

Music Recommendation System with Natural Language Processing

Import Necessary Libraries

```
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
14
```

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

```
↳ [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
```

Load and Explore Dataset

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

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



	Unnamed: 0	artist_name	track_name	release_date	genre	lyrics	len	da
0	0	mukesh	mohabbat bhi jhoothi	1950	pop	hold time feel break feel untrue convince spea...	95	0.00
1	4	frankie laine	i believe	1950	pop	believe drop rain fall grow believe darkest ni...	51	0.03
2	6	johnnie ray	cry	1950	pop	sweetheart 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 × 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 age                      28372 non-null float64
dtypes: float64(23), int64(3), object(5)
memory usage: 6.7+ MB
None

```

```


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

```

```

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

```



	artist_name	track_name	release_date	genre	lyrics	len	dating	viol
0	mukesh	mohabbat bhi jhoothi	1950	pop	hold time feel break feel untrue convince spea...	95	0.000598	0.06
1	frankie laine	i believe	1950	pop	believe drop rain fall grow believe darkest ni...	51	0.035537	0.09
2	johnnie ray	cry	1950	pop	sweetheart send letter goodbye secret feel bet...	24	0.002770	0.00
					kiss lips			

3	pérez prado	patricia	1950	pop	kiss lips want stroll charm mambo chacha merin...	54	0.048249	0.00
4	giorgos papadopoulos	apopse eida oneiro	1950	pop	till darling till matter know till dream live ...	48	0.001350	0.00

5 rows x 30 columns

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

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

```
Data columns (total 30 columns):
```


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

```
dtypes: float64(23), int64(2), object(5)
```

```
memory usage: 6.5+ MB
```

```
None
```

```
1 # Data appears to be mostly numerical types with the majority in float64 typ
2 # Genre, Lyrics, and Topic.
3 # For numerical columns
4 numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
5 print(df[numerical_cols].describe())
```



	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
75%	0.132136	...	0.026622	0.235113	0.032622
max	0.645829	...	0.594459	0.981424	0.958810

	danceability	loudness	acousticness	instrumentalness \
count	28372.000000	28372.000000	2.837200e+04	28372.000000
mean	0.533348	0.665249	3.392347e-01	0.080049
std	0.173218	0.108434	3.267143e-01	0.211245
min	0.005415	0.000000	2.811248e-07	0.000000
25%	0.412975	0.595364	3.423598e-02	0.000000
50%	0.538612	0.679050	2.259028e-01	0.000085
75%	0.656666	0.749026	6.325298e-01	0.009335
max	0.993502	1.000000	1.000000e+00	0.996964

	valence	energy	age
count	28372.000000	28372.000000	28372.000000
mean	0.532864	0.569875	0.425187
std	0.250972	0.244385	0.264107
min	0.000000	0.000000	0.014286
25%	0.329143	0.380361	0.185714
50%	0.539365	0.580567	0.414286
75%	0.738252	0.772766	0.642857
max	1.000000	1.000000	1.000000

[8 rows x 25 columns]

```

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 categories

```

↗ 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
goodbye         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

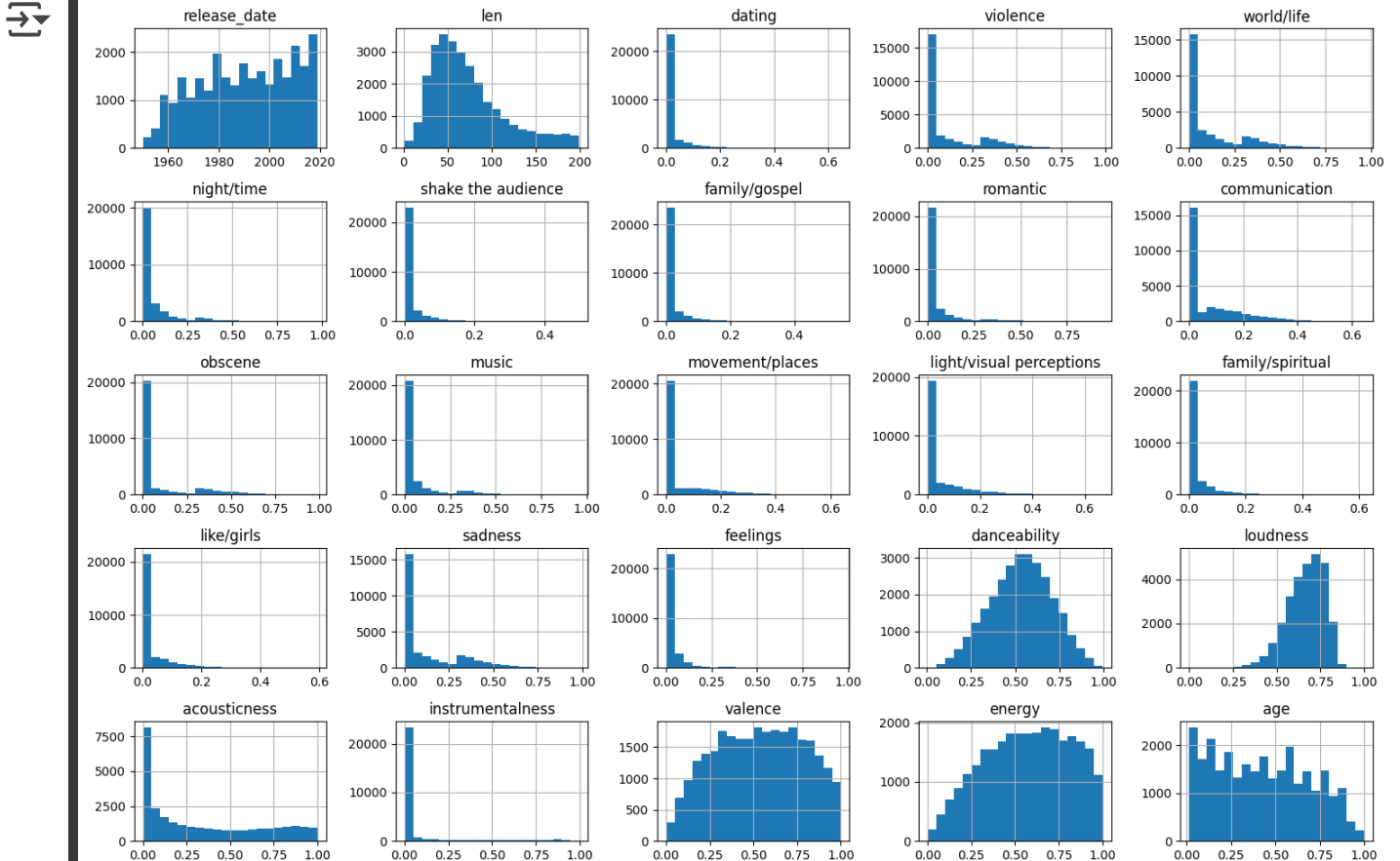
```

Name: count, dtype: int64

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 # Histogram for numerical columns
4 df[numerical_cols].hist(bins=20, figsize=(15, 10))
5 plt.tight_layout()
6 plt.show()

```



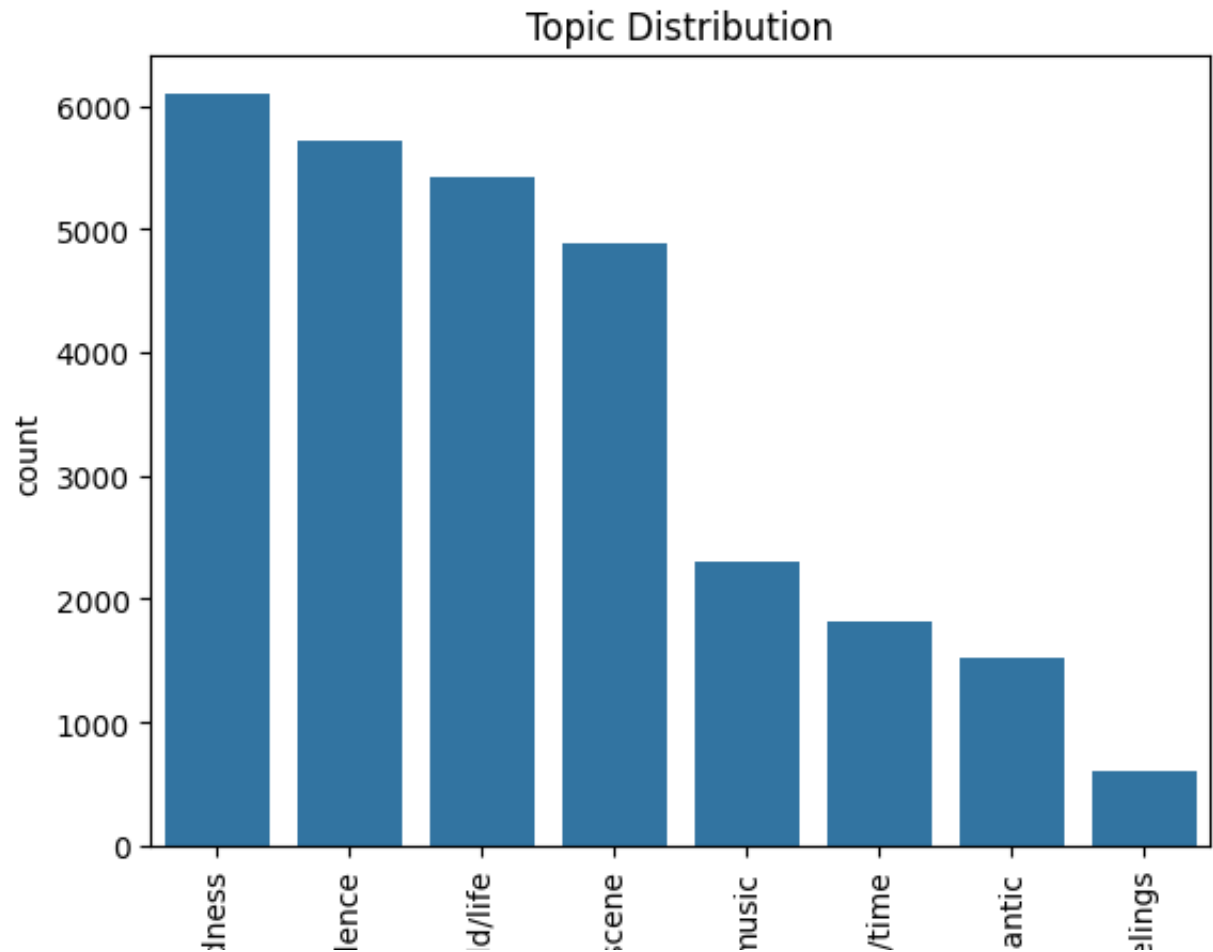
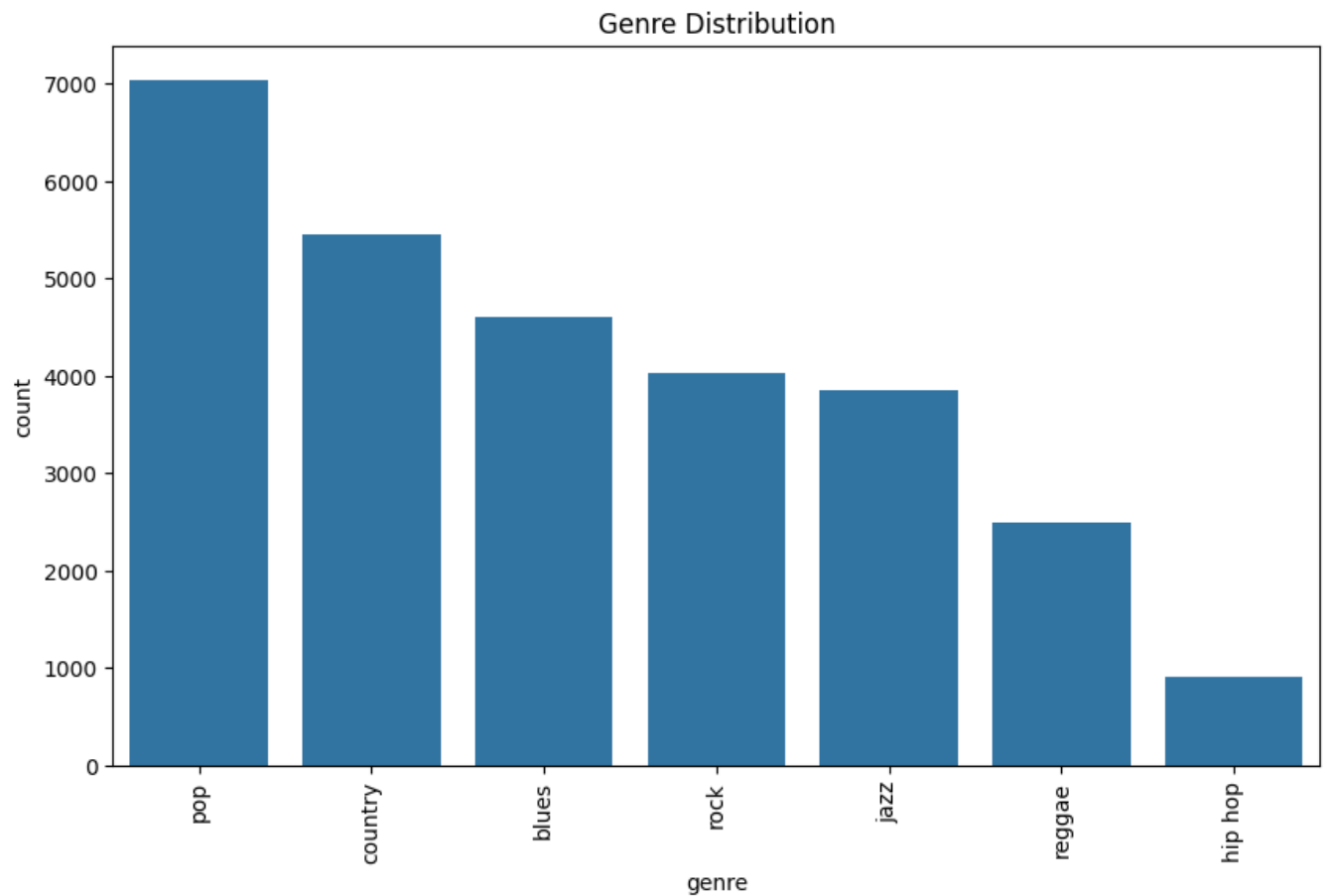
```

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()
```



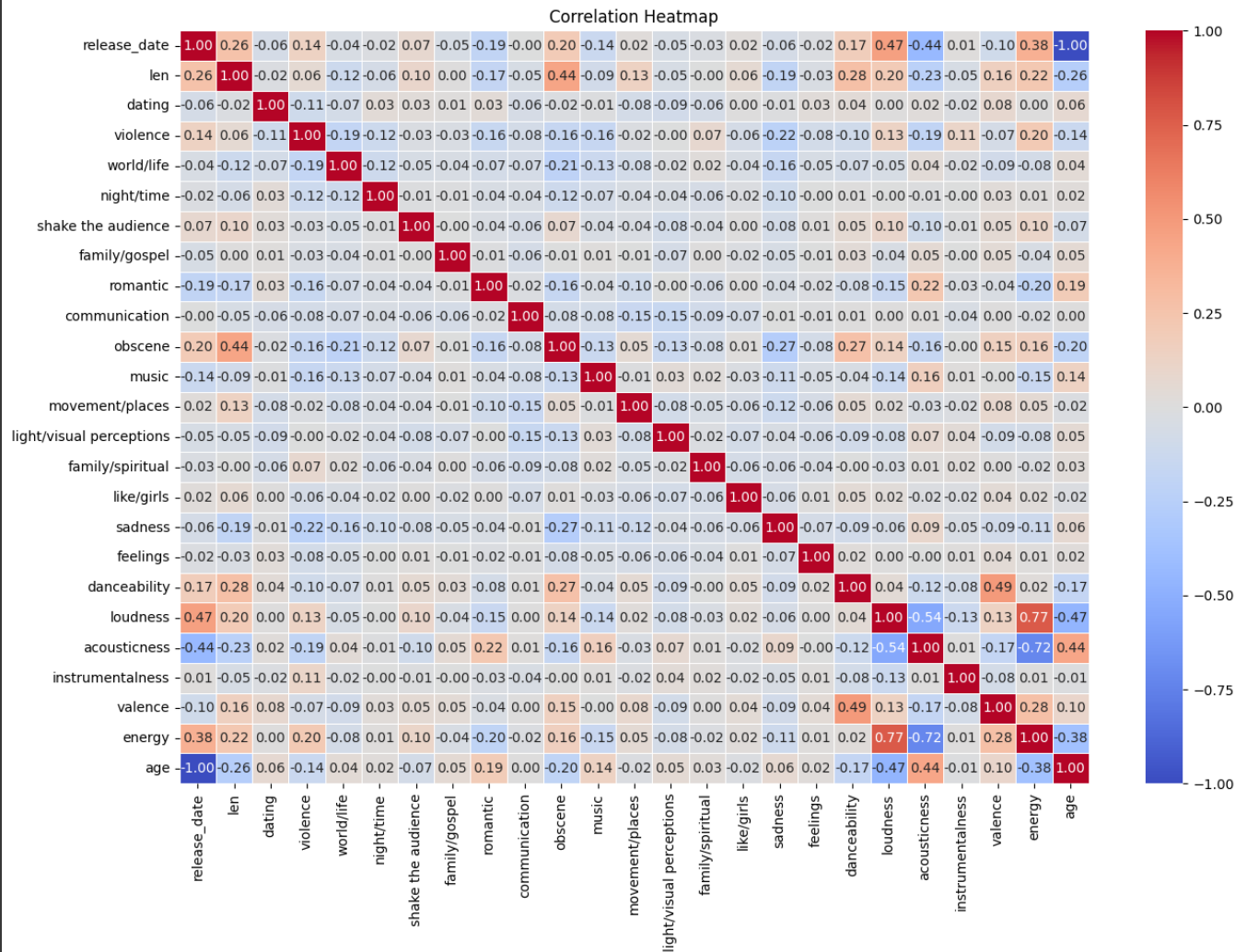


✓ Data Preprocessing

```
1 # Check for missing values
2 missing_data = df.isnull().sum()
3 missing_data_percentage = (missing_data / len(df)) * 100
4 print(missing_data[missing_data > 0])
5 print("Missing Data Percentage:\n", missing_data_percentage[missing_data_per
6
7 # Correlation matrix
8 plt.figure(figsize=(15, 10))
9 corr = df[numerical_cols].corr()
10 sns.heatmap(corr, annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5)
11 plt.title('Correlation Heatmap')
```



```
Series([], dtype: int64)
Missing Data Percentage:
Series([], dtype: float64)
Text(0.5, 1.0, 'Correlation Heatmap')
```



```
1 # Define the pairs for plotting
2 positive_pairs = [('len', 'obscene'), ('loudness', 'energy'), ('energy', 're
3 negative_pairs = [('energy', 'acousticness'), ('age', 'loudness'), ('loudnes
4 # Create a figure and axis for each plot
5 fig, axs = plt.subplots(2, 3, figsize=(15, 10)) # 2 rows, 3 columns for 6 pl
6 fig.suptitle('Scatter Plots Showing Correlation Between Different Features',
7 # Plot positive correlations with line of best fit
```

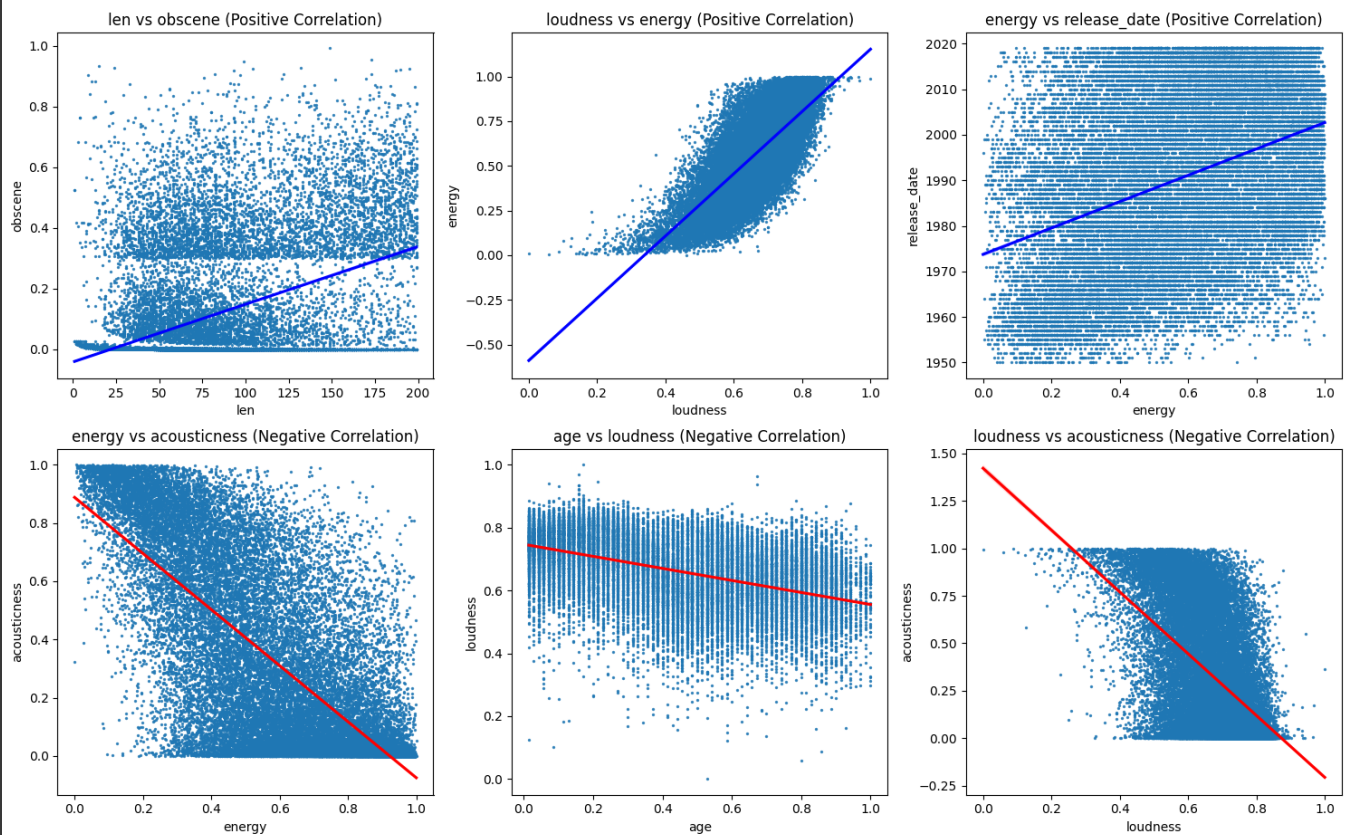
```

8 for i, (x_col, y_col) in enumerate(positive_pairs):
9     sns.regplot(x=x_col, y=y_col, data=df, ax=axes[0, i], scatter_kws={'s': 2},
10    axes[0, i].set_title(f'{x_col} vs {y_col} (Positive Correlation)')
11    axes[0, i].set_xlabel(x_col)
12    axes[0, i].set_ylabel(y_col)
13 # Plot negative correlations with line of best fit
14 for i, (x_col, y_col) in enumerate(negative_pairs):
15     sns.regplot(x=x_col, y=y_col, data=df, ax=axes[1, i], scatter_kws={'s': 2},
16    axes[1, i].set_title(f'{x_col} vs {y_col} (Negative Correlation)')
17    axes[1, i].set_xlabel(x_col)
18    axes[1, i].set_ylabel(y_col)
19 # Adjust the layout
20 plt.tight_layout()
21 plt.subplots_adjust(top=0.9) # Adjust title spacing
22 # Show the plot
23 plt.show()

```



Scatter Plots Showing Correlation Between Different Features






```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

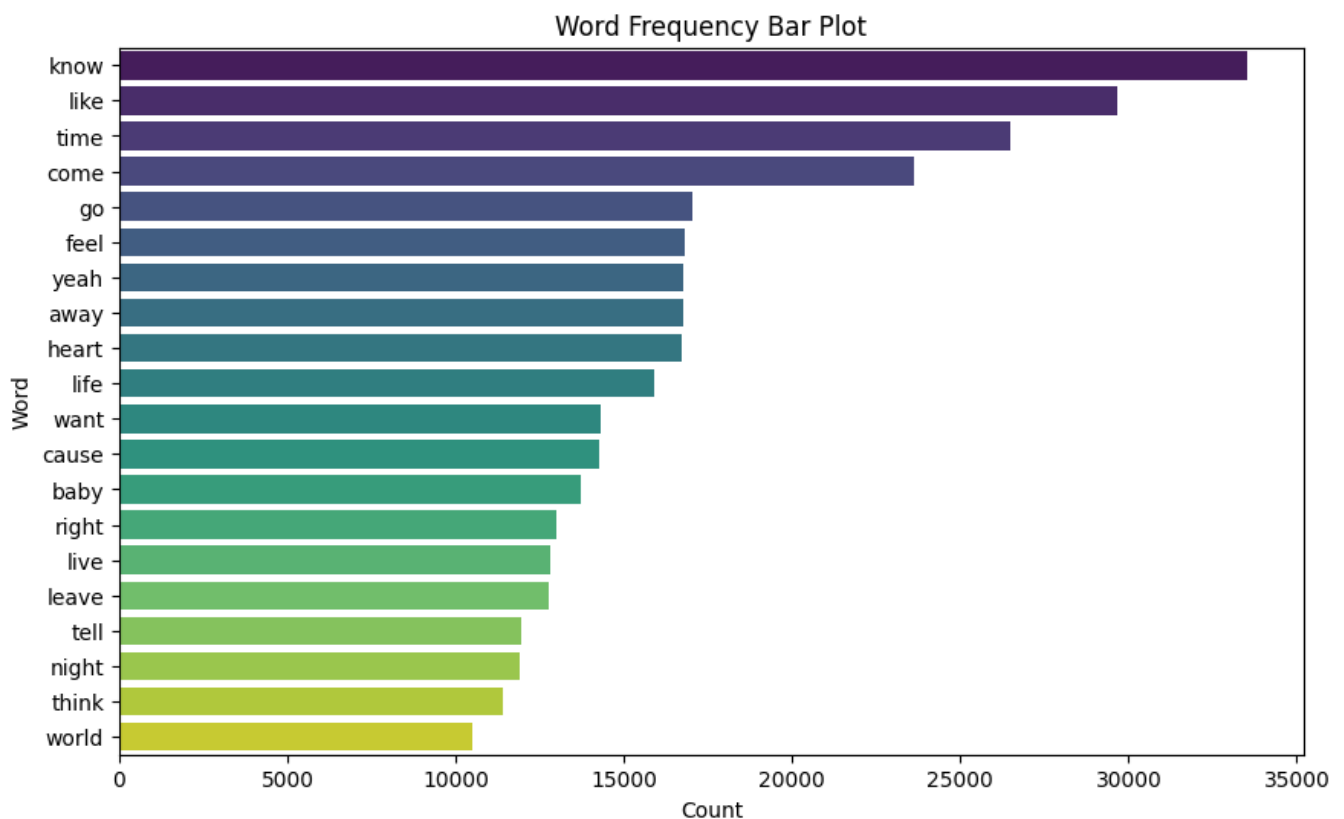
Page 13 of 27

```

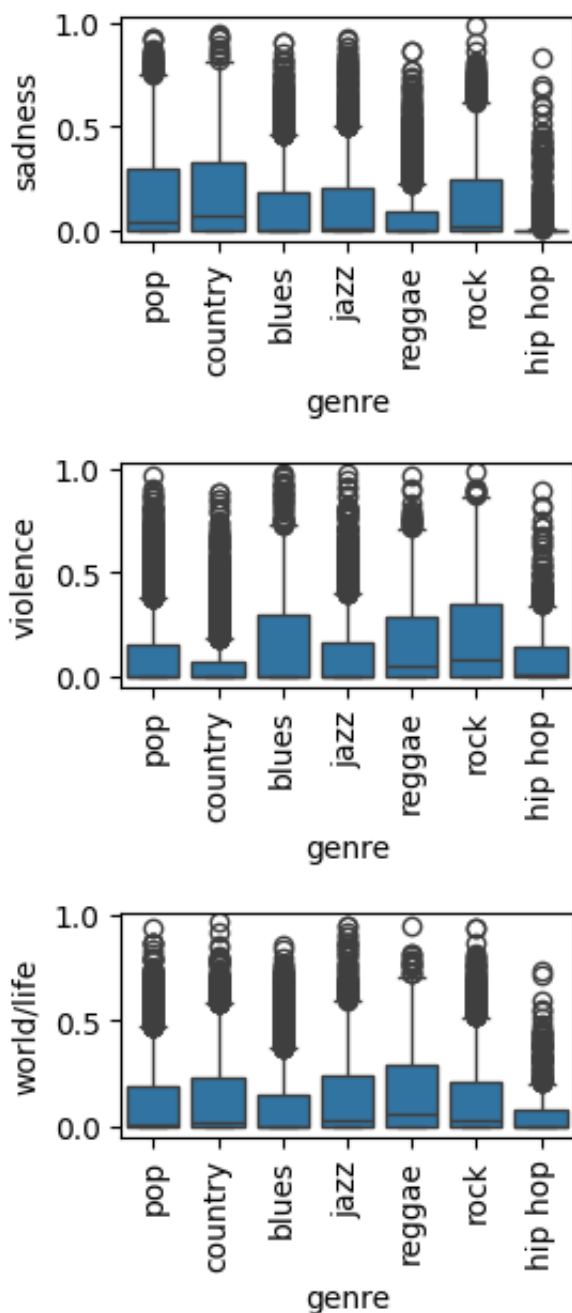
7 df_common_words = pd.DataFrame(common_words, columns=['Word', 'Count'])
8 # Sort by count in descending order for better visualization
9 df_common_words = df_common_words.sort_values(by='Count', ascending=False)
10 # Create the bar plot using seaborn
11 plt.figure(figsize=(10, 6))
12 sns.barplot(x='Count', y='Word', data=df_common_words, palette='viridis')
13 # Set plot labels and title
14 plt.title('Word Frequency Bar Plot')
15 plt.xlabel('Count')
16 plt.ylabel('Word')
17 # Display the plot
18 plt.show()

```

 <ipython-input-23-db81834be142>:12: FutureWarning:
 Passing `palette` without assigning `hue` is deprecated and will be removed
 sns.barplot(x='Count', y='Word', data=df_common_words, palette='viridis')



```
1 # Boxplot for numeric features based on genre
2 plt.figure(figsize=(6, 3))
3 plt.subplot(2,2,1)
4 sns.boxplot(data=df, x='genre', y='sadness') # Example: sadness by genre
5 plt.xticks(rotation=90)
6 plt.show()
7
8 # Boxplot for numeric features based on genre
9 plt.figure(figsize=(6, 3))
10 plt.subplot(2,2,2)
11 sns.boxplot(data=df, x='genre', y='violence') # Example: violence by genre
12 plt.xticks(rotation=90)
13 plt.show()
14
15 # Boxplot for numeric features based on genre
16 plt.figure(figsize=(6, 3))
17 plt.subplot(2,2,3)
18 sns.boxplot(data=df, x='genre', y='world/life')
19 plt.xticks(rotation=90)
20 plt.show()
```



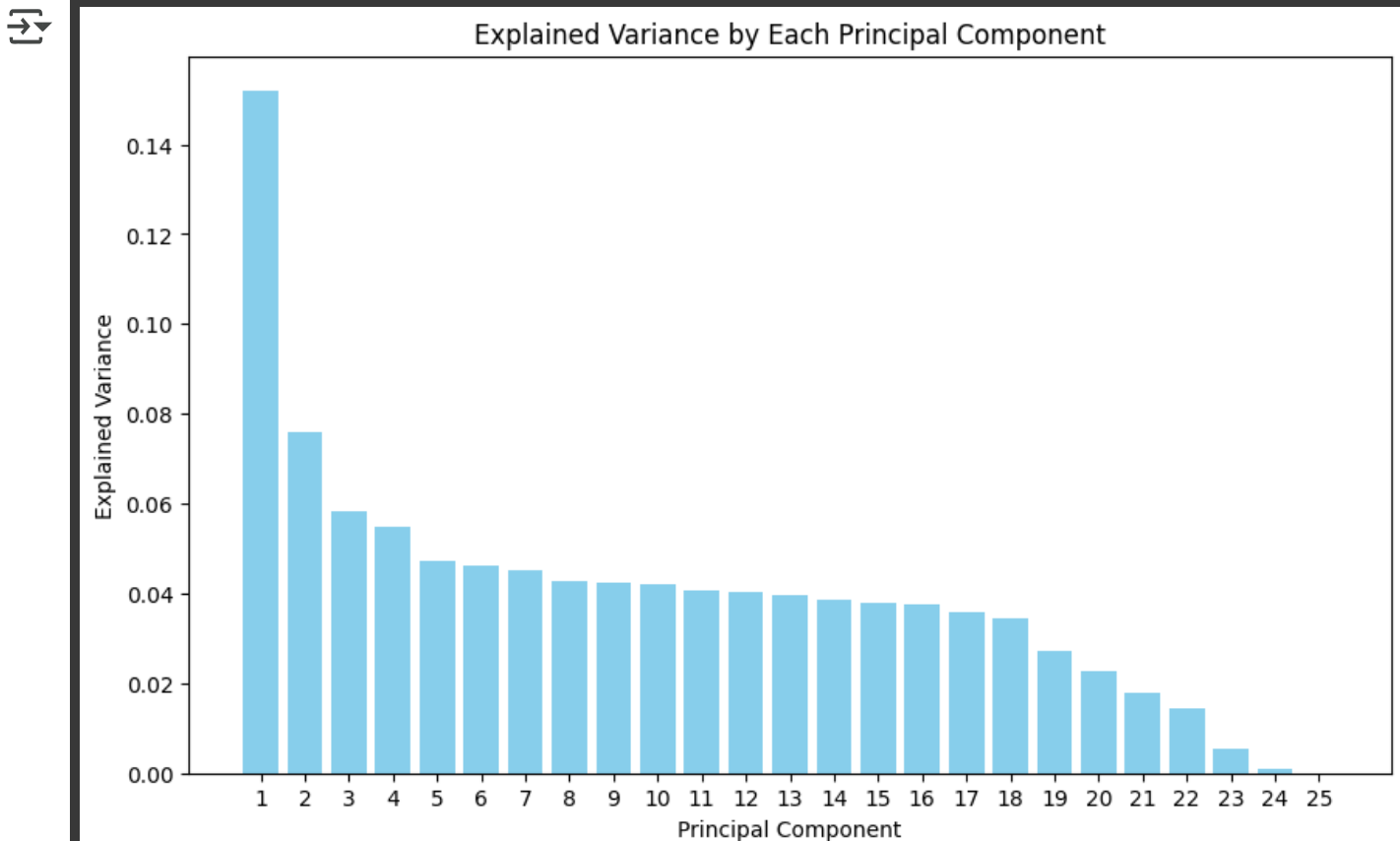
```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.decomposition import PCA
5 from sklearn.preprocessing import StandardScaler
6 # Assuming your DataFrame is called 'df'
7 # Select numerical columns
8 numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
9 numerical_data = df[numerical_cols]
10 # Standardize the data
11 scaler = StandardScaler()
12 numerical_data_scaled = scaler.fit_transform(numerical_data)
13 # Apply PCA

```



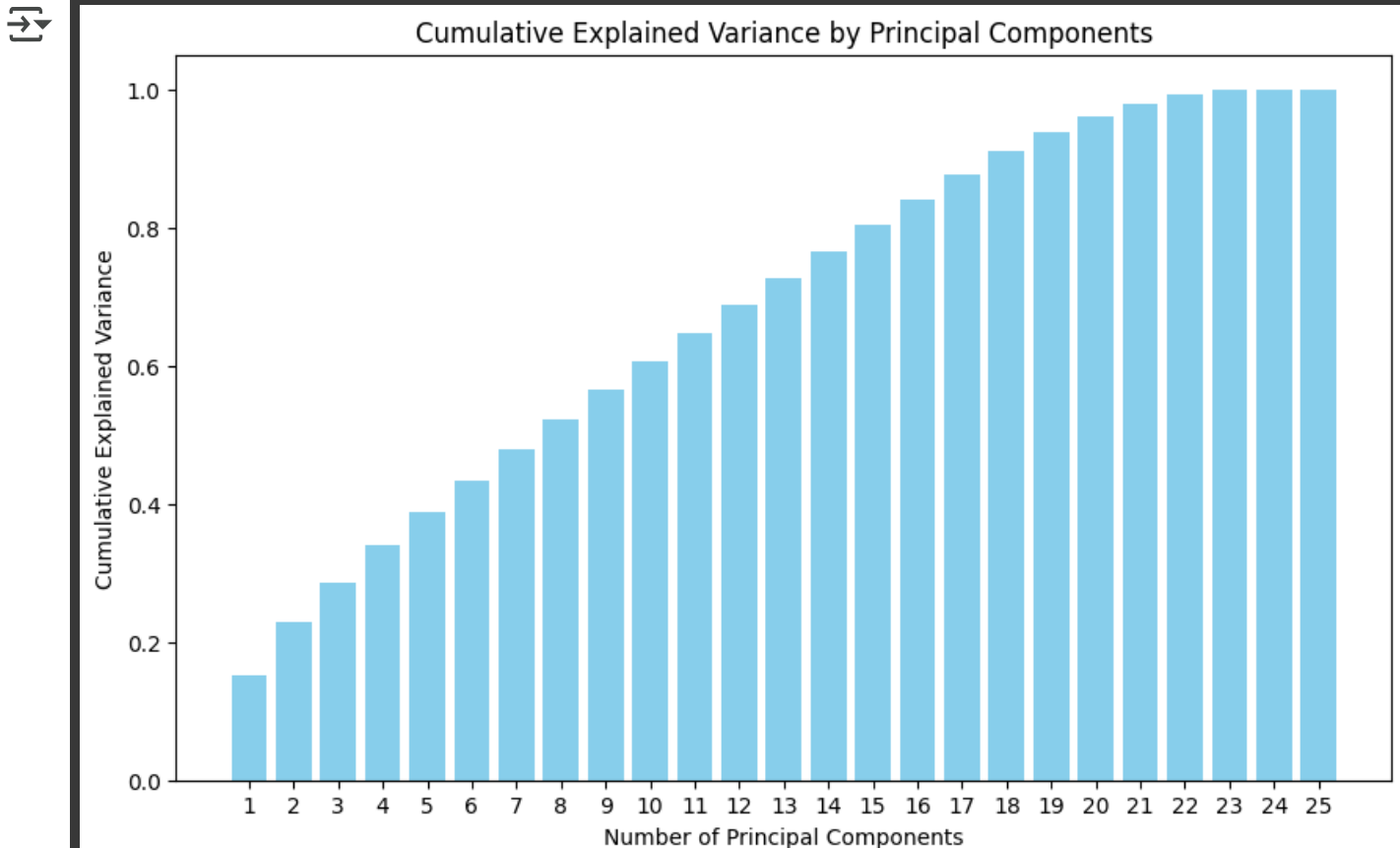
```
14 pca = PCA()
15 pca.fit(numerical_data_scaled)
16 # Explained variance for each component
17 explained_variance = pca.explained_variance_ratio_
18 # Create a bar plot for explained variance
19 plt.figure(figsize=(10, 6))
20 plt.bar(range(1, len(explained_variance) + 1), explained_variance, color='sk
21 plt.title('Explained Variance by Each Principal Component')
22 plt.xlabel('Principal Component')
23 plt.ylabel('Explained Variance')
24 plt.xticks(range(1, len(explained_variance) + 1))
25 plt.show()
```



```

1 cumulative_explained_variance = np.cumsum(explained_variance)
2 # Create a bar plot for cumulative explained variance
3 plt.figure(figsize=(10, 6))
4 plt.bar(range(1, len(cumulative_explained_variance) + 1), cumulative_explain
5 plt.title('Cumulative Explained Variance by Principal Components')
6 plt.xlabel('Number of Principal Components')
7 plt.ylabel('Cumulative Explained Variance')
8 plt.xticks(range(1, len(cumulative_explained_variance) + 1))
9 plt.show()

```



```

1 import pandas as pd
2 import numpy as np
3 from sklearn.decomposition import PCA
4 from sklearn.preprocessing import StandardScaler
5 # Assuming your DataFrame is called 'df'
6 # Select numerical columns
7 numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns

```

```

8 numerical_data = df[numerical_cols]
9 # Standardize the data
10 scaler = StandardScaler()
11 numerical_data_scaled = scaler.fit_transform(numerical_data)
12 # Apply PCA
13 pca = PCA()
14 pca.fit(numerical_data_scaled)
15 # Extract the loadings (components) for the first 3 principal components
16 components_df = pd.DataFrame(pca.components_, columns=numerical_cols)
17 # Print the first 3 components
18 print("First 3 Principal Components:")
19 for i in range(3):
20     print(f"Principal Component {i+1}:")
21     print(components_df.iloc[i])
22     print("\n")

```

↗ First 3 Principal Components:

Principal Component 1:

release_date	0.394588
len	0.254156
dating	-0.020535
violence	0.123877
world/life	-0.062965
night/time	-0.009484
shake the audience	0.090472
family/gospel	-0.026792
romantic	-0.170017
communication	-0.011355
obscene	0.210077
music	-0.121121
movement/places	0.055093
light/visual perceptions	-0.070238
family/spiritual	-0.016020
like/girls	0.026608
sadness	-0.101302
feelings	-0.005960
danceability	0.152552
loudness	0.381933
acousticness	-0.384740
instrumentalness	-0.021810
valence	0.115731
energy	0.396450
age	-0.394588

Name: 0, dtype: float64

Principal Component 2:

release_date	-0.213723
len	0.281174
dating	0.139691
violence	-0.259469
world/life	-0.139391

```

night/time          0.032565
shake the audience  0.091527
family/gospel       0.116300
romantic            0.020838
communication       -0.030348
obscene            0.363784
music              0.073331
movement/places    0.145732
light/visual perceptions -0.156542
family/spiritual   -0.050574
like/girls         0.094260
sadness            -0.156949
feelings           0.044765
danceability       0.461304
loudness           -0.122559
acousticness       0.084466
instrumentalness   -0.123907
valence            0.469775
energy             -0.080393
age                0.213723
Name: 1, dtype: float64

```

```

- . . . . .

```

```

1 # Step 1: 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: Load the dataset
10 file_path = "/content/tcc_ceds_music.csv" # Path to the uploaded file in Cc
11 df = pd.read_csv(file_path)
12
13 # Step 3: Preprocessing
14 # Drop unnecessary columns
15 if 'Unnamed: 0' in df.columns:
16     df.drop(columns=['Unnamed: 0'], inplace=True)
17
18 # Numerical Feature Scaling
19 numeric_features = ['danceability', 'energy', 'acousticness', 'valence', 'lc
20 X_numeric = df[numeric_features]
21
22 # Scale the numerical features
23 scaler = StandardScaler()
24 X_numeric_scaled = scaler.fit_transform(X_numeric)
25
26 # Text Feature Extraction (TF-IDF + SVD)

```

```


27 # Preprocess lyrics: lowercase and keep alphabetic tokens only
28 df['processed_lyrics'] = df['lyrics'].apply(
29     lambda x: ' '.join([word.lower() for word in str(x).split() if word.isal
30 ])
31 )
32 # TF-IDF Vectorization
33 vectorizer = TfidfVectorizer(max_features=5000, stop_words='english')
34 X_text = vectorizer.fit_transform(df['processed_lyrics'])
35
36 # Reduce Dimensionality using SVD
37 svd = TruncatedSVD(n_components=100, random_state=42)
38 X_text_reduced = svd.fit_transform(X_text)
39
40 # Step 4: Combine Numerical and Text Features
41 X_combined = np.hstack((X_numeric_scaled, X_text_reduced))
42
43 # Step 5: Build the Recommendation System
44 knn = NearestNeighbors(n_neighbors=11, metric='cosine') # 11 because input
45 knn.fit(X_combined)
46
47 # Function to recommend songs
48 def recommend_songs(song_name, df, knn_model, combined_features):
49     # Check if song exists
50     if song_name not in df['track_name'].values:
51         print(f"Error: The song '{song_name}' does not exist in the dataset.")
52         return None
53
54     # Find the index of the input song
55     song_index = df[df['track_name'] == song_name].index[0]
56
57     # Get the feature vector for the input song
58     song_features = combined_features[song_index].reshape(1, -1)
59
60     # Find the nearest neighbors
61     distances, indices = knn_model.kneighbors(song_features)
62
63     # Get the indices of the top 10 most similar songs (excluding the input
64     similar_song_indices = indices.flatten()[1:]
65
66     # Return the top 10 most similar songs
67     return df.iloc[similar_song_indices]
68
69 # Step 6: Test the System
70 # Input a song name and get recommendations
71 input_song = "i believe" # Replace with any song name from the dataset
72 recommendations = recommend_songs(input_song, df, knn, X_combined)
73
74 # Display recommendations

```

```

75 if recommendations is not None:
76     print(recommendations[['track_name', 'artist_name', 'danceability', 'ene
77

```



	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
17595	when you're smiling (the whole world smiles wi...	nat king cole
94	september in the rain	the platters
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
17595	0.392397	0.340320	0.860442	0.342539
94	0.286256	0.235211	0.969879	0.241550
152	0.375068	0.287265	0.894578	0.389942

```

1 # Step 1: Check if the input song exists in the dataset
2 input_song = "i believe" # Replace with any song name from the dataset
3 # Check if the song exists in the dataset
4 if input_song not in df['track_name'].values:
5     print(f"Error: The song '{input_song}' does not exist in the dataset.")
6 else:
7     # Step 2: Find the index of the input song
8     song_index = df[df['track_name'] == input_song].index[0]
9
10    # Step 3: Get the feature vector for the input song
11    song_features = X_combined[song_index].reshape(1, -1)
12
13    # Step 4: Find the nearest neighbors
14    distances, indices = knn.kneighbors(song_features)
15
16    # Step 5: Get the indices of the top 10 most similar songs (excluding the
17    similar_song_indices = indices.flatten()[1:]
18
19    # Step 6: Return the top 10 most similar songs
20    recommendations = df.iloc[similar_song_indices]
21    print(recommendations[['track_name', 'artist_name', 'danceability', 'energ

```

	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
17595	when you're smiling (the whole world smiles wi...	nat king cole
94	september in the rain	the platters
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
17595	0.392397	0.340320	0.860442	0.342539
94	0.286256	0.235211	0.969879	0.241550
152	0.375068	0.287265	0.894578	0.389942

✓ Recomender system, Using NLP

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.decomposition import TruncatedSVD
5 from sklearn.preprocessing import StandardScaler
6 from textblob import TextBlob
7 import re
8 import string
9
10 # Step 1: Load the dataset
11 file_path = "/content/tcc_ceds_music.csv" # Path to the uploaded file in Col
12 df = pd.read_csv(file_path)
13 df.drop(columns=['Unnamed: 0'], inplace=True)
14
15 # Step 2: Clean the lyrics
16 def clean_lyrics(text):
17     text = str(text).lower()
18     text = re.sub(f"[{re.escape(string.punctuation)}]", "", text)
19     text = re.sub(r'\d+', '', text)
20     text = re.sub(r'\s+', ' ', text)
21     return text.strip()
22
23 df['clean_lyrics'] = df['lyrics'].apply(clean_lyrics)
24
25 # Step 3: Sentiment analysis
26 df['sentiment'] = df['clean_lyrics'].apply(lambda x: TextBlob(x).sentiment.p
27
28 # Step 4: TF-IDF + SVD
29 vectorizer = TfidfVectorizer(max_features=5000, stop_words='english')
30 X_text = vectorizer.fit_transform(df['clean_lyrics'])
31
32 svd = TruncatedSVD(n_components=100, random_state=42)
33 X_text_reduced = svd.fit_transform(X_text)
34
35 # Step 5: Scale numeric features (including sentiment)
36 numeric_features = ['danceability', 'energy', 'acousticness', 'valence',
37                     'loudness', 'instrumentalness', 'age', 'sentiment']
38 X_numeric = df[numeric_features]
39
40 scaler = StandardScaler()
41 X_numeric_scaled = scaler.fit_transform(X_numeric)
42
43 # Step 6: Combine all features
44 X_combined = np.hstack((X_numeric_scaled, X_text_reduced))

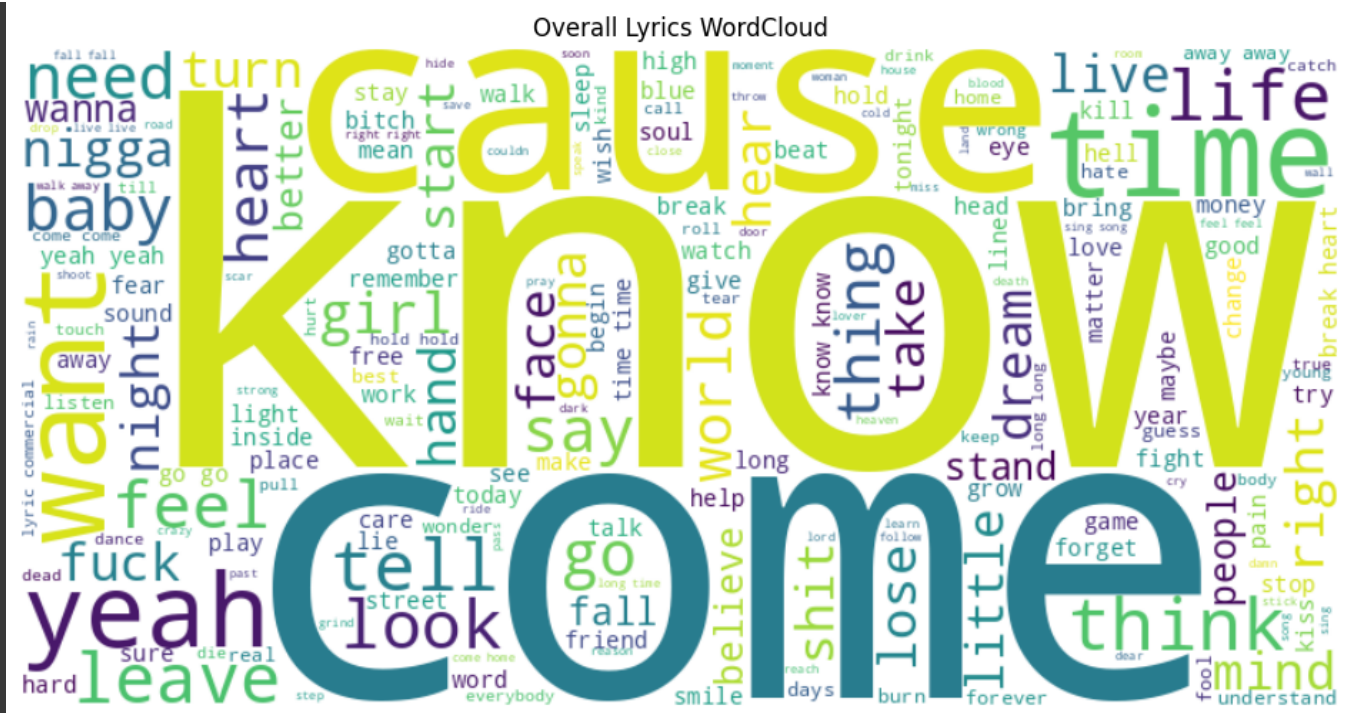
```



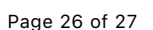
```

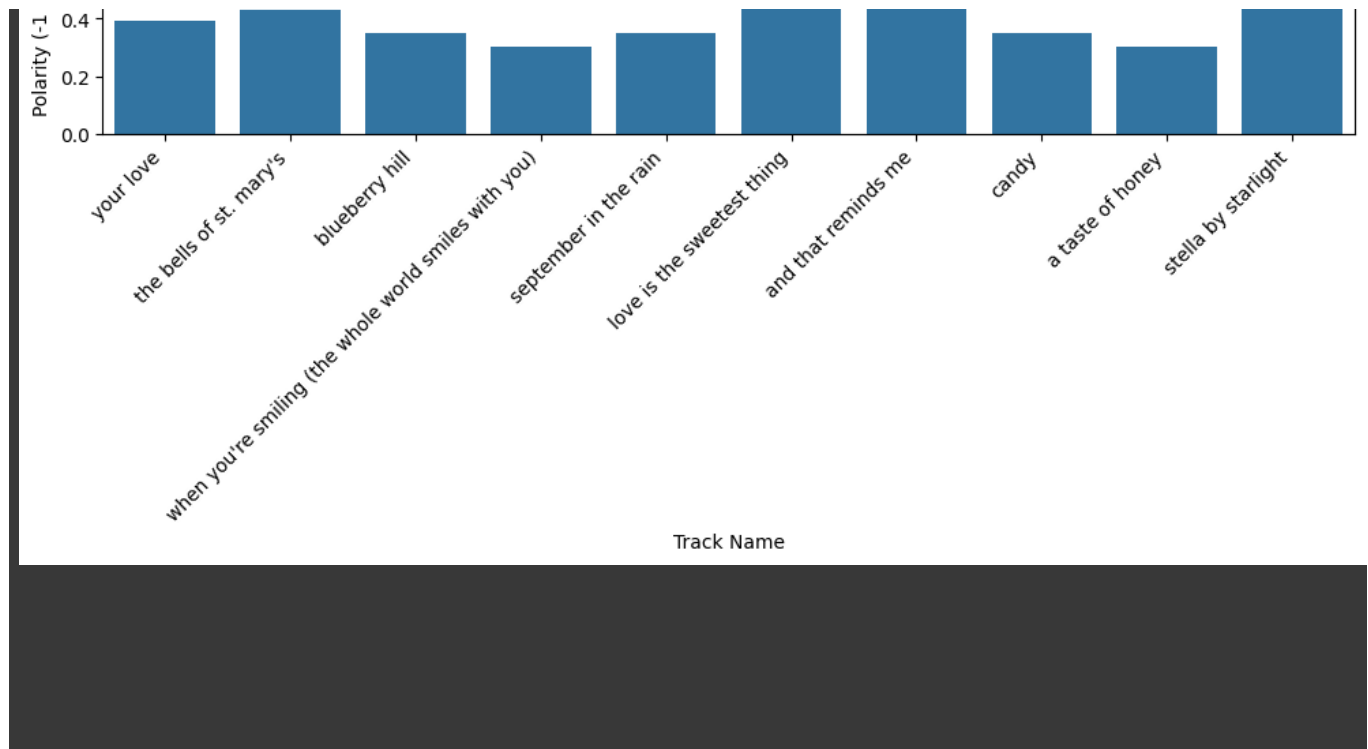
45
46 # Step 7: Build the KNN model
47 knn = NearestNeighbors(n_neighbors=11, metric='cosine')
48 knn.fit(X_combined)
49
50 # Optional: Show WordCloud for general lyrics
51 all_lyrics = ' '.join(df['clean_lyrics'])
52 wordcloud = WordCloud(width=800, height=400, background_color='white').genera
53
54 plt.figure(figsize=(12, 6))
55 plt.imshow(wordcloud, interpolation='bilinear')
56 plt.axis('off')
57 plt.title('Overall Lyrics WordCloud')
58 plt.show()
59
60 # Step 8: Song recommendation function
61 def recommend_songs(song_name, df, knn_model, feature_matrix):
62     if song_name not in df['track_name'].values:
63         print(f"'{song_name}' not found in dataset.")
64         return None
65
66     song_index = df[df['track_name'] == song_name].index[0]
67     song_vector = feature_matrix[song_index].reshape(1, -1)
68
69     distances, indices = knn_model.kneighbors(song_vector)
70     similar_indices = indices.flatten()[1:] # exclude the input song itself
71
72     recommendations = df.iloc[similar_indices].copy()
73     return recommendations
74
75 # Step 9: Test the recommender
76 input_song = "i believe" # Replace with any song from your dataset
77 recommendations = recommend_songs(input_song, df, knn, X_combined)
78
79 if recommendations is not None:
80     print("Top 10 similar songs to:", input_song)
81     display(recommendations[['track_name', 'artist_name', 'danceability', 'en
82
83 # Step 10: Visualize sentiment of recommended songs
84 plt.figure(figsize=(10, 5))
85 sns.barplot(x='track_name', y='sentiment', data=recommendations)
86 plt.xticks(rotation=45, ha='right')
87 plt.title('Sentiment Scores of Recommended Songs')
88 plt.ylabel('Polarity (-1 to 1)')
89 plt.xlabel('Track Name')
90 plt.tight_layout()
91 plt.show()

```



	track_name	artist_name	danceability	energy	valence	sentiment
12900	your love	carla thomas	0.384815	0.310289	0.317807	0.389603
47	the bells of st. mary's	the drifters	0.281924	0.205180	0.231245	0.430000
542	blueberry hill	the lettermen	0.323080	0.366347	0.331204	0.350000
17595	when you're smiling (the whole world smiles wi...	nat king cole	0.392397	0.340320	0.342539	0.300000
94	september in the rain	the platters	0.286256	0.235211	0.241550	0.350000
17121	love is the sweetest thing	ray noble	0.284090	0.309288	0.322960	0.637500
17388	and that reminds me	della reese	0.358822	0.241218	0.309563	0.500000
12742	candy	big maybelle	0.365320	0.240217	0.341509	0.350000
17658	a taste of honey	tony bennett	0.356655	0.237213	0.348722	0.300000
17322	stella by starlight	miles davis	0.425972	0.182157	0.234336	0.491667





Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.