

```

1 import pandas as pd
2 data = pd.read_csv('/content/all_songs_data.csv')
3 print(data.head())
4

```

```

0          Album \
1      Battle of New Orleans
2          That's All
3      "Mr Personality's" 15 Big Hits
4      The Greatest Hits Of Frankie Avalon
5      Paul Anka Sings His Big 15

0          Album URL          Artist \
1  https://genius.com/albums/Johnny-horton/Battle...  Johnny Horton
2  https://genius.com/albums/Bobby-darin/That-s-all  Bobby Darin
3  https://genius.com/albums/Lloyd-price/Mr-perso...  Lloyd Price
4  https://genius.com/albums/Frankie-avalon/The-g...  Frankie Avalon
5  https://genius.com/albums/Paul-anka/Paul-anka-...  Paul Anka

0      Featured Artists          Lyrics \
1      [] [Verse 1] In 1814 we took a little trip Along ...
2      [] Oh the shark, babe Has such teeth, dear And he...
3      [] Over and over I tried to prove my love to you ...
4      [] Hey, Venus! Oh, Venus! Venus, if you will Ple...
5      [] I'm just a lonely boy Lonely and blue I'm all ...

0      Media Rank Release Date \
1  [{'native_uri': 'spotify:track:0dwpdcQkeZqpuoA...  1  1959-04-01
2  [{'native_uri': 'spotify:track:3E5ndy0f06vFDEI...  2      NaN
3  [{'provider': 'youtube', 'start': 0, 'type': '...  3      NaN
4      [] 4      NaN
5      [] 5      NaN

0      Song Title \
1  The Battle Of New Orleans
2      Mack The Knife
3      Personality
4      Venus
5      Lonely Boy

0      Song URL \
1  https://genius.com/Johnny-horton-the-battle-of...
2  https://genius.com/Bobby-darin-mack-the-knife-...
3  https://genius.com/Lloyd-price-personality-lyrics
4  https://genius.com/Frankie-avalon-venus-lyrics
5  https://genius.com/Paul-anka-lonely-boy-lyrics

0      Writers          Year
1  [{'api_path': '/artists/561913', 'header_image...  1959.0
2  [{'api_path': '/artists/218851', 'header_image...  1959.0
3  [{'api_path': '/artists/355804', 'header_image...  1959.0
4  [{'api_path': '/artists/1113175', 'header_imag...  1959.0
5      [] 1959.0

```

```

1 # Descriptive statistics for numerical columns
2 print(data.describe())
3
4 # Descriptive statistics for categorical columns
5 print(data.describe(include='object'))
6

```

```

0      Rank          Year
1  count  6500.000000  6500.000000
2  mean    50.500000  1991.000000
3  std    28.868291  18.763106
4  min     1.000000  1959.000000

```

25%	25.750000	1975.000000
50%	50.500000	1991.000000
75%	75.250000	2007.000000
max	100.000000	2023.000000

	Album	Album URL \
count	6036	6036
unique	4202	4285
top	Greatest Hits	<a href="https://genius.com/albums/Morgan-wallen/One-th...">https://genius.com/albums/Morgan-wallen/One-th...</a>
freq	21	9

	Artist	Featured Artists	Lyrics	Media	Release Date \
count	6500	6384	6384	6384	4563
unique	3181	612	6044	5054	3233
top	Madonna	[]	[Instrumental]	[]	2022-05-06
freq	35	5492	21	1043	10

	Song Title	Song URL	Writers
count	6500	6384	6384
unique	5798	6065	4184
top	Stay	<a href="https://genius.com/Billboard-hot-rap-songs-cha...">https://genius.com/Billboard-hot-rap-songs-cha...</a>	[]
freq	7	5	971

```

1 rank_stats = data['Rank'].describe()
2 rank_frequency = data['Rank'].value_counts().sort_index()
3 top_artists = data['Artist'].value_counts().head(10)
4
5 print("Descriptive Statistics for Song Ranks:")
6 print(rank_stats)
7 print("\nFrequency of Songs by Rank:")
8 print(rank_frequency)
9 print("\nTop 10 Artists by Number of Songs:")
10 print(top_artists)
11

```

➡ Descriptive Statistics for Song Ranks:

```

count    6500.000000
mean      50.500000
std       28.868291
min        1.000000
25%       25.750000
50%       50.500000
75%       75.250000
max       100.000000
Name: Rank, dtype: float64

```

Frequency of Songs by Rank:

```

Rank
1      65
2      65
3      65
4      65
5      65
..
96     65
97     65
98     65
99     65
100    65
Name: count, Length: 100, dtype: int64

```

Top 10 Artists by Number of Songs:

```

Artist
Madonna      35
Mariah Carey  29
Taylor Swift  28
Beatles       27
Elton John    26

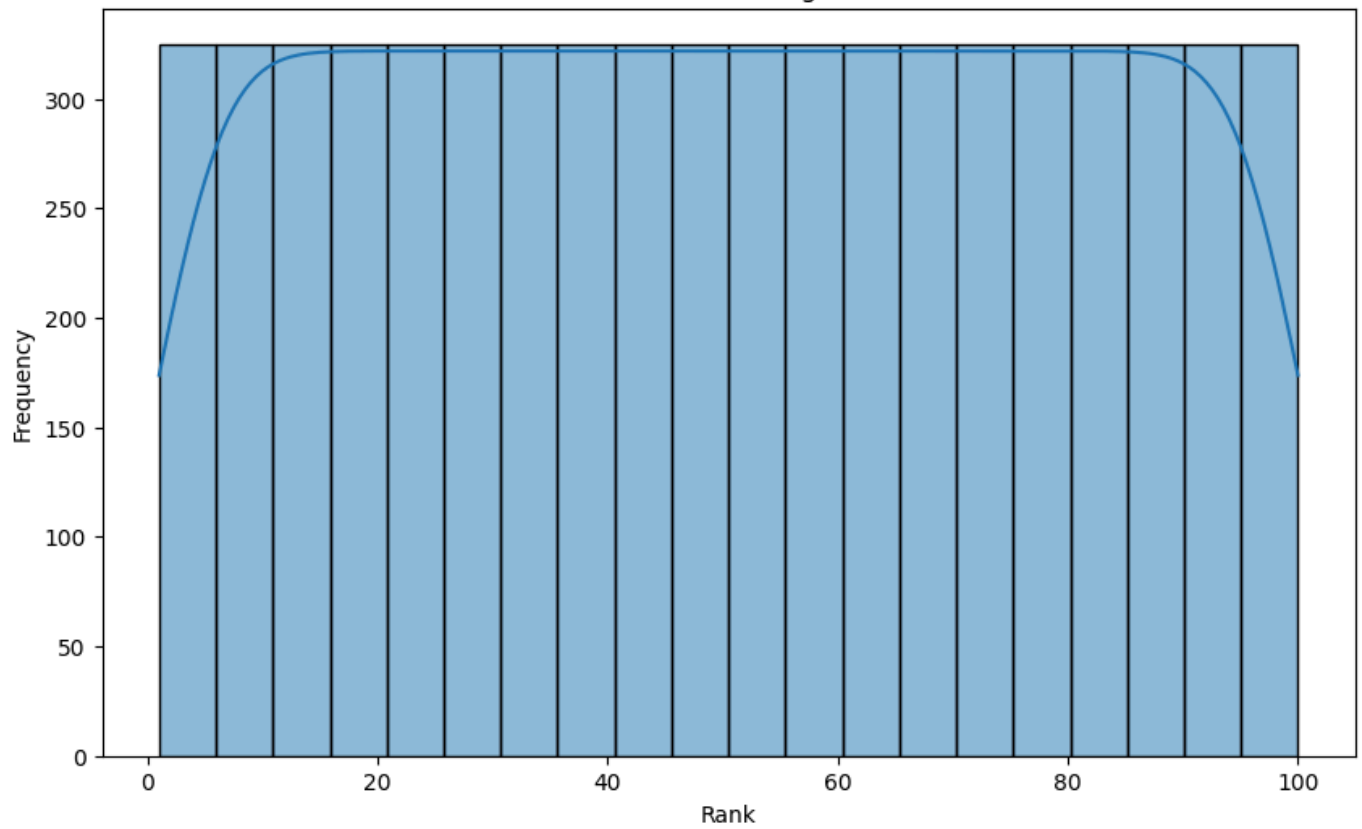
```

```
Stevie Wonder      22
Rihanna            22
Michael Jackson    22
Janet Jackson      21
Whitney Houston    20
Name: count, dtype: int64
```

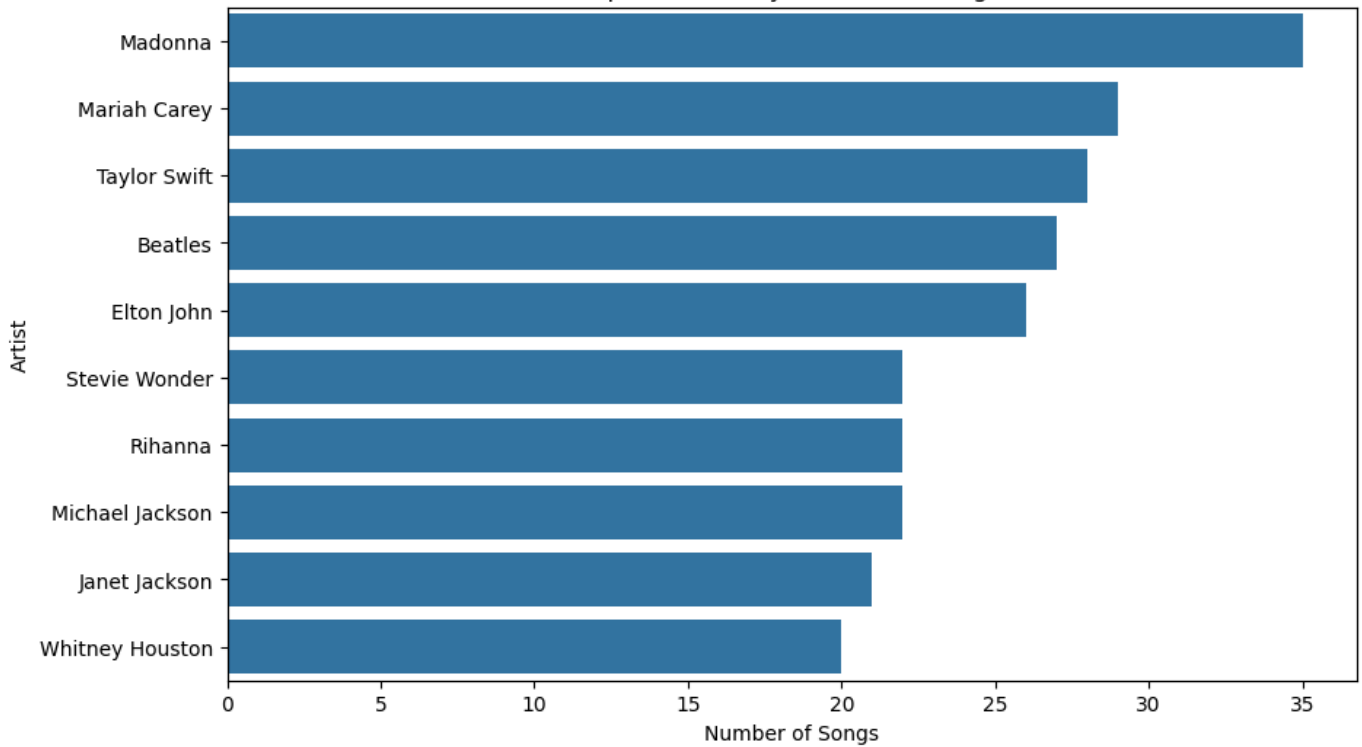
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 plt.figure(figsize=(10, 6))
5 sns.histplot(data['Rank'], bins=20, kde=True)
6 plt.title('Distribution of Song Ranks')
7 plt.xlabel('Rank')
8 plt.ylabel('Frequency')
9 plt.show()
10
11 top_artists = data['Artist'].value_counts().head(10)
12 plt.figure(figsize=(10, 6))
13 sns.barplot(x=top_artists.values, y=top_artists.index)
14 plt.title('Top 10 Artists by Number of Songs')
15 plt.xlabel('Number of Songs')
16 plt.ylabel('Artist')
17 plt.show()
18
```



Distribution of Song Ranks



Top 10 Artists by Number of Songs



```

1 # Converting 'Release Date' to datetime
2 data['Release Date'] = pd.to_datetime(data['Release Date'])
3
4 data['Year'] = data['Release Date'].dt.year
5
6
7 rank_year_corr = data[['Year', 'Rank']].groupby('Year').mean().reset_index()
8 rank_year_corr['Rank'] = rank_year_corr['Rank'].round(2)
9
10 correlation = data['Year'].corr(data['Rank'])
11
12 top_artists = data['Artist'].value_counts().head(10).index
13 rank_distribution_top_artists = data[data['Artist'].isin(top_artists)].groupby('Artist')['Rank'].desc
14
15 print("Average Rank by Year:")
16 print(rank_year_corr)
17 print("\nCorrelation between Year and Rank:")
18 print(correlation)
19 print("\nRank Distribution by Top 10 Artists:")
20 print(rank_distribution_top_artists)
21

```

➔ Average Rank by Year:

	Year	Rank
0	1877.0	26.00
1	1922.0	70.00
2	1955.0	52.00
3	1957.0	77.00
4	1958.0	72.00
...	...	...
66	2020.0	52.36
67	2021.0	45.94
68	2022.0	48.67
69	2023.0	53.90
70	2024.0	82.00

[71 rows x 2 columns]

Correlation between Year and Rank:  
0.04014067434061861

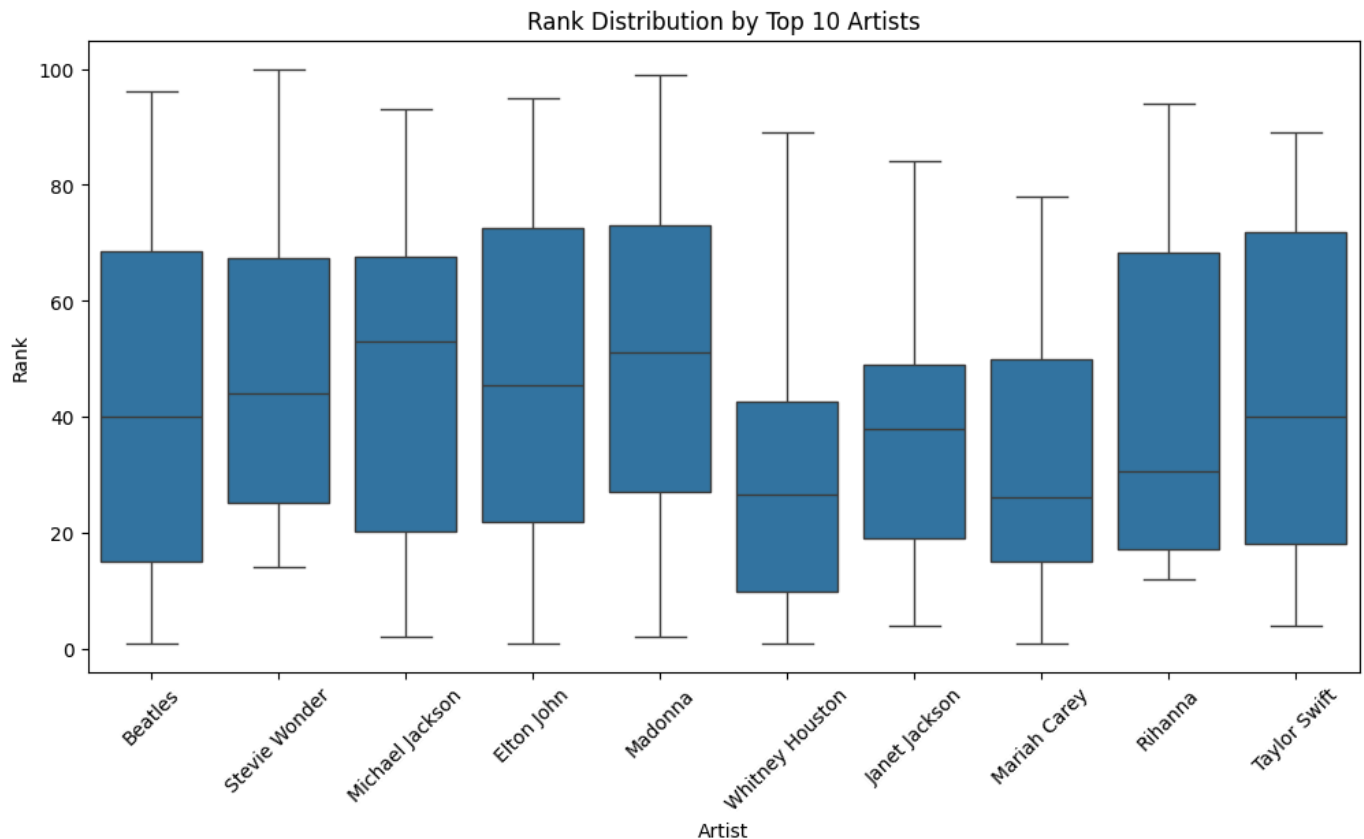
Rank Distribution by Top 10 Artists:

	count	mean	std	min	25%	50%	75%	max
Artist								
Beatles	27.0	42.777778	31.624398	1.0	15.00	40.0	68.50	96.0
Elton John	26.0	47.038462	29.200659	1.0	21.75	45.5	72.50	95.0
Janet Jackson	21.0	38.619048	21.973794	4.0	19.00	38.0	49.00	84.0
Madonna	35.0	49.285714	28.678308	2.0	27.00	51.0	73.00	99.0
Mariah Carey	29.0	33.724138	24.335301	1.0	15.00	26.0	50.00	78.0
Michael Jackson	22.0	47.409091	30.316833	2.0	20.25	53.0	67.50	93.0
Rihanna	22.0	42.090909	29.012163	12.0	17.25	30.5	68.25	94.0
Stevie Wonder	22.0	48.318182	26.141349	14.0	25.25	44.0	67.25	100.0
Taylor Swift	28.0	43.642857	27.214123	4.0	18.00	40.0	71.75	89.0
Whitney Houston	20.0	29.200000	23.625588	1.0	9.75	26.5	42.50	89.0

```

1 top_artists = data['Artist'].value_counts().head(10).index
2 plt.figure(figsize=(12, 6))
3 sns.boxplot(x='Artist', y='Rank', data=data[data['Artist'].isin(top_artists)])
4 plt.title('Rank Distribution by Top 10 Artists')
5 plt.xlabel('Artist')
6 plt.ylabel('Rank')
7 plt.xticks(rotation=45)
8 plt.show()
9

```



At this stage, I am moving on to exploring the categorical data and cleaning/preprocessing the text

```

1 # Checking for missing values
2 print(f"Missing Lyrics: {data['Lyrics'].isnull().sum()}")
3 empty_lyrics = data[data['Lyrics'].apply(lambda x: isinstance(x, str) and len(x) == 0)]
4 print(f"Empty Lyrics: {empty_lyrics.shape[0]}")
5 invalid_lyrics = data[data['Lyrics'].apply(lambda x: not isinstance(x, str))]
6 print(f"Invalid Lyrics (Non-string values): {invalid_lyrics.shape[0]}")
7

```



```

Missing Lyrics: 116
Empty Lyrics: 0
Invalid Lyrics (Non-string values): 116

```

```

1 # Removing missing values
2 data_cleaned = data.dropna(subset=['Lyrics'])
3
4 data_cleaned = data_cleaned[data_cleaned['Lyrics'].apply(lambda x: isinstance(x, str))]
5
6 print(f"Data after cleaning: {data_cleaned.shape[0]} rows")
7 print(data_cleaned['Lyrics'].isnull().sum()) # Check if there are any missing lyrics
8

```



```

Data after cleaning: 6384 rows
0

```

```

1 !pip install spacy
2 !python -m spacy download en_core_web_sm

```

 Show hidden output


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

Next, I am loading the necessary tools to preprocess the lyrical data

```

1 import spacy
2 import re
3 from nltk.corpus import stopwords
4 from nltk.stem import WordNetLemmatizer
5
6 # Implementing custom stopwords
7 custom_stopwords = set(stopwords.words('english')).union({'chorus', 'verse', 'bridge', 'hook', 'intro
8 nlp = spacy.load("en_core_web_sm")
9 lemmatizer = WordNetLemmatizer()
10
11 #in the first round, i noticed contractions were not removed properly so I used a dictionary to ensure
12 def expand_contractions(text):
13     contractions_dict = {
14         "don't": "do not", "can't": "cannot", "won't": "will not", "didn't": "did not",
15         "isn't": "is not", "aren't": "are not", "wasn't": "was not", "weren't": "were not",
16         "hasn't": "has not", "haven't": "have not", "hadn't": "had not", "doesn't": "does not",
17         "didn't": "did not", "couldn't": "could not", "shouldn't": "should not", "mightn't": "might not",
18         "mustn't": "must not", "let's": "let us", "i'm": "i am", "you're": "you are", "he's": "he is",
19         "she's": "she is", "it's": "it is", "we're": "we are", "they're": "they are", "that's": "that",
20         "what's": "what is", "who's": "who is", "where's": "where is", "how's": "how is"
21     }
22     for word, expansion in contractions_dict.items():
23         text = re.sub(r'\b' + word + r'\b', expansion, text)
24     return text
25
26 def preprocess_text(text):
27     text = expand_contractions(text)
28
29     # Converting to lowercase
30     text = text.lower()
31
32     # Removing non-alphanumeric characters (keeping spaces and words)
33     text = re.sub(r'[^a-zA-Z\s]', '', text)
34
35     # Tokenizing the text using spaCy
36     doc = nlp(text)
37
38     words = [token.text for token in doc if token.text not in custom_stopwords and not token.is_punct]
39
40     # Lemmatizing words
41     words = [lemmatizer.lemmatize(word) for word in words]
42
43     # Rejoining words
44     cleaned_text = ' '.join(words)
45
46     return cleaned_text
47 data_cleaned['cleaned_lyrics'] = data_cleaned['Lyrics'].apply(preprocess_text)
48
49 print(data_cleaned[['Song Title', 'Lyrics', 'cleaned_lyrics']].head())
50

```



	Song Title	Lyrics	cleaned_lyrics
0	The Battle Of New Orleans		
1	Mack The Knife		
2	Personality		

```

3          Venus
4      Lonely Boy

```

```

                                Lyrics \
0 [Verse 1] In 1814 we took a little trip Along ...
1 Oh the shark, babe Has such teeth, dear And he...
2 Over and over I tried to prove my love to you ...
3 Hey, Venus! Oh, Venus! Venus, if you will Ple...
4 I'm just a lonely boy Lonely and blue I'm all ...

```

```

                                cleaned_lyrics
0      took little trip along colonel jackson mig...
1 oh shark babe teeth dear show pearly white jac...
2 tried prove love friend say fool ill fool ...
3 hey venus oh venus venus please send little ...
4 lonely boy lonely blue alone nothin got ever...

```

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

Now that the data is preprocessed, we move on to vectorizing the text data to fit the LDA model

```

1 import pandas as pd
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.decomposition import LatentDirichletAllocation
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 texts = data_cleaned['cleaned_lyrics'].dropna()
8
9 vectorizer = CountVectorizer(stop_words='english', max_features=5000)
10
11 X = vectorizer.fit_transform(texts)
12
13 lda = LatentDirichletAllocation(n_components=5, random_state=42)
14 lda.fit(X)
15
16 # Retrieving the words for each topic
17 words = vectorizer.get_feature_names_out()
18
19 #display/interpret top words/topic
20 def print_top_words(model, feature_names, n_top_words=10):
21     topic_labels = {}
22     for topic_idx, topic in enumerate(model.components_):
23         top_words_idx = topic.argsort()[::-n_top_words - 1:-1]
24         top_words = [feature_names[i] for i in top_words_idx]
25
26
27         print(f"Topic {topic_idx + 1}:")
28         print(" ".join(top_words))
29
30     #Labels/category interpretations were created by human coder after first iteration retrieved t
31     if topic_idx == 0:
32         topic_labels[topic_idx] = "Reflections and Life Experiences"
33     elif topic_idx == 1:
34         topic_labels[topic_idx] = "Communication and Family Dynamics"
35     elif topic_idx == 2:
36         topic_labels[topic_idx] = "Love and Emotional Longing"
37     elif topic_idx == 3:
38         topic_labels[topic_idx] = "Street Culture and Raw Emotions"
39     elif topic_idx == 4:
40         topic_labels[topic_idx] = "Romantic Desire and Affection"
41
42     return topic_labels

```



```
43
44 topic_labels = print_top_words(lda, words)
45
46 # Adding the dominant topic for each song
47 topic_probabilities = lda.transform(X)
48 dominant_topic = topic_probabilities.argmax(axis=1)
49 data_cleaned['dominant_topic'] = dominant_topic
50
51 # Mapping topic labels to the songs
52 data_cleaned['topic_label'] = data_cleaned['dominant_topic'].map(topic_labels)
53
54 print(data_cleaned[['Song Title', 'dominant_topic', 'topic_label']].head())
55
56 # Visualizing topic distribution across songs
57 plt.figure(figsize=(10, 6))
58 sns.countplot(x='topic_label', data=data_cleaned, palette='viridis')
59 plt.title('Distribution of Topics Across Songs')
60 plt.xlabel('Topic Label')
61 plt.ylabel('Number of Songs')
62 plt.xticks(rotation=45)
63 plt.show()
64
```

```

Topic 1:
said like man day time old know hand make say
Topic 2:
say tell mr know mam like boy come dad want
Topic 3:
love time night heart day away ill know life let
Topic 4:
like nigga got la bitch shit know love ai fuck
Topic 5:
baby yeah love oh know got like na want girl

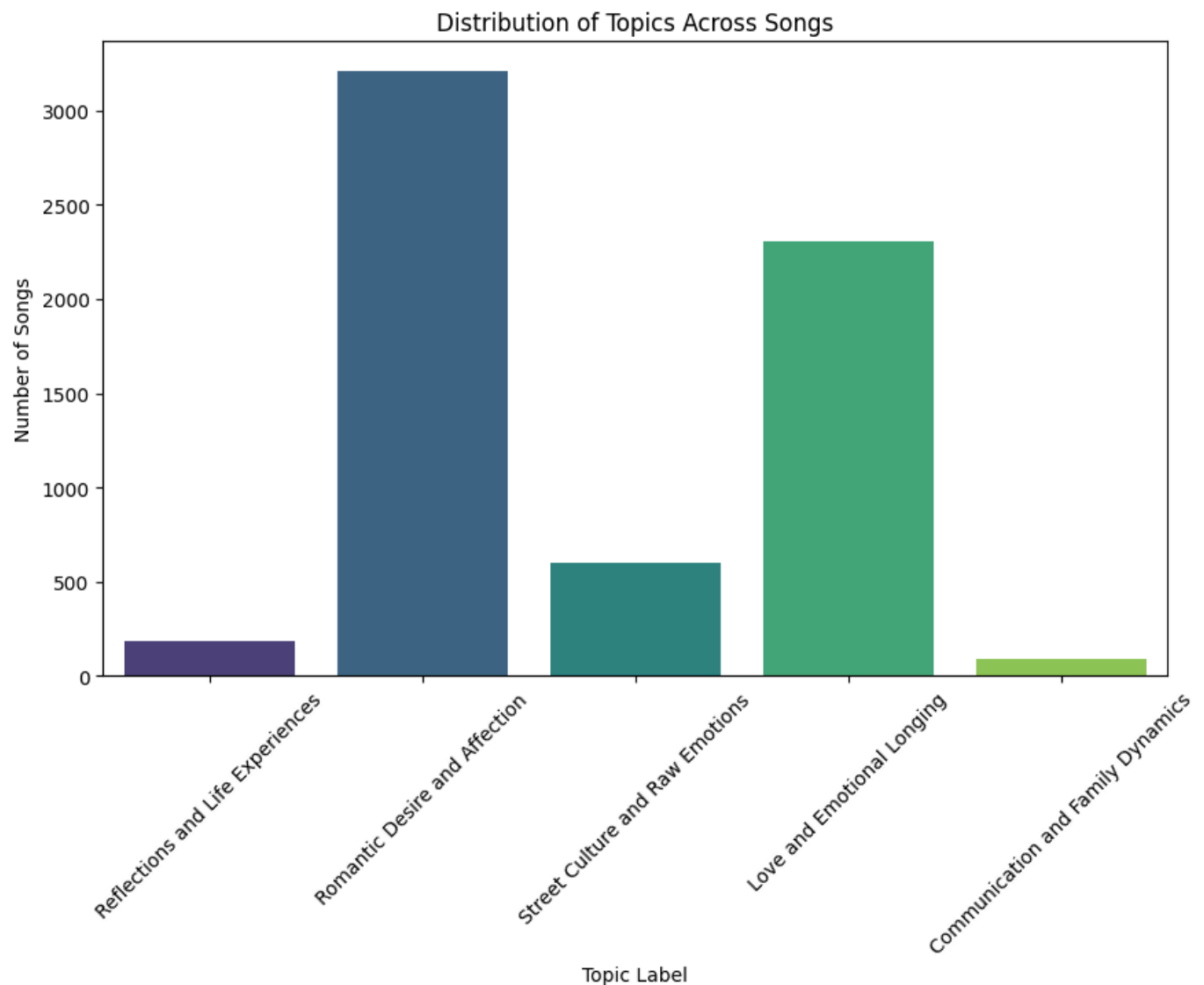
```

	Song Title	dominant_topic	topic_label
0	The Battle Of New Orleans	0	Reflections and Life Experiences
1	Mack The Knife	4	Romantic Desire and Affection
2	Personality	3	Street Culture and Raw Emotions
3	Venus	0	Reflections and Life Experiences
4	Lonely Boy	2	Love and Emotional Longing

```
<ipython-input-40-152ba0ce7419>:58: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `

```
sns.countplot(x='topic_label', data=data_cleaned, palette='viridis')
```



Testing the coherence score to confirm quality of topics

```

1 from gensim.models import CoherenceModel
2 from gensim.corpora import Dictionary
3
4 texts = [text.split() for text in data_cleaned['cleaned_lyrics']]
5
6 dictionary = Dictionary(texts)
7 corpus = [dictionary.doc2bow(text) for text in texts]
8
9 from gensim.models import LdaModel
10 lda_gensim = LdaModel(corpus, num_topics=5, id2word=dictionary)
11
12 coherence_model_lda = CoherenceModel(model=lda_gensim, texts=texts, dictionary=dictionary, coherence=
13 coherence_score = coherence_model_lda.get_coherence()
14 print(f"Coherence Score: {coherence_score}")
15

```

⚠ WARNING:gensim.models.ldamodel:too few updates, training might not converge; consider increasing the Coherence Score: 0.5125489353605521

0.51 coherence score represents moderate coherence

Comparing ML models for classifying the dominant topic of song lyrics

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.svm import SVC
5 from sklearn.metrics import classification_report, accuracy_score
6 from sklearn.feature_extraction.text import CountVectorizer
7
8 vectorizer = CountVectorizer(stop_words='english', max_features=5000)
9 X = vectorizer.fit_transform(data_cleaned['cleaned_lyrics']).dropna()
10 y = data_cleaned['dominant_topic']
11
12 #Split data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
14
15 # Logistic Regression
16 logreg = LogisticRegression(max_iter=1000)
17 logreg.fit(X_train, y_train)
18 logreg_predictions = logreg.predict(X_test)
19
20 # Random Forest Classifier
21 rf = RandomForestClassifier(random_state=42)
22 rf.fit(X_train, y_train)
23 rf_predictions = rf.predict(X_test)
24
25 # Support Vector Machine (SVM)
26 svm = SVC(kernel='linear', random_state=42)
27 svm.fit(X_train, y_train)
28 svm_predictions = svm.predict(X_test)
29
30 # Evaluation
31 print("Logistic Regression Performance:")
32 print(classification_report(y_test, logreg_predictions))
33 print("Accuracy:", accuracy_score(y_test, logreg_predictions))
34
35 print("\nRandom Forest Performance:")
36 print(classification_report(y_test, rf_predictions))
37 print("Accuracy:", accuracy_score(y_test, rf_predictions))
38

```

```

39 print("\nSVM Performance:")
40 print(classification_report(y_test, svm_predictions))
41 print("Accuracy:", accuracy_score(y_test, svm_predictions))
42

```

Logistic Regression Performance:

	precision	recall	f1-score	support
0	0.68	0.41	0.51	37
1	0.45	0.26	0.33	19
2	0.83	0.89	0.86	460
3	0.90	0.70	0.79	125
4	0.89	0.91	0.90	636
accuracy			0.86	1277
macro avg	0.75	0.64	0.68	1277
weighted avg	0.86	0.86	0.86	1277

Accuracy: 0.860610806577917

Random Forest Performance:

	precision	recall	f1-score	support
0	0.46	0.16	0.24	37
1	0.67	0.11	0.18	19
2	0.80	0.71	0.75	460
3	0.94	0.36	0.52	125
4	0.74	0.94	0.83	636
accuracy			0.76	1277
macro avg	0.72	0.45	0.50	1277
weighted avg	0.77	0.76	0.74	1277

Accuracy: 0.7635082223962412

SVM Performance:

	precision	recall	f1-score	support
0	0.31	0.38	0.34	37
1	0.26	0.26	0.26	19
2	0.81	0.84	0.82	460
3	0.83	0.70	0.76	125
4	0.89	0.88	0.89	636
accuracy			0.82	1277
macro avg	0.62	0.61	0.61	1277
weighted avg	0.83	0.82	0.82	1277

Accuracy: 0.8238057948316366

1 #I have selected Logistic Regression to continue topic predictions in my analysis

```

1 vectorizer = CountVectorizer(stop_words='english', max_features=5000)
2 X = vectorizer.fit_transform(data_cleaned['cleaned_lyrics']).dropna()
3 y = data_cleaned['dominant_topic']
4
5 # Retraining Logistic Regression model
6 logreg = LogisticRegression(max_iter=1000)
7 logreg.fit(X, y)
8
9 # Get predictions
10 data_cleaned['predicted_topic'] = logreg.predict(X)
11
12 # Map the predicted topic to a label
13 topic_labels = {
14     0: "Reflections and Life Experiences",

```

```

15     1: "Communication and Family Dynamics",
16     2: "Love and Emotional Longing",
17     3: "Street Culture and Raw Emotions",
18     4: "Romantic Desire and Affection"
19 }
20 data_cleaned['predicted_topic_label'] = data_cleaned['predicted_topic'].map(topic_labels)
21

```

Using the logistic regression model, continue visualizations

```


1 #grouping by decade
2 data_cleaned['Decade'] = (data_cleaned['Year'] // 10) * 10
3 decade_topic_distribution = data_cleaned.groupby(['Decade', 'predicted_topic']).size().unstack(fill_value=0)
4 decade_topic_distribution_normalized = decade_topic_distribution.div(decade_topic_distribution.sum(axis=1))
5

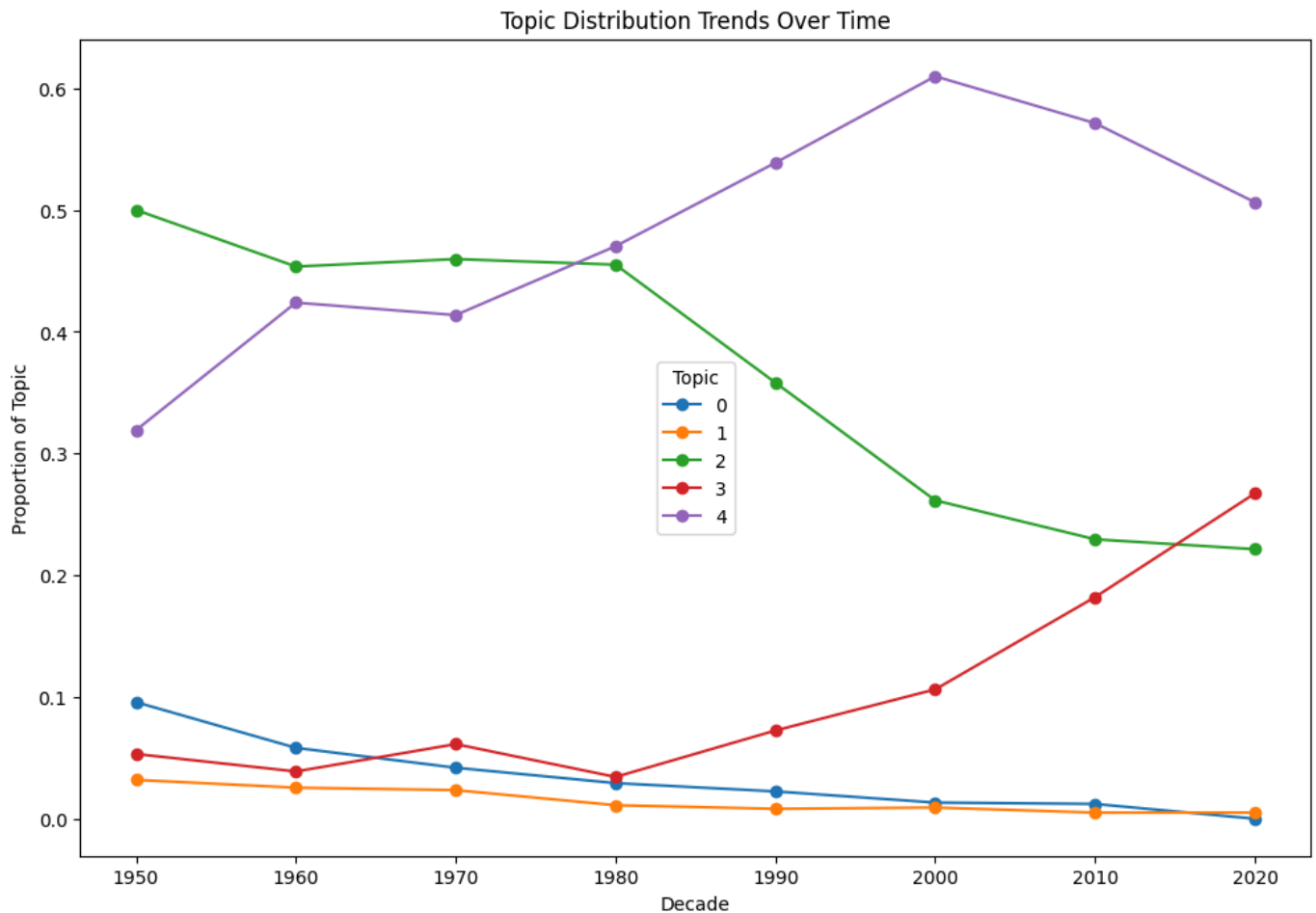
```

```

1 # Plotting topic distribution trends over time
2 plt.figure(figsize=(12, 8))
3 decade_topic_distribution_normalized.plot(kind='line', marker='o', figsize=(12, 8))
4 plt.title('Topic Distribution Trends Over Time')
5 plt.xlabel('Decade')
6 plt.ylabel('Proportion of Topic')
7 plt.legend(title='Topic')
8 plt.show()
9

```

 <Figure size 1200x800 with 0 Axes>

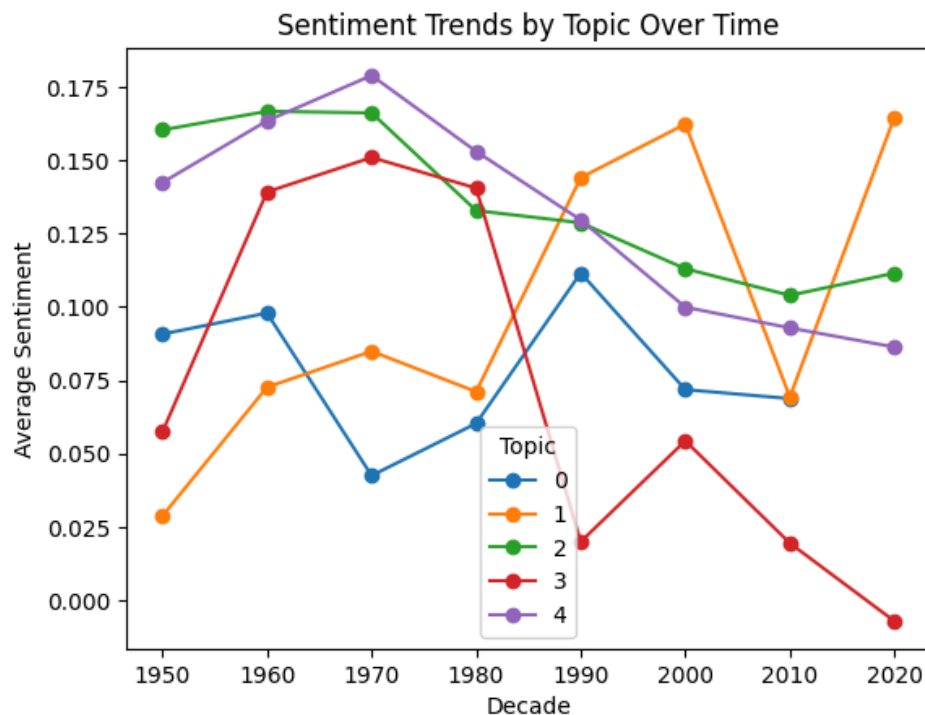


```

1 from textblob import TextBlob
2
3 #sentiment polarity
4 def calculate_sentiment(text):
5     return TextBlob(text).sentiment.polarity
6
7 data_cleaned['sentiment'] = data_cleaned['cleaned_lyrics'].apply(calculate_sentiment)
8
9 topic_sentiment_by_decade = data_cleaned.groupby(['Decade', 'dominant_topic'])['sentiment'].mean().un:
10
11 # Visualizing sentiment trends by decade
12 plt.figure(figsize=(12, 8))
13 topic_sentiment_by_decade.plot(kind='line', marker='o')
14 plt.title('Sentiment Trends by Topic Over Time')
15 plt.xlabel('Decade')
16 plt.ylabel('Average Sentiment')
17 plt.legend(title='Topic')
18 plt.show()
19

```

<Figure size 1200x800 with 0 Axes>



```

1 # Sentiment summary stats
2 sentiment_stats = data_cleaned['sentiment'].describe()
3 print("Sentiment Summary Statistics:")
4 print(sentiment_stats)
5

```

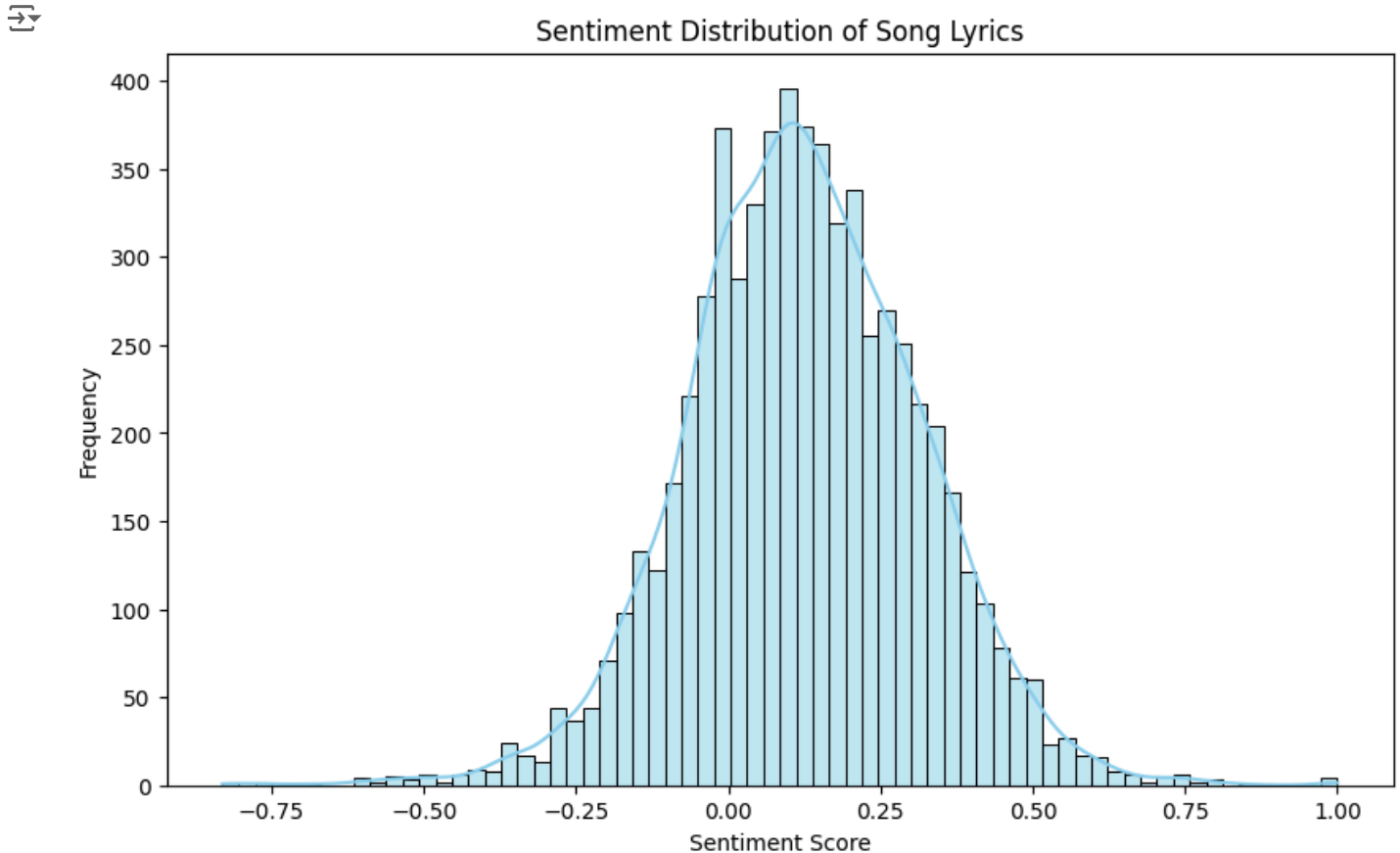
Sentiment Summary Statistics:

count	6384.000000
mean	0.123441
std	0.192460
min	-0.831108
25%	0.000000
50%	0.119649
75%	0.250969
max	1.000000
Name: sentiment, dtype: float64	

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Plot the distribution of sentiment
5 plt.figure(figsize=(10, 6))
6 sns.histplot(data_cleaned['sentiment'], kde=True, color='skyblue')
7 plt.title('Sentiment Distribution of Song Lyrics')
8 plt.xlabel('Sentiment Score')
9 plt.ylabel('Frequency')
10 plt.show()
11

```



```

1 decade_sentiment = data_cleaned.groupby('Decade')['sentiment'].mean()
2
3 # Plot sentiment over decades
4 plt.figure(figsize=(12, 6))
5 decade_sentiment.plot(kind='line', marker='o', color='green')
6 plt.title('Average Sentiment of Song Lyrics by Decade')
7 plt.xlabel('Decade')
8 plt.ylabel('Average Sentiment')
9 plt.grid(True)
10 plt.show()
11

```

