

Capstone Project 3

Course#3 Machine Learning

Santanu Chandra

email: santanu.chandra@gmail.com

Submission date - 08-09-2025

Creating Cohorts of Songs

Problem Scenario:

The customer always looks forward to specialized treatment, whether shopping on an e-commerce website or watching Netflix. The customer desires content that aligns with their preferences. To maintain customer engagement, companies must consistently provide the most relevant information.

Starting with Spotify, a Swedish audio streaming and media service provider, boasts over 456 million active monthly users, including more than 195 million paid subscribers as of September 2022. The company aims to create cohorts of different songs to enhance song recommendations. These cohorts will be based on various relevant features, ensuring that each group contains similar types of songs.

Problem Objective:

As a data scientist, you should perform exploratory data analysis and cluster analysis to create cohorts of songs. The goal is to better understand the various factors that create a cohort of songs.

Data Description:

The dataset comprises information from Spotify's API regarding all albums by the Rolling Stones available on Spotify. It's crucial to highlight that each song possesses a unique ID.

Steps to Perform:

1. Initial data inspection and data cleaning:

- a. Examine the data initially to identify duplicates, missing values, irrelevant entries,

```
import pandas as pd
import numpy as np
import scipy.stats
import seaborn as sns
from sklearn.metrics.pairwise import cosine_similarity
import matplotlib.pyplot as plt
import statistics
import operator
```

1.a. Load the data

```
# load the data from 'rolling_stones_spotify' into a DataFrame named 'spotify'.
spotify_df = pd.read_csv('rolling_stones_spotify.csv')
# Display the first five rows of the DataFrame using the 'head()' method to provide a quick
spotify_df.head(5)
```

		Unnamed: 0	name	album	release_date	track_number			
		0	Concert Intro Music - Live	Licked Live In NYC			1	2IEkywLJ4ykbhi1yRQvmsT	sp
		1	Fighting Man - Live	Licked Live In NYC	2022-06-10		2	6GVgVJBKkGJoRfarYRvGTU	spot
		2	Start Me Up - Live	Licked Live In NYC	2022-06-10		3	1Lu761pZ0dBTGpzaQoZNW	spoti
		3	If You Can't Rock Me - Live	Licked Live In NYC	2022-06-10		4	1agTQzOTUnGNggycxEqiDH	spot
		4	Don't Stop - Live	Licked Live In NYC	2022-06-10		5	7piGJR8YndQBQWVXv6KtQw	spotif

1.b Summary of the Dataframe

```
# Use the 'info()' method to output a concise summary of the DataFrame, including the number
spotify_df.info()
```

```
##identify type of dataset
#spotify_df.dtypes
```

```
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1610 entries, 0 to 1609
Data columns (total 18 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Unnamed: 0        1610 non-null    int64  
 1   name              1610 non-null    object  
 2   album              1610 non-null    object  
 3   release_date      1610 non-null    object  
 4   track_number      1610 non-null    int64  
 5   id                1610 non-null    object  
 6   uri               1610 non-null    object  
 7   acousticness      1610 non-null    float64 
 8   danceability      1610 non-null    float64 
 9   energy             1610 non-null    float64 
 10  instrumentalness  1610 non-null    float64 
 11  liveness          1610 non-null    float64 
 12  loudness          1610 non-null    float64 
 13  speechiness       1610 non-null    float64 
 14  tempo              1610 non-null    float64 
 15  valence            1610 non-null    float64 
 16  popularity         1610 non-null    int64  
 17  duration_ms        1610 non-null    int64  
dtypes: float64(9), int64(4), object(5)
memory usage: 226.5+ KB
```

Summary

- Type- object 5
- Type integer 4
- Type float 9

Total parameters = 18

1.c Identify missing values

```
# Calculates the total number of missing (null) values in each column using the 'isnull().sum()' method
spotify_df.isnull().sum()
```

	0
Unnamed: 0	0
name	0
album	0
release_date	0
track_number	0
id	0
uri	0
acousticness	0
danceability	0
energy	0
instrumentalness	0
liveness	0
loudness	0
speechiness	0
tempo	0
valence	0
popularity	0
duration_ms	0

dtype: int64

Observations:

- There are 18 columns of data for each of 1610 entries
- No missing values in any of the columns

```
# In case there were missing values
```

```
# Drop rows with missing values in key features
features = ['acousticness','danceability','energy','instrumentalness','liveness','loudness',
spotify_df_clean = spotify_df.dropna(subset=features)
```

```
# Generate descriptive statistics for the numerical columns of the 'animes' DataFrame
spotify_df.describe()
```

	Unnamed: 0	track_number	acousticness	danceability	energy	instrumentalness
count	1610.000000	1610.000000	1610.000000	1610.000000	1610.000000	1610.000000
mean	804.500000	8.613665	0.250475	0.468860	0.792352	0.164175
std	464.911282	6.560220	0.227397	0.141775	0.179886	0.276200
min	0.000000	1.000000	0.000009	0.104000	0.141000	0.000000
25%	402.250000	4.000000	0.058350	0.362250	0.674000	0.000000
50%	804.500000	7.000000	0.183000	0.458000	0.848500	0.013750
75%	1206.750000	11.000000	0.403750	0.578000	0.945000	0.179000
max	1609.000000	47.000000	0.994000	0.887000	0.999000	0.996000

find unique values

```
# Calculate the number of unique albums in the spotify_df dataset
len(spotify_df['album'].unique())
```

→ 90

```
# Calculate the number of unique name in the spotify_df dataset
len(spotify_df['name'].unique())
```

→ 954

```
# Calculate the number of unique release data in the spotify_df dataset
len(spotify_df['release_date'].unique())
```

→ 57

```
# Calculate the number of unique id in the spotify_df dataset
len(spotify_df['id'].unique())
```

→ 1610

```
# Calculate the number of unique uri in the spotify_df dataset
len(spotify_df['uri'].unique())
```

→ 1610

```
len(spotify_df['Unnamed: 0'].unique())
```

→ 1610

Observations:

- The Unnamed: 0 is the number count of the songs , counts from 0 to 1609 - total of 1610 songs
- The id and and uri are unique to each song
- There are 90 unique albums in the list
- There are unique 954 name in the list ? What is the definition of the name ?
- There are unique 57 release dates, that means multiple albums might have released on the same date

```
import statistics

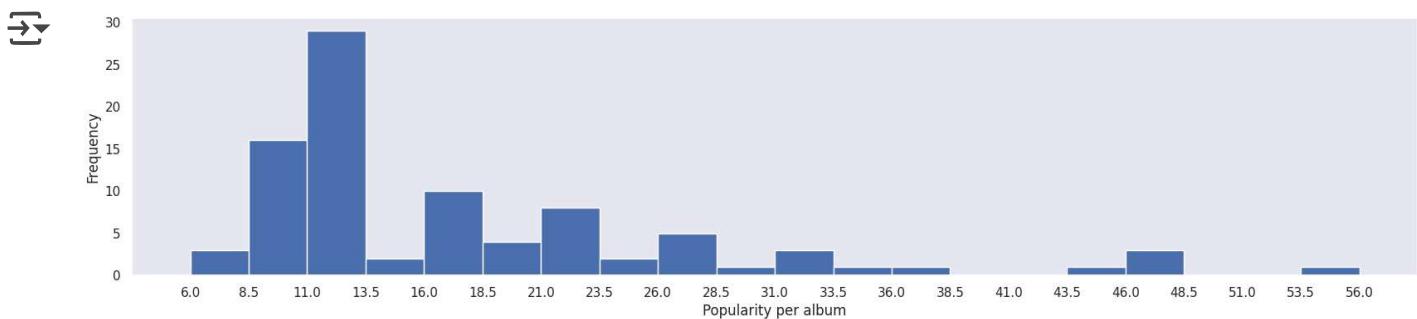
# Calculate the popularity per album
popularity_per_album = spotify_df.groupby('album')['popularity'].count()

# Calculate the mean popularity per album
mean_popularity_per_album = popularity_per_album.mean()

# Print the results
mean_popularity_per_album
```

→ np.float64(17.8888888888889)

```
# Plot the histogram and capture the return values
plt.figure(figsize=(20, 4))
counts, bins, patches = plt.hist(popularity_per_album, bins=20)
# Customize the x-axis labels to show bin edges
plt.xticks(bins) # Set x-axis ticks to match bin edges
plt.xlabel('Popularity per album')
plt.ylabel('Frequency')
plt.show()
```



Observations:

- Majority of the albums have popularity between 10-16 with only a few with popularity more than 50.
- The mean popularity is 20.7
- The max is 80 (dont see it in the histogram)
- The mean popularity per album is 17.8

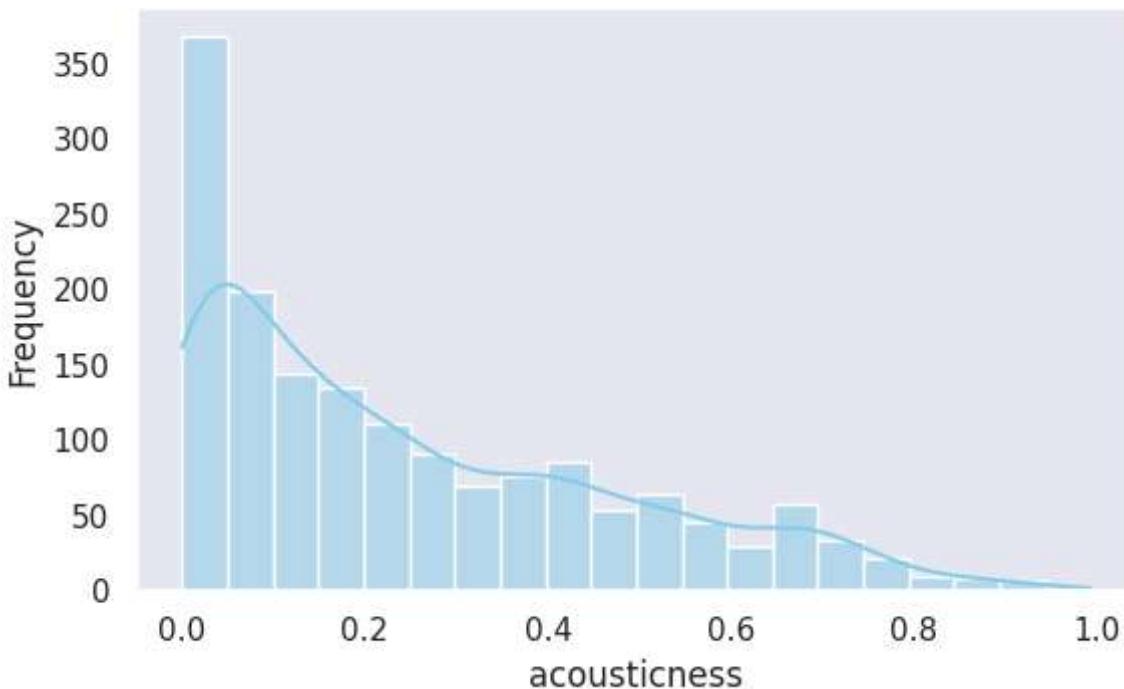
1.c PPlot Distributions

```
# Set up the plotting style
sns.set(style="dark")

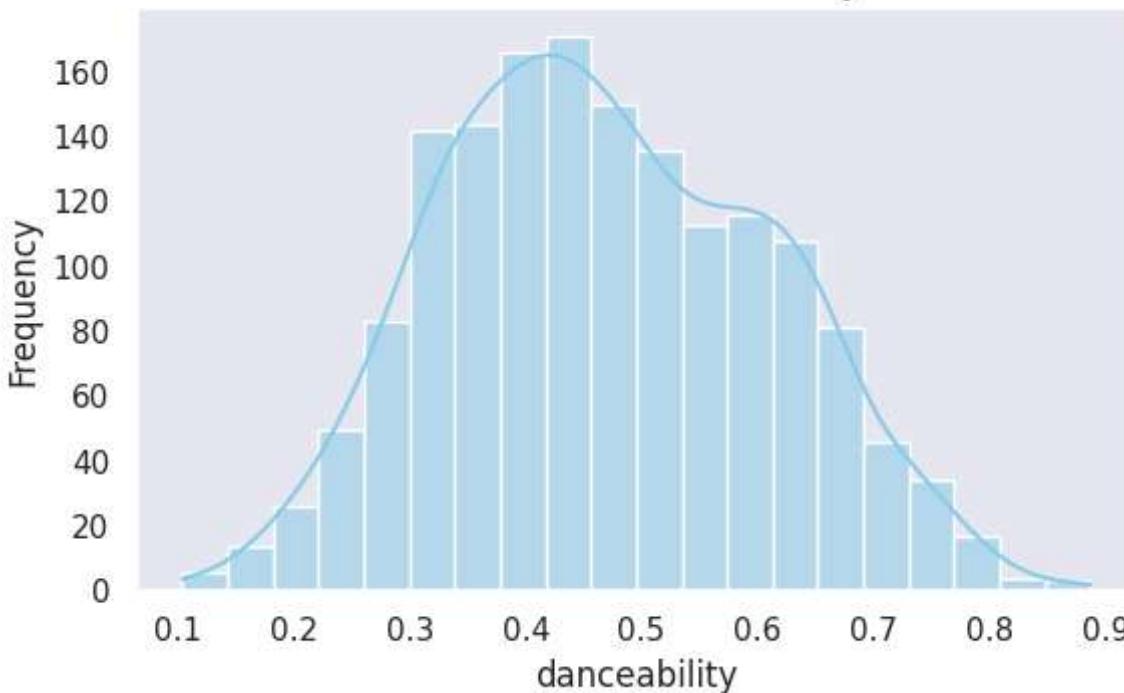
# Histograms for each feature
for feature in features:
    plt.figure(figsize=(6, 4))
    sns.histplot(spotify_df[feature], kde=True, bins=20, color='skyblue')
    plt.title(f'Distribution of {feature}')
    plt.xlabel(feature)
    plt.ylabel('Frequency')
    plt.tight_layout()
    plt.savefig(f"{feature}_histogram.png")
```



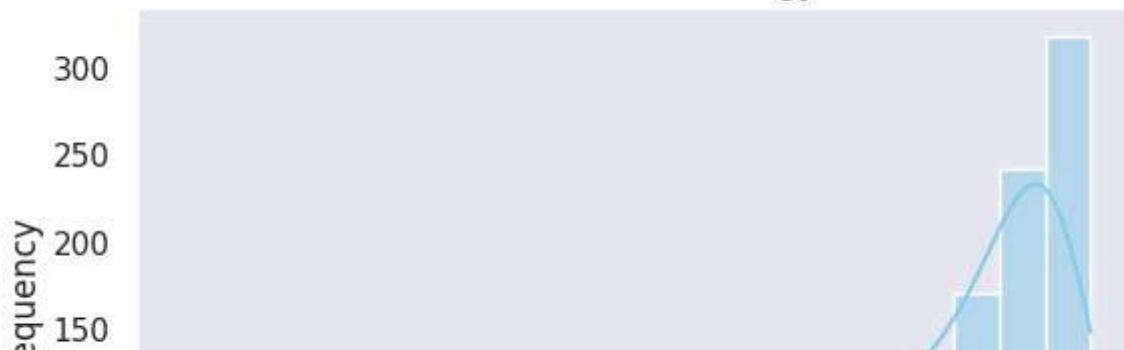
Distribution of acousticness

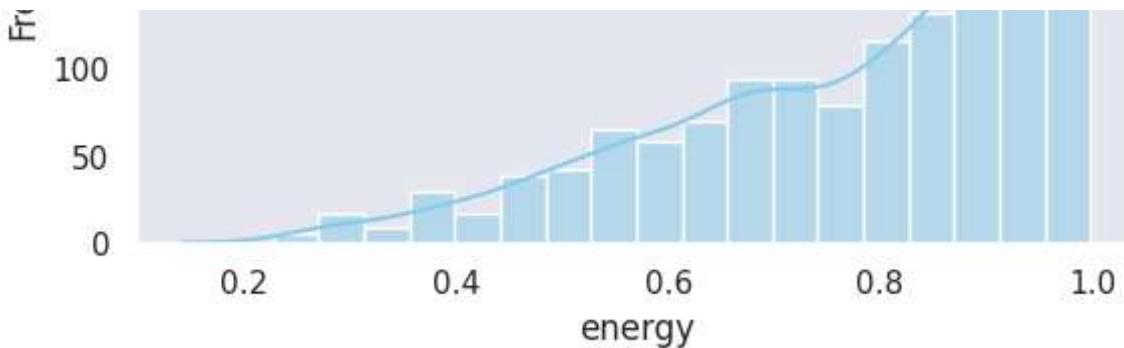


Distribution of danceability

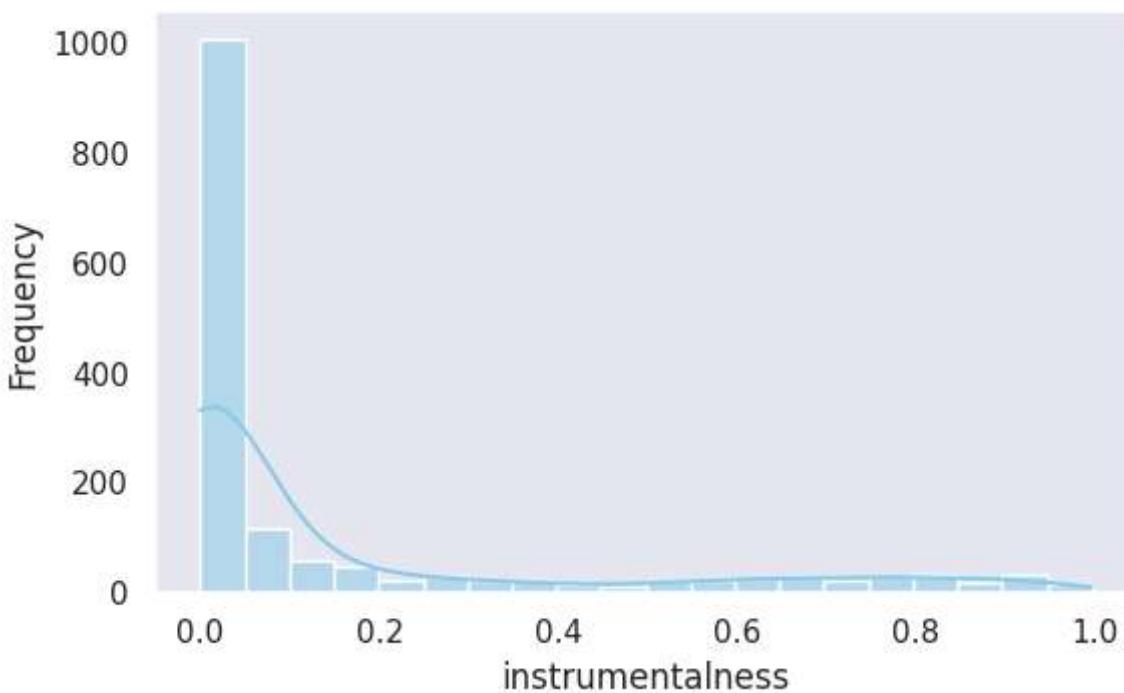


Distribution of energy

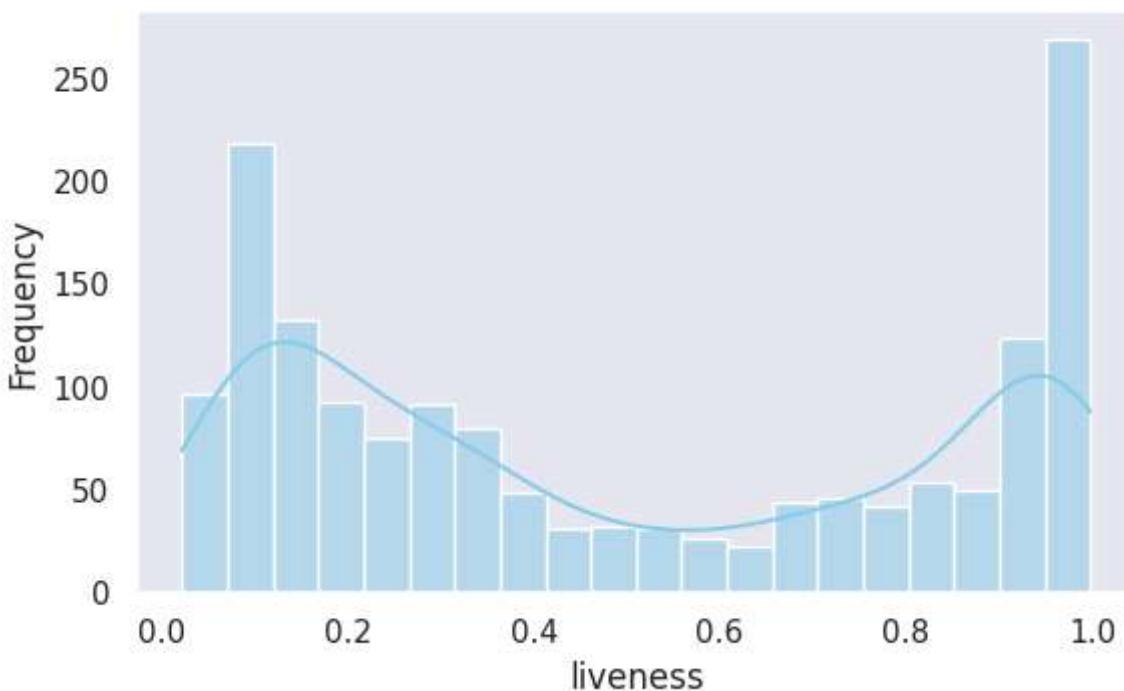




Distribution of instrumentalness



Distribution of liveness



Distribution of loudness

Observations

- Energy: Most songs are highly energetic.
- Danceability: Moderate to high danceability is common.
- Valence: Emotional positivity varies widely
- Loudness: Songs tend to be loud, typical of rock music
- Tempo: A broad range of tempos, with some clustering around 120 BPM.

100

150

200

250

300

350

400

450

500

550

600

650

700

750

800

850

900

950

1000

1050

1100

1150

1200

1250

1300

1350

1400

1450

1500

1550

1600

1650

1700

1750

1800

1850

1900

1950

2000

2050

2100

2150

2200

2250

2300

2350

2400

2450

2500

2550

2600

2650

2700

2750

2800

2850

2900

2950

3000

3050

3100

3150

3200

3250

3300

3350

3400

3450

3500

3550

3600

3650

3700

3750

3800

3850

3900

3950

4000

4050

4100

4150

4200

4250

4300

4350

4400

4450

4500

4550

4600

4650

4700

4750

4800

4850

4900

4950

5000

5050

5100

5150

5200

5250

5300

5350

5400

5450

5500

5550

5600

5650

5700

5750

5800

5850

5900

5950

6000

6050

6100

6150

6200

6250

6300

6350

6400

6450

6500

6550

6600

6650

6700

6750

6800

6850

6900

6950

7000

7050

7100

7150

7200

7250

7300

7350

7400

7450

7500

7550

7600

7650

7700

7750

7800

7850

7900

7950

8000

8050

8100

8150

8200

8250

8300

8350

8400

8450

8500

8550

8600

8650

8700

8750

8800

8850

8900

8950

9000

9050

9100

9150

9200

9250

9300

9350

9400

9450

9500

9550

9600

9650

9700

9750

9800

9850

9900

9950

10000

10050

10100

10150

10200

10250

10300

10350

10400

10450

10500

10550

10600

10650

10700

10750

10800

10850

10900

10950

11000

11050

11100

11150

11200

11250

11300

11350

11400

11450

11500

11550

11600

11650

11700

11750

11800

11850

11900

11950

12000

12050

12100

12150

12200

12250

12300

12350

12400

12450

12500

12550

12600

12650

12700

12750

12800

12850

12900

12950

13000

13050

13100

13150

13200

13250

13300

13350

13400

13450

13500

13550

13600

13650

13700

13750

13800

13850

13900

13950

14000

14050

14100

14150

14200

14250

14300

14350

14400

14450

14500

14550

14600

14650

14700

14750

14800

14850

14900

14950

15000

15050

15100

15150

15200

15250

15300

15350

15400

15450

15500

15550

15600

15650

15700

15750

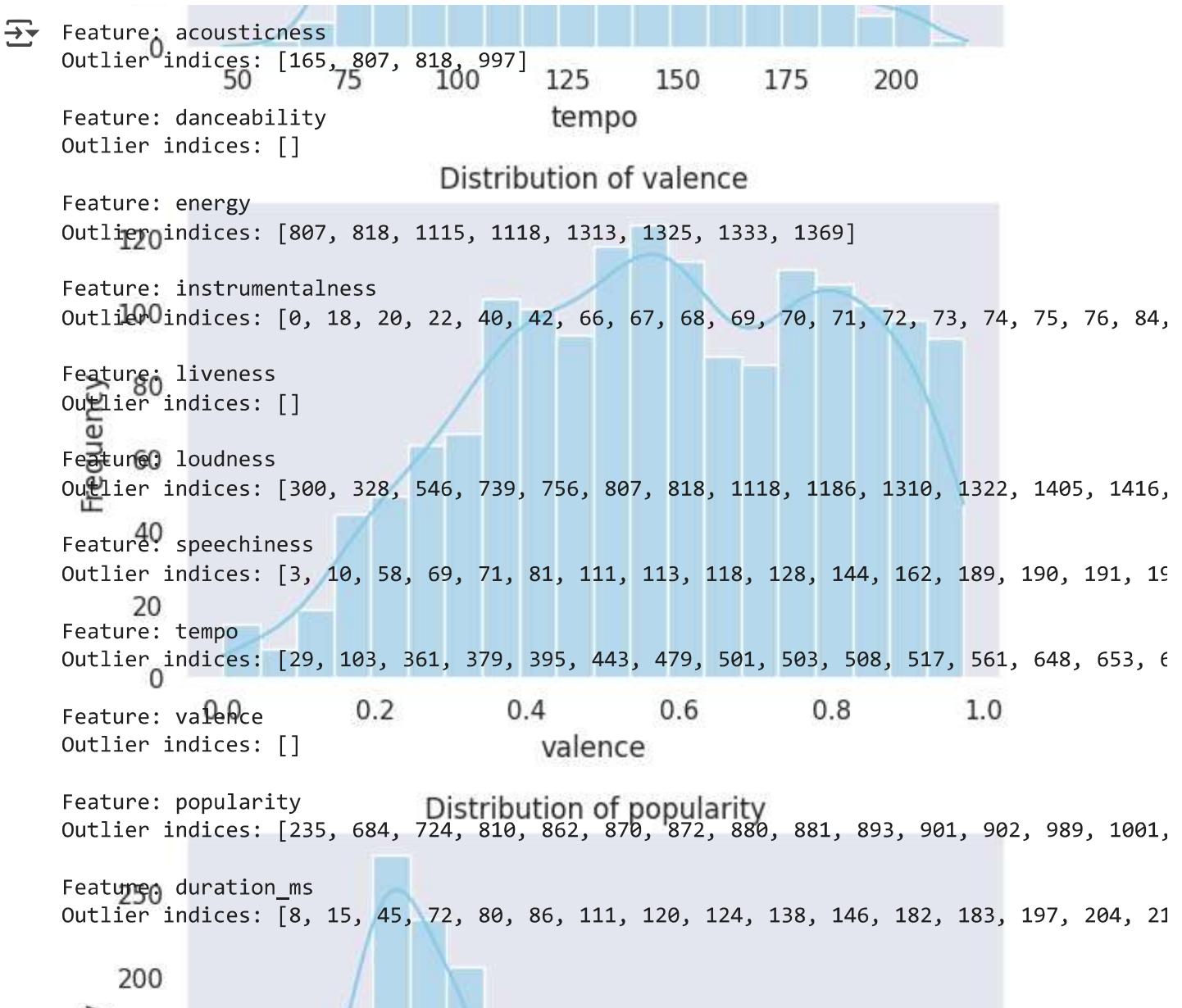
15800

15850

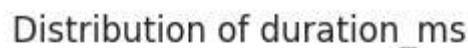
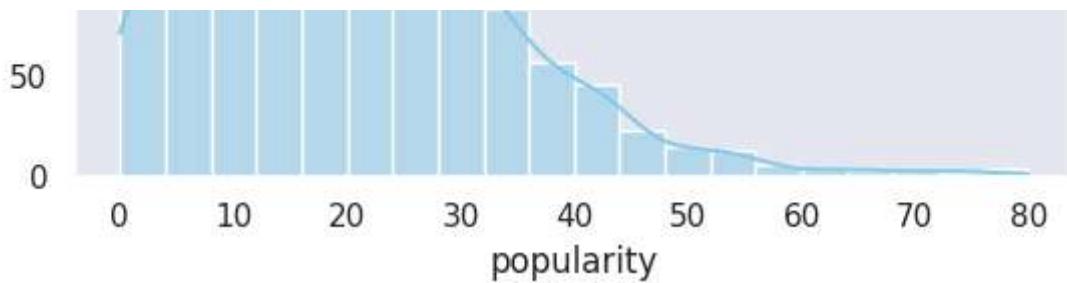
15900

15950

16000



```
#groupby the data by albums  
album_df = spotify_df.groupby(['album'])  
album_df.head()
```



	Unnamed: 0	name	album	release_date	track_number	id
0	0	Concert Intro Music - Live	Licked Live In NYC	2022-06-10	1	2IEkywLJ4ykahi1yRQvmsT
1	100	Street Fighting Man - Live	Licked Live In NYC	2022-06-10	2	6GVgVJBKkGJoRfarYRvGTU
2	2	Start Me Up - Live	Licked Live In NYC	2022-06-10	3	1Lu761pZ0dBTGpzaQoZNW
3	3	If You Can't Rock Me - Live	Licked Live In NYC	2022-06-10	4	1agTQzOTUnGNggycEqiDH
4	4	Don't Stop - Live	Licked Live In NYC	2022-06-10	5	7piGJR8YndQBQWVXv6KtQw
...
1598	1598	Route 66	The Rolling Stones	1964-04-16	1	1FRP8d6I2jm3DS5f78ZrhK
1599	1599	I Just Want To Make Love To You - Mono Version	The Rolling Stones	1964-04-16	2	7j96wehhMtN0fkVvFhD8Ix
1600	1600	Honest I Do	The Rolling Stones	1964-04-16	3	22Bvku5X3odiXj2wbtgY4T
1601	1601	Mona (I Need You Baby)	The Rolling Stones	1964-04-16	4	1o2wsWx1RkkNuVp6Z21HC
1602	1602	Now I've Got A Witness	The Rolling Stones	1964-04-16	5	6QswKScFURI1x3gZ9isXK7

450 rows × 18 columns

3. Perform exploratory Data Analysis

```
# Need the unique name of each album in spotify_df
album_names = spotify_df['album'].unique()

# Convert to DataFrame
album_table = pd.DataFrame(album_names, columns=['Album Name'])

# Print as table
print(album_table)
```

```
→          Album Name
0      Licked Live In NYC
1      Live At The El Mocambo
2  Tattoo You (Super Deluxe)
3          Tattoo You
4    A Bigger Bang (Live)
..
85          ...
86          ...
87  England's Newest Hitmakers
88  England's Newest Hit Makers
89      The Rolling Stones

[90 rows x 1 columns]
```

```
# Access album names from rows 85 and 86
album_85 = album_table.iloc[85]['Album Name']
album_86 = album_table.iloc[86]['Album Name']

# Print both names
print("Album 85:", album_85)
print("Album 86:", album_86)
```

```
→ Album 85: 12 x 5
Album 86: 12 X 5
```

```
album_85
```

```
→ '12 x 5'
```

```
album_names
type(album_names)

→ numpy.ndarray

album_names.shape

→ (90,)
```

1.e Identify Duplicates

```
# Levenshtein distance function
def levenshtein_distance(s1, s2):
    if len(s1) < len(s2):
        return levenshtein_distance(s2, s1)

    previous_row = list(range(len(s2) + 1))
    for i, c1 in enumerate(s1):
        current_row = [i + 1]
        for j, c2 in enumerate(s2):
            insertions = previous_row[j + 1] + 1
            deletions = current_row[j] + 1
            substitutions = previous_row[j] + (c1 != c2)
            current_row.append(min(insertions, deletions, substitutions))
        previous_row = current_row

    return previous_row[-1]

# Calculate distances
results = []
for i in range(len(album_names)):
    for j in range(i + 1, len(album_names)):
        dist = levenshtein_distance(album_names[i], album_names[j])
        results.append({
            "Album 1": album_names[i],
            "Album 2": album_names[j],
            "Levenshtein Distance": dist
        })

# Convert to DataFrame
album_name_compare_df = pd.DataFrame(results)
print(album_name_compare_df)
```

	Album 1	Album 2 \
0	Licked Live In NYC	Live At The El Mocambo
1	Licked Live In NYC	Tattoo You (Super Deluxe)
2	Licked Live In NYC	Tattoo You
3	Licked Live In NYC	A Bigger Bang (Live)
4	Licked Live In NYC	Steel Wheels Live

```

...
4000           12 X 5   England's Newest Hit Makers
4001           12 X 5   The Rolling Stones
4002 England's Newest Hitmakers  England's Newest Hit Makers
4003 England's Newest Hitmakers   The Rolling Stones
4004 England's Newest Hit Makers   The Rolling Stones

```

Levenshtein Distance

0	17
1	22
2	17
3	16
4	15
...	...
4000	25
4001	16
4002	3
4003	22
4004	23

[4005 rows x 3 columns]

```
# Show Album1 and Album2 names when the Levenshtein distance is less than 5
album_name_compare_df[album_name_compare_df['Levenshtein Distance'] < 6]
```

	Album 1	Album 2	Levenshtein Distance
462	Steel Wheels Live	Steel Wheels	5
1074	On Air	12 x 5	5
1075	On Air	12 X 5	5
2970	Love You Live (Remastered)	Love You Live (Remastered 2009)	5
3787	Flowers	Now!	5
3852	got LIVE if you want it!	Got Live if you want it!	4
3914	December's Children (And Everybody's)	December's Children (and everybody's)	4
3941	Out Of Our Heads	Out Of Our Heads (UK)	5
3950	Out Of Our Heads (US Sequence)	Out Of Our Heads (UK Sequence)	1
3984	The Rolling Stones, Now!	The Rolling Stones No. 2	4
3995	12 x 5	12 X 5	1

```
#show all the songs from Steel Wheels LIve
#spotify_df[spotify_df['album'] == '12 X 5']
```

```
#show all the songs from Steel Wheels LIve
#spotify_df[spotify_df['album'] == 'Out Of Our Heads (US Sequence)']
```



```
#show all the songs from Steel Wheels LIve
#spotify_df[spotify_df['album'] == 'Out Of Our Heads (UK Sequence)']
```



```
#show all the songs from Steel Wheels LIve
#spotify_df[spotify_df['album'] == 'Out Of Our Heads']
```



```
#show all the songs from Steel Wheels LIve
#spotify_df[spotify_df['album'] == 'Out Of Our Heads (UK)']
```



```
#show all the songs from December's Children (And Everybody's)
#spotify_df[spotify_df['album'] == 'December's Children (And Everybody's)']
```



```
#show all the songs from December's Children (And Everybody's)
#spotify_df[spotify_df['album'] == 'December's Children (and everybody's)']
```



```
#show all the songs from December's Children (And Everybody's)
#spotify_df[spotify_df['album'] == 'got LIVE if you want it!']
```



```
#show all the songs from December's Children (And Everybody's)
#spotify_df[spotify_df['album'] == 'Got Live if you want it!']
```



```
#show all the songs from December's Children (And Everybody's)
#spotify_df[spotify_df['album'] == 'Love You Live (Remastered)']
```



```
#show all the songs from December's Children (And Everybody's)
#spotify_df[spotify_df['album'] == 'Love You Live (Remastered 2009)']
```



```
#show all the songs from Steel Wheels LIve
#spotify_df[spotify_df['album'] == 'Steel Wheels Live']
```



```
#show all the songs from Steel Wheels
#spotify_df[spotify_df['album'] == 'Steel Wheels']
```

2. Refine data

```
# drop all the songs of the album 'Steel Wheel'  
spotify_df = spotify_df[spotify_df['album'] != 'Out Of Our Heads (UK)']  
spotify_df = spotify_df[spotify_df['album'] != 'Out Of Our Heads (UK Sequence)']  
spotify_df = spotify_df[spotify_df['album'] != 'Out Of Our Heads (US Sequence)']  
spotify_df = spotify_df[spotify_df['album'] != 'Love You Live (Remastered 2009)']
```

```
# Calculate the number of unique albums in the spotify_df dataset  
len(spotify_df['album'].unique())
```

→ 86

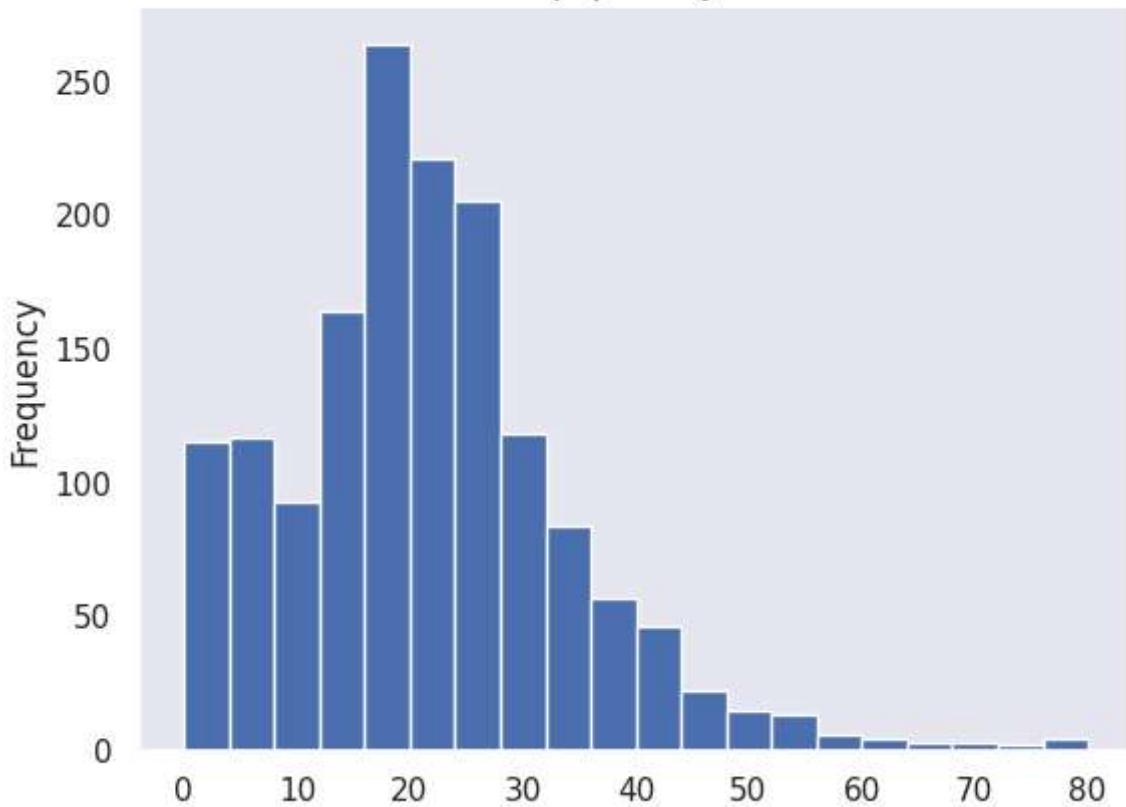
▼ popularity

```
# @title popularity
```

```
spotify_df['popularity'].plot(kind='hist', bins=20, title='popularity')  
plt.gca().spines[['top', 'right',]].set_visible(False)
```

→

popularity



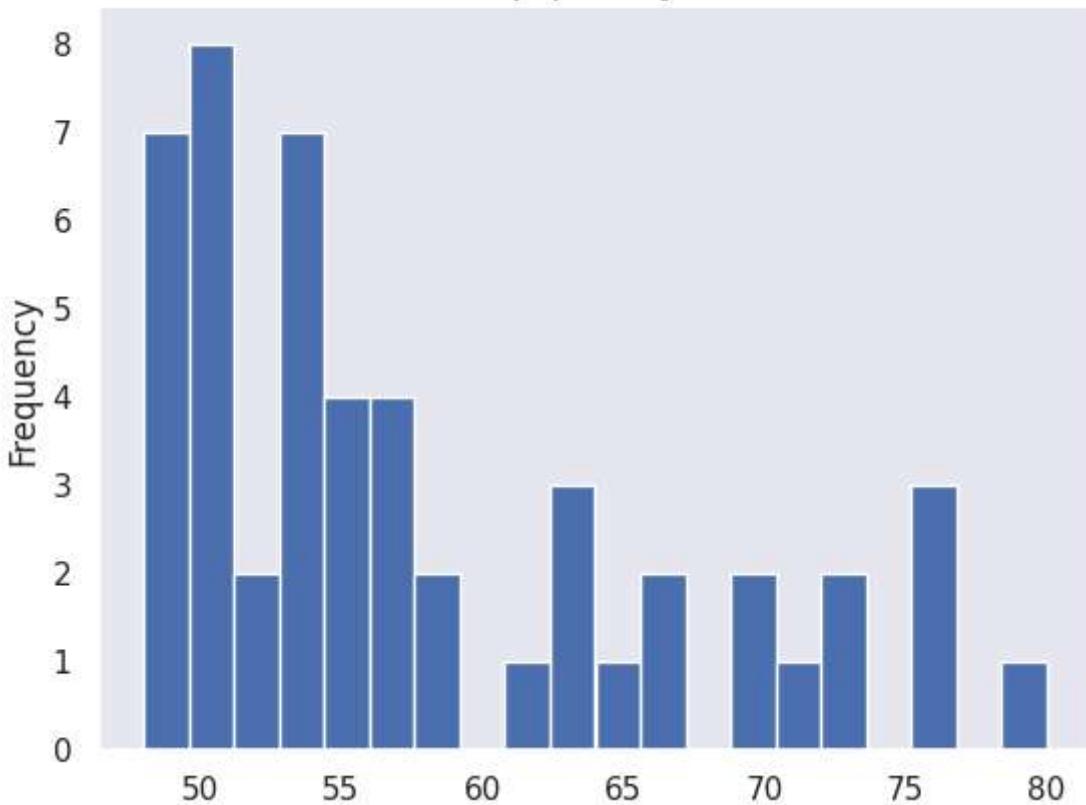
```
# create group by popularity ascending order and take the top 50  
top_50_popularity = spotify_df.sort_values(by='popularity', ascending=False).head(50)  
top_50_popularity.shape
```

→ (50, 18)

```
top_50_popularity['popularity'].plot(kind='hist', bins=20, title='popularity')
plt.gca().spines[['top', 'right', ]].set_visible(False)
```

→

popularity



3.a. Utilize suitable visualizations to identify the two albums that should be recommended to anyone based on the number of popular songs in each album

```
# Find out which album has the most songs in the list
top_50_popularity_album = top_50_popularity['album'].value_counts()

# Print the top 2 albums with highest popularity
top_50_popularity_album.head(2)
```



count

album

Sticky Fingers (Remastered)	6
Let It Bleed	5

dtype: int64

top_50_popularity.shape

(50, 18)

3.c Examine the relationship between songs popularity and various factors

```
# Define popular songs as those with a popularity score above 80
popular_songs = spotify_df[spotify_df['popularity'] > 60]

# Count the number of popular songs per album
popular_counts = popular_songs['album'].value_counts()

# Identify the top 2 albums with the most popular songs
top_2_albums = popular_counts.head(2)

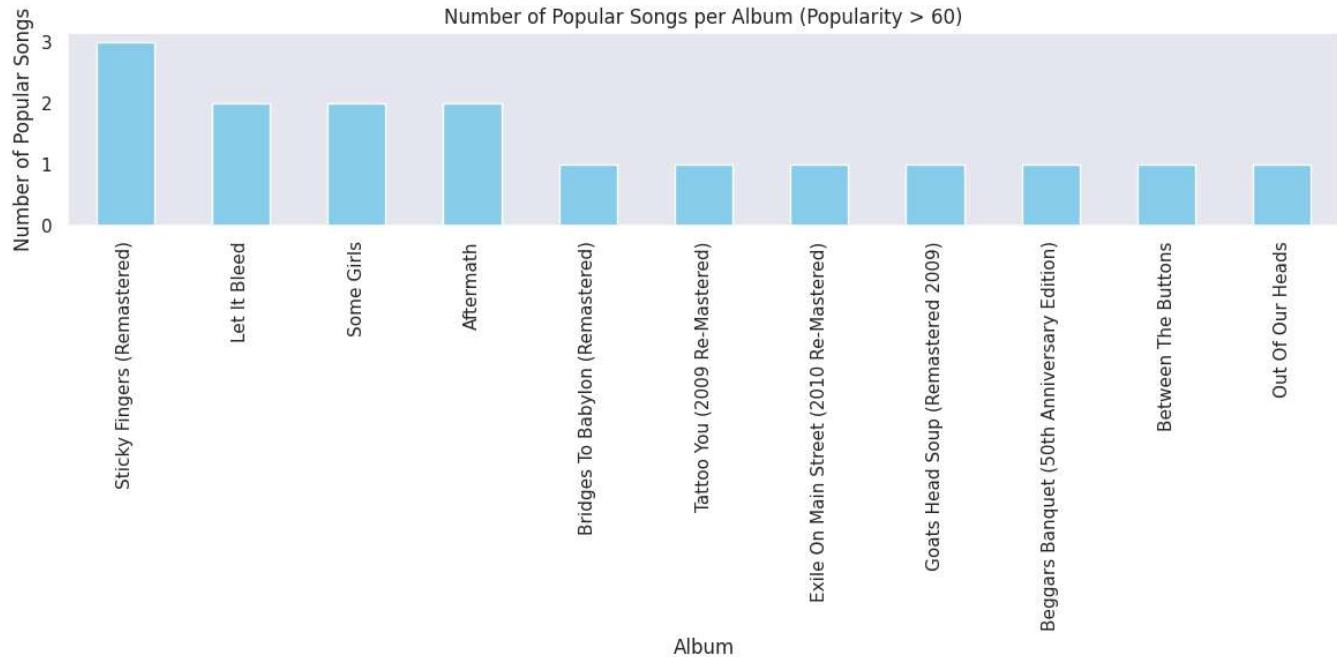
# Plot the number of popular songs per album
plt.figure(figsize=(12, 6))
popular_counts.plot(kind='bar', color='skyblue')
plt.title('Number of Popular Songs per Album (Popularity > 60)')
plt.xlabel('Album')
plt.ylabel('Number of Popular Songs')
plt.xticks(rotation=90)
plt.tight_layout()
plt.savefig("popular_songs_per_album.png")

# Display the top 2 albums
print("Top 2 albums with the most popular songs:")
print(top_2_albums)
```

→ Top 2 albums with the most popular songs:

album

Sticky Fingers (Remastered)	3
Let It Bleed	2
Name: count, dtype: int64	



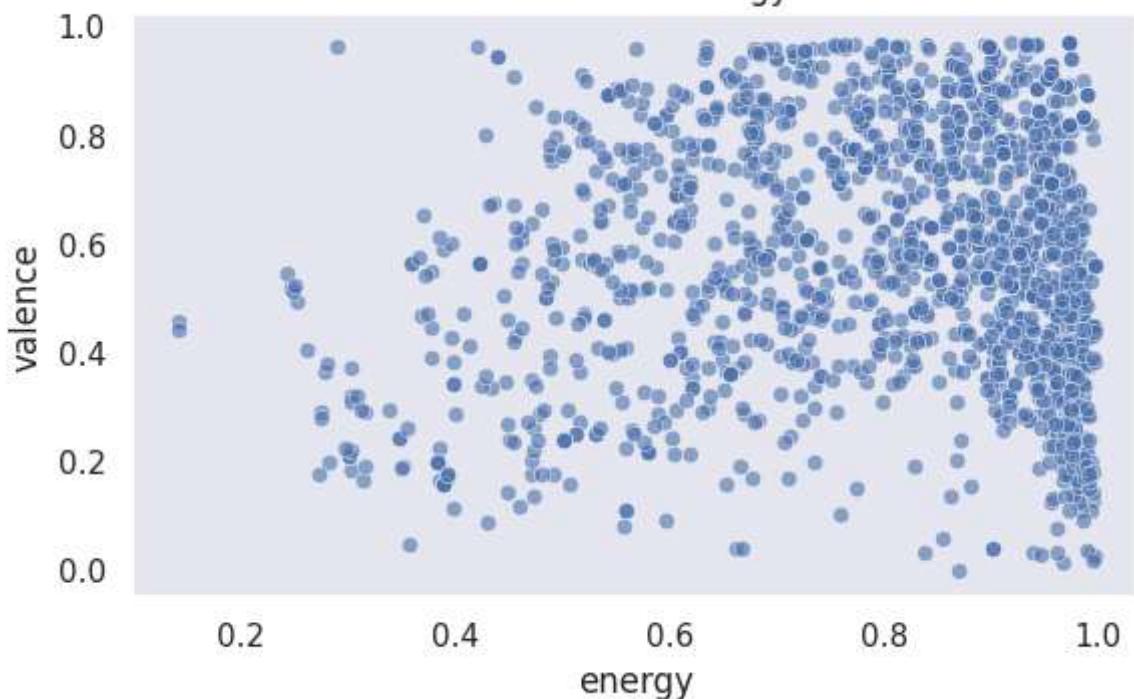
#popular_songs

```
# Scatter plots to explore relationships
scatter_pairs = [
    ('energy', 'valence'),
    ('danceability', 'valence'),
    ('loudness', 'energy'),
    ('tempo', 'danceability')
]
for x, y in scatter_pairs:
    plt.figure(figsize=(6, 4))
    sns.scatterplot(data=spotify_df, x=x, y=y, alpha=0.6)
    plt.title(f'{y} vs {x}')
```

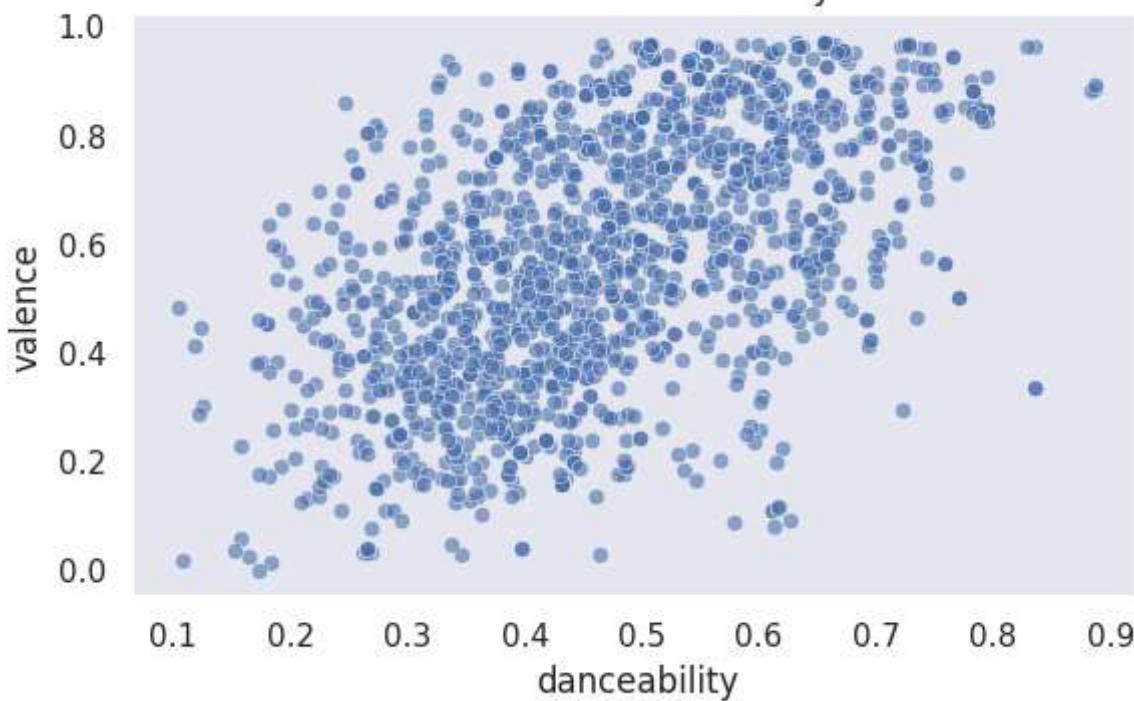
```
plt.xlabel(x)
plt.ylabel(y)
plt.tight_layout()
plt.savefig(f"{y}_vs_{x}_scatter.png")
```



valence vs energy



valence vs danceability



energy vs loudness



0.4

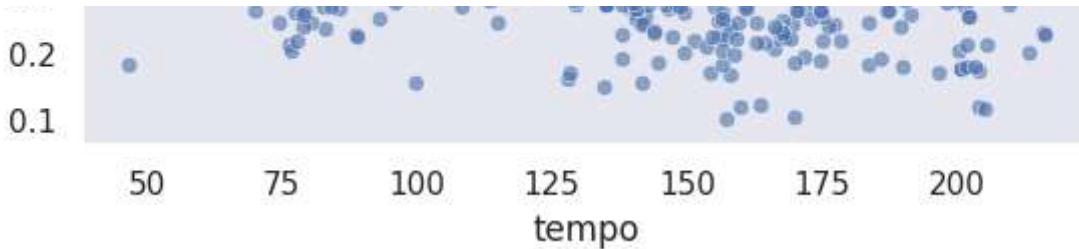
Feature Relationships

- Energy vs. Valence: Slight positive correlation—more energetic songs tend to be more positive.
- Danceability vs. Valence: Positive correlation—danceable songs are often more cheerful.
- Energy vs. Loudness: Strong correlation—louder songs are more energetic.
- Danceability vs. Tempo: No strong pattern, but some clustering..

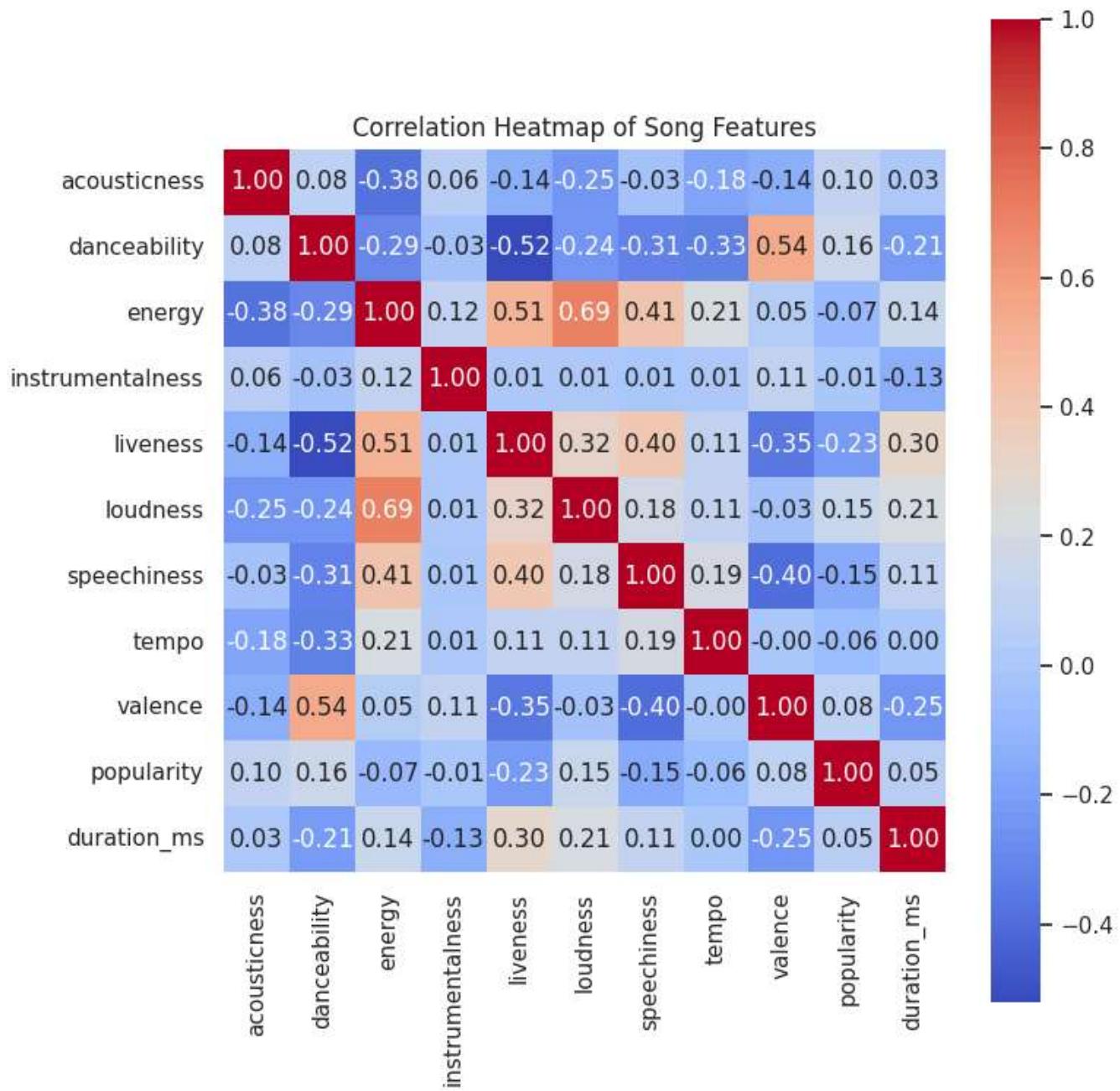
Correlation heatmap

```
plt.figure(figsize=(8, 8))
corr_matrix = spotify_df[features].corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f", square=True)
plt.title('Correlation Heatmap of Song Features')
plt.tight_layout()
plt.savefig("correlation_heatmap.png")
```

```
print("Exploratory data analysis completed. Visualizations saved as PNG files.")
```



➡ Exploratory data analysis completed. Visualizations saved as PNG files.



Observations

- The maximum positive correlation coefficient is 0.69 (loudness vs energy)
- The maximum negative correlation is -0.52 (liveness vs danceability)
- None of the factors are super correlated - CANNOT Drop any factor

3.d Perform Principal Component Analysis

```
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

# Standardize the features
scaler = StandardScaler()
scaled_features = scaler.fit_transform(spotify_df[features])

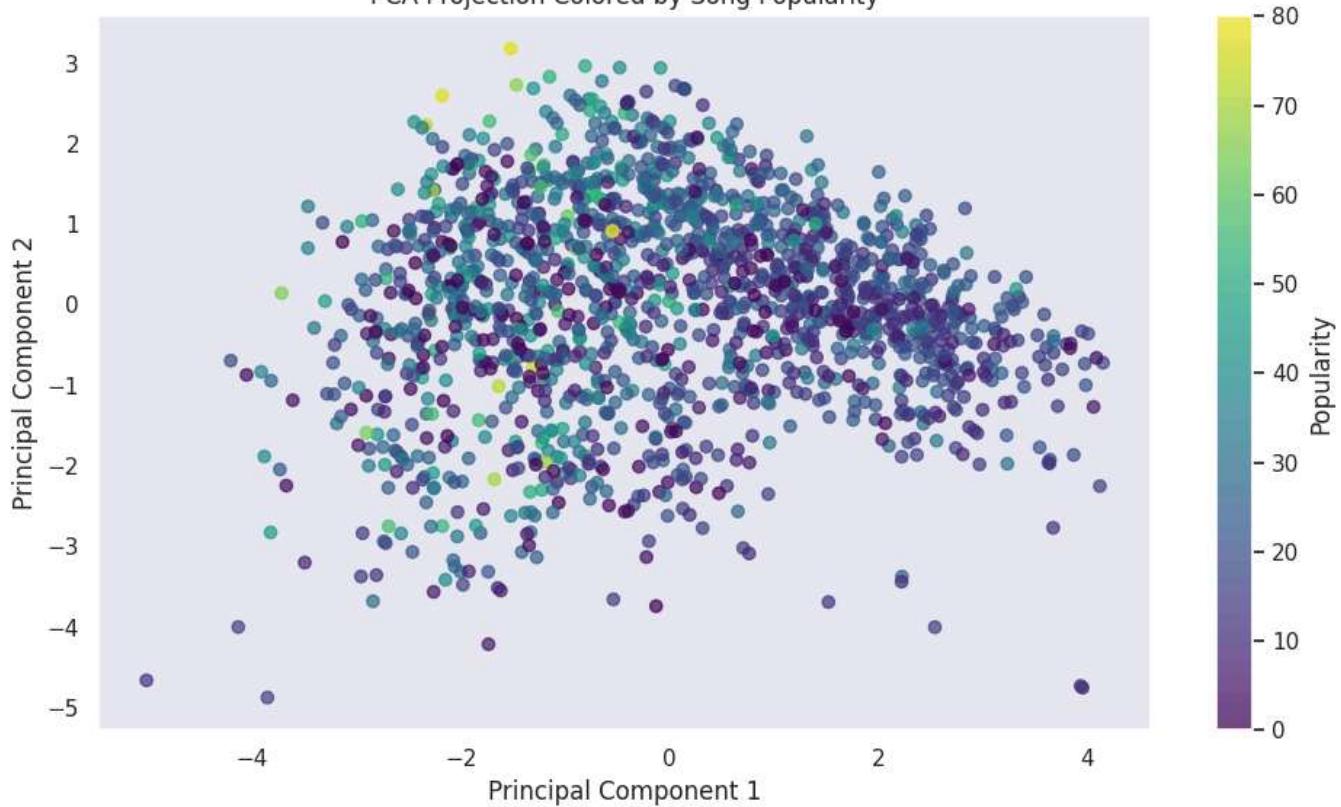
# Apply PCA
pca = PCA(n_components=2)
pca_components = pca.fit_transform(scaled_features)

# Create a DataFrame with PCA components and popularity
pca_df = pd.DataFrame(data=pca_components, columns=['PC1', 'PC2'])
pca_df['popularity'] = spotify_df['popularity'].values

# Plot PCA components colored by popularity
plt.figure(figsize=(10, 6))
scatter = plt.scatter(pca_df['PC1'], pca_df['PC2'], c=pca_df['popularity'], cmap='viridis',
plt.colorbar(scatter, label='Popularity')
plt.title('PCA Projection Colored by Song Popularity')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.tight_layout()
plt.savefig("pca_popularity_projection.png")
```



PCA Projection Colored by Song Popularity



```
# Define popularity threshold
popularity_threshold = 30

# Create a new column to categorize songs
spotify_df['popularity_group'] = spotify_df['popularity'].apply(lambda x: 'Popular' if x >=
    popularity_threshold else 'Not Popular')

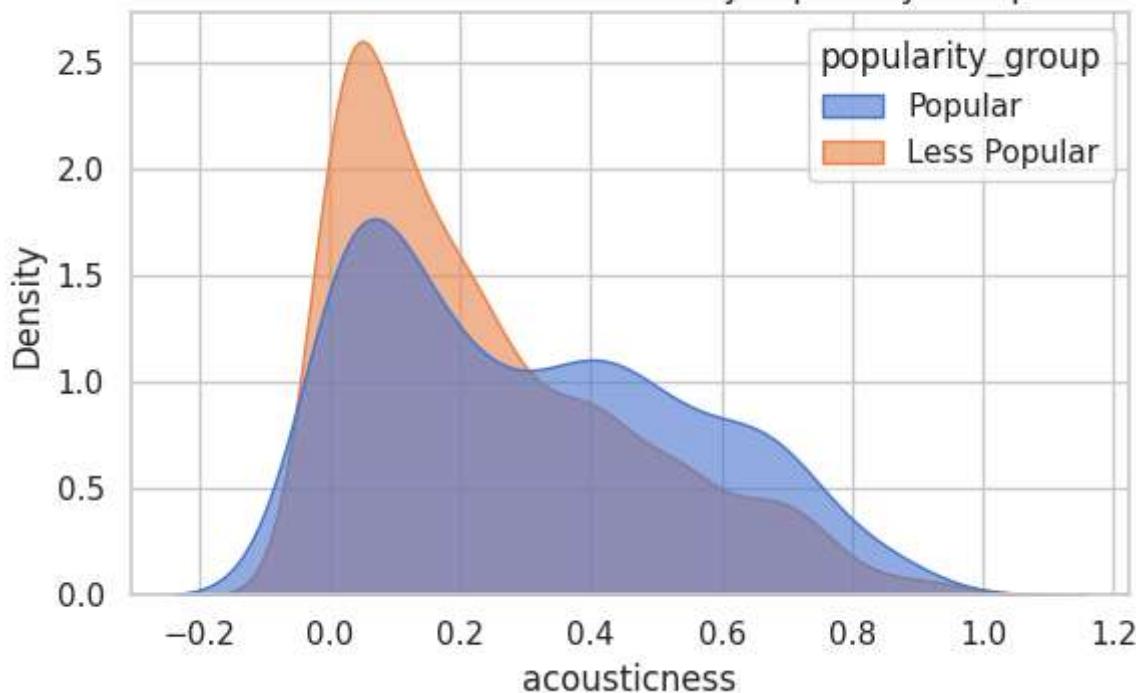
# Set seaborn style
sns.set(style="whitegrid")

# Generate KDE plots for each feature
for feature in features:
    plt.figure(figsize=(6, 4))
    sns.kdeplot(data=spotify_df, x=feature, hue='popularity_group', fill=True, common_norm=False)
    plt.title(f'Distribution of {feature} by Popularity Group')
    plt.xlabel(feature)
    plt.ylabel('Density')
    plt.tight_layout()
```

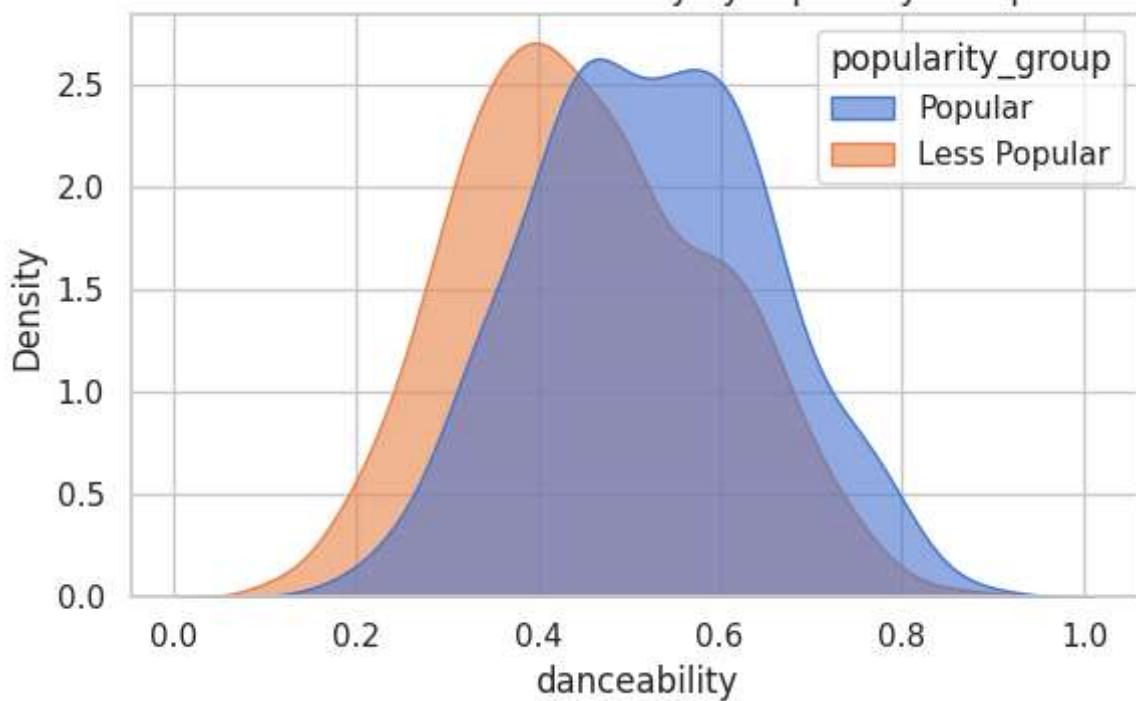
```
plt.savefig(f"{feature}_popularity_kde.png")
```



Distribution of acousticness by Popularity Group

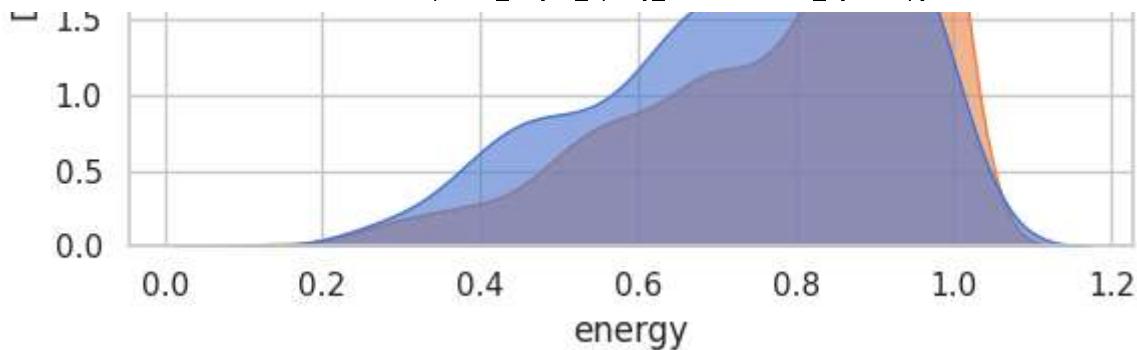


Distribution of danceability by Popularity Group

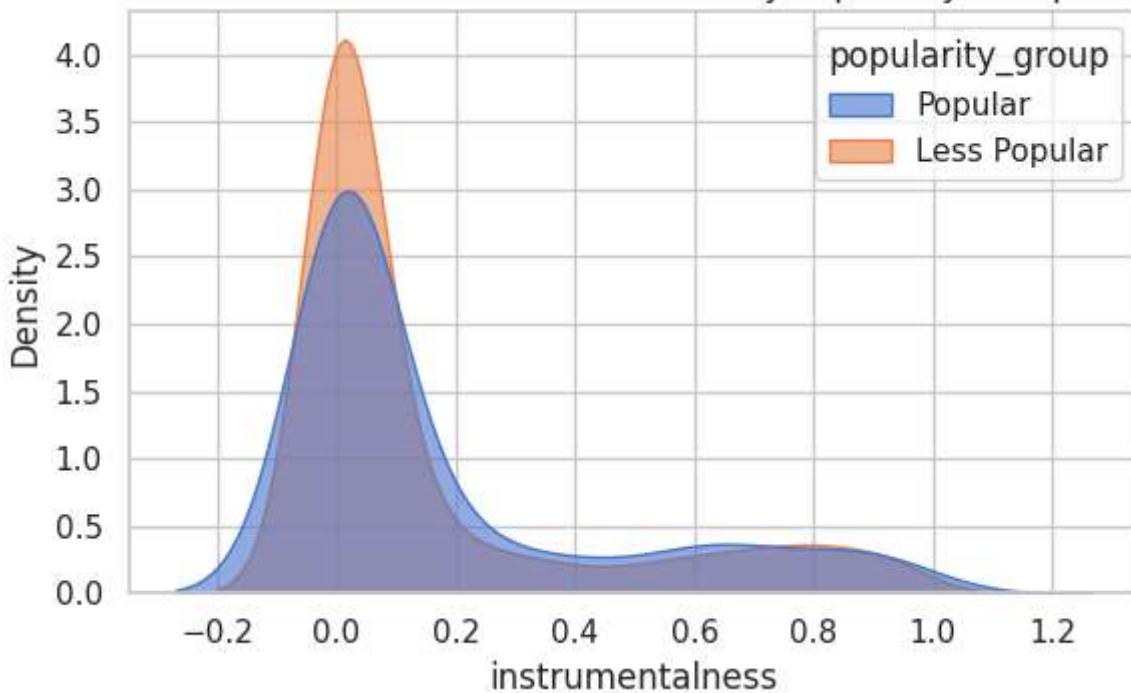


Distribution of energy by Popularity Group

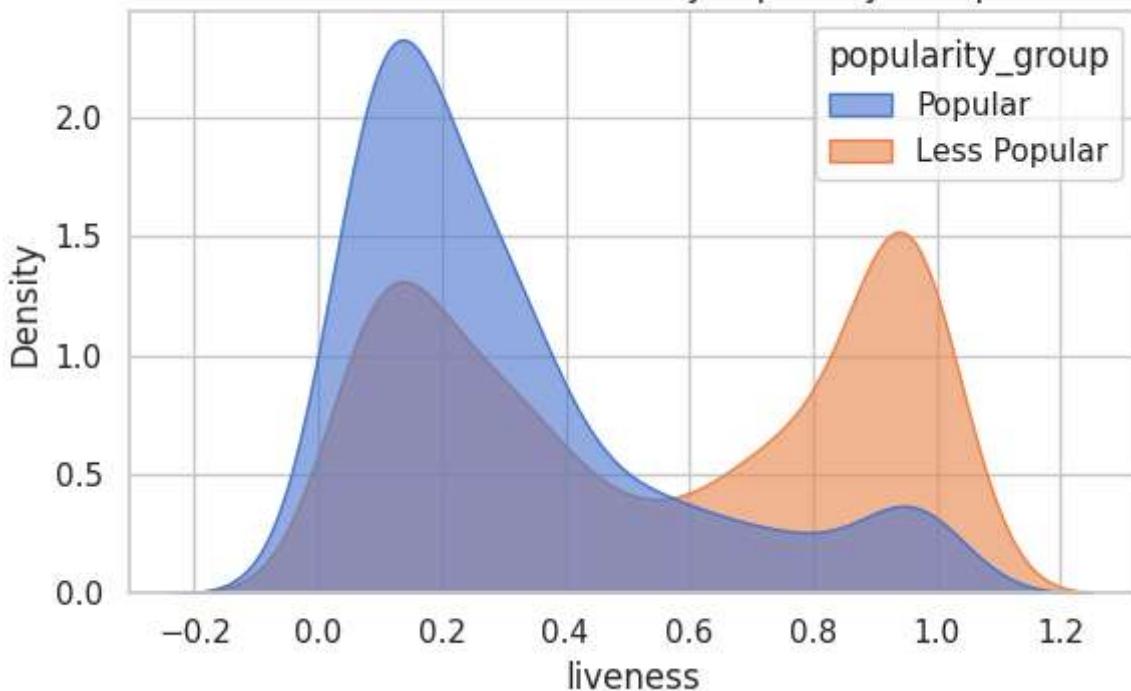




Distribution of instrumentalness by Popularity Group



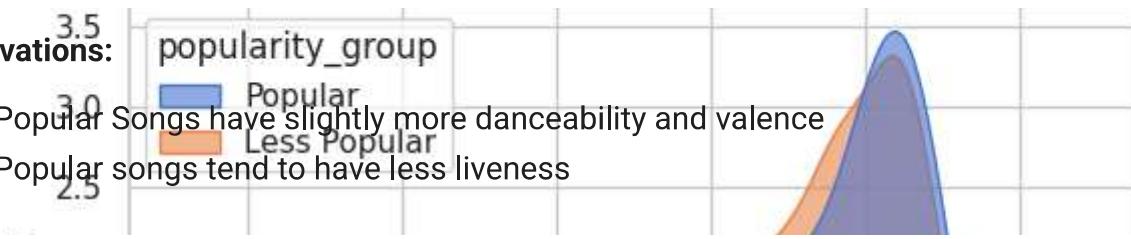
Distribution of liveness by Popularity Group



Distribution of loudness by Popularity Group

Observations: popularity_group

- Popular Songs have slightly more danceability and valence
- Popular songs tend to have less liveness



```

import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

# Standardize the features
scaler = StandardScaler()
scaled_features = scaler.fit_transform(spotify_df[features])

# Identify the optimal number of clusters using the elbow method
inertia = []
for k in range(1, 11):
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(scaled_features)
    inertia.append(kmeans.inertia_)

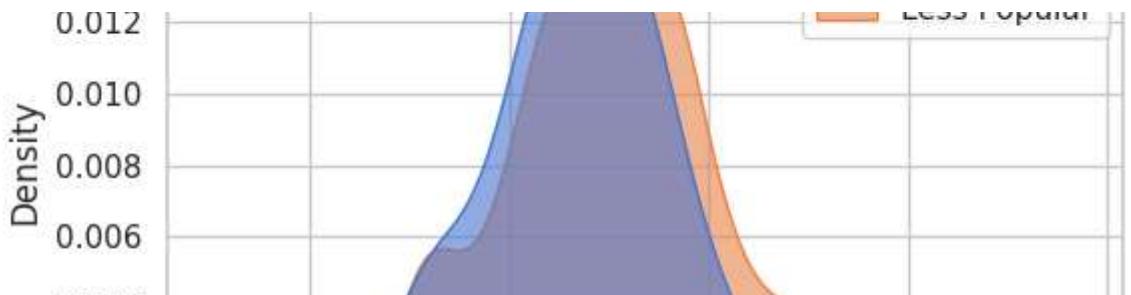
# Plot the elbow curve
plt.figure(figsize=(8, 5))
plt.plot(range(1, 11), inertia, marker='o')
plt.title('Elbow Method for Optimal Number of Clusters')
plt.xlabel('Number of Clusters')
plt.ylabel('Inertia')
plt.tight_layout()
plt.savefig("elbow_curve.png")

# Apply KMeans clustering with optimal number of clusters (e.g., 3 based on elbow curve)
optimal_k = 3
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
spotify_df['cluster'] = kmeans.fit_predict(scaled_features)

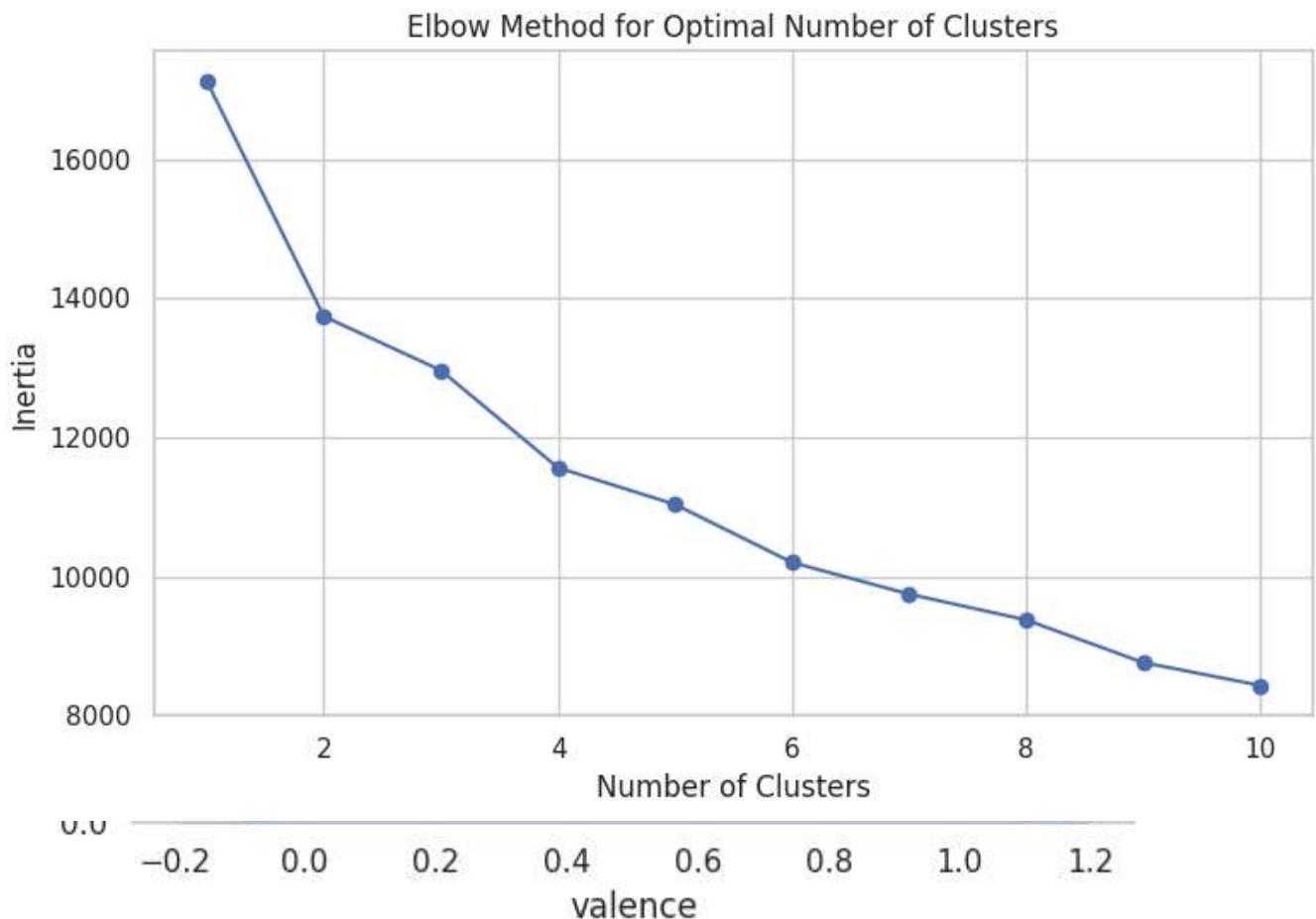
# Define each cluster based on feature averages
cluster_summary = spotify_df.groupby('cluster')[features].mean()
cluster_summary.to_csv("cluster_summary.csv")

print("Cluster analysis completed.")
print("Cluster summary saved to 'cluster_summary.csv'.")

```



0.004
Cluster analysis completed.
Cluster summary saved to 'cluster_summary.csv'.



Distribution of popularity by Popularity Group

```
# Apply KMeans clustering with 3 clusters
kmeans = KMeans(n_clusters=3, random_state=42)
clusters = kmeans.fit_predict(scaled_features)

# Create a DataFrame for visualization
pca_df = pd.DataFrame(data=pca_components, columns=['PC1', 'PC2'])
pca_df['cluster'] = clusters

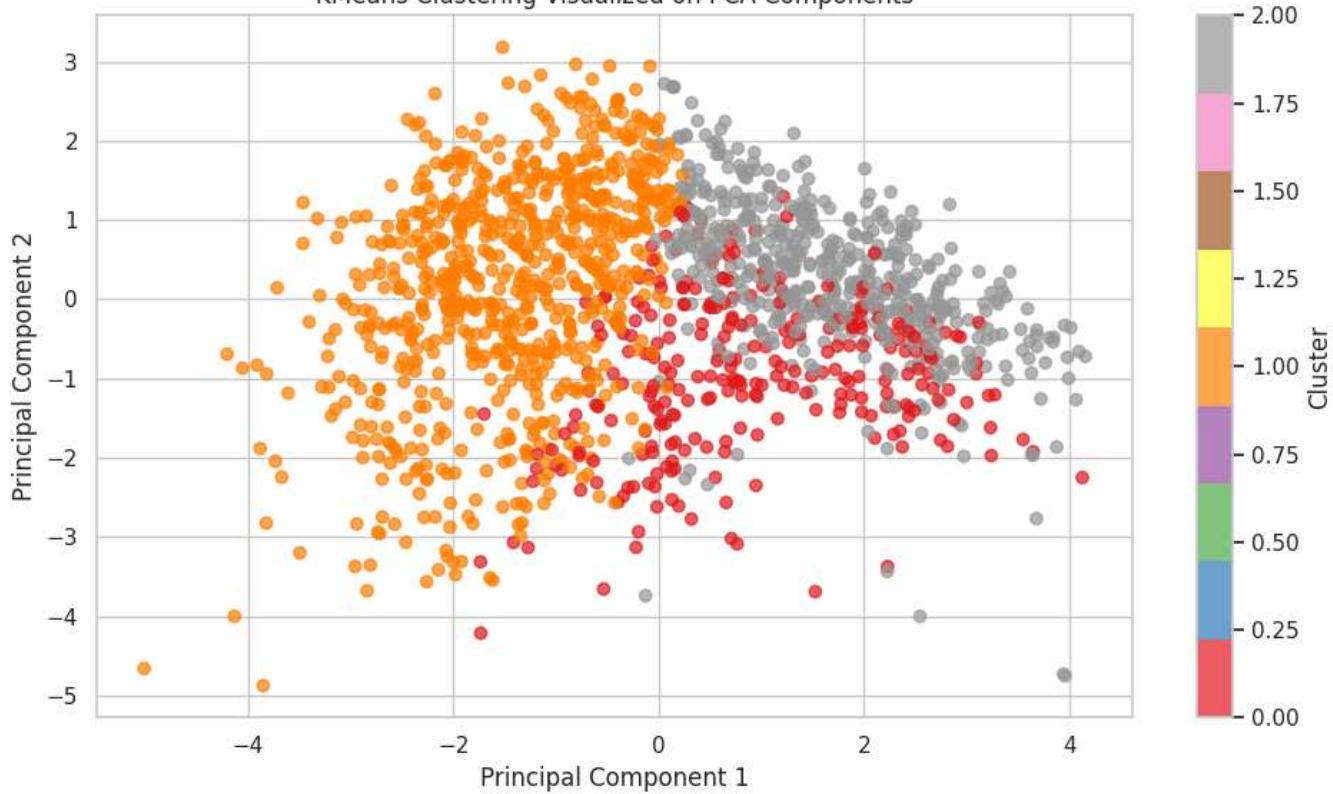
# Plot the clusters
plt.figure(figsize=(10, 6))
scatter = plt.scatter(pca_df['PC1'], pca_df['PC2'], c=pca_df['cluster'], cmap='Set1', alpha=0.8)
plt.title('KMeans Clustering Visualized on PCA Components')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(scatter, label='Cluster')
plt.tight_layout()
plt.savefig("cluster_pca_visualization_3.png")
```

1e-6 Distribution of duration_ms by Popularity Group





KMeans Clustering Visualized on PCA Components



Start coding or [generate](#) with AI.

```
# Apply KMeans clustering with 3 clusters
kmeans = KMeans(n_clusters=5, random_state=42)
clusters = kmeans.fit_predict(scaled_features)

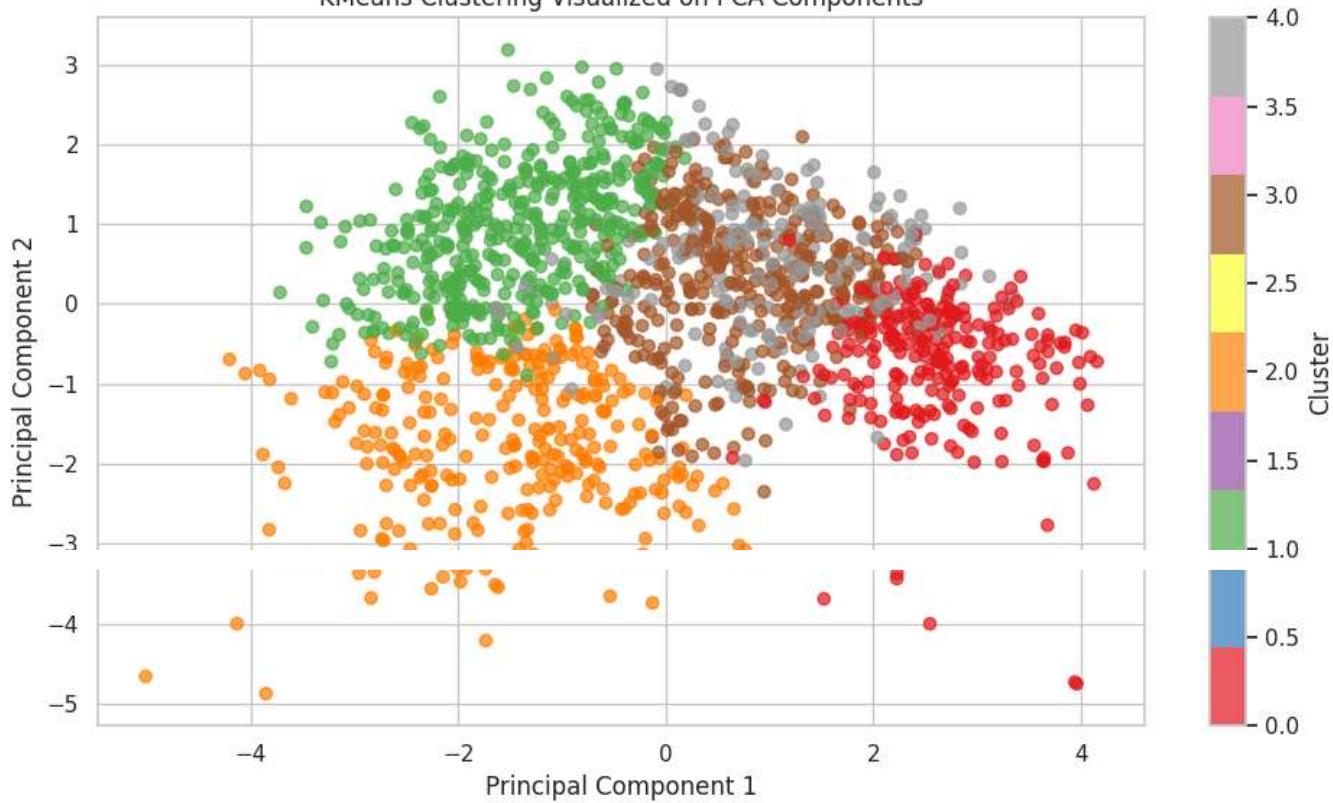
# Create a DataFrame for visualization
pca_df = pd.DataFrame(data=pca_components, columns=['PC1', 'PC2'])
pca_df['cluster'] = clusters

# Plot the clusters
plt.figure(figsize=(10, 6))
scatter = plt.scatter(pca_df['PC1'], pca_df['PC2'], c=pca_df['cluster'], cmap='Set1', alpha=0.5)
plt.title('KMeans Clustering Visualized on PCA Components')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(scatter, label='Cluster')
plt.tight_layout()
```

```
plt.savefig("cluster_pca_visualization_5.png")
```



KMeans Clustering Visualized on PCA Components



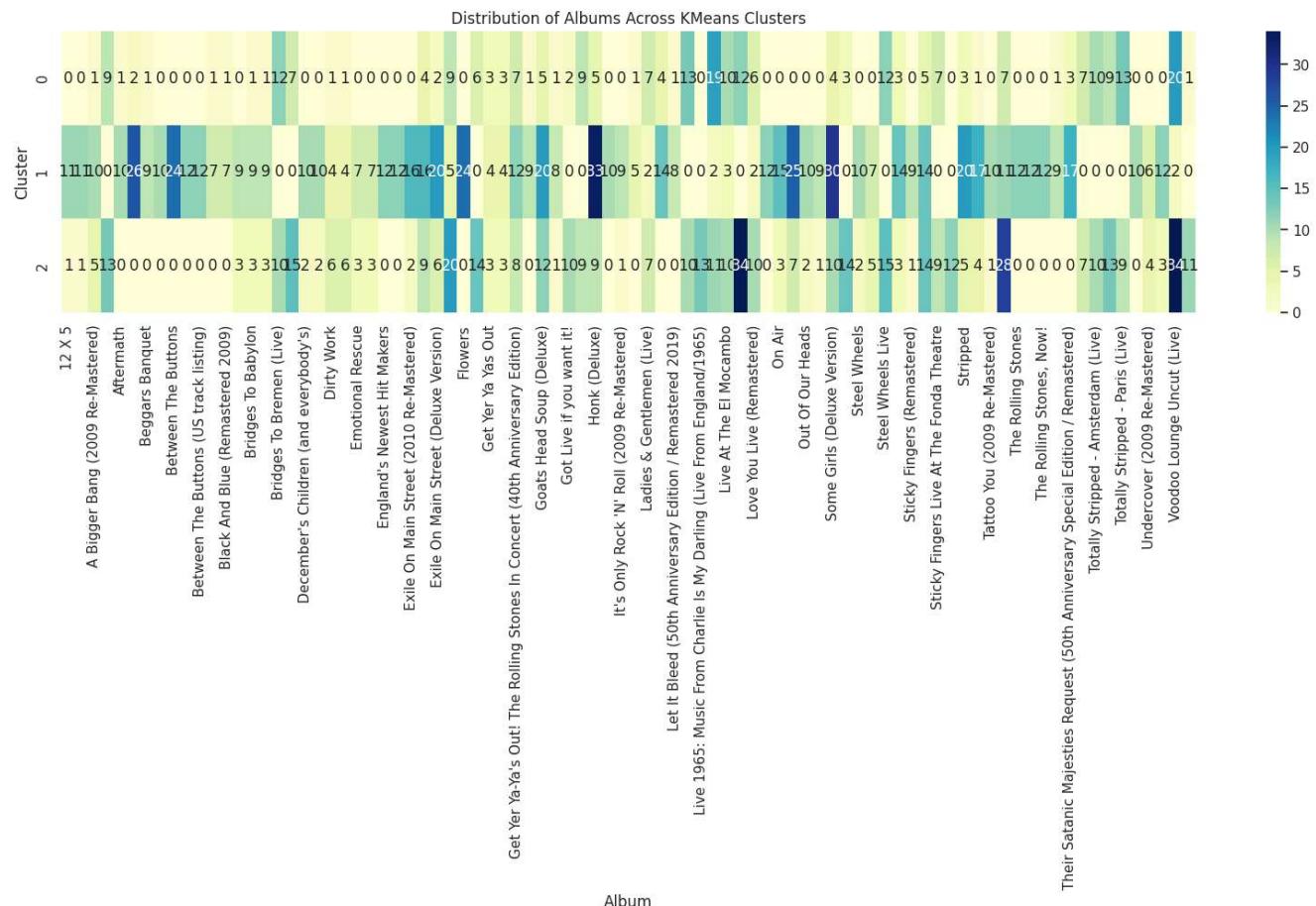
```
# Apply KMeans clustering (assuming 3 clusters based on previous analysis)
kmeans = KMeans(n_clusters=3, random_state=42)
spotify_df['cluster'] = kmeans.fit_predict(scaled_features)

# Group by cluster and album to count songs
cluster_album_counts = spotify_df.groupby(['cluster', 'album']).size().unstack(fill_value=0)

# Plot heatmap of album distribution across clusters
plt.figure(figsize=(20, 4))
sns.heatmap(cluster_album_counts, cmap='YlGnBu', annot=True, fmt='d')
plt.title('Distribution of Albums Across KMeans Clusters')
plt.xlabel('Album')
plt.ylabel('Cluster')
plt.xticks(rotation=90)
plt.tight_layout()
```

```
plt.savefig("album_cluster_heatmap.png")
```

→ /tmp/ipython-input-3466032068.py:15: UserWarning: Tight layout not applied. The bottom
plt.tight_layout()



```
# Apply KMeans clustering (assuming 3 clusters based on previous analysis)
kmeans = KMeans(n_clusters=3, random_state=42)
spotify_df['cluster'] = kmeans.fit_predict(scaled_features)

# Compute average popularity per cluster
popularity_by_cluster = spotify_df.groupby('cluster')['popularity'].mean()

# Plot the average popularity per cluster
plt.figure(figsize=(8, 5))
popularity_by_cluster.plot(kind='bar', color='coral')
plt.title('Average Popularity per KMeans Cluster')
plt.xlabel('Cluster')
plt.ylabel('Average Popularity')
plt.xticks(rotation=0)
plt.tight_layout()
plt.savefig("cluster_popularity_bar.png")

# Display the average popularity values
print("Average popularity per cluster:")
print(popularity_by_cluster)
```

→ Average popularity per cluster:
cluster
0 18.766187