PROJECT REPORT

UE19CS332-Algorithms of the Intelligent Web and Information Retrieval
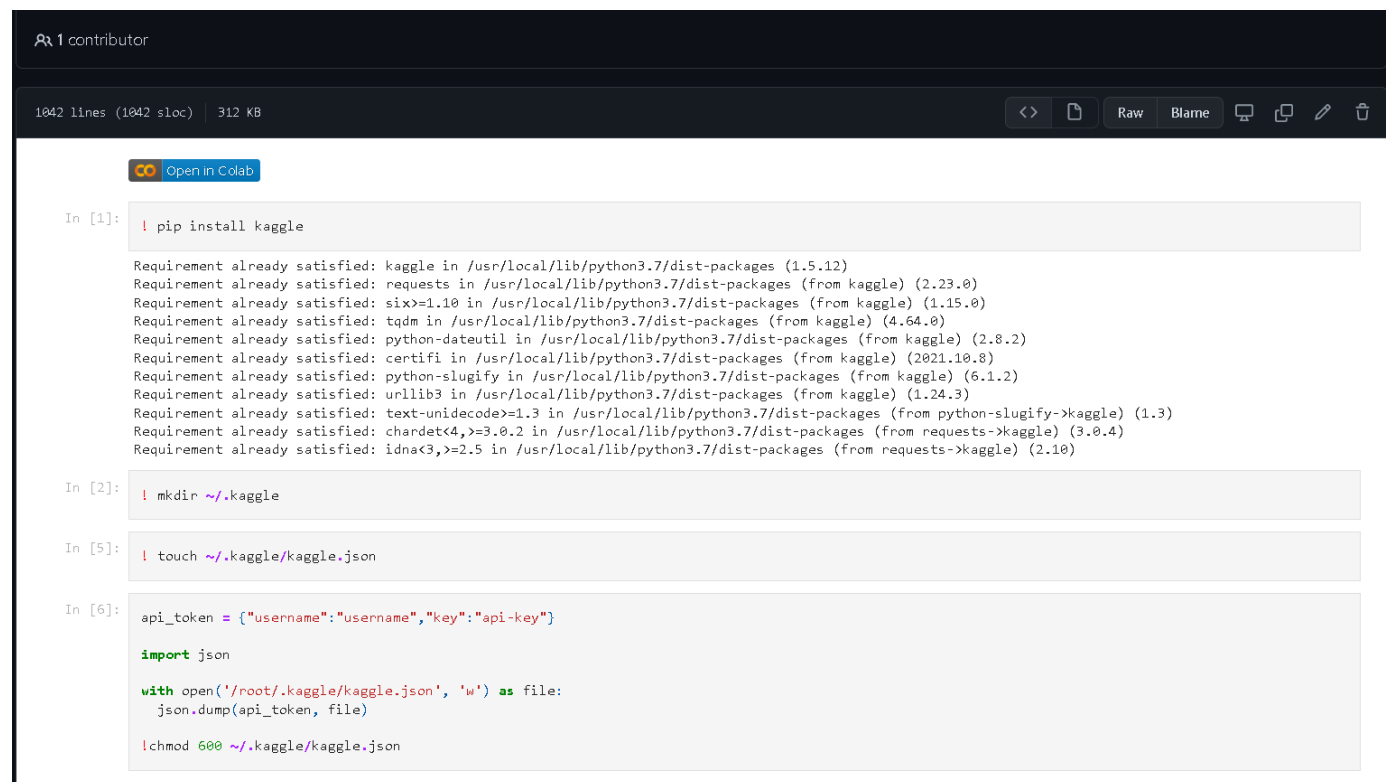
Project Title: Spotify Recommendation System

Team: Sapna Singh – PES2UG19CS366

   Siddharth Chattar – PES2UG19CS388

   Sneha Sujit Saha – PES2UG19CS393

Code and the Output:

```
!kaggle datasets download -d vatsalmavani/spotify-dataset
```

```
Downloading spotify-dataset.zip to /content
 91% 15.0M/16.5M [00:00<00:00, 74.7MB/s]
100% 16.5M/16.5M [00:00<00:00, 74.6MB/s]
```

```
! unzip spotify-dataset
```

```
Archive:  spotify-dataset.zip
  inflating: data/data.csv
  inflating: data/data_by_artist.csv
  inflating: data/data_by_genres.csv
  inflating: data/data_by_year.csv
  inflating: data/data_w_genres.csv
```

```python
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")
```

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

```python
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
```

```python
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
```

```
<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 [ ]:
```python
from yellowbrick.target import FeatureCorrelation

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

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()
```
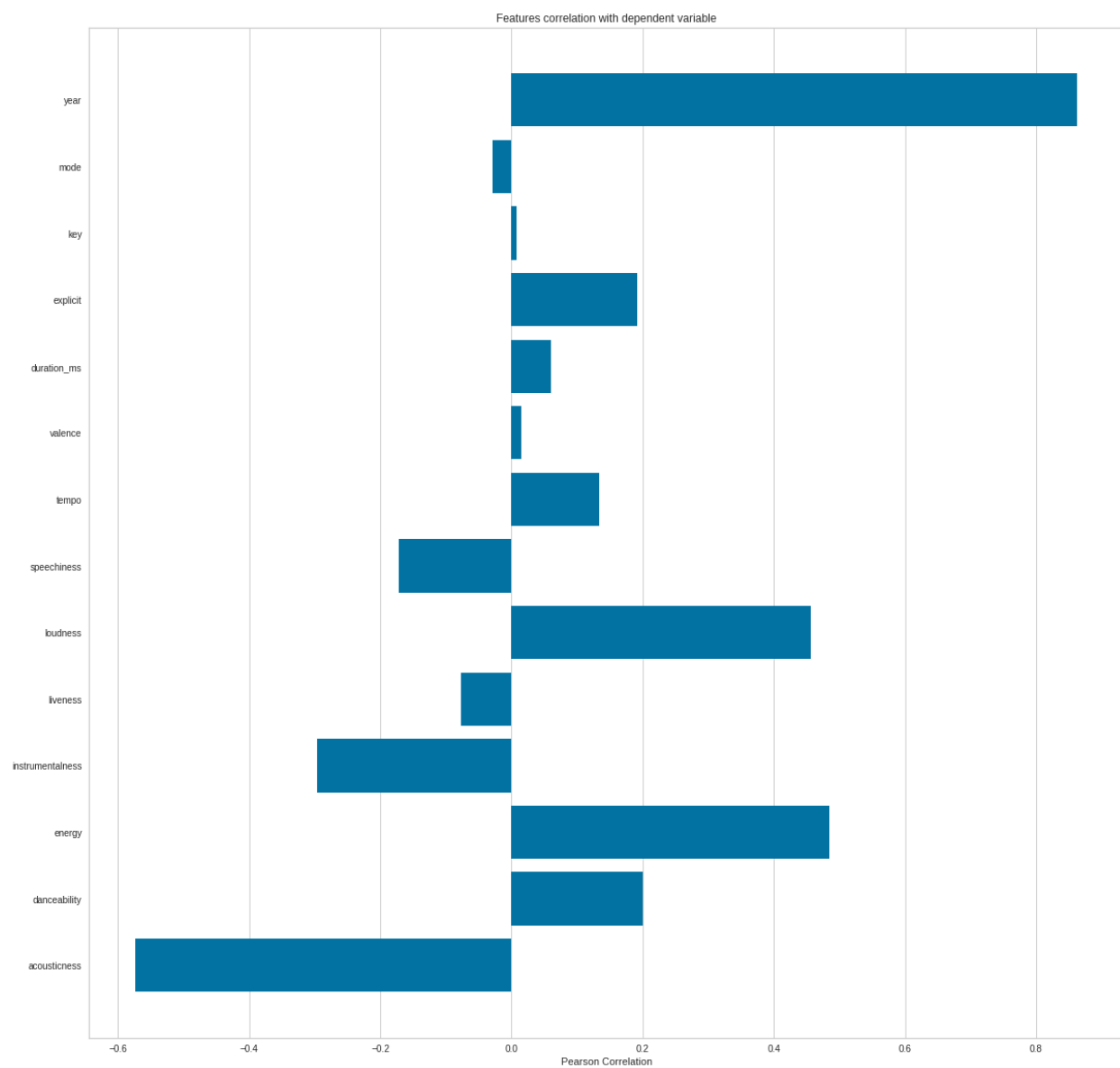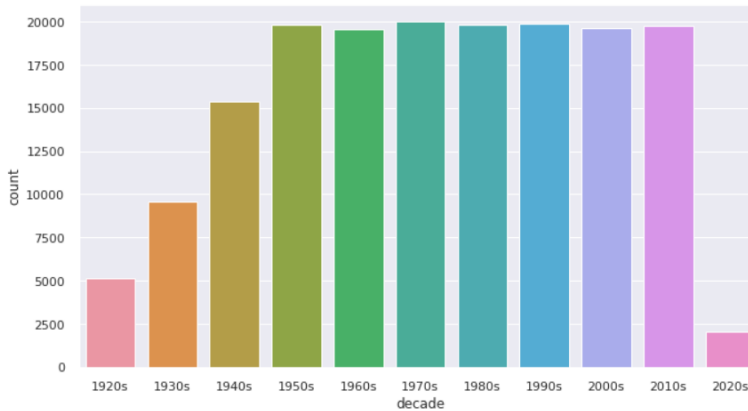
Features correlation with dependent variable

Features correlation with dependent variable

In [ ]:

```python
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[ ]: `<matplotlib.axes._subplots.AxesSubplot at 0x7f5798d61a10>`



In [ ]:

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

In [ ]:

```python
top10_genres = genre_data.nlargest(10, 'popularity')

fig = px.bar(top10_genres, x='genres', y=['valence', 'energy', 'danceability', 'acousticness'], barmode='group')
fig.show()
```

In [ ]:

```python
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 [ ]:

```python
from sklearn.manifold import TSNE

tsne_pipeline = Pipeline([('scaler', StandardScaler()), ('tsne', TSNE(n_components=2, verbose=1))])
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 0.006s...
[t-SNE] Computed neighbors for 2973 samples in 0.476s...
[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: 76.116959
[t-SNE] KL divergence after 1000 iterations: 1.388681
```

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

```
Collecting spotipy
  Downloading spotipy-2.19.0-py3-none-any.whl (27 kB)
Requirement already satisfied: six>=1.15.0 in /usr/local/lib/python3.7/dist-packages (from spotipy) (1.15.0)
Collecting requests>=2.25.0
  Downloading requests-2.27.1-py2.py3-none-any.whl (63 kB)
     |████████████████████████████████| 63 kB 1.5 MB/s
Collecting urllib3>=1.26.0
  Downloading urllib3-1.26.9-py2.py3-none-any.whl (138 kB)
     |████████████████████████████████| 138 kB 18.2 MB/s
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.7/dist-packages (from requests>=2.25.0->spotipy) (2.10)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from requests>=2.25.0->spotipy) (2021.10.8)
Requirement already satisfied: charset-normalizer~=2.0.0 in /usr/local/lib/python3.7/dist-packages (from requests>=2.25.0->spotipy) (2.0.12)
Installing collected packages: urllib3, requests, spotipy
  Attempting uninstall: urllib3
    Found existing installation: urllib3 1.24.3
    Uninstalling urllib3-1.24.3:
      Successfully uninstalled urllib3-1.24.3
  Attempting uninstall: requests
    Found existing installation: requests 2.23.0
    Uninstalling requests-2.23.0:
      Successfully uninstalled requests-2.23.0
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the followin
g dependency conflicts.
google-colab 1.0.0 requires requests~=2.23.0, but you have requests 2.27.1 which is incompatible.
datascience 0.10.6 requires folium==0.2.1, but you have folium 0.8.3 which is incompatible.
Successfully installed requests-2.27.1 spotipy-2.19.0 urllib3-1.26.9
```

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

          sp = spotipy.Spotify(auth_manager=SpotifyClientCredentials(client_id='ae86770e0e2946afaf50f25db5ba1a33',
                                                                    client_secret='7cab13d55eeb4dd296fdc99fd177bc08'))
          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 [ ]:   find_song('Kabira', '2013')
```

Out[ ]:

| | name | year | explicit | duration_ms | popularity | danceability | energy | key | loudness | mode | ... | instrumentalness | liveness | valence | tempo | type | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Kabira | 2013 | 0 | 223460 | 69 | 0.59 | 0.555 | 2 | -6.861 | 1 | ... | 0 | 0.192 | 0.191 | 84.027 | audio_features | 4bD9z9qa4qg9BhryvYWB |

1 rows × 22 columns

```python
from collections import 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', 'explicit',
 'instrumentalness', 'key', 'liveness', 'loudness', 'mode', 'popularity', 'speechiness', 'tempo']
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)
```

```python
        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')
```

```python
recommend_songs([{'name': 'Kabira', 'year':2013}],  data)
```

```
[{'artists': "['Lee Brice']", 'name': "I Don't Dance", 'year': 2014},
 {'artists': "['Lady A']", 'name': 'American Honey', 'year': 2010},
 {'artists': "['James Blunt']", 'name': 'Carry You Home', 'year': 2007},
 {'artists': "['Joji', 'BENEE']", 'name': 'Afterthought', 'year': 2020},
 {'artists': "['Zara Larsson']", 'name': 'Uncover', 'year': 2012},
 {'artists': "['ZAYN']", 'name': 'Good Years', 'year': 2017},
 {'artists': "['The Head and the Heart']",
  'name': 'People Need A Melody',
  'year': 2019},
 {'artists': "['Gustavo Cerati']", 'name': 'Crimen', 'year': 2006},
 {'artists': "['Kodaline']", 'name': 'High Hopes', 'year': 2013}]
```