

Build a Recommendation System Music

<https://www.kaggle.com/datasets/vatsalmavani/spotify-dataset>

```
In [ ]: pip install spotipy
```

```
In [ ]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import spotipy
import os
%matplotlib inline
```

```
In [ ]: spotify_data = pd.read_csv('desktop/python/data/data.csv')
genre_data = pd.read_csv('desktop/python/data/data_by_genres.csv')
data_by_year = pd.read_csv('desktop/python/data/data_by_year.csv')
```

```
In [ ]: import plotly.express as px
sound_features = ['acousticness', 'danceability', 'energy', 'instrumentalness', 'liveness', 'valence']
fig = px.line(data_by_year, x='year', y=sound_features)
fig.show()
```

```
In [ ]: fig = px.line(data_by_year, x='year', y='tempo')
fig.show()
```

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

```
In [ ]: 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, n_jobs=-1))])
X = genre_data.select_dtypes(np.number)
cluster_pipeline.fit(X)
genre_data['cluster'] = cluster_pipeline.predict(X)
```

```
In [ ]: from sklearn.manifold import TSNE
tsne_pipeline = Pipeline([('scaler', StandardScaler()), ('tsne', TSNE(n_components=2, verbose=2))])
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']
```

```
In [ ]: import plotly.express as px
fig = px.scatter(
    projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'genres'])
fig.show()
```

```
In [ ]: song_cluster_pipeline = Pipeline([('scaler', StandardScaler()),
                                       ('kmeans', KMeans(n_clusters=20,
                                                         verbose=2, n_jobs=4)), verbose=True])
X = spotify_data.select_dtypes(np.number)
number_cols = list(X.columns)
song_cluster_pipeline.fit(X)
song_cluster_labels = song_cluster_pipeline.predict(X)
spotify_data['cluster_label'] = song_cluster_labels
```

```
In [ ]: 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'] = spotify_data['name']
projection['cluster'] = spotify_data['cluster_label']
```

```
In [ ]: import plotly.express as px
fig = px.scatter(projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'title'])
fig.show()
```

```
In [ ]: from spotipy.oauth2 import SpotifyClientCredentials
from collections import defaultdict

sp = spotipy.Spotify(auth_manager=SpotifyClientCredentials(client_id=os.environ["SPOTIFY_CLIENT_ID"],
                                                         client_secret=os.environ["SPOTIFY_CLIENT_SECRET"]))
```

```

def find_song(name, year):
    """
    This function returns a dataframe with data for a song given the name and release year.
    The function uses Spotify to fetch audio features and metadata for the specified song.
    """

    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 [ ]: from collections import defaultdict
        from scipy.spatial.distance import cdist
        import difflib

number_cols = ['valence', 'year', 'acousticness', 'danceability', 'duration_ms', 'energy', 'explicit',
               'instrumentalness', 'key', 'liveness', 'loudness', 'mode', 'popularity', 'speechiness', 'tempo']

def get_song_data(song, spotify_data):
    """
    Gets the song data for a specific song. The song argument takes the form of a dictionary with
    key-value pairs for the name and release year of the song.
    """

    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):
    """
    Gets the mean vector for a list of songs.
    """

    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):
    """
    Utility function for flattening a list of dictionaries.
    """

    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):
    """
    Recommends songs based on a list of previous songs that a user has listened to.
    """

    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 [ ]: 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}], spotify_data)

```

```

In [ ]: recommend_songs([{'name':'Beat It', 'year': 1982},
                        {'name': 'Billie Jean', 'year': 1988},
                        {'name': 'Thriller', 'year': 1982}], spotify_data)

```