

✓ Music Recommendation System with Natural Language Processing

✓ Introduction.

We designed a music recommendation system using a dataset of over 28,000 Spotify tracks containing both audio features and full lyrics. Our goal was to build two systems: one based on traditional content-based filtering, and another enhanced by Natural Language Processing (NLP) to capture lyrical sentiment.

For our project, we chose to use a Spotify data set from Kaggle and create two recommender systems. One of these recommender systems used simpler content-based filtering while the other was based around our own sentiment analysis of the song lyrics. Our data set had just over 28,000 observations and 31 total variables, including 'energy', 'acousticness', and other measures of song characteristics that we used to generate similarity scores. We chose this data set for our project because we are all passionate about music, and we were curious to see if a recommender system like this could help us to discover new music.

Data Description.

The chosen data set for our project is a data set of songs released from 1950-2019, from Kaggle (\url{<https://www.kaggle.com/saurabhshahane/music-dataset-1950-to-2019>}). The data has 28,372 total observations and 31 columns/variables. Here are all of them listed:

- Unnamed: 0
- artist_name
- track_name
- release_date
- genre
- lyrics
- len
- dating

- violence
- world/life
- night/time
- shake the audience
- family/gospel
- romantic
- communication
- obscene
- music
- movement/places
- light/visual perceptions
- family/spiritual
- like/girls
- sadness
- feelings
- danceability
- loudness
- acousticness
- instrumentalness
- valence
- energy
- topic
- age.

All of the variables in this data set will be predictors in some way, given that we will be creating a recommendation system based on all inputs, although we will focus more specifically on 'lyrics' for creating our sentiment analysis output. Many of these variables are numeric, but there are also string and longer text (lyrics) outputs as well. It appears that there is no missing data, which means that our cleaning process should be fairly easy and this data set is nearly ready to be analyzed.

✓ Background Information

There are two primary types of recommender systems:

- **Content-based** systems recommend items by analyzing the properties or attributes of the items themselves (such as audio features and lyrics in our case). These systems recommend new items by comparing item feature vectors using similarity metrics such as cosine similarity.

Example: Songs with similar energy, danceability, and lyrical content are recommended to a user based on previously liked songs.

- **Collaborative filtering** systems, by contrast, rely on user interaction data (such as song ratings or listening history) to find patterns across users or items. These methods compute similarity either between users or between items based on shared behavior.
 - In a **user-based** system, the algorithm compares the listening histories of different users to find similar users.
 - In an **item-based** system, items (e.g., songs) are compared based on how similar users have rated them.

Cosine Similarity

Across both types of systems, the metric of **cosine similarity** is commonly used to measure how similar two vectors are.

In the context of music, these vectors could represent:

- Audio features (e.g., danceability, energy)
- TF-IDF word embeddings of lyrics
- Combined feature sets including sentiment

Interpretation of values:

- **+1**: Perfect similarity
- **0**: No similarity (orthogonal)
- **-1**: Complete opposition

$$\Rightarrow S_C(A, B) := \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \cdot \sqrt{\sum_{i=1}^n B_i^2}}$$

In our project, we developed two content-based recommender systems using Spotify song metadata and lyrics:

Uses standardized numeric audio features like `danceability`, `energy`, `valence`, etc., along with reduced lyric vectors from TF-IDF + Truncated SVD.

Cosine similarity identifies the 10 most similar songs to the input track.

Incorporates **sentiment analysis** using TextBlob to score lyrics from -1 to +1. These scores are combined with audio and textual features. Word clouds were generated to visualize word distributions for positive and negative sentiment songs.

In both systems, **cosine similarity** serves as the backbone for identifying emotional and musical similarity between tracks.

https://colab.research.google.com/drive/1zZmPfa8vrBMxsabNPY_MgdXB9MFKTOa8#scrollTo=Dcpl_6WhUIVd

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import nltk
7 from nltk.tokenize import word_tokenize
8 from nltk.corpus import stopwords
9 from sklearn.feature_extraction.text import TfidfVectorizer
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.neighbors import NearestNeighbors
12 from sklearn.decomposition import TruncatedSVD
13 from wordcloud import WordCloud

```

```

1 # Ensure necessary NLTK resources are downloaded
2 nltk.download('punkt')
3 nltk.download('stopwords')

```

```

[Out] [nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True

```

```

1 file_path = "/content/tcc_ceds_music.csv" # Path to the uploaded file in C
2 df = pd.read_csv(file_path)

```

```

1 # Display first few rows and dataset info
2 display(df.head())
3 display(df.info())

```

```

[Out]
   Unnamed: 0  artist_name  track_name  release_date  genre  lyrics  len  da
0           0      mukesh    mohabbat      1950      pop  hold time  95  0.00
              feel break
              feel untrue
              convince
              spea...

1           4  frankie laine    i believe      1950      pop  believe  51  0.00
              drop rain
              fall grow
              believe
              darkest
              ni...

              sweetheart

```

2	6	johnnie ray	cry	1950	pop	send letter goodbye secret feel bet...	24	0.00
3	10	pérez prado	patricia	1950	pop	kiss lips want stroll charm mambo chacha merin...	54	0.04
4	12	giorgos papadopoulos	apopse eida oneiro	1950	pop	till darling till matter know till dream live ...	48	0.00

5 rows x 31 columns

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 28372 entries, 0 to 28371
```

```
Data columns (total 31 columns):
```

#	Column	Non-Null Count		Dtype
---	-----	-----	-----	-----
0	Unnamed: 0	28372	non-null	int64
1	artist_name	28372	non-null	object
2	track_name	28372	non-null	object
3	release_date	28372	non-null	int64
4	genre	28372	non-null	object
5	lyrics	28372	non-null	object
6	len	28372	non-null	int64
7	dating	28372	non-null	float64
8	violence	28372	non-null	float64
9	world/life	28372	non-null	float64
10	night/time	28372	non-null	float64
11	shake the audience	28372	non-null	float64
12	family/gospel	28372	non-null	float64
13	romantic	28372	non-null	float64
14	communication	28372	non-null	float64
15	obscene	28372	non-null	float64
16	music	28372	non-null	float64
17	movement/places	28372	non-null	float64
18	light/visual perceptions	28372	non-null	float64
19	family/spiritual	28372	non-null	float64
20	like/girls	28372	non-null	float64
21	sadness	28372	non-null	float64
22	feelings	28372	non-null	float64
23	danceability	28372	non-null	float64
24	loudness	28372	non-null	float64
25	acousticness	28372	non-null	float64
26	instrumentalness	28372	non-null	float64
27	valence	28372	non-null	float64
28	energy	28372	non-null	float64
29	topic	28372	non-null	object
30	---	28372	non-null	float64

```

30 age                28372 non-null float64
dtypes: float64(23), int64(3), object(5)
memory usage: 6.7+ MB
None

```

As can be seen, we have 30 columns and 28,372 rows. There are no missing values in any column which makes our analysis easier and we do not have to do any imputation. Most of the variables are float64 data types with the exception being artist_name, track_name, genre, lyrics, and age which are objects. release_date and len are int64 variables which are still numeric. This means we have 25 numeric variables and 5 categorical ones.

✓ EDA.

We will begin by examining the dataset, by viewing the head to get a sense of the column layouts and their values. We will then use the info command to check the data types of each variable as well as the non-null count.

```

1 # Drop unnecessary columns
2 df.drop(columns=['Unnamed: 0'], inplace=True)

```

```

1 # For categorical columns
2 categorical_cols = df.select_dtypes(include=['object']).columns
3 categorical_cols = categorical_cols[categorical_cols != 'lyrics']
4 for col in categorical_cols:
5     print(f'{col} unique values: {df[col].nunique()}')
6     print(df[col].value_counts().head()) # Show top 5 most frequent categ
7
8 # Data appears to be mostly numerical types with the majority in float64 ty
9 # Genre, Lyrics, and Topic.
10 # For numerical columns
11 numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
12 print(df[numerical_cols].describe())

```

```

↗ artist_name unique values: 5426
artist_name
johnny cash      190
ella fitzgerald  188
dean martin      146
willie nelson    131
george jones     107
Name: count, dtype: int64
track_name unique values: 23689

```

```

track_name
tonight      17
stay         15
hold on      15
without you  14
home         13
Name: count, dtype: int64
genre unique values: 7
genre
pop          7042
country      5445
blues        4604
rock         4034
jazz         3845
Name: count, dtype: int64
topic unique values: 8
topic
sadness      6096
violence     5710
world/life   5420
obscene      4882
music        2303
Name: count, dtype: int64

```

	release_date	len	dating	violence	world/life
count	28372.000000	28372.000000	28372.000000	28372.000000	28372.000000
mean	1990.236888	73.028444	0.021112	0.118396	0.120973
std	18.487463	41.829831	0.052370	0.178684	0.172200
min	1950.000000	1.000000	0.000291	0.000284	0.000291
25%	1975.000000	42.000000	0.000923	0.001120	0.001170
50%	1991.000000	63.000000	0.001462	0.002506	0.006579
75%	2007.000000	93.000000	0.004049	0.192608	0.197793
max	2019.000000	199.000000	0.647706	0.981781	0.962105

	night/time	shake	the audience	family/gospel	romantic \
count	28372.000000		28372.000000	28372.000000	28372.000000
mean	0.057387		0.017422	0.017045	0.048681
std	0.111923		0.040670	0.041966	0.106095
min	0.000289		0.000284	0.000289	0.000284
25%	0.001032		0.000993	0.000923	0.000975
50%	0.001949		0.001595	0.001504	0.001754
75%	0.065842		0.010002	0.004785	0.042301
max	0.973684		0.497463	0.545303	0.940789

	communication	...	like/girls	sadness	feelings \
count	28372.000000	...	28372.000000	28372.000000	28372.000000
mean	0.076680	...	0.028057	0.129389	0.030996
std	0.109538	...	0.058473	0.181143	0.071652
min	0.000291	...	0.000284	0.000284	0.000289
25%	0.001144	...	0.000975	0.001144	0.000993
50%	0.002632	...	0.001595	0.005263	0.001754

We will now look at the Top 5 most common values for each categorical variable. We can see Johnny Cash is the most occurring artist, sadness is the most popular topic, and pop is the most popular genre.

```

1 # For categorical columns
2 categorical_cols = df.select_dtypes(include=['object']).columns
3 categorical_cols = categorical_cols[categorical_cols != 'lyrics']
4 for col in categorical_cols:
5     print(f'{col} unique values: {df[col].nunique()}')
6     print(df[col].value_counts().head()) # Show top 5 most frequent categ
7
8 # Data appears to be mostly numerical types with the majority in float64 ty
9 # Genre, Lyrics, and Topic.
10 # For numerical columns
11 numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
12 print(df[numerical_cols].describe())

```

⇒ artist_name unique values: 5426

artist_name

johnny cash 190

ella fitzgerald 188

dean martin 146

willie nelson 131

george jones 107

Name: count, dtype: int64

track_name unique values: 23689

track_name

tonight 17

stay 15

hold on 15

without you 14

home 13

Name: count, dtype: int64

genre unique values: 7

genre

pop 7042

country 5445

blues 4604

rock 4034

jazz 3845

Name: count, dtype: int64

topic unique values: 8

topic

sadness 6096

violence 5710

world/life 5420

obscene 4882

music 2303

Name: count, dtype: int64

	release_date	len	dating	violence	world/life
count	28372.000000	28372.000000	28372.000000	28372.000000	28372.000000
mean	1990.236888	73.028444	0.021112	0.118396	0.120973
std	18.487463	41.829831	0.052370	0.178684	0.172200
min	1950.000000	1.000000	0.000291	0.000284	0.000291
25%	1975.000000	42.000000	0.000923	0.001120	0.001170
50%	1991.000000	63.000000	0.001462	0.002506	0.006579
75%	2007.000000	93.000000	0.004049	0.192608	0.197793
max	2019.000000	199.000000	0.647706	0.981781	0.962105

	night/time	shake the audience	family/gospel	romantic \
count	28372.000000	28372.000000	28372.000000	28372.000000
mean	0.057387	0.017422	0.017045	0.048681
std	0.111923	0.040670	0.041966	0.106095
min	0.000289	0.000284	0.000289	0.000284
25%	0.001032	0.000993	0.000923	0.000975
50%	0.001949	0.001595	0.001504	0.001754
75%	0.065842	0.010002	0.004785	0.042301
max	0.973684	0.497463	0.545303	0.940789

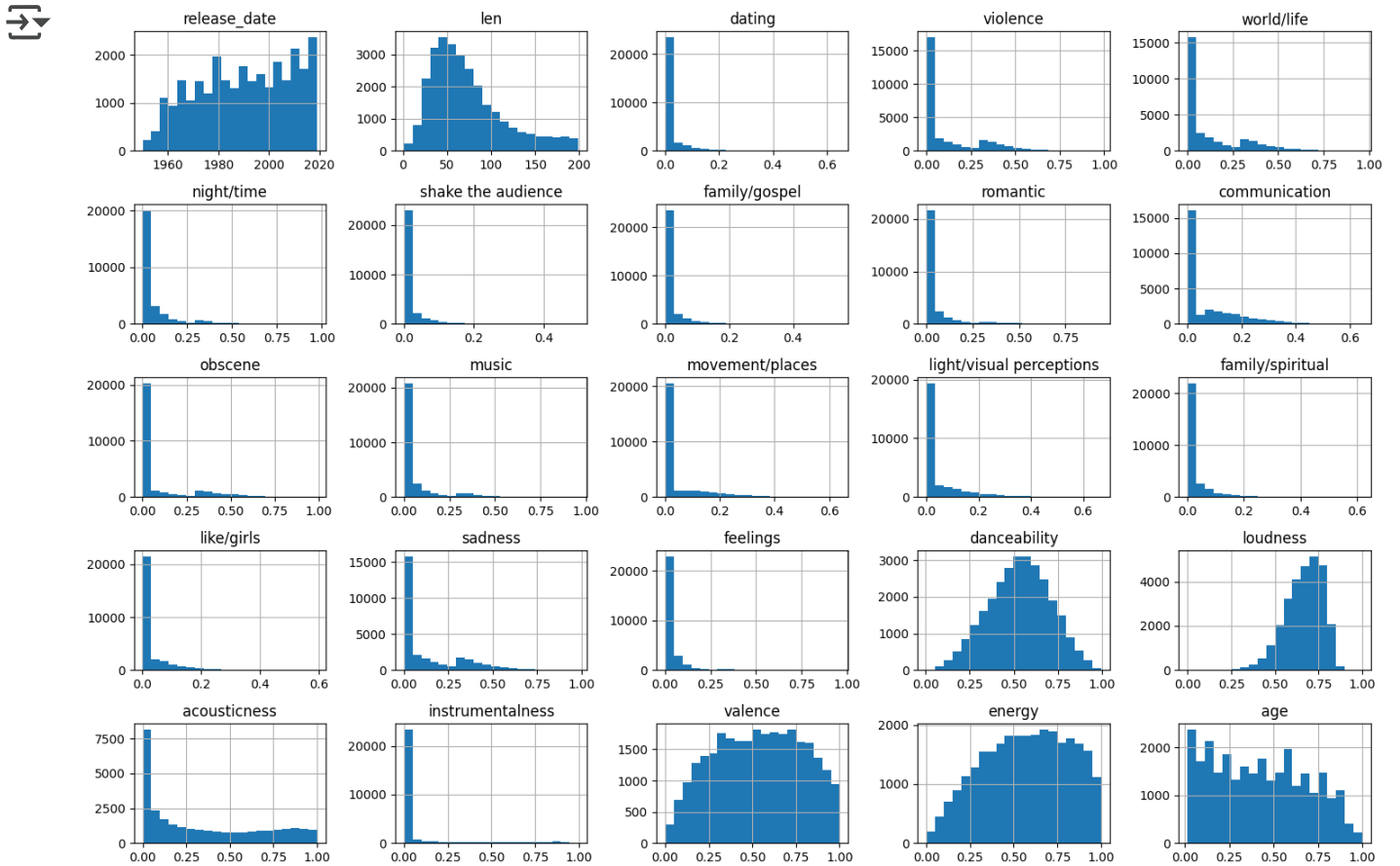
	communication ...	like/girls	sadness	feelings \
count	28372.000000	28372.000000	28372.000000	28372.000000
mean	0.076680	0.028057	0.129389	0.030996
std	0.109538	0.058473	0.181143	0.071652
min	0.000291	0.000284	0.000284	0.000289
25%	0.001144	0.000975	0.001144	0.000993
50%	0.002637	0.001505	0.005763	0.001754

We will now look at histograms of our 25 numerical variables to get a sense of how they are distributed. These histograms display count values so we can see if certain patterns emerge and this will help us find correlations between variables.

```

1 # Histogram for numerical columbns
2 df[numerical_cols].hist(bins=20, figsize=(15, 10))
3 plt.tight_layout()
4 plt.show()

```



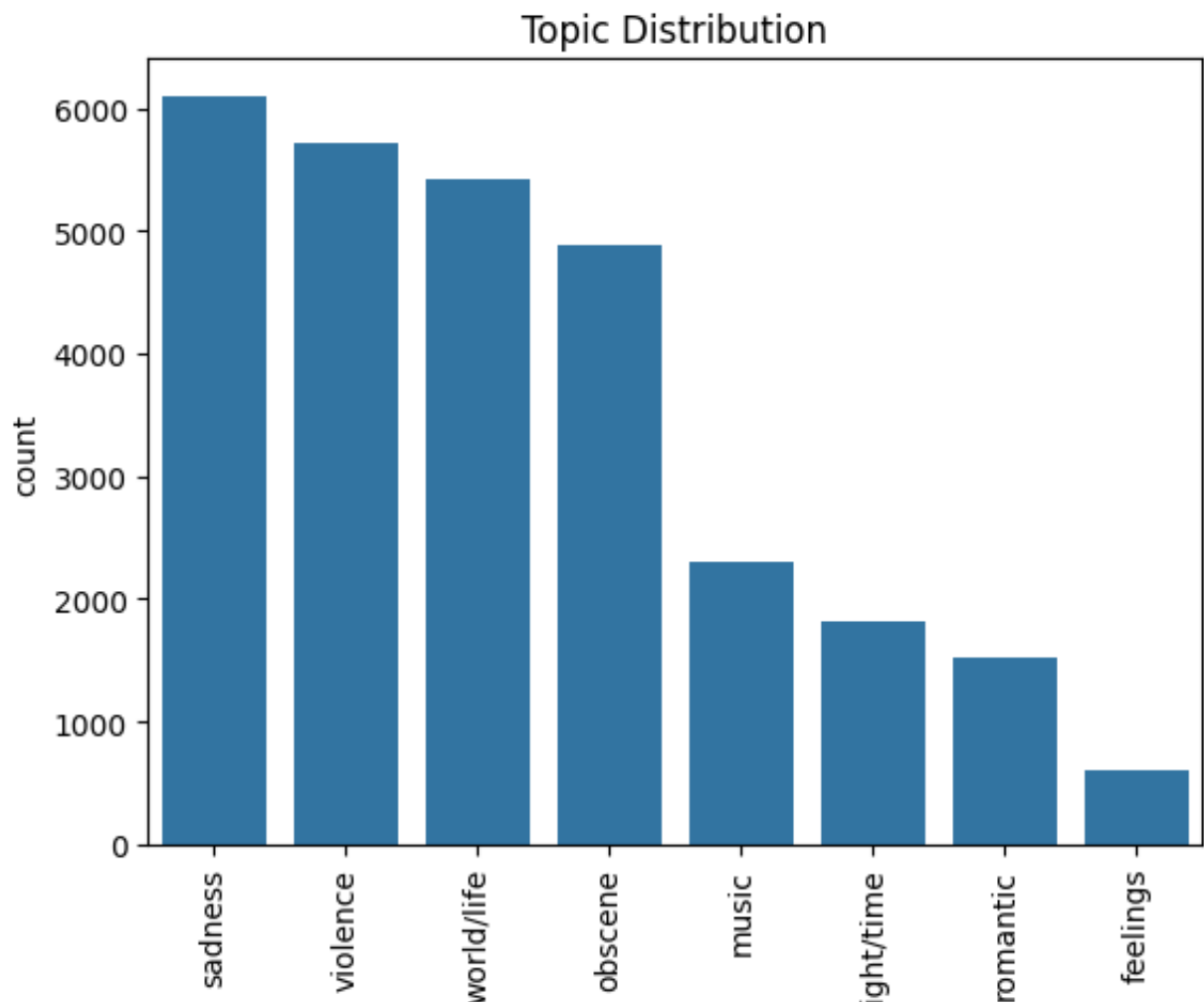
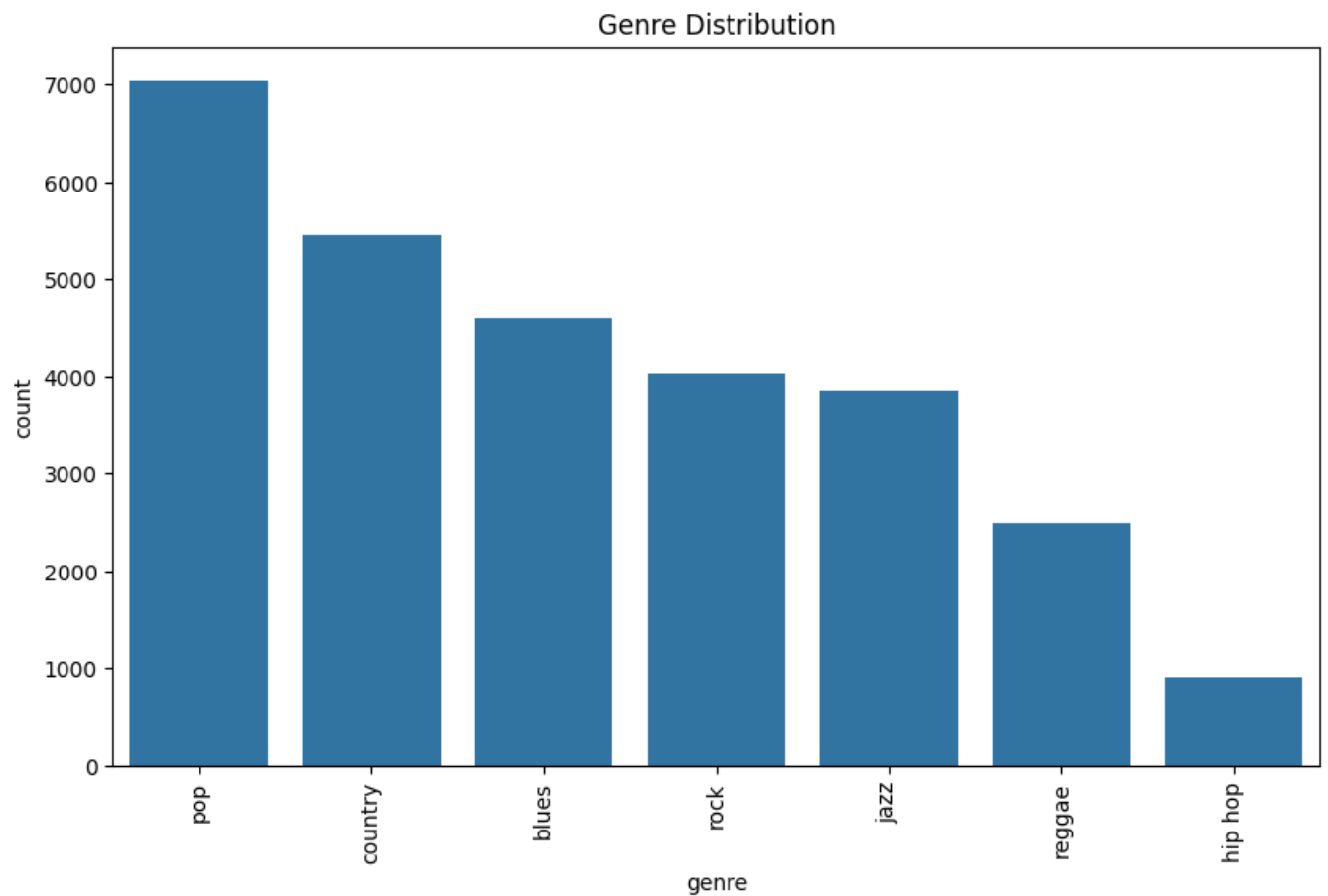
A few key patterns emerge from these histograms. Length is right tailed with a mean of about 50, and release date appears to have more songs from newer years. Danceability, valence, and energy have means of about 0.5 and have even tails on both sides. The bulk of the variables that contain emotional connotations of songs such as world/life, romantic, and sadness have very high means at just above 0, and have few values outside of this range.

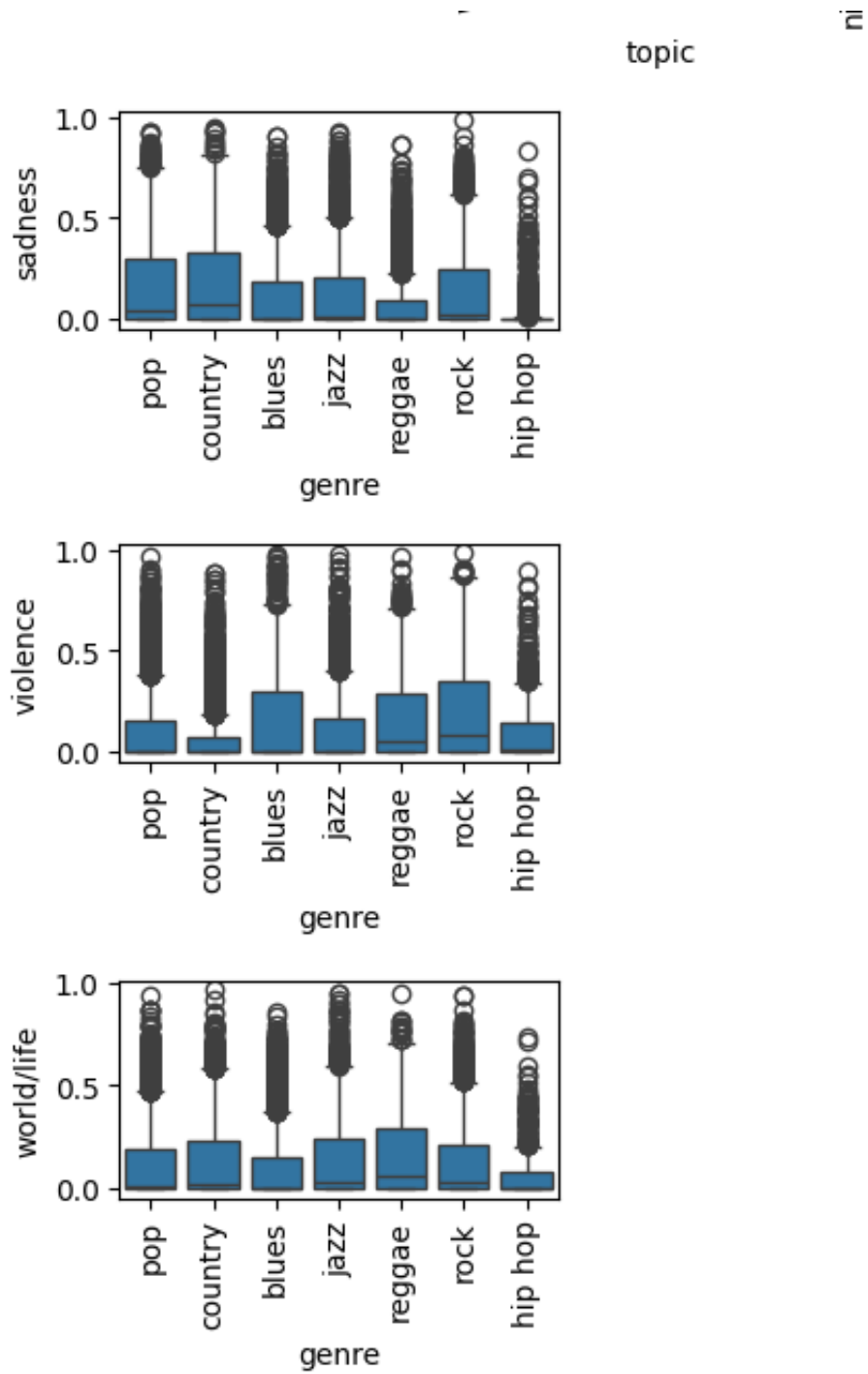
We will now look at barplots of the counts of each genre and song topic. Both genre and topic are expected to categorize similar songs together as differences between musical genres are quite large and users are predicted to generally like similar genre songs.

```

1 # Plot the distribution of genres in the dataset
2 plt.figure(figsize=(10, 6))
3 sns.countplot(data=df, x='genre', order=df['genre'].value_counts().index)
4 plt.title('Genre Distribution') # Set plot title
5 plt.xticks(rotation=90) # Rotate x-axis labels for better readability
6 plt.show()
7
8 # Plot the distribution of topics in the dataset
9 sns.countplot(data=df, x='topic', order=df['topic'].value_counts().index)
10 plt.title('Topic Distribution') # Set plot title
11 plt.xticks(rotation=90) # Rotate x-axis labels for better readability
12 plt.show()
13
14 # Create a boxplot to visualize the distribution of 'sadness' scores across
15 plt.figure(figsize=(6, 3))
16 plt.subplot(2, 2, 1) # Specify subplot location (1st plot in a 2x2 layout)
17 sns.boxplot(data=df, x='genre', y='sadness') # Boxplot for sadness by genre
18 plt.xticks(rotation=90) # Rotate x-axis labels
19 plt.show()
20
21 # Create a boxplot to visualize the distribution of 'violence' scores across
22 plt.figure(figsize=(6, 3))
23 plt.subplot(2, 2, 2) # Specify subplot location (2nd plot in a 2x2 layout)
24 sns.boxplot(data=df, x='genre', y='violence') # Boxplot for violence by genre
25 plt.xticks(rotation=90) # Rotate x-axis labels
26 plt.show()
27
28 # Create a boxplot to visualize the distribution of 'world/life' scores across
29 plt.figure(figsize=(6, 3))
30 plt.subplot(2, 2, 3) # Specify subplot location (3rd plot in a 2x2 layout)
31 sns.boxplot(data=df, x='genre', y='world/life') # Boxplot for world/life by genre
32 plt.xticks(rotation=90) # Rotate x-axis labels
33 plt.show()
34

```

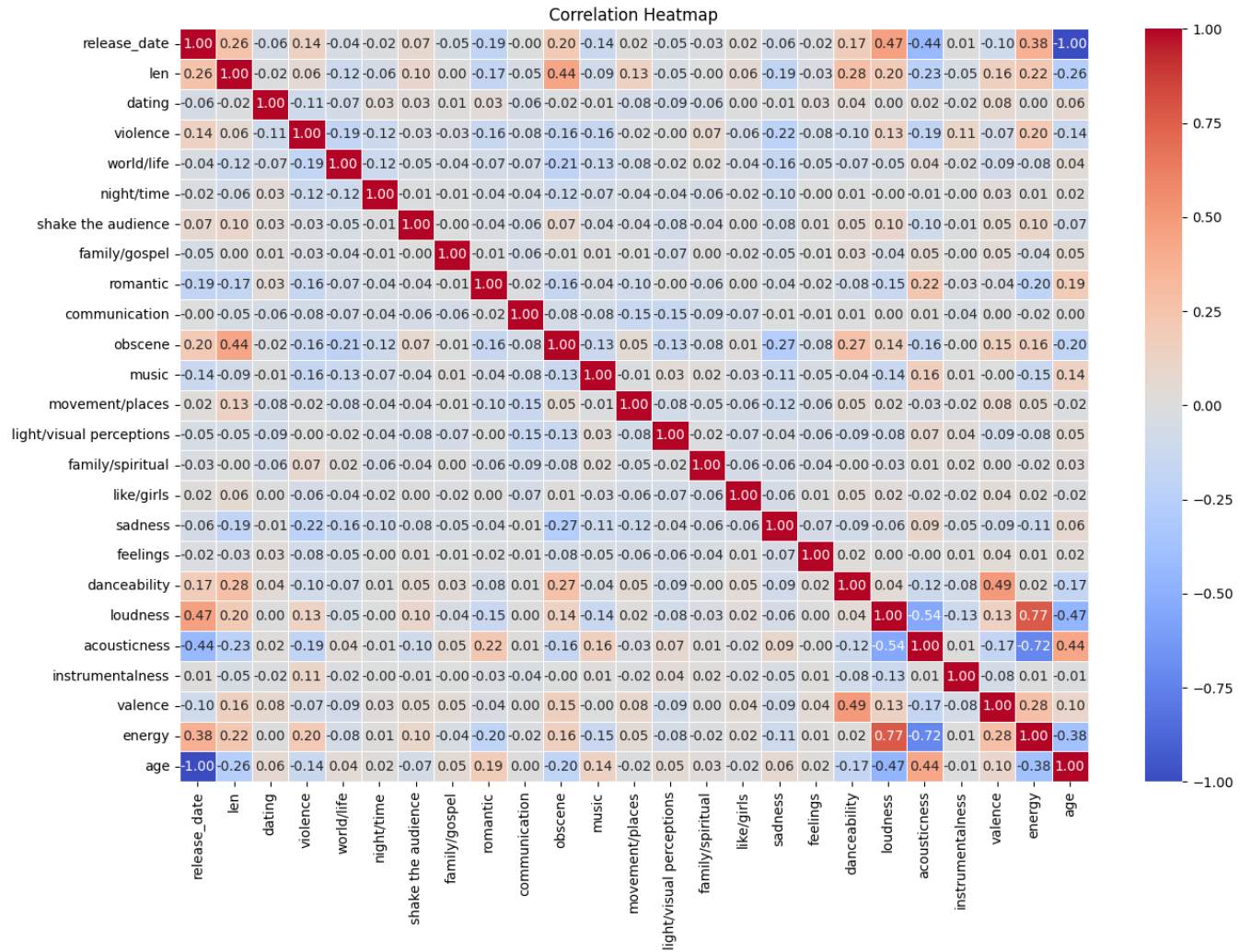




As can be seen, pop is the top category with a difference of about 1500 from the second place country. Blues, Jazz, and Rock have similar counts and reggae and hip-hop have much less. In terms of topic there is a big difference in the values of the top 4 topics and the bottom 4. The top 4: sadness, violence, world/life, and obscene have values around 4800-6000 and the bottom 4 are all below 2500. We can generally expect to see these top genres and topics being recommended more by our NLP algorithm.

We will now look at a correlation matrix of all numeric values in our dataset. We will look for variables for strong positive and negative correlations as these relationships will be impactful in interpreting the results of our song interpreter.

```
1 # Correlation matrix
2 plt.figure(figsize=(15, 10))
3 corr = df[numerical_cols].corr()
4 sns.heatmap(corr, annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5)
5 plt.title('Correlation Heatmap')
6 plt.show()
```



We can see that the largest Positive Correlations are between: Len and Obscene, Loudness and Energy, and Danceability and Valence We can see that the largest Negative Correlations are between: Energy and Acousticness, Age and Loudness, and Loudness and Acousticness Intuitively many of these make sense, loud songs are more likely to have more energy and acousticness would lead to quiet songs which would have a negative correlation with loudness. An interesting unexpected correlation is that longer songs seem more likely to be more obscene and also that the valence (positivity or negativity) of a song is correlated with its danceability.

We will now look at scatterplots with a best fit linear regression line of the top 3 positive and negative relationships based on our correlation matrix. We want to visualize how strong these relationships are as these will provide good insight into what factors we can expect to be most influential in our algorithm.

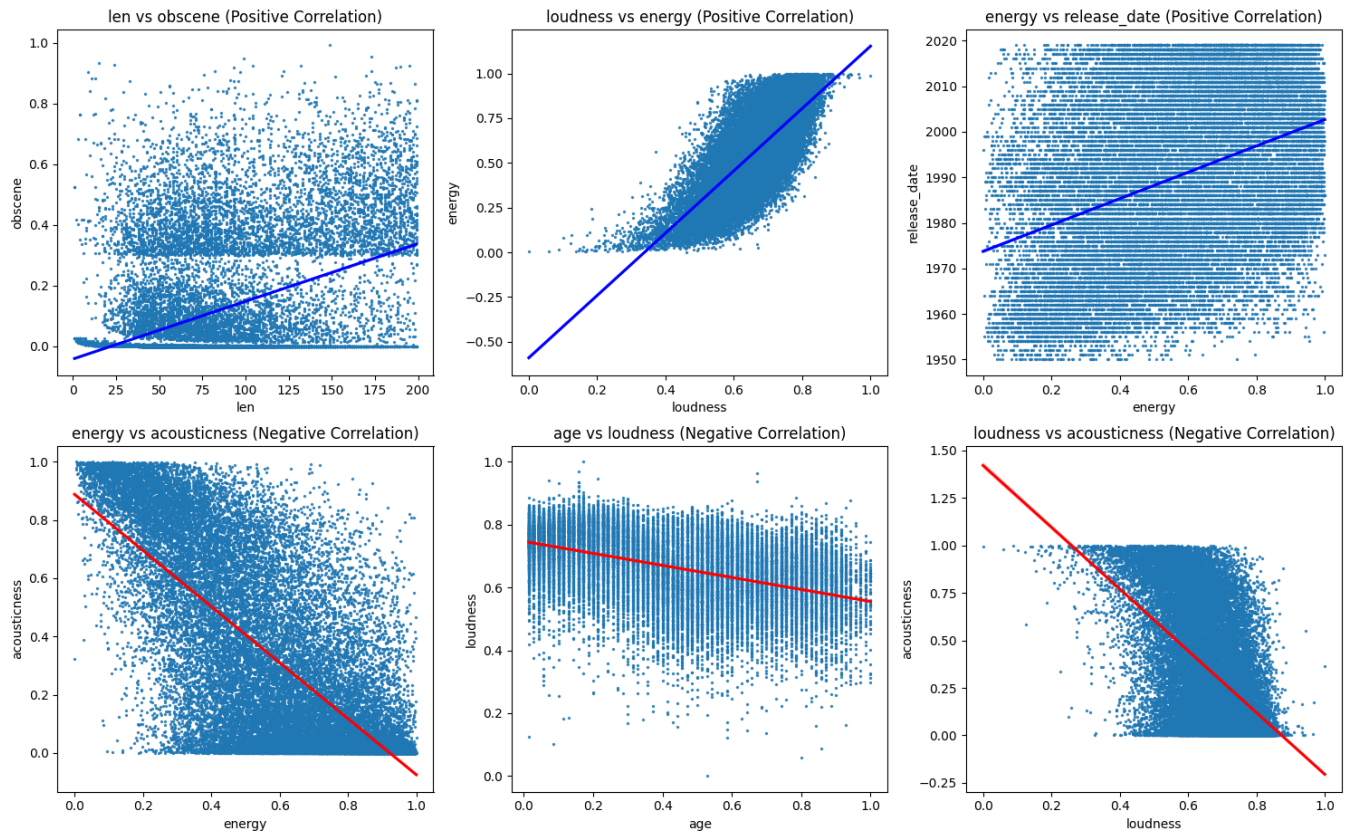
```

1 # Define the pairs for plotting
2 positive_pairs = [('len', 'obscene'), ('loudness', 'energy'), ('energy', 'r
3 negative_pairs = [('energy', 'acousticness'), ('age', 'loudness'), ('loudne
4
5 # Create a figure and axis for each plot
6 fig, axs = plt.subplots(2, 3, figsize=(15, 10)) # 2 rows, 3 columns for 6
7 fig.suptitle('Scatter Plots Showing Correlation Between Different Features'
8
9 # Plot positive correlations with line of best fit
10 for i, (x_col, y_col) in enumerate(positive_pairs):
11     sns.regplot(x=x_col, y=y_col, data=df, ax=axs[0, i], scatter_kws={'s':
12     axs[0, i].set_title(f'{x_col} vs {y_col} (Positive Correlation)')
13     axs[0, i].set_xlabel(x_col)
14     axs[0, i].set_ylabel(y_col)
15
16 # Plot negative correlations with line of best fit
17 for i, (x_col, y_col) in enumerate(negative_pairs):
18     sns.regplot(x=x_col, y=y_col, data=df, ax=axs[1, i], scatter_kws={'s':
19     axs[1, i].set_title(f'{x_col} vs {y_col} (Negative Correlation)')
20     axs[1, i].set_xlabel(x_col)
21     axs[1, i].set_ylabel(y_col)
22
23 # Adjust the layout
24 plt.tight_layout()
25 plt.subplots_adjust(top=0.9) # Adjust title spacing
26
27 # Show the plot
28 plt.show()

```



Scatter Plots Showing Correlation Between Different Features



We can immediately see that loudness and energy has a very clear relationship and the regression line indicates a strong positive correlation. The negative relationships have very strong trend lines and pattern of points appear to support their accuracy. Loudness vs. Acousticness has a very strong negative correlation which makes sense as loudness, energy, and acousticness all are very closely related. (Len and Obscene) and (Energy Release_Date) have strong positive relationships but the pattern isn't quite as apparent without the best fit line.

We will now plot and examine a word cloud before doing our NLP process so we can visualize the most popular words in each song's lyrics. We expect songs with similar vocabularies to be recommended together.



We will now partition each word, remove stop words and plot a barplot of the most frequently occurring words.

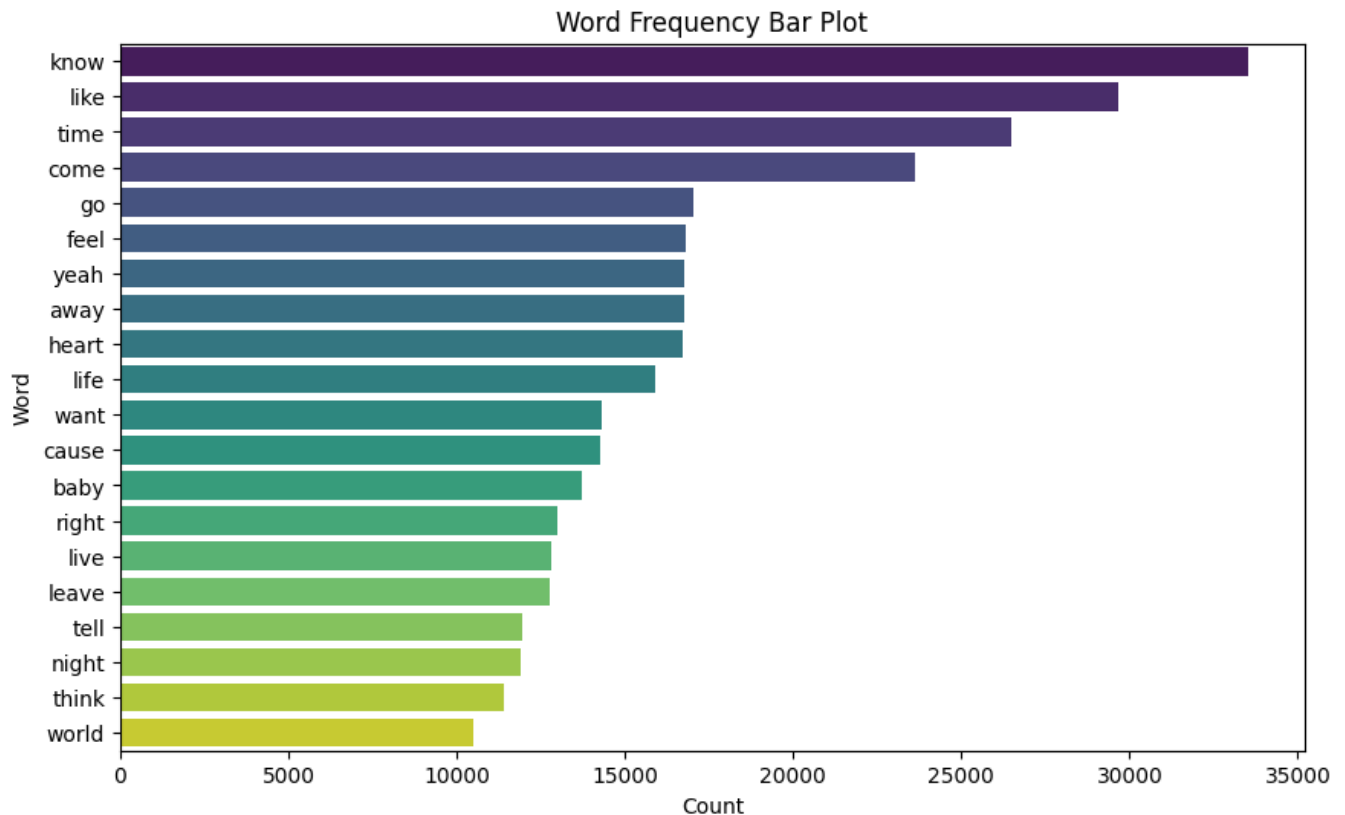
Page 20 of 37

```
3
4 # Tokenizing and removing stopwords
5 tokens = [word for word in all_lyrics.split() if word.lower() not in stop_w
6
7 # Get the most common words
8 common_words = Counter(tokens).most_common(20)
9 # Example common_words list
10
11 # Convert the list of tuples into a DataFrame
12 df_common_words = pd.DataFrame(common_words, columns=['Word', 'Count'])
13
14 # Sort by count in descending order for better visualization
15 df_common_words = df_common_words.sort_values(by='Count', ascending=False)
16
17 # Create the bar plot using seaborn
18 plt.figure(figsize=(10, 6))
19 sns.barplot(x='Count', y='Word', data=df_common_words, palette='viridis')
20
21 # Set plot labels and title
22 plt.title('Word Frequency Bar Plot')
23 plt.xlabel('Count')
24 plt.ylabel('Word')
25
26 # Display the plot
27 plt.show()
```

```
<ipython-input-13-7721d693af20>:19: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed

```
sns.barplot(x='Count', y='Word', data=df_common_words, palette='viridis')
```



We can see that the words know, like, time, and come have by far the most occurrences. We can see visualizing inspecting the words that these all seem like very common words musicians use in pop and country music in particular as well as music in general.

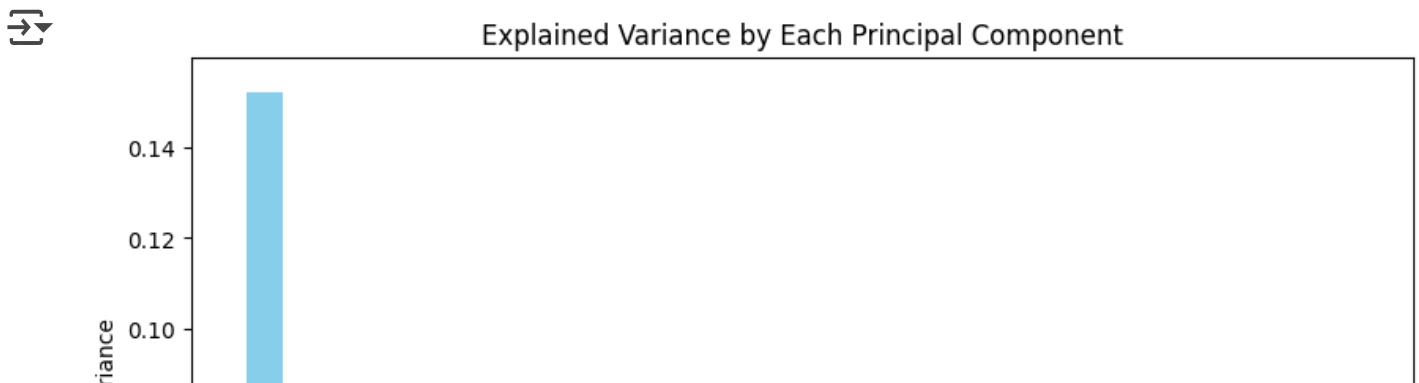
We will now perform Principle Component Analysis to examine which variables have the highest loadings on each principle component axis and how much each variable contributes to the cumulative variance.

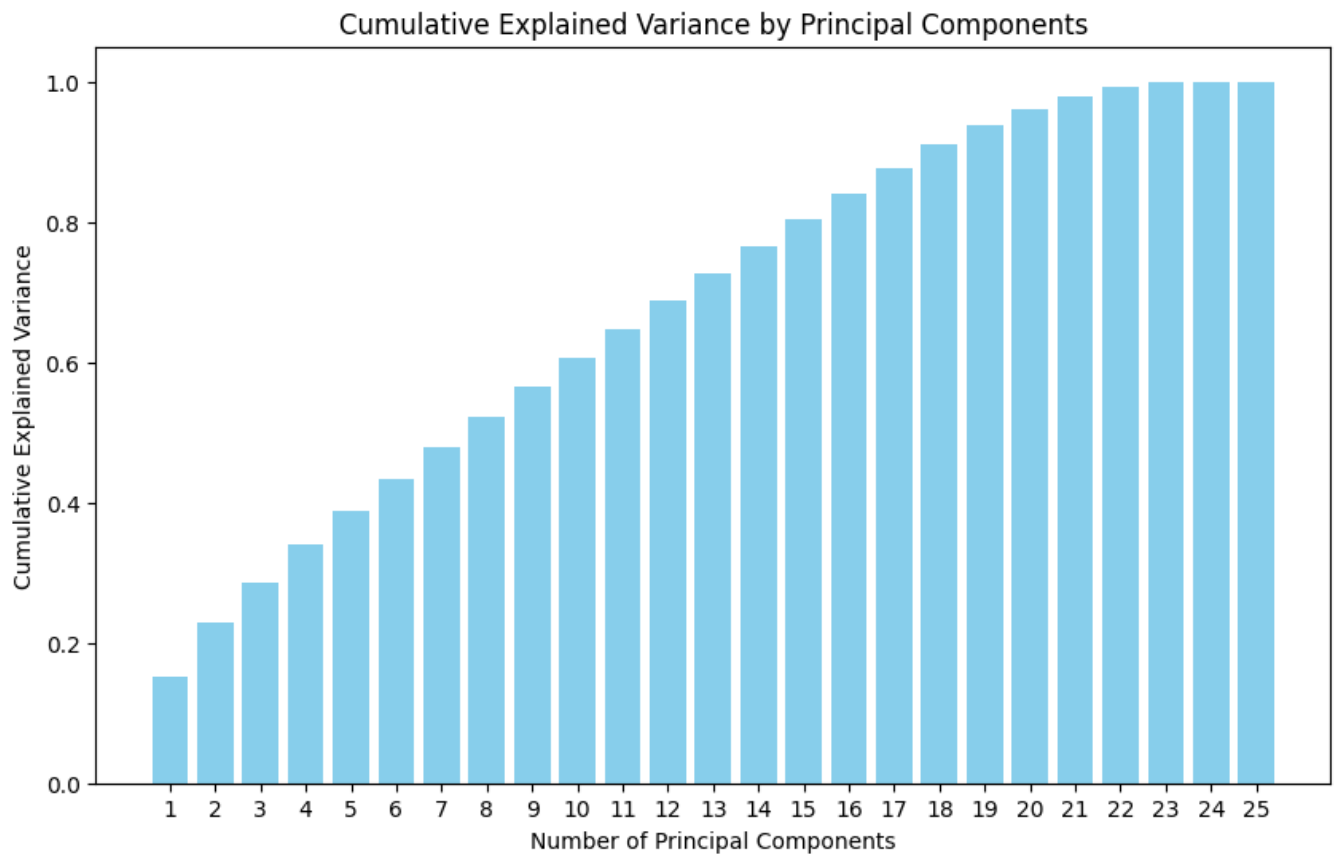
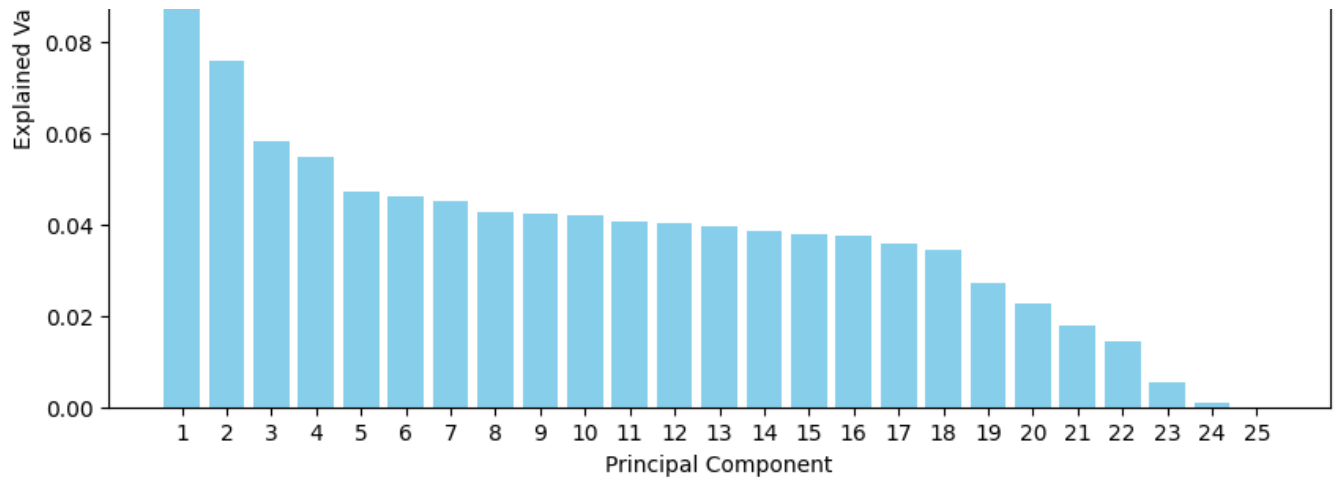
```
1 import matplotlib.pyplot as plt
```

```

2 from sklearn.decomposition import PCA
3 from sklearn.preprocessing import StandardScaler
4
5 # Select numerical columns
6 numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
7 numerical_data = df[numerical_cols]
8
9 # Standardize the data
10 scaler = StandardScaler()
11 numerical_data_scaled = scaler.fit_transform(numerical_data)
12
13 # Apply PCA
14 pca = PCA()
15 pca.fit(numerical_data_scaled)
16
17 # Explained variance for each component
18 explained_variance = pca.explained_variance_ratio_
19
20 # Create a bar plot for explained variance
21 plt.figure(figsize=(10, 6))
22 plt.bar(range(1, len(explained_variance) + 1), explained_variance, color='s')
23 plt.title('Explained Variance by Each Principal Component')
24 plt.xlabel('Principal Component')
25 plt.ylabel('Explained Variance')
26 plt.xticks(range(1, len(explained_variance) + 1))
27 plt.show()
28
29 # Cumulative explained variance (optional)
30 cumulative_explained_variance = np.cumsum(explained_variance)
31
32 # Create a bar plot for cumulative explained variance
33 plt.figure(figsize=(10, 6))
34 plt.bar(range(1, len(cumulative_explained_variance) + 1), cumulative_explai
35 plt.title('Cumulative Explained Variance by Principal Components')
36 plt.xlabel('Number of Principal Components')
37 plt.ylabel('Cumulative Explained Variance')
38 plt.xticks(range(1, len(cumulative_explained_variance) + 1))
39 plt.show()

```





We can see firstly based on the first plot that the first principle contributes by far the most to the dataset. PC1 has an explained variance of about 0.16 compared to just 0.08 for PC2. It also seems that PC3 - PC18 each contribute a similar amount of explained variance and linearly increase the total proportion explained. This indicates that to reach an explained cumulative variance of 0.9 we would still need about 18 PC components to do so.

```
1 # Select numerical columns
2 numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
3 numerical_data = df[numerical_cols]
4
5 # Standardize the data
6 scaler = StandardScaler()
7 numerical_data_scaled = scaler.fit_transform(numerical_data)
8
9 # Apply PCA
10 pca = PCA()
11 pca.fit(numerical_data_scaled)
12
13 # Extract the loadings (components) for the first 3 principal components
14 components_df = pd.DataFrame(pca.components_, columns=numerical_cols)
15
16 # Print the first 3 components
17 print("First 3 Principal Components:")
18 for i in range(3):
19     print(f"Principal Component {i+1}:")
20     top_loadings = components_df.iloc[i].abs().nlargest(6).index
21     for feature in top_loadings:
22         print(f"{feature}: {components_df.iloc[i][feature]:.4f}")
23     print("\n")
```

↔ First 3 Principal Components:

Principal Component 1:

energy: 0.3964
 release_date: 0.3946
 age: -0.3946
 acousticness: -0.3847
 loudness: 0.3819
 len: 0.2542

Principal Component 2:

valence: 0.4698
 danceability: 0.4613
 obscene: 0.3638
 len: 0.2812
 violence: -0.2595
 age: 0.2137

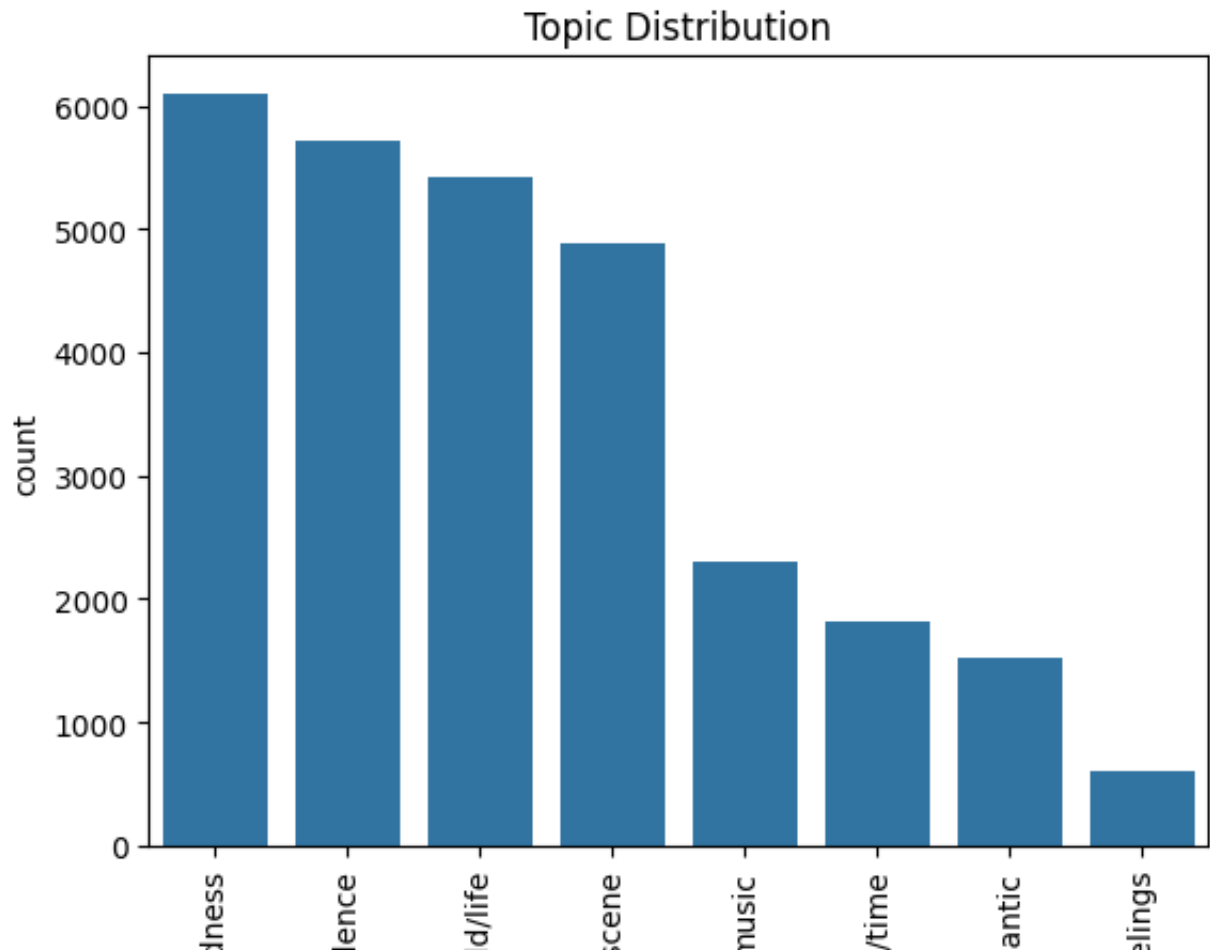
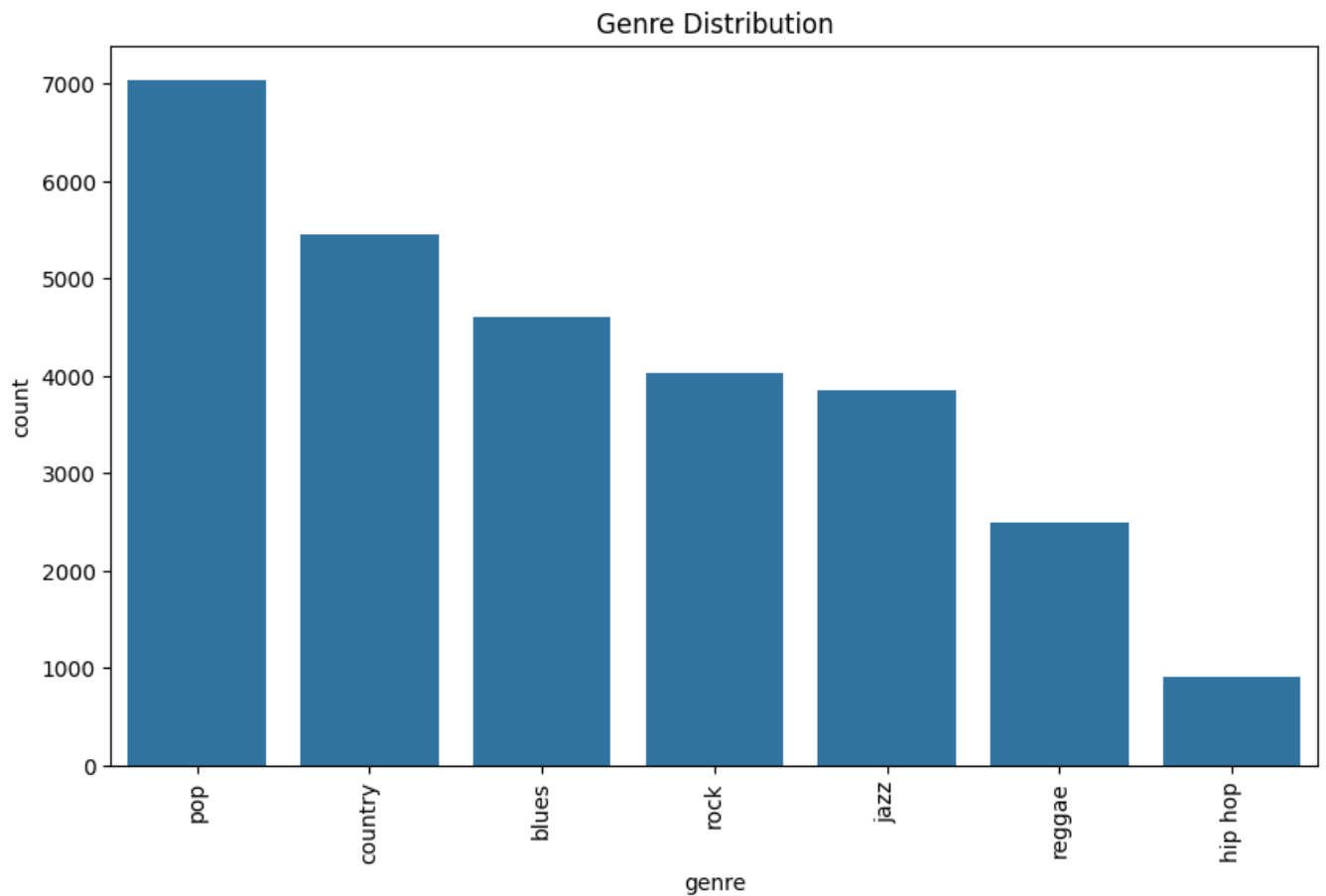
Principal Component 3:

sadness: 0.3171
 len: -0.2765
 dating: 0.2704
 communication: 0.2703
 movement/places: -0.2572
 obscene: -0.2534

In the table above, we calculated the 5 highest absolute value loadings for the first 3 PC components and also displayed the variables alongside. We can see that for PC1 the first four variables energy, age, release_date, and acousticness have very similar loading scores and there is a drop off of about 0.13 to the 5th greatest loading score variable, len. Valence and Danceability have the highest loadings for PC2 and Sadness has by far the greatest loading for PC3. These relationships between PC loadings are similar in many ways to some of the relationships in the correlation matrix above and indicate significant relationships between several predictor variables.

```
1 plt.figure(figsize=(10, 6))
2 sns.countplot(data=df, x='genre', order=df['genre'].value_counts().index)
3 plt.title('Genre Distribution')
4 plt.xticks(rotation=90)
5 plt.show()
6 # Plot for other categorical columns if needed (e.g., artist_name, topic)
7 sns.countplot(data=df, x='topic', order=df['topic'].value_counts().index)
8 plt.title('Topic Distribution')
9 plt.xticks(rotation=90)
```

```
10 plt.show()
```



sac
vio
worl
obs
r
night
rom
fet
topic

✓ Recommender System

To give users generalized recommendations of the most popular and most relevant songs, we will explore the world of recommender systems. We want to start off with a generalized recommender system that gives users the option to input a song in the Spotify music system and provides generalized recommendations of the most similar songs to the given input. This recommender system follows the concept of content based methods in recommender systems, where we recommend new music songs to users based on the similarity of songs that the user has liked or previewed before.

```

1 # Step 1: We import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from sklearn.decomposition import TruncatedSVD
7 from sklearn.neighbors import NearestNeighbors
8
9 # Step 2: We load the dataset
10 file_path = "/content/tcc_ceds_music.csv" # Path to the uploaded file in C
11 df = pd.read_csv(file_path)
12
13 # Step 3: Preprocessing for our system
14
15 if 'Unnamed: 0' in df.columns:
16     df.drop(columns=['Unnamed: 0'], inplace=True)
17
18 numeric_features = ['danceability', 'energy', 'acousticness', 'valence', 'l
19 X_numeric = df[numeric_features]
20

```

```

21 scaler = StandardScaler()
22 X_numeric_scaled = scaler.fit_transform(X_numeric)
23
24 df['processed_lyrics'] = df['lyrics'].apply(
25     lambda x: ' '.join([word.lower() for word in str(x).split() if word.isa
26 ])
27
28 vectorizer = TfidfVectorizer(max_features=5000, stop_words='english')
29 X_text = vectorizer.fit_transform(df['processed_lyrics'])
30
31 svd = TruncatedSVD(n_components=100, random_state=42)
32 X_text_reduced = svd.fit_transform(X_text)
33
34 # Step 4: We combine Numerical and Text Features
35 X_combined = np.hstack((X_numeric_scaled, X_text_reduced))
36
37 # Step 5: We build the Recommendation System
38 knn = NearestNeighbors(n_neighbors=11, metric='cosine') # 11 because input
39 knn.fit(X_combined)
40
41
42 def recommend_songs(song_name, df, knn_model, combined_features):
43
44     if song_name not in df['track_name'].values:
45         print(f"Error: The song '{song_name}' does not exist in the dataset")
46         return None
47
48     song_index = df[df['track_name'] == song_name].index[0]
49
50     song_features = combined_features[song_index].reshape(1, -1)
51
52     distances, indices = knn_model.kneighbors(song_features)
53
54     similar_song_indices = indices.flatten()[1:]
55
56     return df.iloc[similar_song_indices]
57
58 # Step 6: We test our Recommendation System
59 input_song = "i believe"
60 recommendations = recommend_songs(input_song, df, knn, X_combined)
61
62 if recommendations is not None:
63     print(recommendations[['track_name', 'artist_name', 'danceability', 'en
64

```

		track_name	artist_name
23549		i believe in the man in the sky	elvis presley
70		temptation	the platters
17243		pretend	nat king cole
17121		love is the sweetest thing	ray noble
17223		snow	rosemary clooney
200		try a little tenderness	the platters
7292		young love	sonny james
94		september in the rain	the platters
17595	when you're smiling (the whole world smiles wi...		nat king cole
152		birth of the blues	perry como

	danceability	energy	acousticness	valence
23549	0.421640	0.293271	0.865462	0.415705
70	0.384815	0.322301	0.931727	0.389942
17243	0.333911	0.275253	0.930723	0.293075
17121	0.284090	0.309288	0.968875	0.322960
17223	0.338243	0.301279	0.873494	0.332234
200	0.367486	0.297275	0.774096	0.367271
7292	0.412975	0.350330	0.788152	0.325021
94	0.286256	0.235211	0.969879	0.241550
17595	0.392397	0.340320	0.860442	0.342539
152	0.375068	0.287265	0.894578	0.389942

In order to implement the recommender system, we first perform numerical feature scaling in order to obtain various TF-IDF values. Next, we reduce dimensionality by applying Singular Value Decomposition before finally building the recommender system. This ensures we have a working recommender system where we are able to input a song name and get the top ten most recommended most similar songs.

✓ Recommender System, Using NLP

We have successfully built a recommender system using simpler content-based methods. Now, we want to see how incorporating NLP and sentiment analysis will enhance our recommendations.

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from sklearn.decomposition import TruncatedSVD
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.neighbors import NearestNeighbors

```

```

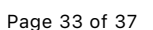
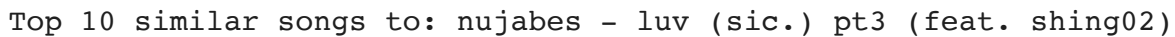
9 from wordcloud import WordCloud
10 from textblob import TextBlob
11 import re
12 import string
13
14 # Step 1: Load the dataset
15 file_path = "/content/tcc_ceds_music.csv" # Path to the uploaded file in
16 df = pd.read_csv(file_path)
17 df.drop(columns=['Unnamed: 0'], inplace=True)
18
19 # Create unique song identifier: artist - track
20 df['song_id'] = df['artist_name'] + " - " + df['track_name']
21
22 # Step 2: Clean the lyrics
23 def clean_lyrics(text):
24     text = str(text).lower()
25     text = re.sub(f"[{re.escape(string.punctuation)}]", "", text)
26     text = re.sub(r'\d+', '', text)
27     text = re.sub(r'\s+', ' ', text)
28     return text.strip()
29
30 df['clean_lyrics'] = df['lyrics'].apply(clean_lyrics)
31
32 # Step 3: Sentiment analysis
33 df['sentiment'] = df['clean_lyrics'].apply(lambda x: TextBlob(x).sentiment)
34
35 # Show WordCloud for positive/negative lyrics
36 positive_lyrics = ' '.join(df[df['sentiment'] > 0.4]['clean_lyrics'])
37 negative_lyrics = ' '.join(df[df['sentiment'] < -0.4]['clean_lyrics'])
38
39 fig, ax = plt.subplots(1, 2, figsize=(16, 6))
40 wordcloud_pos = WordCloud(width=800, height=400, background_color='white')
41 wordcloud_neg = WordCloud(width=800, height=400, background_color='black',
42
43 ax[0].imshow(wordcloud_pos, interpolation='bilinear')
44 ax[0].axis('off')
45 ax[0].set_title("Positive Lyrics WordCloud")
46
47 ax[1].imshow(wordcloud_neg, interpolation='bilinear')
48 ax[1].axis('off')
49 ax[1].set_title("Negative Lyrics WordCloud")
50
51 plt.tight_layout()
52 plt.show()
53
54 # Step 4: TF-IDF + SVD
55 vectorizer = TfidfVectorizer(max_features=5000, stop_words='english')
56 X_text = vectorizer.fit_transform(df['clean_lyrics'])

```

```

57
58 svd = TruncatedSVD(n_components=100, random_state=42)
59 X_text_reduced = svd.fit_transform(X_text)
60
61 # Step 5: Scale numeric features (including sentiment)
62 numeric_features = ['danceability', 'energy', 'acousticness', 'valence',
63                     'loudness', 'instrumentalness', 'age', 'sentiment']
64 X_numeric = df[numeric_features]
65
66 scaler = StandardScaler()
67 X_numeric_scaled = scaler.fit_transform(X_numeric)
68
69 # Step 6: Combine all features
70 X_combined = np.hstack((X_numeric_scaled, X_text_reduced))
71
72 # Step 7: Build the KNN model
73 knn = NearestNeighbors(n_neighbors=11, metric='cosine')
74 knn.fit(X_combined)
75
76 # Step 8: Song recommendation function using 'song_id'
77 def recommend_songs(song_id, df, knn_model, feature_matrix):
78     if song_id not in df['song_id'].values:
79         print(f"'{song_id}' not found in dataset.")
80         return None
81
82     song_index = df[df['song_id'] == song_id].index[0]
83     song_vector = feature_matrix[song_index].reshape(1, -1)
84
85     distances, indices = knn_model.kneighbors(song_vector)
86     similar_indices = indices.flatten()[1:] # exclude the input song itse
87
88     recommendations = df.iloc[similar_indices].copy()
89     return recommendations
90
91 # Step 9: Test the recommender
92 input_song_id = "nujabes - luv (sic.) pt3 (feat. shing02)" # Replace with
93 recommendations = recommend_songs(input_song_id, df, knn, X_combined)
94
95 if recommendations is not None:
96     print("Top 10 similar songs to:", input_song_id)
97     display(recommendations[['song_id', 'danceability', 'energy', 'valence'
98
99     # Step 10: Visualize sentiment of recommended songs
100     plt.figure(figsize=(10, 5))
101     sns.barplot(x='song_id', y='sentiment', data=recommendations)
102     plt.xticks(rotation=45, ha='right')
103     plt.title('Sentiment Scores of Recommended Songs')
104     plt.ylabel('Polarity (-1 to 1)')

```

easton corbin

jermaine dupri - let's

1.

di-

walker hay

black rob - pd w-

Song (Artist - Track)

Each song lyric is first assigned a sentiment polarity score ranging from -1 (very negative) to +1 (very positive). To capture the most meaningful language features, the top 5,000 important words from each lyric are extracted using TF-IDF. These high-dimensional features are then reduced to 100 components using Truncated SVD, allowing for efficient processing. In addition to lyrical data, numeric audio features such as danceability, energy, and valence are scaled and combined with the lyric embeddings to form a unified feature vector, referred to as $X_{combined}$. Using this combined representation, a K-Nearest Neighbors (KNN) model with cosine similarity is applied to find the 10 most similar songs based on both textual and acoustic characteristics, including track name, artist, danceability, energy, valence, and sentiment. Word clouds are generated to visualize word frequency patterns—one showing overall common words across all lyrics, and two others contrasting the vocabulary of positively vs. negatively scored lyrics. Finally, the system outputs the top 10 most similar tracks along with key metadata, and presents a bar chart of their sentiment scores to help assess the emotional similarity between the input song and the recommendations.

Methods/Results

The first step for our project was exploring and understanding the data set. We first got rid of any unnecessary columns, along with any NA data points of which there weren't many. After this, we created some histogram visuals to examine the distributions of our variables and look for any possible outliers that may need to be dealt with. We also created a correlation matrix to help us understand the hidden correlations within the data. This was incredibly

helpful, as it allowed us to see that our data was not just randomly generated but rather it contains meaningful patterns that mirror what we'd expect based on our previous knowledge. For example, we were able to find and plot the negative relationship between acousticness and energy, which is exactly what we'd expect based on the contents of these variables. We then prepared our NLP by looking at a wordcloud of our data's 'lyrics' column and then creating a barplot of the most common words, after tokenizing and removing stop words. The last step before building our recommenders was performing a PCA analysis on the data to understand the influence of specific variables and further strengthen our argument of statistically significant relationships being present in the data. We found that age, energy, release_date, and acousticness contributed most to our PC1, putting them among our most influential variables and once again suggesting a strong presence of correlation within the data.

Our project used two recommendation systems. The first recommendation system that we implemented was the simpler content based recommendation system. For this recommendation system, we apply the concept of content-based filtering. Content-based filtering refers to recommendation systems that are built based on the content being investigated. This is achieved by calculating the similarity between various items and comparable items that the user has liked before. For example, if the user had previously shown a preference for item A, and item A and item B are similar, then based on content-based filtering, the recommendation system recommends item B to the user. In order to implement the content-based filtering recommender system for our Spotify Music Recommender project, we provide recommendations for a specific song based on features including danceability, energy, acousticness, and valence. To start, we first perform numerical feature scaling in order to obtain various TF-IDF values. Next, we reduce dimensionality by applying Singular Value Decomposition before we combine the numerical and text features together. Finally, we build the function for our recommender system, with the function taking in inputs including "song_name", "df", "knn_model", and "combined_features". In order to test the validity of our recommender system, we input a song from the Spotify music list called "i believe", and we obtain the top ten most similar songs to "i believe" with values for each specific song based on the features of danceability, energy, acousticness, and valence. We now have a working recommender system where we are able to input a song name and get the top ten most recommended most similar songs.

✓ Conclusion

In this project, we successfully designed and implemented two **content-based music recommender systems** using a comprehensive dataset of over 28,000 Spotify tracks. Our objective was to generate meaningful song recommendations by analyzing both **audio-based features** and **lyrical content**.

The first system utilized standardized numeric features such as `danceability`, `energy`, `valence`, and others, combined with reduced TF-IDF vectors from lyrics. By applying **cosine similarity** through a K-Nearest Neighbors model, we generated recommendations based on musical and lyrical proximity. The second system expanded this foundation by incorporating **sentiment scores** derived from lyrics using **TextBlob**, enabling the model to account for emotional tone in the recommendation process.

Key takeaways include:

- **EDA (Exploratory Data Analysis)** revealed important trends: *Pop* was the dominant genre, and topics like *sadness* and *violence* were most prevalent. Strong correlations such as `energy ~ loudness` and `acousticness ~ loudness` aligned with musical intuition.
- **PCA** showed that dimensionality reduction was both necessary and effective. While the first principal component explained around 16% of the variance, nearly 18 components were needed to reach 90% cumulative variance.
- The **word cloud and frequency plots** confirmed that common emotional and expressive words dominate lyrical themes, with “know”, “like”, and “time” being the most frequent.
- Our recommendation system demonstrated strong practical value, effectively recommending emotionally and sonically similar tracks – as shown in the example for “*I Believe*” by Frankie Laine, which returned songs with similar sentiment and style from artists like *Nat King Cole* and *Elvis Presley*.

This project illustrates the power of combining **structured numerical data** with **unstructured text-based NLP** to build intelligent recommender systems. Moving forward, incorporating collaborative filtering or hybrid models could further improve personalization. Additionally, real-time user feedback could be integrated to dynamically tune recommendations.

