2021.10.8)
Requirement already satisfied: idna<4,>=2.5 in c:\users\chari\anaconda3\lib\site-packages (from requests>=2.25.0->spotipy) (3.3)
```

In [16]:

```
import spotipy
from spotipy.oauth2 import SpotifyClientCredentials
from collections import defaultdict

sp = spotipy.Spotify(auth_manager=SpotifyClientCredentials(client_id="39cf48e7258b4fa397c3b
                                                         client_secret="41fb0dfa090249e8a

def find_song(name, year):
    song_data = defaultdict()
    results = sp.search(q= 'track: {} year: {}'.format(name,year), limit=1)
    if results['tracks']['items'] == []:
        return None

    results = results['tracks']['items'][0]
    track_id = results['id']
    audio_features = sp.audio_features(track_id)[0]

    song_data['name'] = [name]
    song_data['year'] = [year]
    song_data['explicit'] = [int(results['explicit'])]
    song_data['duration_ms'] = [results['duration_ms']]
    song_data['popularity'] = [results['popularity']]

    for key, value in audio_features.items():
        song_data[key] = value

    return pd.DataFrame(song_data)
```

In [17]:

```

from collections import defaultdict
from sklearn.metrics import euclidean_distances
from scipy.spatial.distance import cdist
import difflib

number_cols = ['valence', 'year', 'acousticness', 'danceability', 'duration_ms', 'energy',
               'instrumentalness', 'key', 'liveness', 'loudness', 'mode', 'popularity', 'speechiness', 't

def get_song_data(song, spotify_data):
    try:
        song_data = spotify_data[(spotify_data['name'] == song['name'])
                                & (spotify_data['year'] == song['year'])].iloc[0]
        return song_data

    except IndexError:
        return find_song(song['name'], song['year'])

def get_mean_vector(song_list, spotify_data):

    song_vectors = []

    for song in song_list:
        song_data = get_song_data(song, spotify_data)
        if song_data is None:
            print('Warning: {} does not exist in Spotify or in database'.format(song['name']))
            continue
        song_vector = song_data[number_cols].values
        song_vectors.append(song_vector)
        song_matrix = np.array(list(song_vectors))
    return np.mean(song_matrix, axis=0)

def flatten_dict_list(dict_list):

    flattened_dict = defaultdict()
    for key in dict_list[0].keys():
        flattened_dict[key] = []

    for dictionary in dict_list:
        for key, value in dictionary.items():
            flattened_dict[key].append(value)

    return flattened_dict

def recommend_songs( song_list, spotify_data, n_songs=10):

    metadata_cols = ['name', 'year', 'artists']
    song_dict = flatten_dict_list(song_list)

    song_center = get_mean_vector(song_list, spotify_data)
    scaler = song_cluster_pipeline.steps[0][1]
    scaled_data = scaler.transform(spotify_data[number_cols])
    scaled_song_center = scaler.transform(song_center.reshape(1, -1))
    distances = cdist(scaled_song_center, scaled_data, 'cosine')
    index = list(np.argsort(distances)[: , :n_songs][0])

    rec_songs = spotify_data.iloc[index]
    rec_songs = rec_songs[~rec_songs['name'].isin(song_dict['name'])]

```

```
return rec_songs[metadata_cols].to_dict(orient='records')
```

In [18]:

```
recommend_songs([{'name': 'Night changes', 'year': 2014},  
                 {'name': 'Friends', 'year': 2018},  
                 {'name': 'Shape of you', 'year': 2017},  
                 {'name': 'Just the way you are ', 'year': 2010},  
                 {'name': 'Baby', 'year': 2010}], data)
```

Out[18]:

```
[{'name': 'Dynamite', 'year': 2020, 'artists': "['BTS']"},  
 {'name': "What's Love Got to Do with It",  
  'year': 2020,  
  'artists': "['Kygo', 'Tina Turner']"},  
 {'name': 'Mi Niña',  
  'year': 2020,  
  'artists': "['Wisin', 'Myke Towers', 'Los Legendarios']"},  
 {'name': 'Breaking Me', 'year': 2019, 'artists': "['Topic', 'A7S']"},  
 {'name': 'Dynamite', 'year': 2020, 'artists': "['BTS']"},  
 {'name': 'Telepathy', 'year': 2020, 'artists': "['BTS']"},  
 {'name': 'Stay', 'year': 2017, 'artists': "['Zedd', 'Alessia Cara']"},  
 {'name': 'Con Calma', 'year': 2019, 'artists': "['Daddy Yankee', 'Snow']"},  
 {'name': 'Jangueo', 'year': 2019, 'artists': "['Alex Rose', 'Rafa Pabö  
n']"},  
 {'name': 'Be Honest (feat. Burna Boy)',  
  'year': 2019,  
  'artists': "['Jorja Smith', 'Burna Boy']"}]
```

In [ ]: