

# Buidling Music Recommendation System using Spotify Dataset

Hello and welcome to my kernel. In this kernel, I have created Music Recommendation System using Spotify Dataset. To do this, I presented some of the visualization processes to understand data and some done EDA(Exploratory Data Analysis) so we can select features that are relevant to create a Recommendation System.

Import Libraries

```
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
import plotly offline

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")
genre_data = pd.read_csv("data_by_genre.csv")
year_data = pd.read_csv("data_by_year.csv")
artist_data = pd.read_csv("data_by_artist.csv")

In [3]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 176653 entries, 0 to 176652
Data columns (total 15 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   valence     176653 non-null  float64
 1   year        176653 non-null  int64
 2   acousticness 176653 non-null  float64
 3   artists     176653 non-null  object
 4   danceability 176653 non-null  float64
 5   duration_ms 176653 non-null  int64
 6   energy      176653 non-null  float64
 7   explicit    176653 non-null  int64
 8   id          176653 non-null  object
 9   instrumentalness 176653 non-null  float64
10   key         176653 non-null  int64
11   liveliness  176653 non-null  float64
12   loudness    176653 non-null  float64
13   mode       176653 non-null  int64
14   name        176653 non-null  object
15   popularity  176653 non-null  int64
16   release_date 176653 non-null  object
17   speechiness 176653 non-null  float64
18   tempo       176653 non-null  float64
dtypes: float64(9), int64(6), object(4)
memory usage: 24.7+ MB

In [4]: 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   liveliness  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   year        2973 non-null  int64
13   key         2973 non-null  int64
dtypes: float64(11), int64(2), object(1)
memory usage: 325.3+ KB

In [5]: year_data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 180 entries, 0 to 99
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   mode       180 non-null  int64
 1   year       180 non-null  int64
 2   acousticness 180 non-null  float64
 3   danceability 180 non-null  float64
 4   duration_ms 180 non-null  float64
 5   energy      180 non-null  float64
 6   instrumentalness 180 non-null  float64
 7   liveliness  180 non-null  float64
 8   loudness    180 non-null  float64
 9   speechiness 180 non-null  float64
10   tempo       180 non-null  float64
11   valence     180 non-null  float64
12   popularity  180 non-null  int64
13   key         180 non-null  int64
dtypes: float64(11), int64(3)
memory usage: 11.1 KB

In [6]: artist_data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2880 entries, 0 to 2879
Data columns (total 15 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   mode       2880 non-null  int64
 1   count      2880 non-null  int64
 2   acousticness 2880 non-null  float64
 3   artists     2880 non-null  object
 4   danceability 2880 non-null  float64
 5   duration_ms 2880 non-null  float64
 6   energy      2880 non-null  float64
 7   instrumentalness 2880 non-null  float64
 8   liveliness  2880 non-null  float64
 9   loudness    2880 non-null  float64
10   loudness    2880 non-null  float64
11   popularity  2880 non-null  int64
12   release_date 2880 non-null  object
13   popularity  2880 non-null  int64
14   key         2880 non-null  int64
dtypes: float64(11), int64(3), object(1)
memory usage: 3.3+ MB

The datasets are downloaded by a trusted and genuine source which is kaggle.com

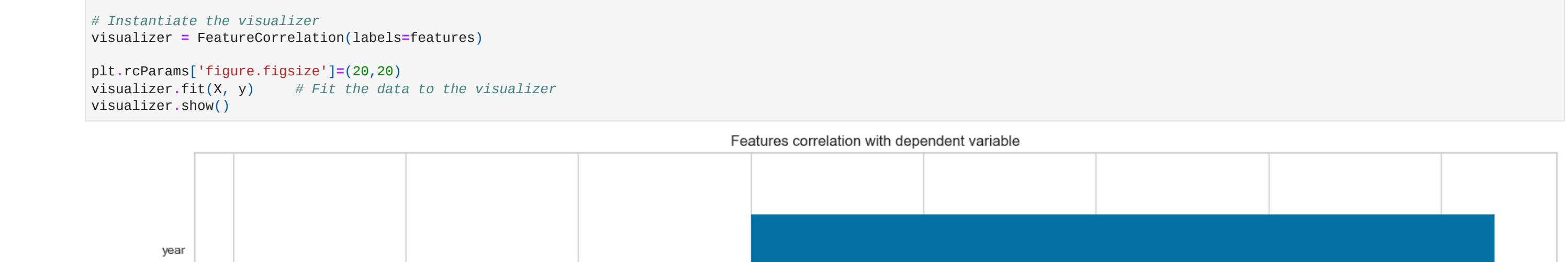
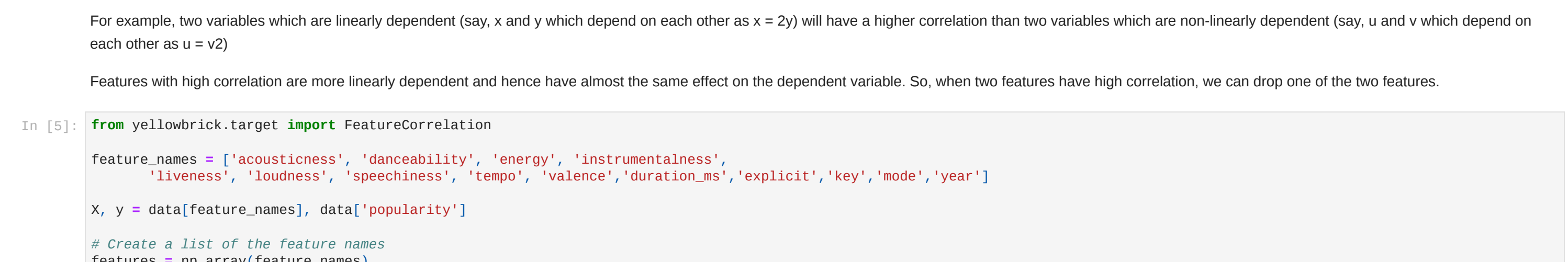
We are going to check for the analysis with the target as popularity. Before going to do that let's check for the Feature Correlation by considering a few features and for that, I'm going to use the yellowbrick package.

Yellowbrick extends the Scikit-Learn API to make model selection and hyperparameter tuning easier. Under the hood, it's using Matplotlib.

Correlation is a statistical term which in common usage refers to how close two variables are to having a linear relationship with each other.

For example, two variables which are linearly dependent (say, x and y which depend on each other as x = 2y) will have a higher correlation than two variables which are non-linearly dependent (say, u and v which depend on each other as u = v^2).

Features with high correlation are more linearly dependent and hence have almost the same effect on the dependent variable. So, when two features have high correlation, we can drop one of the two features.
```



Conclusions from EDA

Most of the songs range between 1950s-2010. Energy in songs has increased over the time. Acousticness in songs has reduced greatly over the decades. We can clearly see that loudness has dominantly increased over the decades and is at its peak in 2020. In top 15 genres we can see that energy and danceability are most noticeable features.

Clustering

Clustering is a data science technique in machine learning that groups similar rows in a data set. After running a clustering technique, a new column appears in the data set to indicate the group each row of data fits into best. Since rows of data, or data points, often represent people, physical transactions, documents or other important entities, these groups tend to form clusters of similar entities that have several kinds of real-world applications.

Clustering is sometimes referred to as unsupervised machine learning. To perform clustering, labels for past known outcomes – a dependent, y target or label variable – are generally unnecessary. For example, when applying a clustering method to a mortgage loan application process, it's not necessary to know whether the applicants made their past mortgage payments. Rather, you need demographic, psychographic, behavioral, geographic or other information about the applicants in a mortgage portfolio. A clustering method will then attempt to group the applicants based on that information. This method stands in contrast to supervised learning, in which mortgage default risk for new applicants, for example, can be predicted based on patterns in data labeled with past default outcomes.

Since clustering generally groups input data, it can be very creative and flexible. Clustering can be used for data exploration and preprocessing, as well as specific applications. From a technical perspective, common applications of clustering include:

- Data visualization. Data often contains natural groups or segments, and clustering is useful to large find them. Visualizing clusters can be a highly informative data analysis approach. Prototypes are data points that represent many other points and help explain data and models. If a cluster represents a large market segment, then the data point at the cluster center – or cluster centroid – is the prototypical member of that market segment. Sampling. Since clustering can define groups in the data, clusters can be used to create different types of data samples. Drawing an equal number of data points from each cluster in a data set, for example, can create a balanced sample of the population represented by that data set. Segments for models. Sometimes the predictive performance of supervised models – regression, decision trees and neural networks, for example – can be improved by using the information learned from unsupervised approaches such as clusters. Data scientists might create clusters as inputs to other models or build separate models for each cluster. For business applications, clustering is a battle-tested tool in market segmentation and fraud detection. Clustering is also useful for categorizing documents, making product recommendations and in other applications where grouping entities makes sense.

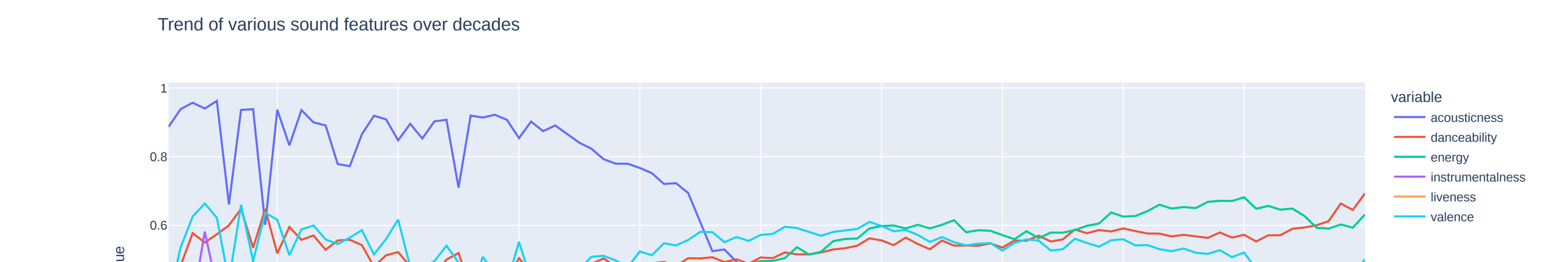
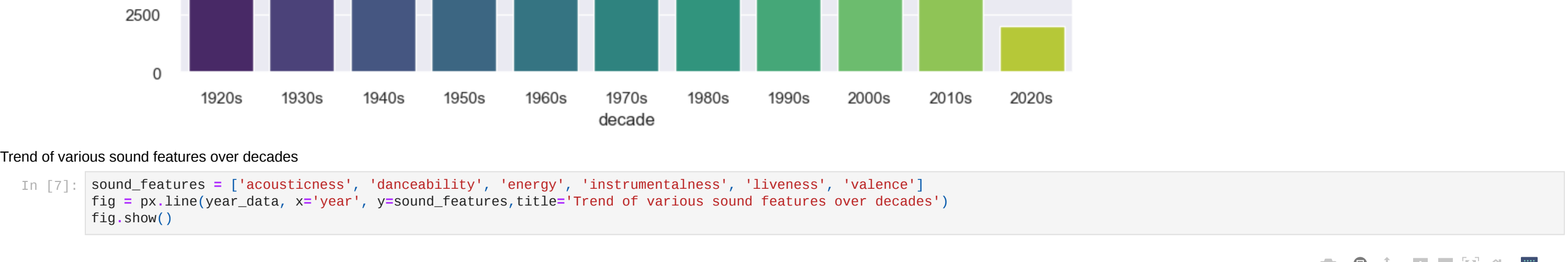
We performed k-means clustering on datasets to divide the dataset based on genre and songs.

In [10]:

```
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline

cluster_pipeline = Pipeline([
    ('scaler', StandardScaler()),
    ('kmeans', KMeans(n_clusters=20, verbose=False)) # Remove n_jobs parameter here
])

X = genre_data.select_dtypes(include=[np.number])
cluster_pipeline.fit(X)
genre_data['cluster'] = cluster_pipeline.named_steps['kmeans'].predict(X)
```



Model Buiding

In [40]:

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

new SPOTIFY_CLIENT_ID= 05f87c2b3ca747aba558fb77796582
new SPOTIFY_CLIENT_SECRET= 01b07f0130104025b9a93bf89793999

sp = spotipy.Spotify(auth_manager=SpotifyClientCredentials(client_id=new SPOTIFY_CLIENT_ID,
                                                         client_secret=new SPOTIFY_CLIENT_SECRET))

def find_song(name, year):
    song_data = defaultdict()
    results = sp.search(q=f'track: {year} {name}', 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)

env = SPOTIFY_CLIENT_ID=05f87c2b3ca747aba558fb77796582
env = SPOTIFY_CLIENT_SECRET=01b07f0130104025b9a93bf89793999
```

In [41]:

```
# Visualizing the Clusters with PCA

from sklearn.decomposition import PCA

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

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

In [42]:

```
def get_song_data(song, spotify_data):
    try:
        song_data = spotify_data[spotify_data['name'] == song['name']]
        s = (spotify_data['year'] == song['year']) && (spotify_data['name'] == song['name'])
        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(f'Warning: {song} does not exist in Spotify or in database'.format(song['name']))
            continue
        song_vectors.append(song_data)
    song_vectors = np.array(list(song_vectors))
    return np.mean(song_vectors, 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:
            flattened_dict[key].append(dictionary[key])
    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)
    scaled_data = song_cluster_pipeline.steps[0][1].transform(song_data)
    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:-1])
    rec_songs = spotify_data.iloc[index]
    rec_songs = rec_songs[['name', 'year', 'artists']]
    return rec_songs[metadata_cols].to_dict(orient='records')
```

In [43]:

```
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(f'Warning: {song} does not exist in Spotify or in the database'.format(song['name']))
            continue
        song_vector = song_data[number_cols].values
        if len(song_vector.shape) == 1:
            song_vector = song_vector[0] # Extract vectors as a 2D array
            print(f'Song: {song["name"]}, Vector: {song_vector}')
        song_vectors.append(song_vector)
    song_matrix = np.array(list(song_vectors))
    return np.mean(song_matrix, axis=0)
```

Testing

In [43]:

```
recommend_songs([{'name': 'Come As You Are', 'year': 1991},
                  {'name': 'Smells Like Teen Spirit', 'year': 1991},
                  {'name': 'Lithium', 'year': 1992},
                  {'name': 'All Apologies', 'year': 1993},
                  {'name': 'Stay Away', 'year': 1993}], data)
```

Couldn't read cache at: cache

Song: Come As You Are, Vector: [0.539 1991 0.00016 0.5 2.09202 0.8240000000000001 0 0.00161 4 0.0916

0.5 0.72 0.008 120.125]

Song: Smells Like Teen Spirit, Vector: [0.72 1991 2.556e-05 0.562 301920 0.932 0 0.000173 1 0.106 -4.556 1 75

0.064 116.761]

Song: Lithium, Vector: [0.388 1992 0.000203 0.688 250093 0.599 0 0.0 0.0782 -0.176 1 32 0.038

123.29]

Song: All Apologies, Vector: [0.359 1993 0.0793 0.446 233173 0.632 0 0.00266 1 0.0881 -12.197 1 69

0.0 0.174]

Couldn't write token to cache at: cache

Couldn't read cache at: cache

Couldn't write token to cache at: cache

Song: Stay Away, Vector: [1.52080e-01 1.99300e+03 2.20080e-02 0.67000e-01 1.51891e+05

-5.65900e+00 1.80000e+00 5.80000e+01 8.20000e-02 1.58058e+02

['name': 'My Wish', 'year': 2006, 'artists': ['Rascal Flatts']],

['name': 'Shivers Dressed Man - 2006 Remaster',

'year': 1989,

'artists': ['ZZ Top']],

['name': 'Nuestro Amor',

'year': 2005,

'artists': ['RBD', 'Anahí', 'Dulce María', 'Maite Perroni', 'Christian Chávez', 'Christopher von Uckermann', 'Alfonso Herrera']],

['name': 'Believe', 'year': 1998, 'artists': ['Cher']],

['name': 'Far Away', 'year': 2005, 'artists': ['Nickelback']],

['name': 'Psychosocial', 'year': 2008, 'artists': ['Slipknot']],

['name': 'Heart of Glass (Live From The JMWart Festival)',

'year': 2020,

'artists': ['Miley Cyrus']],

['name': 'What 2've Done', 'year': 2007, 'artists': ['Linkin Park']],

['name': 'In the End', 'year': 2001, 'artists': ['Nine Inch Nails']],

['name': 'Hanging By A Moment', 'year': 2008, 'artists': ['LilRequiem']]

In [ ]: