

```
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
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import accuracy_score, classification_report

nltk.download('stopwords')
nltk.download('punkt')

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

train_path = "/content/train_data.txt"
train_data = pd.read_csv(train_path, sep=':::', names=['Title', 'Genre', 'Description'], engine='python')

print(train_data.describe())

              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



print(train_data.info())

<class 'pandas.core.frame.DataFrame'>
Int64Index: 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

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

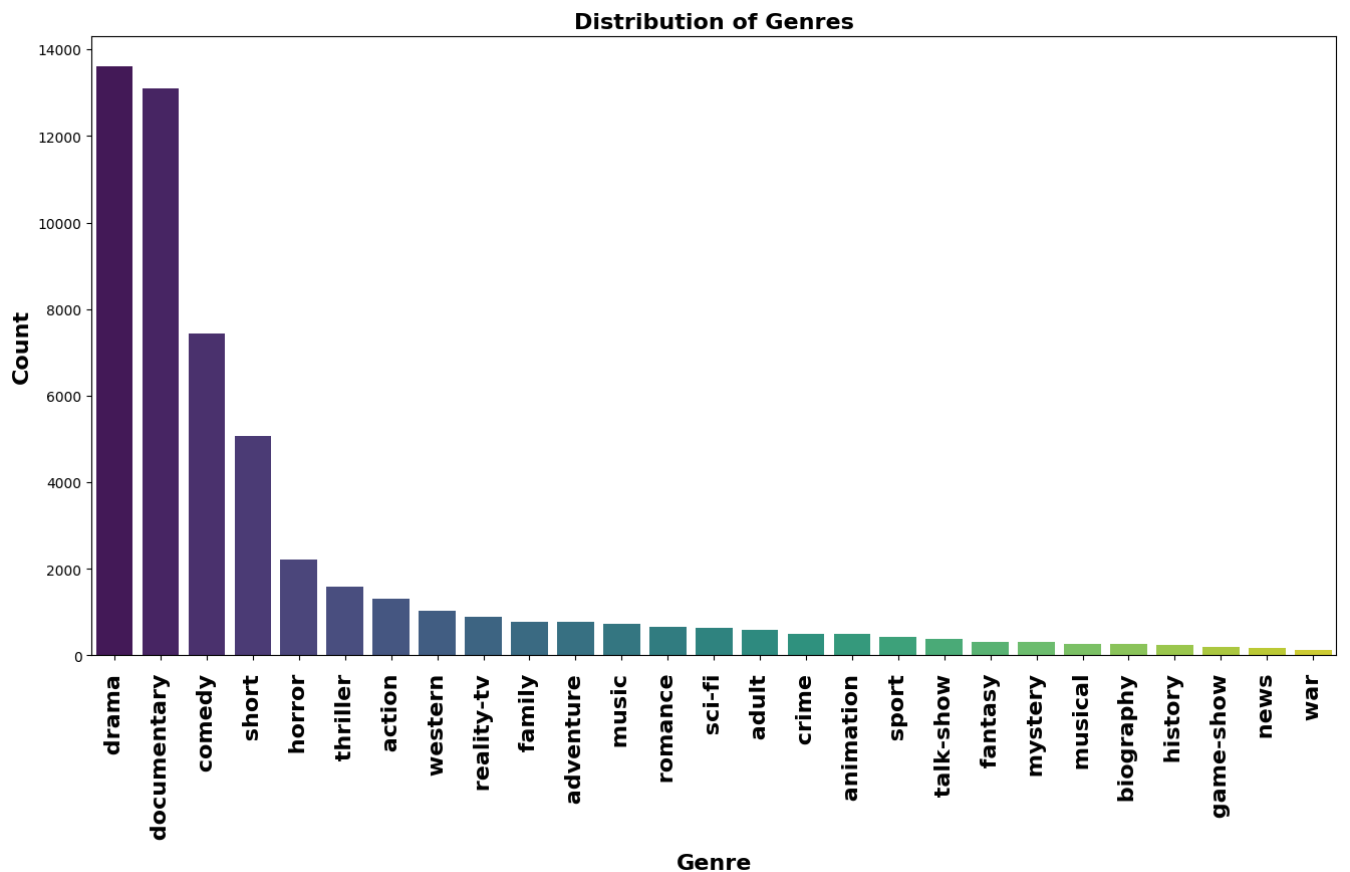
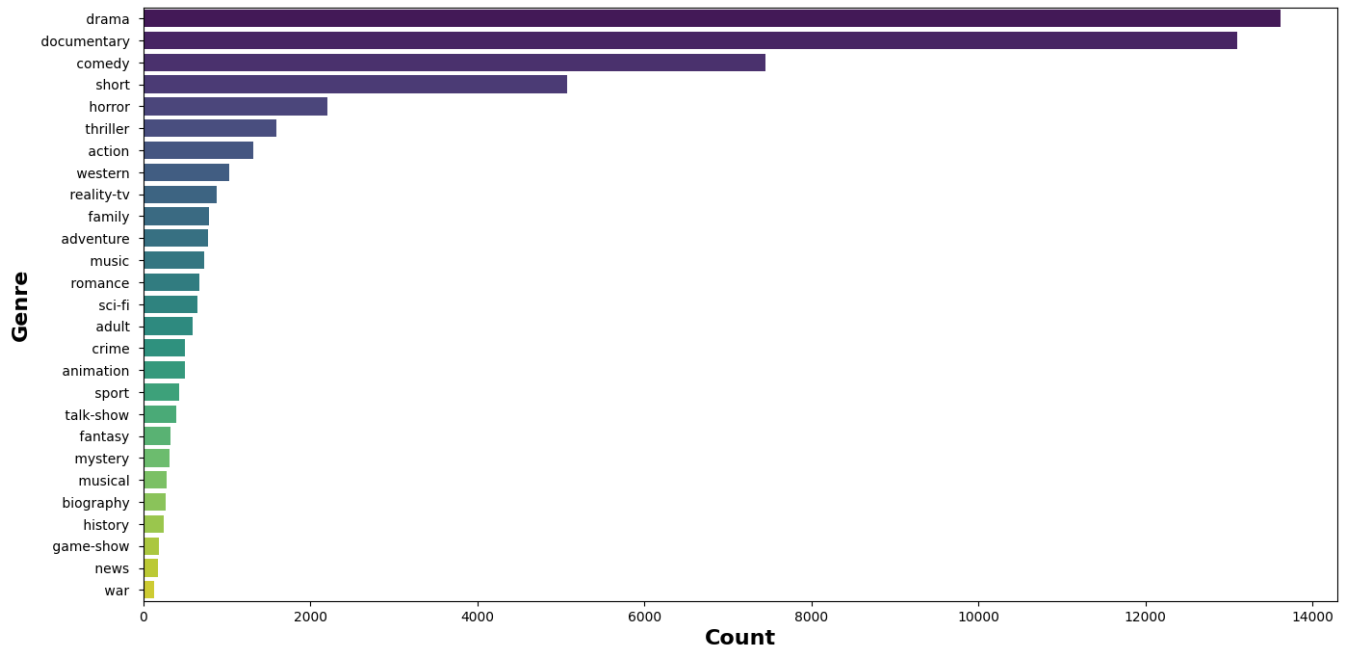
Title      0
Genre      0
Description 0
dtype: int64

test_path = "/content/test_data.txt"
test_data = pd.read_csv(test_path, sep=':::', names=['Id', 'Title', 'Description'], engine='python')
test_data.head()
```

		Id	Title	Description	
0	1		Edgar's Lunch (1998)	L.R. Brane loves his life - his car, his apar...	
1	2		La guerra de papá (1977)	Spain, March 1964: Quico is a very naughty ch...	
2	3		Off the Beaten Track (2010)	One year in the life of Albin and his family ...	
3	4		Meu Amigo Hindu (2015)	His father has died, he hasn't spoken with hi...	
4	5		Er nu zhai (1955)	Before he was known internationally as a mart...	

```
# Plot the distribution of genres in the training data
plt.figure(figsize=(16, 8))
sns.countplot(data=train_data, y='Genre', order=train_data['Genre'].value_counts().index, palette='viridis')
plt.xlabel('Count', fontsize=16, fontweight='bold')
plt.ylabel('Genre', fontsize=16, fontweight='bold')

# Plot the distribution of genres using a bar plot
plt.figure(figsize=(16, 8))
counts = train_data['Genre'].value_counts()
sns.barplot(x=counts.index, y=counts, palette='viridis')
plt.xlabel('Genre', fontsize=16, fontweight='bold')
plt.ylabel('Count', fontsize=16, fontweight='bold')
plt.title('Distribution of Genres', fontsize=16, fontweight='bold')
plt.xticks(rotation=90, fontsize=16, fontweight='bold')
plt.show()
```



```

# 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
train_data['Text_cleaning'] = train_data['Description'].apply(clean_text)
test_data['Text_cleaning'] = test_data['Description'].apply(clean_text)

# 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()

# Initialize the TF-IDF vectorizer
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'])

# Split the data into training and validation sets
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)

# Make predictions on the validation set
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))

```

```

Validation Accuracy: 0.44526422576777647
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

```

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

```

```

# Use the trained model to make predictions on the test data
X_test_predictions = classifier.predict(X_test)
test_data['Predicted_Genre'] = X_test_predictions

```

```

# Save the test_data DataFrame with predicted genres to a CSV file
test_data.to_csv('predicted_genres.csv', index=False)

```

```

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

```

	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 Twist (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]
```