

1 CaseCraft: The Analytics Sprint – Project 3

1.0.1 Spotify Playlist Recommender

Subheading: Building a content-based recommender system using audio features to suggest similar tracks within a playlist.

1.0.2 Key Objectives

- Load and explore Spotify track metadata
- Engineer features for similarity comparison
- Build a content-based recommender using cosine similarity
- Visualize audio feature distributions and clusters
- Generate recommendations for a given track
- Summarize insights and potential improvements

1.0.3 Requirements

- pandas
- numpy
- sklearn
- matplotlib
- seaborn

1.0.4 Dataset Description

Simulated Spotify playlist dataset with ~500 tracks containing:

- `track_name`: Name of the song

- artist: Artist name
- genre: Genre label
- danceability, energy, valence, tempo, acousticness, instrumentalness: Audio features
- duration_ms: Track length in milliseconds

```
[7]: import pandas as pd
import numpy as np
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
import seaborn as sns

# Simulate dataset
np.random.seed(42)
genres = ['Pop', 'Rock', 'Hip-Hop', 'Jazz', 'EDM']
artists = [f'Artist_{i}' for i in range(1, 21)]
tracks = [f'Track_{i}' for i in range(1, 501)]

data = []
for track in tracks:
    artist = np.random.choice(artists)
    genre = np.random.choice(genres)
    features = np.random.rand(6) # Changed from 7 to 6
    duration = np.random.randint(180000, 300000)
    data.append([track, artist, genre, *features, duration])

columns = ['track_name', 'artist', 'genre', 'danceability', 'energy', 'valence',
           'tempo', 'acousticness', 'instrumentalness', 'duration_ms'] # 10
# Matches 10 values

df = pd.DataFrame(data, columns=columns)
df.head()
```

```
[7]: track_name  artist  genre  danceability  energy  valence  tempo \
0    Track_1  Artist_7   Jazz    0.950714  0.731994  0.598658  0.156019
1    Track_2  Artist_4  Hip-Hop    0.020584  0.969910  0.832443  0.212339
2    Track_3  Artist_12  Pop      0.291229  0.611853  0.139494  0.292145
3    Track_4  Artist_19  Jazz     0.514234  0.592415  0.046450  0.607545
4    Track_5  Artist_14  Rock     0.385417  0.015966  0.230894  0.241025

acousticness  instrumentalness  duration_ms
0      0.155995           0.058084      292727
1      0.181825           0.183405      233707
2      0.366362           0.456070      264654
```

3	0.170524	0.065052	190627
4	0.683264	0.609997	197159

1.0.5 Feature Engineering and Scaling

Select the audio features and scale them for the recommender system.

```
[8]: # Select audio features for recommendation
audio_features = ['danceability', 'energy', 'valence', 'tempo', 'acousticness', 'instrumentalness', 'duration_ms']
df_features = df[audio_features]

# Scale the features
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df_features)

# Convert scaled features back to a DataFrame for easier handling (optional)
df_scaled = pd.DataFrame(df_scaled, columns=audio_features)
df_scaled.head()
```

```
[8]:  danceability    energy    valence    tempo    acousticness    instrumentalness \
0      1.569121    0.779174    0.378873   -1.155217      -1.219452      -1.485185
1     -1.631279    1.620691    1.198801   -0.955120      -1.128224      -1.047884
2     -0.700042    0.354231   -1.231510   -0.671584      -0.476475      -0.096430
3      0.067277    0.285477   -1.557832    0.448981      -1.168136      -1.460871
4     -0.375961   -1.753442   -0.910952   -0.853202       0.642760       0.440691

    duration_ms
0      1.527138
1     -0.182120
2      0.714126
3     -1.429745
4     -1.240574
```

1.0.6 Distribution of Danceability Across Genres

Visualizing how danceability varies by genre.

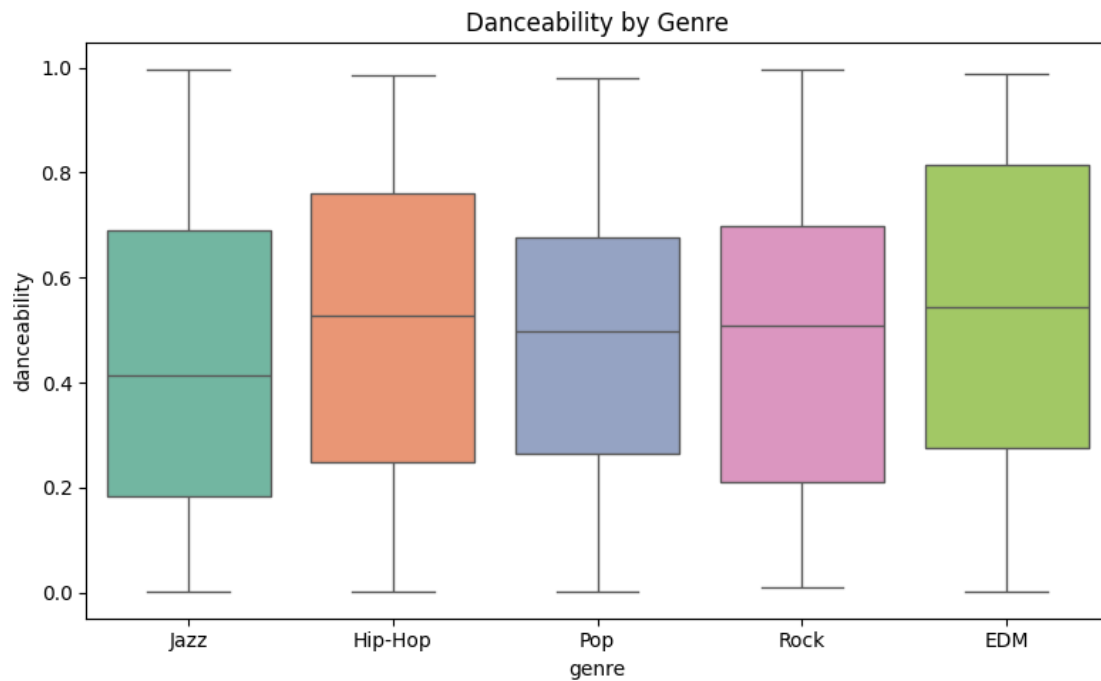
```
[9]: plt.figure(figsize=(8, 5))
sns.boxplot(data=df, x='genre', y='danceability', palette='Set2')
plt.title("Danceability by Genre")
plt.tight_layout()
plt.show()
```

/tmp/ipython-input-1254055711.py:2: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same

effect.

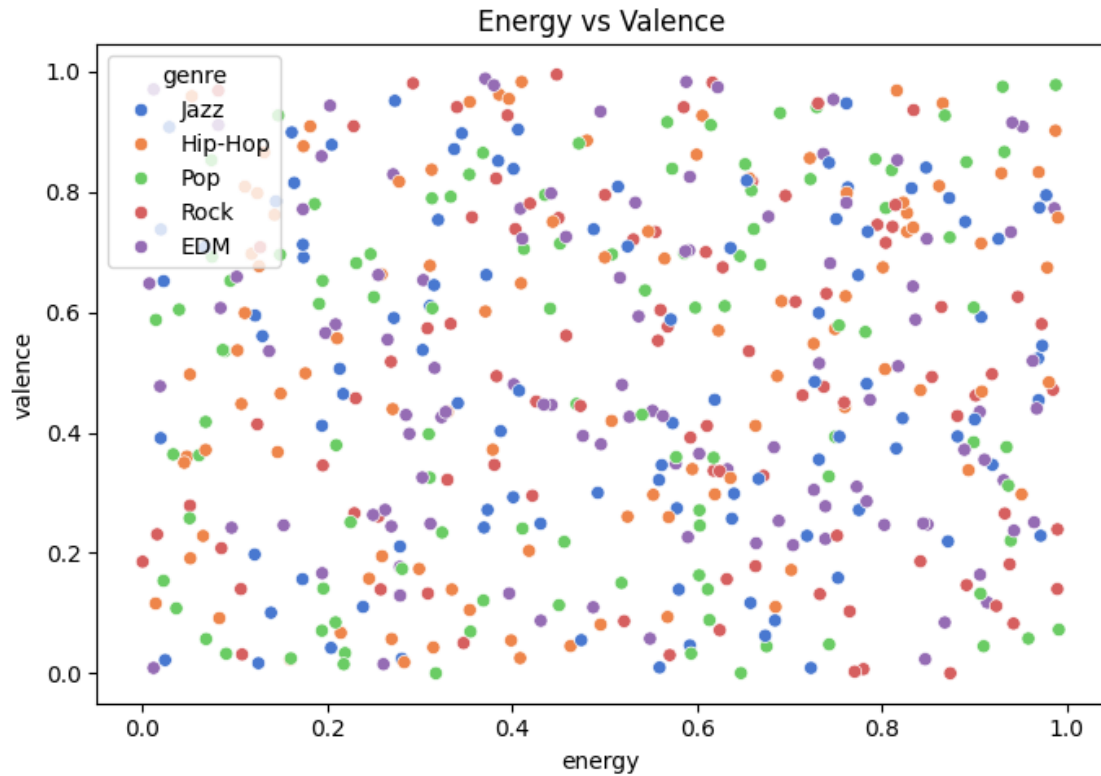
```
sns.boxplot(data=df, x='genre', y='danceability', palette='Set2')
```



1.0.7 Energy vs Valence Scatter Plot

Exploring emotional tone and intensity of tracks.

```
[10]: plt.figure(figsize=(7, 5))
sns.scatterplot(data=df, x='energy', y='valence', hue='genre', palette='muted')
plt.title("Energy vs Valence")
plt.tight_layout()
plt.show()
```

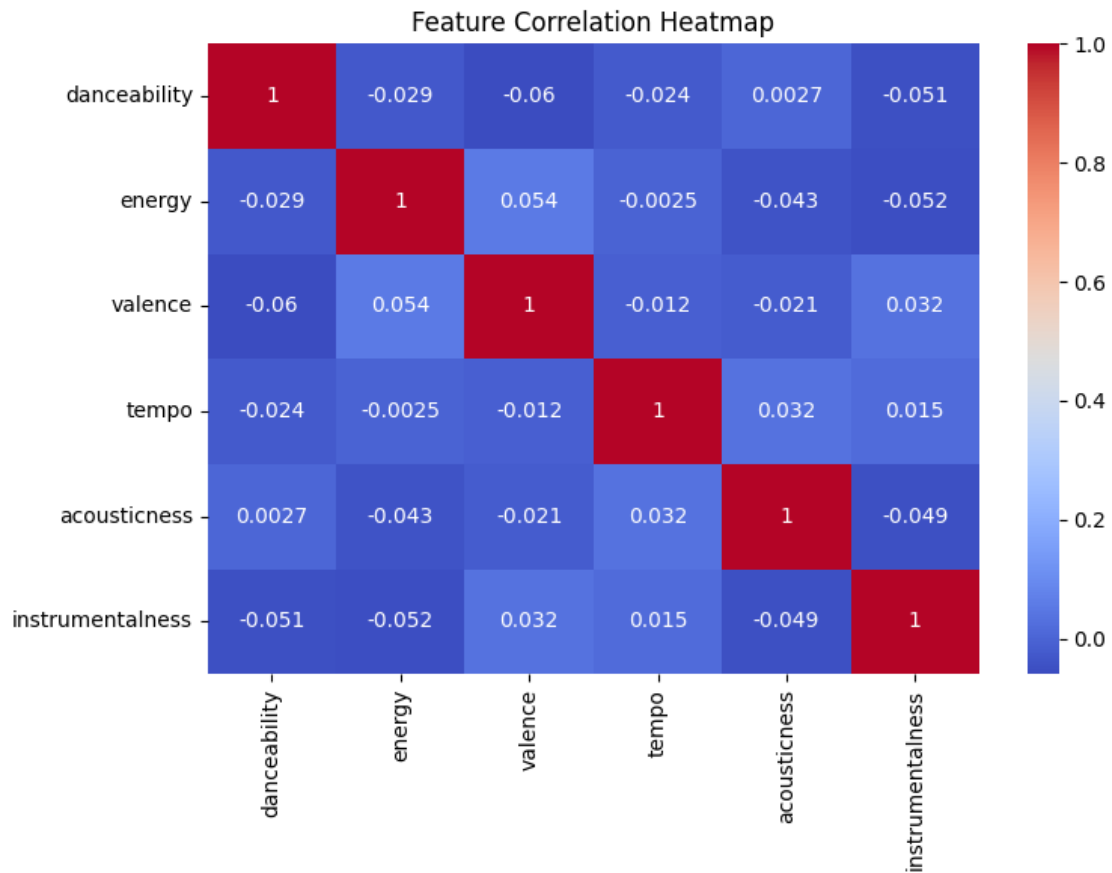


1.0.8 Correlation Heatmap of Audio Features

Understanding relationships between features.

```
[11]: features = ['danceability', 'energy', 'valence', 'tempo', 'acousticness',
↳ 'instrumentalness']
corr = df[features].corr()

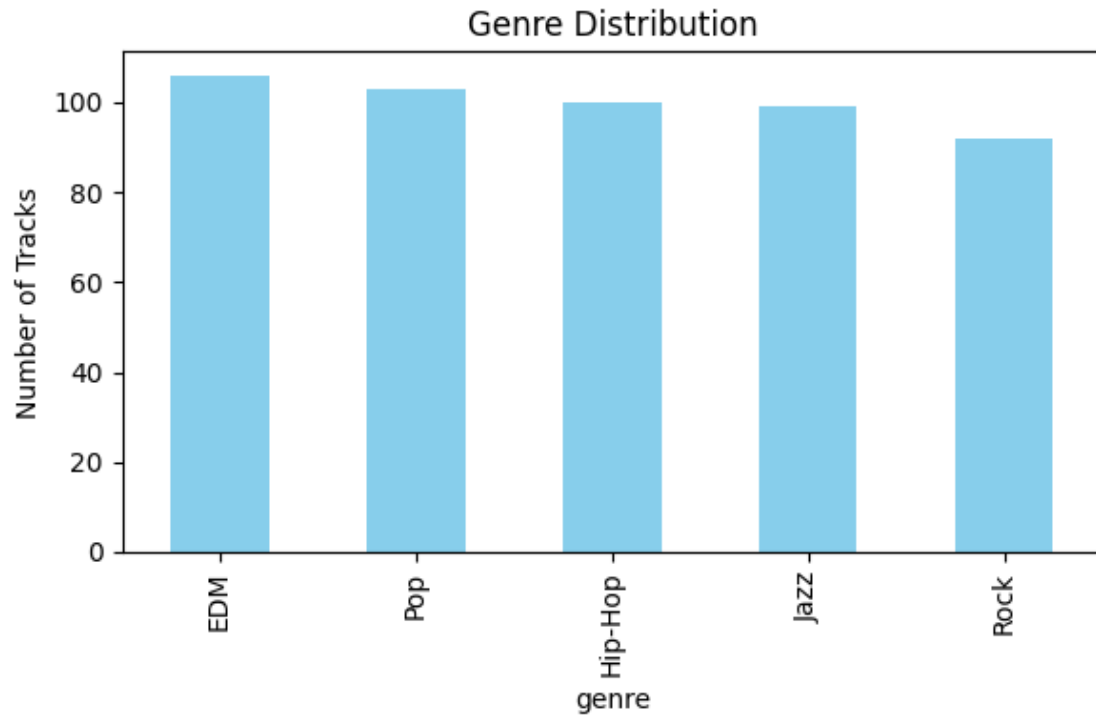
plt.figure(figsize=(8, 6))
sns.heatmap(corr, annot=True, cmap='coolwarm')
plt.title("Feature Correlation Heatmap")
plt.tight_layout()
plt.show()
```



1.0.9 Genre Distribution

Bar chart showing genre frequency in the playlist.

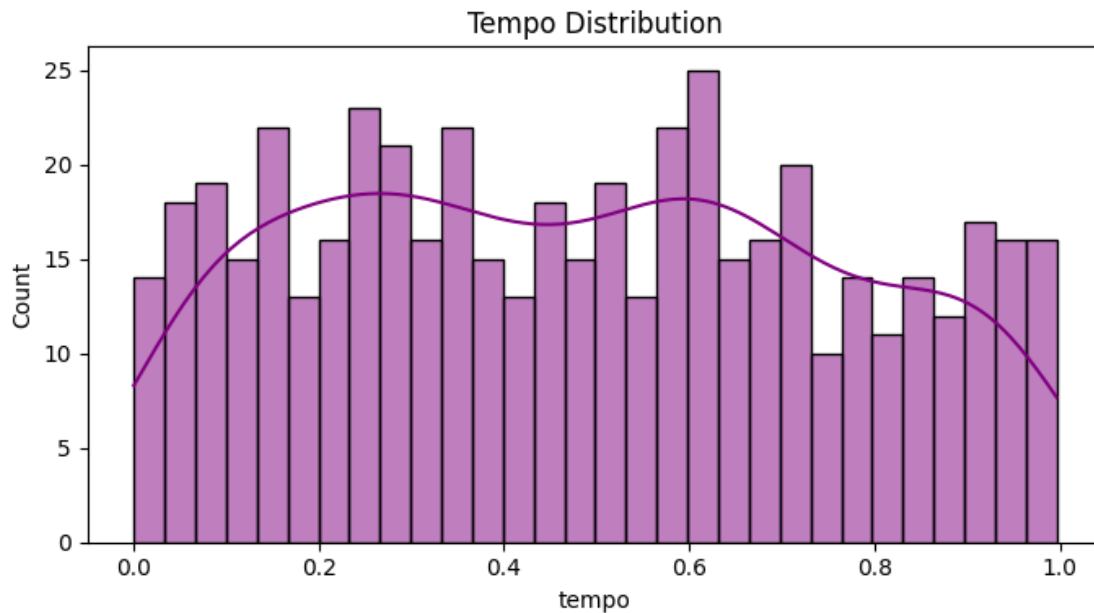
```
[12]: plt.figure(figsize=(6, 4))
df['genre'].value_counts().plot(kind='bar', color='skyblue')
plt.title("Genre Distribution")
plt.ylabel("Number of Tracks")
plt.tight_layout()
plt.show()
```



1.0.10 Tempo Distribution Histogram

Visualizing tempo spread across tracks.

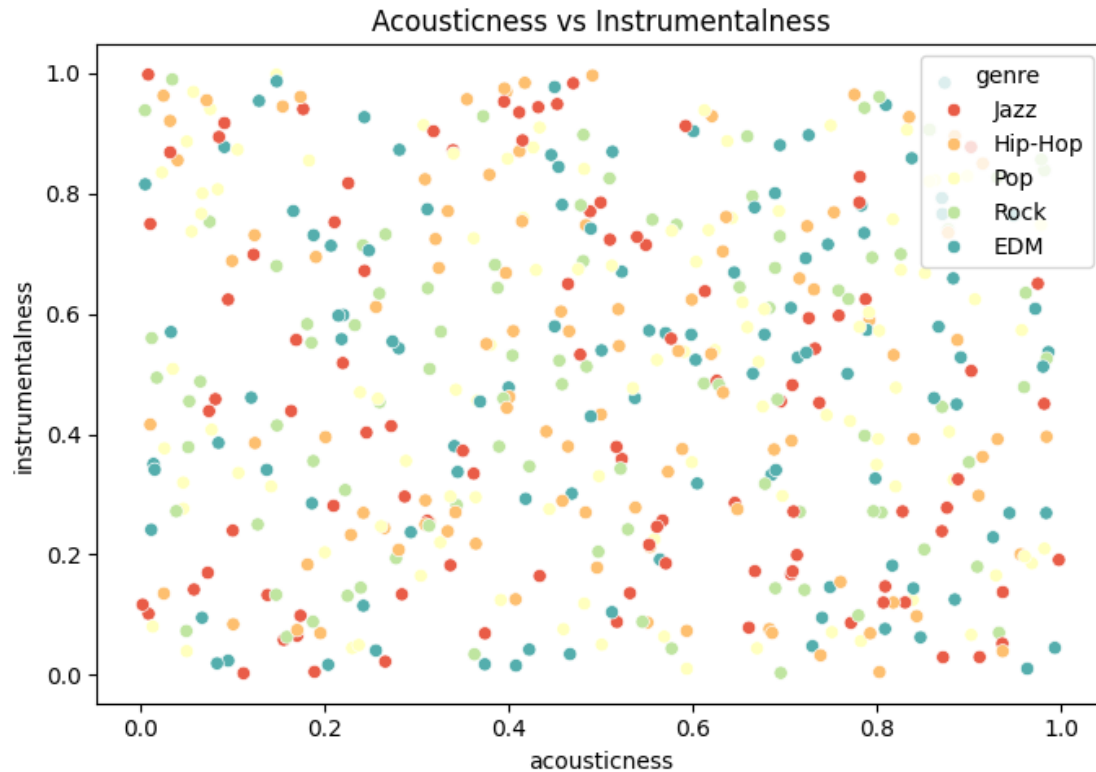
```
[13]: plt.figure(figsize=(7, 4))
sns.histplot(df['tempo'], bins=30, kde=True, color='purple')
plt.title("Tempo Distribution")
plt.tight_layout()
plt.show()
```



1.0.11 Acousticness vs Instrumentalness

Scatter plot to explore ambient vs instrumental nature.

```
[14]: plt.figure(figsize=(7, 5))
sns.scatterplot(data=df, x='acousticness', y='instrumentalness', hue='genre',
               palette='Spectral')
plt.title("Acousticness vs Instrumentalness")
plt.tight_layout()
plt.show()
```

1.0.12 Recommender System Logic

```
[15]: # Select features for similarity
feature_cols = ['danceability', 'energy', 'valence', 'tempo', 'acousticness', 'instrumentalness']

scaler = StandardScaler()
scaled_features = scaler.fit_transform(df[feature_cols])

# Compute cosine similarity
similarity_matrix = cosine_similarity(scaled_features)

# Recommendation function
def recommend(track_name, top_n=5):
    idx = df[df['track_name'] == track_name].index[0]
    sim_scores = list(enumerate(similarity_matrix[idx]))
    sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)[1:top_n+1]
    recommendations = df.iloc[[i[0] for i in sim_scores]][['track_name', 'artist', 'genre']]
    return recommendations

# Example usage
```

```
recommend('Track_42')
```

```
[15]:
```

	track_name	artist	genre
295	Track_296	Artist_5	Jazz
43	Track_44	Artist_4	Rock
176	Track_177	Artist_20	EDM
408	Track_409	Artist_17	Pop
323	Track_324	Artist_12	Hip-Hop

1.0.13 Report Summary and Conclusion

1. Feature Insights

- Danceability and energy vary significantly across genres
- Valence and energy show moderate correlation, suggesting emotional tone clusters
- Tempo distribution is centered around 120–140 BPM, typical for pop and EDM

2. Recommender Logic

- Content-based filtering using cosine similarity on scaled audio features
- Tracks with similar acoustic profiles are recommended regardless of artist or genre

3. Sample Output

Track_42 returns 5 similar tracks based on audio features, enabling playlist expansion

4. Business Implications

- Can be used to auto-curate playlists based on mood or genre
- Enhances user engagement through personalized discovery
- Scalable to larger datasets with real Spotify API integration