```python
import numpy as np  # NumPy: Numerical computing library for arrays and matrices.
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import re
import nltk
import string
from nltk.corpus import stopwords
from nltk.stem import LancasterStemmer
# Import necessary modules
from sklearn.preprocessing import OneHotEncoder
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import accuracy_score, classification_report
# Import the necessary module
from sklearn.feature_extraction.text import TfidfVectorizer


# Load the training data
train_path = "train_data.txt"
train_data = pd.read_csv(train_path, sep=':::', names=['Title', 'Genre', 'Description'], engine='python')


print(train_data.describe())


print(train_data.info())


print(train_data.isnull().sum())


# Load the test data
test_path = "test_data.txt"
test_data = pd.read_csv(test_path, sep=':::', names=['Id', 'Title', 'Description'], engine='python')
test_data.head()
#Count each genre value
train_data.Genre.value_counts()


# Plot count plot
plt.figure(figsize=(12,8))
counts = train_data.Genre.value_counts()
sns.barplot(x=counts.index, y=counts, color='blue')
plt.xlabel('Genre' ,fontsize=14, fontweight='bold')
plt.ylabel('Count', fontsize=14, fontweight='bold')
plt.title('Distribution of Genres', fontsize=16, fontweight='bold')
plt.xticks(rotation=90, fontsize=14, fontweight='bold');


# Plot the distribution of genres in the training data
plt.figure(figsize=(14, 7))
sns.countplot(data=train_data, y='Genre', order=train_data['Genre'].value_counts().index, palette='viridis')
plt.xlabel('Count', fontsize=14, fontweight='bold')
plt.ylabel('Genre', fontsize=14, fontweight='bold')


# Plot the distribution of genres using a bar plot
plt.figure(figsize=(14, 7))
counts = train_data['Genre'].value_counts()
sns.barplot(x=counts.index, y=counts, palette='viridis')
plt.xlabel('Genre', fontsize=14, fontweight='bold')
plt.ylabel('Count', fontsize=14, fontweight='bold')
plt.title('Distribution of Genres', fontsize=16, fontweight='bold')
plt.xticks(rotation=90, fontsize=14, fontweight='bold')
plt.show()
train_data.info()
#Finda any null value
train_data.isnull().sum()
# Initialize the stemmer and stop words
stemmer = LancasterStemmer()
stop_words = set(stopwords.words('english'))


# Define the clean_text function
def clean_text(text):
    text = text.lower()  # Lowercase all characters
    text = re.sub(r'@\S+', '', text)  # Remove Twitter handles
    text = re.sub(r'http\S+', '', text)  # Remove URLs
    text = re.sub(r'pic.\S+', '', text)
    text = re.sub(r"[^a-zA-Z+']", ' ', text)  # Keep only characters
    text = re.sub(r'\s+[a-zA-Z]\s+', ' ', text + ' ')  # Keep words with length > 1 only
    text = "".join([i for i in text if i not in string.punctuation])
    words = nltk.word_tokenize(text)
    stopwords = nltk.corpus.stopwords.words('english')  # Remove stopwords
    text = " ".join([i for i in words if i not in stopwords and len(i) > 2])
    text = re.sub("\s[\s]+", " ", text).strip()  # Remove repeated/leading/trailing spaces
    return text


# Apply the clean_text function to the 'Description' column in the training and test data
```

```python
train_data['Text_cleaning'] = train_data['Description'].apply(clean_text)
test_data['Text_cleaning'] = test_data['Description'].apply(clean_text)
# Droping the redundant data
print("shape before drop nulls",train_data.shape)
train_data = train_data.drop_duplicates()
print("shape after drop nulls",train_data.shape)
# Calculate the length of cleaned text
train_data['length_Text_cleaning'] = train_data['Text_cleaning'].apply(len)
# Visualize the distribution of text lengths
plt.figure(figsize=(8, 7))
sns.histplot(data=train_data, x='length_Text_cleaning', bins=20, kde=True, color='blue')
plt.xlabel('Length', fontsize=14, fontweight='bold')
plt.ylabel('Frequency', fontsize=14, fontweight='bold')
plt.title('Distribution of Lengths', fontsize=16, fontweight='bold')
plt.show()

# Set up the figure with two subplots
plt.figure(figsize=(12, 6))

# Subplot 1: Original text length distribution
plt.subplot(1, 2, 1)
original_lengths = train_data['Description'].apply(len)
plt.hist(original_lengths, bins=range(0, max(original_lengths) + 100, 100), color='blue', alpha=0.7)
plt.title('Original Text Length')
plt.xlabel('Text Length')
plt.ylabel('Frequency')

# Subplot 2: Cleaned text length distribution
plt.subplot(1, 2, 2)
cleaned_lengths = train_data['Text_cleaning'].apply(len)

plt.hist(cleaned_lengths, bins=range(0, max(cleaned_lengths) + 100, 100), color='green', alpha=0.7)
plt.title('Cleaned Text Length')
plt.xlabel('Text Length')
plt.ylabel('Frequency')

# Adjust layout and display the plots
plt.tight_layout()
plt.show()

tfidf_vectorizer = TfidfVectorizer()

# Fit and transform the training data
X_train = tfidf_vectorizer.fit_transform(train_data['Text_cleaning'])

# Transform the test data
X_test = tfidf_vectorizer.transform(test_data['Text_cleaning'])

X = X_train
y = train_data['Genre']
X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize and train a Multinomial Naive Bayes classifier
classifier = MultinomialNB()
classifier.fit(X_train, y_train)

y_pred = classifier.predict(X_val)

# Evaluate the performance of the model
accuracy = accuracy_score(y_val, y_pred)
print("Validation Accuracy:", accuracy)
print(classification_report(y_val, y_pred))

X_test_predictions = classifier.predict(X_test)
test_data['Predicted_Genre'] = X_test_predictions

test_data.to_csv('predicted_genres.csv', index=False)

# Display the 'test_data' DataFrame with predicted genres
print(test_data)
```

```
                              Title    Genre  \
count                         54214    54214
unique                        54214       27
top      Oscar et la dame rose (2009)    drama
freq                              1    13613

                                       Description
count                                        54214
unique                                       54086
top        Grammy - music award of the American academy ...
freq                                            12
<class 'pandas.core.frame.DataFrame'>
Index: 54214 entries, 1 to 54214
Data columns (total 3 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   Title        54214 non-null  object
 1   Genre        54214 non-null  object
 2   Description  54214 non-null  object
dtypes: object(3)
memory usage: 1.7+ MB
None
Title          0
Genre          0
Description    0
dtype: int64
<ipython-input-50-3206035593e8>:46: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `le

  sns.countplot(data=train_data, y='Genre', order=train_data['Genre'].value_counts().index, palette='viridis')
<ipython-input-50-3206035593e8>:53: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `le

  sns.barplot(x=counts.index, y=counts, palette='viridis')
```
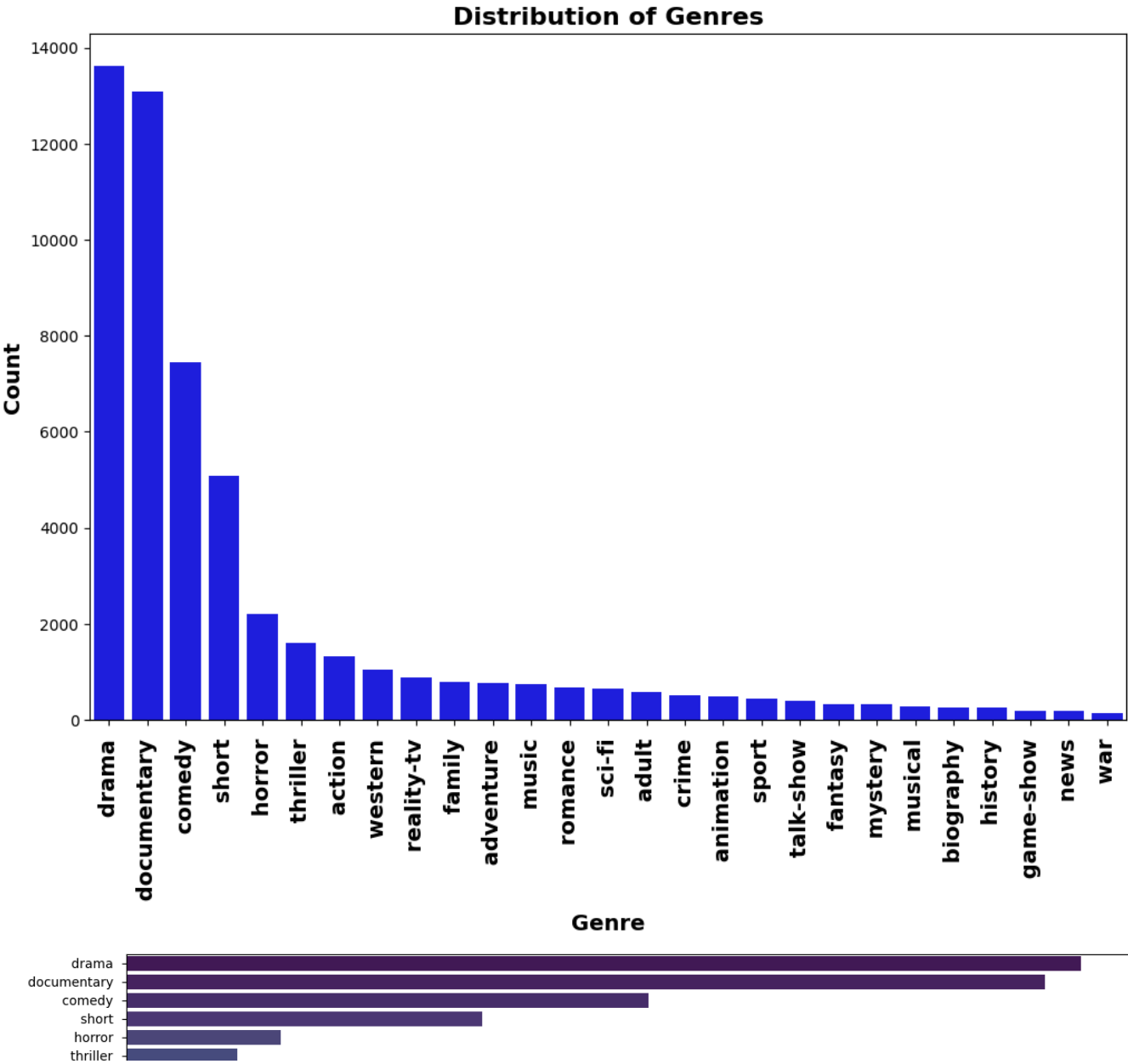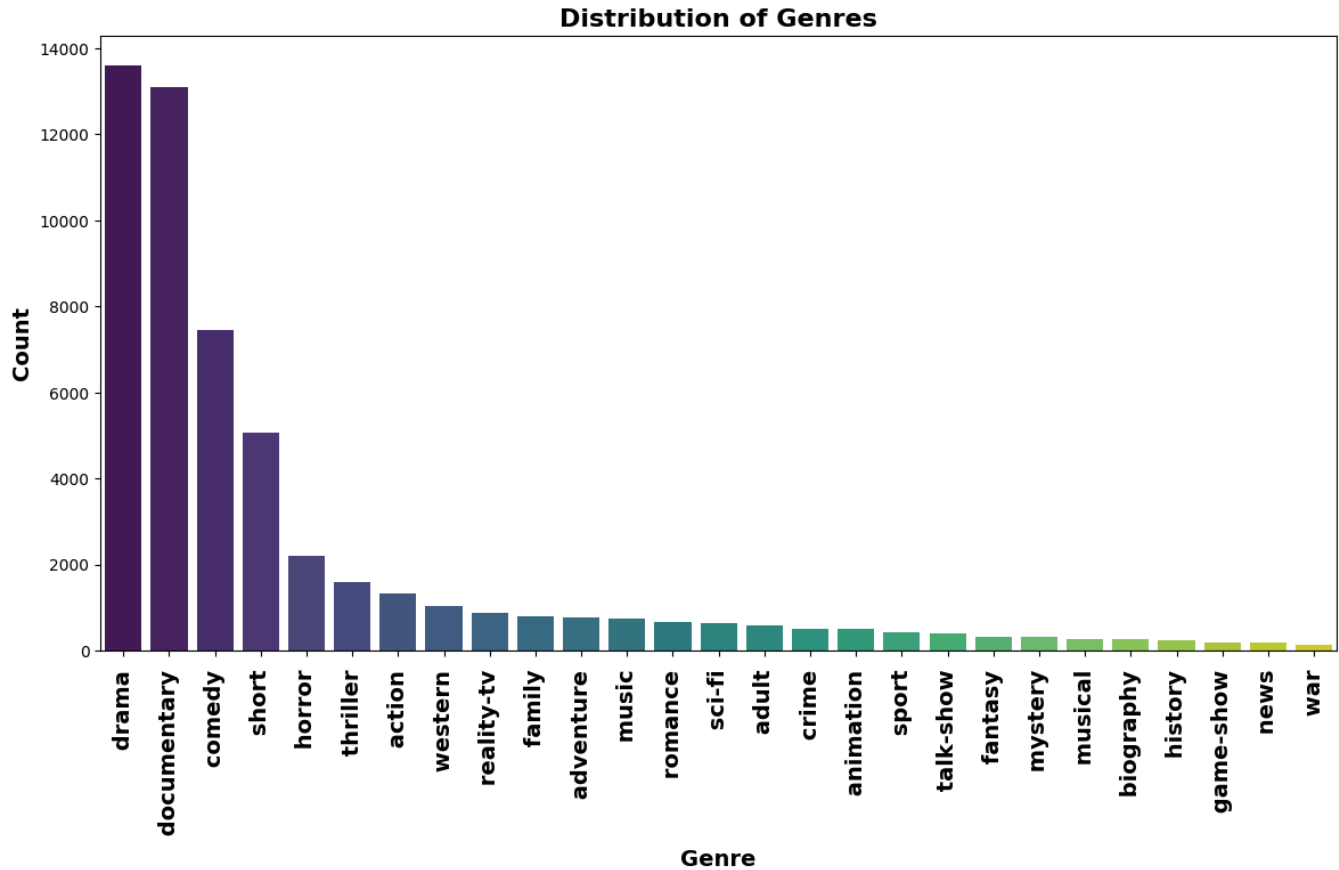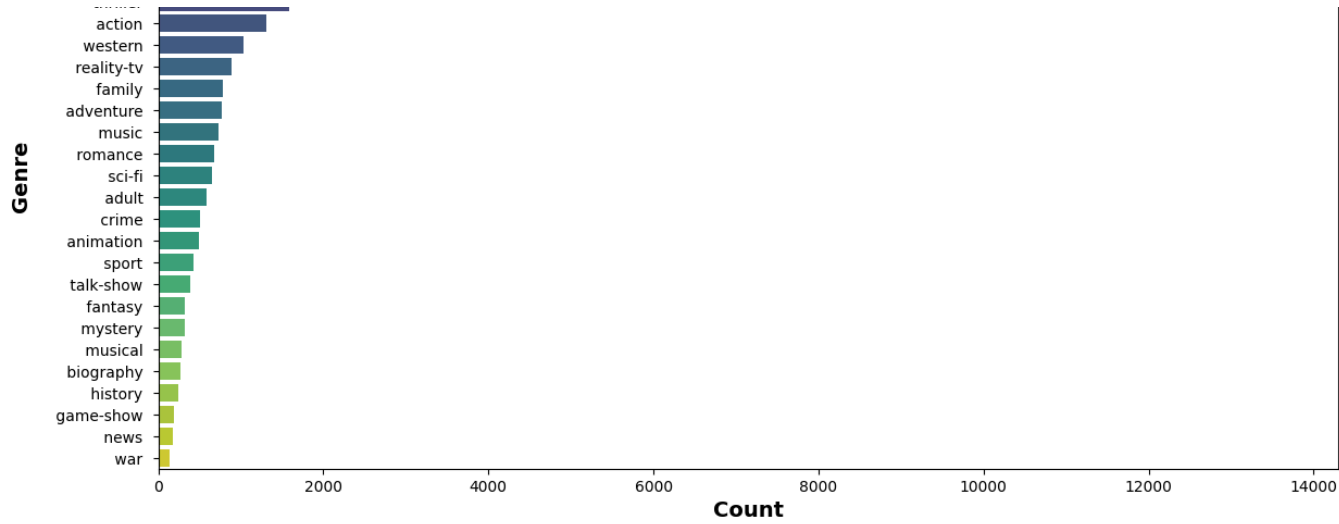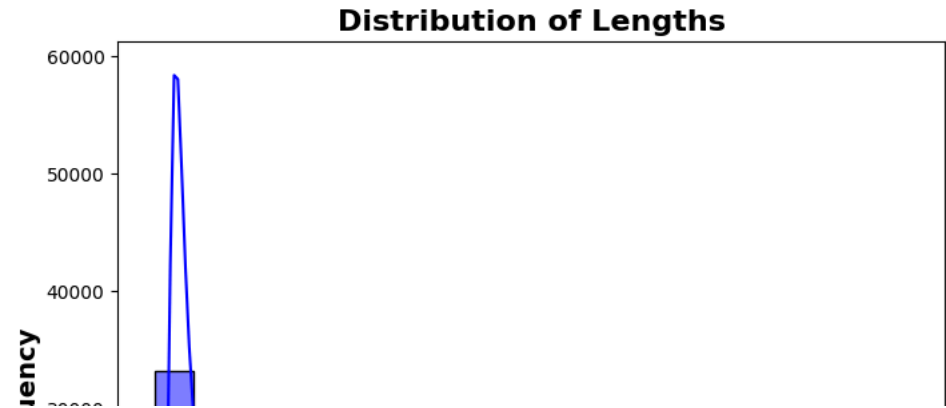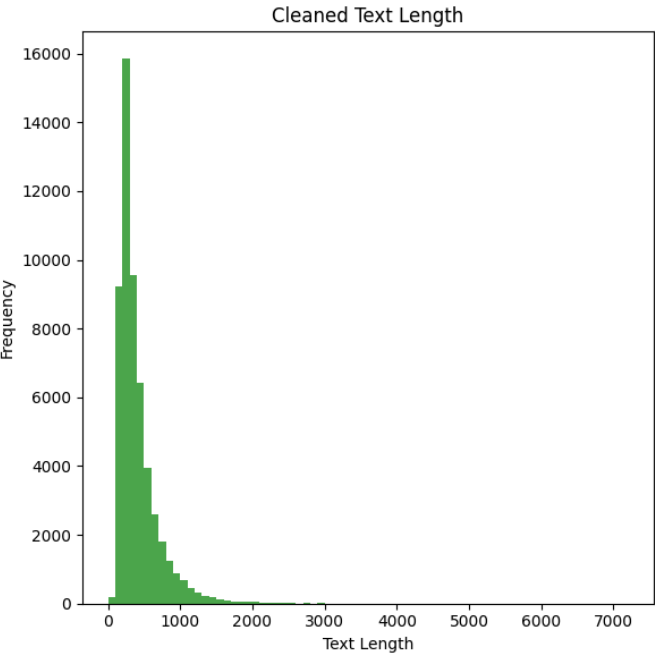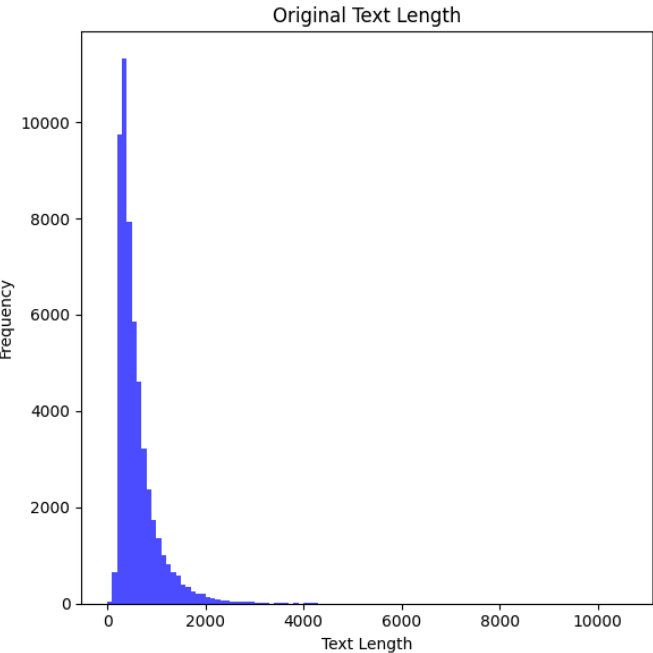
## Distribution of Genres

**Distribution of Genres**



```
<class 'pandas.core.frame.DataFrame'>
Index: 54214 entries, 1 to 54214
Data columns (total 3 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   Title        54214 non-null  object
 1   Genre        54214 non-null  object
 2   Description  54214 non-null  object
dtypes: object(3)
memory usage: 1.7+ MB
shape before drop nulls (54214, 4)
shape after drop nulls (54214, 4)
```

**Distribution of Lengths**

Original Text Length



Cleaned Text Length



Validation Accuracy: 0.44526422576777647
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are i
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are i
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are i
  _warn_prf(average, modifier, msg_start, len(result))

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| action       | 0.00      | 0.00   | 0.00     | 263     |
| adult        | 0.00      | 0.00   | 0.00     | 112     |
| adventure    | 0.00      | 0.00   | 0.00     | 139     |
| animation    | 0.00      | 0.00   | 0.00     | 104     |
| biography    | 0.00      | 0.00   | 0.00     | 61      |
| comedy       | 0.61      | 0.04   | 0.07     | 1443    |
| crime        | 0.00      | 0.00   | 0.00     | 107     |
| documentary  | 0.54      | 0.90   | 0.67     | 2659    |
| drama        | 0.38      | 0.88   | 0.53     | 2697    |
| family       | 0.00      | 0.00   | 0.00     | 150     |
| fantasy      | 0.00      | 0.00   | 0.00     | 74      |
| game-show    | 0.00      | 0.00   | 0.00     | 40      |
| history      | 0.00      | 0.00   | 0.00     | 45      |
| horror       | 0.00      | 0.00   | 0.00     | 431     |
| music        | 0.00      | 0.00   | 0.00     | 144     |
| musical      | 0.00      | 0.00   | 0.00     | 50      |
| mystery      | 0.00      | 0.00   | 0.00     | 56      |
| news         | 0.00      | 0.00   | 0.00     | 34      |
| reality-tv   | 0.00      | 0.00   | 0.00     | 192     |
| romance      | 0.00      | 0.00   | 0.00     | 151     |
| sci-fi       | 0.00      | 0.00   | 0.00     | 143     |
| short        | 0.50      | 0.00   | 0.00     | 1045    |
| sport        | 0.00      | 0.00   | 0.00     | 93      |
| talk-show    | 0.00      | 0.00   | 0.00     | 81      |
| thriller     | 0.00      | 0.00   | 0.00     | 309     |
| war          | 0.00      | 0.00   | 0.00     | 20      |
| western      | 0.00      | 0.00   | 0.00     | 200     |
|              |           |        |          |         |
| accuracy     |           |        | 0.45     | 10843   |
| macro avg    | 0.08      | 0.07   | 0.05     | 10843   |
| weighted avg | 0.36      | 0.45   | 0.31     | 10843   |

```
        Id                    Title  \
0        1          Edgar's Lunch (1998)
```

```
1          2        La guerra de papá (1977)
2          3       Off the Beaten Track (2010)
3          4          Meu Amigo Hindu (2015)
4          5               Er nu zhai (1955)
...       ...                            ...
54195   54196   "Tales of Light & Dark" (2013)
54196   54197      Der letzte Mohikaner (1965)
54197   54198              Oliver Twink (2007)
54198   54199                Slipstream (1973)
54199   54200       Curitiba Zero Grau (2010)


                                Description  \
0       L.R. Brane loves his life - his car, his apar...
1       Spain, March 1964: Quico is a very naughty ch...
2       One year in the life of Albin and his family ...
3       His father has died, he hasn't spoken with hi...
4       Before he was known internationally as a mart...
...                                                   ...
54195   Covering multiple genres, Tales of Light & Da...
54196   As Alice and Cora Munro attempt to find their...
54197   A movie 169 years in the making. Oliver Twist...
54198   Popular, but mysterious rock D.J Mike Mallard...
54199   Curitiba is a city in movement, with rhythms ...


                                Text_cleaning Predicted_Genre
0       brane loves life car apartment job especially ...          drama
1       spain march quico naughty child three belongin...          drama
2       one year life albin family shepherds north tra...    documentary
3       father died hasnt spoken brother years serious...          drama
4       known internationally martial arts superstar b...          drama
...                                                   ...            ...
54195   covering multiple genres tales light dark anth...          drama
54196   alice cora munro attempt find father british o...          drama
54197   movie years making oliver twist artful dodger ...          drama
54198   popular mysterious rock mike mallard askew bro...          drama
54199   curitiba city movement rhythms different pulsa...    documentary

[54200 rows x 5 columns]
```