

```
In [1]: import numpy as np
import pandas as pd

import seaborn as sns
import matplotlib.pyplot as plt

from collections import defaultdict
from scipy.spatial.distance import pdist
from sklearn.preprocessing import MinMaxScaler, StandardScaler

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.sample(5)

Out[3]:
```

	valence	year	acousticness	artists	danceability	duration_ms	energy	explicit	id	instrumentalness	key	liveness	loudness	mode	name	pop
14274	0.167	1978	0.506	[Little River Band]	0.511	262947	0.323	0	5ZpRH2XO7MaeCtHwUkZq	0.000114	5	0.3200	-10.028	1	Santay's Sole - 2010 Digital Remaster	1
59927	0.338	1944	0.963	[Richard Strauss; Franz Schubert; Lisa del...	0.374	753333	0.316	0	6tHtBvMyShYsRBNdAS	0.000000	10	0.8780	-17.502	0	Arabia/ Op. 78 / TV 263 / Act 2 - Ein Feigling...	0
52588	0.296	1996	0.107	[Coke]	0.436	201613	0.579	0	1pRtXvDCJ8XECBxP0wYd	0.000730	9	0.0996	-9.795	0	Friend is a Four Letter Word	0
15728	0.744	1956	0.872	[J.J. Johnson; Winding]	0.648	175227	0.481	0	1PsrRtEPYMytB5SXWCest	0.871000	0	0.1030	-10.219	1	The Continental (You Kiss While You're Dancing)	1
132557	0.634	1978	0.542	[Emmylou Harris]	0.475	189240	0.386	0	5VEUvXvUC4AhYsYcYCDMUkK	0.000000	3	0.0783	-11.423	1	Easy from New On - 2003 Remaster	1

```
In [4]: genre_data.sample(5)

Out[4]:
```

	mode	genres	acousticness	danceability	duration_ms	energy	instrumentalness	liveness	loudness	speechiness	tempo	valence	popularity	key
378	1	british modern classical	0.876342	0.288750	234667.338217	0.184769	0.417125	0.160951	-21.579563	0.051742	106.918516	0.221183	13.468956	2
8	1	deta blues	0.771977	0.579439	200740.886489	0.368764	0.113124	0.179875	-12.740045	0.076983	117.211747	0.599151	24.429914	10
2044	1	old/funk hip hop	0.124070	0.729840	236709.893915	0.628074	0.002893	0.164853	-7.680831	0.221095	107.575201	0.540111	47.764114	10
1139	0	gaming edm	0.089090	0.563561	245632.186650	0.810934	0.263486	0.200551	-5.112201	0.072846	132.444843	0.326599	58.918980	7
831	1	deep&nb	0.155000	0.426000	182707.000000	0.953000	0.000000	0.132000	-3.262000	0.106000	169.930000	0.678000	51.000000	9

```
In [5]: year_data.sample(5)

Out[5]:
```

	mode	year	acousticness	danceability	duration_ms	energy	instrumentalness	liveness	loudness	speechiness	tempo	valence	popularity	key
74	1	1995	0.302512	0.559046	246624.091500	0.578580	0.132269	0.201374	-10.119450	0.063754	119.281978	0.530247	44.801000	0
14	1	1935	0.778386	0.555869	220124.247036	0.246367	0.225873	0.229300	-15.414750	0.353912	110.607957	0.545876	1.501318	7
13	1	1934	0.981149	0.528706	189356.126298	0.262131	0.276382	0.213453	-14.756875	0.162457	115.390846	0.588805	1.257785	0
31	1	1952	0.874301	0.457032	229434.727000	0.253026	0.208413	0.232972	-15.956866	0.137362	108.232280	0.443135	3.338000	0
12	1	1933	0.899898	0.570290	196219.257598	0.279699	0.183949	0.209072	-13.069009	0.091123	112.522000	0.599410	6.896698	7

```
In [6]: artist_data.sample(5)

Out[6]:
```

	mode	count	acousticness	artists	danceability	duration_ms	energy	instrumentalness	liveness	loudness	speechiness	tempo	valence	popularity	key
7122	0	2	0.170000	Sandra De Sá	0.842000	267267.000000	0.549000	0.001220	0.080400	-10.857000	0.044700	96.566000	0.967000	51.000000	5
2190	1	39	0.949974	Eino-Pekka Salonen	0.274154	198321.179467	0.17727	0.834154	0.233482	-23.138128	0.048528	99.85359	0.133895	7.307692	0
5905	0	2	0.953000	Dave Thomas Junior	0.365000	330484.000000	0.158000	0.039200	0.096100	-14.019000	0.032000	123.531000	0.074000	58.000000	3
243	1	12	0.000521	ASIAN KUNG-FU GENERATION	0.417000	244089.000000	0.91850	0.092434	0.136000	-4.071500	0.065660	156.35350	0.489167	57.000000	4
14215	1	2	0.507000	Lauren O'Connell	0.415000	185263.000000	0.34300	0.818000	0.349000	-13.526000	0.032100	170.90400	0.485000	53.000000	7

```
In [7]: datasets = [{"data": data}, {"genre_data": genre_data}, {"year_data": year_data}, {"artist_data": artist_data}]

In [8]: data['release_date'] = pd.to_datetime(data['release_date'])
year_data['year'] = pd.to_datetime(year_data['format']).year
artist_data['year'] = pd.to_datetime(artist_data['year']).year
```

Checking the Data

```
In [9]: for name, df in datasets:
print("Info about the dataset: (name)")
print("-" * 30)
print(df.info())
print()

Info about the dataset: data
-----
<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 datetime64[ns]
 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 object
 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 datetime64[ns]
 17 speechiness 170653 non-null float64
 18 tempo 170653 non-null float64
 dtypes: datetime64[ns](2), float64(9), int64(5), object(3)
memory usage: 24.7+ MB
None

Info about the dataset: genre_data
-----
<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 object
 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

Info about the dataset: year_data
-----
<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 datetime64[ns]
 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: datetime64[ns](1), float64(11), int64(2)
memory usage: 11.1 KB
None

Info about the dataset: artist_data
-----
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 28680 entries, 0 to 28679
Data columns (total 15 columns):
 # Column Non-Null Count Dtype
---
 0 count 28680 non-null int64
 1 acousticness 28680 non-null float64
 2 artists 28680 non-null object
 3 danceability 28680 non-null float64
 4 duration_ms 28680 non-null float64
 5 energy 28680 non-null float64
 6 instrumentalness 28680 non-null float64
 7 liveness 28680 non-null float64
 8 loudness 28680 non-null float64
 9 popularity 28680 non-null int64
 10 release_date 28680 non-null datetime64[ns]
 11 speechiness 28680 non-null float64
 12 tempo 28680 non-null float64
 13 valence 28680 non-null float64
 14 key 28680 non-null int64
 dtypes: float64(11), int64(3), object(1)
memory usage: 3.3+ MB
None

In [10]: for name, df in datasets:
print("Missing Values in: (name)")
print("-" * 30)
print(df.isnull().sum())
print()

Missing Values in: data
-----
valence 0
year 0
acousticness 0
artists 0
danceability 0
duration_ms 0
energy 0
explicit 0
id 0
instrumentalness 0
key 0
liveness 0
loudness 0
mode 0
name 0
popularity 0
release_date 0
speechiness 0
tempo 0
dtype: int64

Missing Values in: genre_data
-----
mode 0
genres 0
acousticness 0
danceability 0
duration_ms 0
energy 0
instrumentalness 0
liveness 0
loudness 0
popularity 0
tempo 0
valence 0
dtype: int64

Missing Values in: year_data
-----
mode 0
year 0
acousticness 0
danceability 0
duration_ms 0
energy 0
instrumentalness 0
liveness 0
loudness 0
popularity 0
tempo 0
valence 0
dtype: int64

Missing Values in: artist_data
-----
mode 0
count 0
acousticness 0
artists 0
danceability 0
duration_ms 0
energy 0
instrumentalness 0
liveness 0
loudness 0
popularity 0
release_date 0
speechiness 0
tempo 0
valence 0
dtype: int64

In [11]: for name, df in datasets:
print("Duplicates in the dataset: (name)")
print("-" * 30)
print(df.duplicated(keep=False).sum())
print()

Duplicates in the dataset: data
-----
0

Duplicates in the dataset: genre_data
-----
0

Duplicates in the dataset: year_data
-----
0

Duplicates in the dataset: artist_data
-----
0

In [12]: for name, df in datasets:
print("Unique Values in: (name)")
print("-" * 30)
print(df.nunique())
print()

Unique Values in: data
-----
valence 1733
year 100
acousticness 4689
artists 34988
danceability 1246
duration_ms 51755
energy 2332
explicit 2
id 170653
instrumentalness 5481
key 12
liveness 7
loudness 25410
mode 2
name 133338
popularity 189
release_date 10988
speechiness 1526
tempo 84694
dtype: int64

Unique Values in: genre_data
-----
mode 2
genres 2973
acousticness 2798
danceability 2725
duration_ms 2872
energy 2778
instrumentalness 2731
liveness 2789
loudness 2873
speechiness 2797
tempo 2872
valence 2745
popularity 2188
key 12
dtype: int64

Unique Values in: year_data
-----
mode 1
year 100
acousticness 100
danceability 100
duration_ms 100
energy 100
instrumentalness 100
liveness 100
loudness 100
popularity 100
tempo 100
valence 100
dtype: int64

Unique Values in: artist_data
-----
mode 2
count 378
acousticness 14127
artists 28680
danceability 18613
duration_ms 23960
energy 12126
instrumentalness 15517
loudness 23862
speechiness 10950
tempo 24891
valence 11882
popularity 4663
key 12
dtype: int64
```

Clustering

```
In [13]: from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline

cluster_pipeline = Pipeline([('scaler', StandardScaler()), ('kmeans', KMeans(n_clusters=12))])
X = genre_data.select_dtypes(np.number)
cluster_pipeline.fit(X)
genre_data['cluster'] = cluster_pipeline.predict(X)

In [14]: 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'], title='Clusters of genres')
fig.show()

[TSNE] Computing 91 nearest neighbors...
[TSNE] Indexed 2973 samples in 0.322s...
[TSNE] Computed neighbors for 2973 samples in 0.412s...
[TSNE] Computed conditional probabilities for sample 1600 / 2973
[TSNE] Computed conditional probabilities for sample 2000 / 2973
[TSNE] Computed conditional probabilities for sample 2973 / 2973
[TSNE] Mean sigma: 0.777516
[TSNE] KL divergence after 250 iterations with early exaggeration: 76.166277
[TSNE] KL divergence after 1000 iterations: 1.392958

Clusters of genres
```

```
In [15]: song_cluster_pipeline = Pipeline([('scaler', StandardScaler()), ('kmeans', KMeans(n_clusters=25,
                                             verbose=False))])
X = data.select_dtypes(np.number)
song_cluster_pipeline.fit(X)
song_cluster_labels = song_cluster_pipeline.predict(X)
data['cluster_label'] = song_cluster_labels

In [16]: 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'], title='Clusters of songs')
fig.show()
```

Building the recommender system

```
In [17]: data['year'] = data['year'].year

In [18]: number_cols = ['valence', 'year', 'acousticness', 'danceability', 'duration_ms', 'energy', 'explicit', 'year',
                    'instrumentalness', 'key', 'liveness', 'loudness', 'mode', 'popularity', 'speechiness', 'tempo']

In [19]: def get_song_data(name, data):
    try:
        return data[data['name'].str.lower() == name].iloc[0]
    except IndexError:
        return None

In [20]: def get_mean_vector(song_list, data):
    song_vectors = []
    for song in song_list:
        song_data = get_song_data(song['name'], data)
        if song_data is None:
            print("Warning: {} does not exist in the dataset".format(song['name']))
            return None
        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)

In [21]: def flatten_dict_list(dict_list):
    flattened_dict = defaultdict(list)
    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

In [22]: min_max_scaler = MinMaxScaler()
normalized_data = min_max_scaler.fit_transform(data[number_cols])
standard_scaler = StandardScaler()
scaled_normalized_data = standard_scaler.fit_transform(normalized_data)

In [23]: def recommend_songs(seed_songs, data, n_recommendations=10):
    metadata_cols = ['name', 'artists', 'year']
    song_center = get_mean_vector(seed_songs, data)

    if song_center is None:
        return []

    normalized_song_center = min_max_scaler.transform([song_center])
    scaled_normalized_song_center = standard_scaler.transform(normalized_song_center)
    distances = cdist(scaled_normalized_song_center, scaled_normalized_data, 'euclidean')
    indices = np.argsort(distances)[0]

    rec_songs = []
    for i in index:
        song_name = data.iloc[i]['name']
        if song_name not in [song['name'] for song in seed_songs] and song_name not in [song['name'] for song in rec_songs]:
            rec_songs.append(song_name)
            if len(rec_songs) == n_recommendations:
                break
    return pd.DataFrame(rec_songs)[metadata_cols].to_dict(orient='records')

In [24]: my_songs = [
    'Stumble',
    'Your Song',
]
my_songs = [{"name": name} for name in my_songs]
n_recommendations = 20

recommended_songs = recommend_songs(my_songs, data, n_recommendations)
recommended_df = pd.DataFrame(recommended_songs)

for idx, song in enumerate(recommended_songs, start=1):
    print("{}({}). [song['name']] by [song['artists']] ([song['year']])")

1. Can't Take My Eyes Off You by ['Erkin Koray'] (1968)
2. Candle in the Wind - Remastered 2014 by ['Elton John'] (1973)
3. New Kid in Town - 2013 Remaster by ['Eagles'] (1976)
4. Reasons by ['Earth, Wind & Fire'] (1975)
5. Skyline Pigeon - Piano Version by ['Elton John'] (1973)
6. Fly Me to the Moon (in other words) by ['Bobby Darin'] (1964)
7. It's All Over Now, Baby Blue (Feat. Van Morrison) by ['The Beatles', 'Van Morrison'] (1966)
8. Seni Her Gövdüme by ['Erkin Koray'] (1974)
9. Goodbye - Low Budget - Remastered 2014 by ['Elton John'] (1973)
10. Kissed by the Light - Single edit by ['Manfred Mann's Earth Band'] (1976)
11. Voicé les clés by ['Gerard Lenorman'] (1976)
12. I Never Loved a Man (The Way I Love You) by ['Aretha Franklin'] (1967)
13. Beautiful Boy (Darling Boy) - Stripped Down Mix, 2010 by ['John Lennon'] (1980)
14. I Can See Clearly Now - Edit by ['Johnny Nash'] (1972)
15. I Can See Clearly Now by ['Johnny Nash'] (1972)
16. Get Together by ['The Youngbloods'] (1967)
17. Metamorphose Ambulante by ['Raül Seixas'] (1973)
18. Stay by ['Oingo Boingo'] (1985)
19. Que Te Voya Bonito by ['Vicente Fernández'] (1972)
20. Maggie May by ['Rod Stewart'] (1971)
```