

```

1 from google.colab import files
2
3 # This will prompt you to upload a file
4 uploaded = files.upload()

```



Choose Files 3 files

- description.txt(text/plain) - 366 bytes, last modified: 6/18/2021 - 100% done
- test\_data.txt(text/plain) - 34780789 bytes, last modified: 6/18/2021 - 100% done
- train\_data.txt(text/plain) - 35435236 bytes, last modified: 6/18/2021 - 100% done

Saving description.txt to description.txt  
 Saving test\_data.txt to test\_data.txt  
 Saving train\_data.txt to train\_data.txt

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import accuracy_score, classification_report
6 import re

```

```

1 # Load the training data
2 train_df = pd.read_csv('train_data.txt', delimiter=':::', engine='python', names=['ID', 'Title', 'Genre', 'Description'])
3
4 # Load the test data
5 test_df = pd.read_csv('test_data.txt', delimiter=':::', engine='python', names=['ID', 'Title', 'Description'])
6
7 # Check the first few rows of the data
8 print(train_df.head())
9 print(test_df.head())

```



ID	Title	Genre
0 1	Oscar et la dame rose (2009)	drama
1 2	Cupid (1997)	thriller
2 3	Young, Wild and Wonderful (1980)	adult
3 4	The Secret Sin (1915)	drama
4 5	The Unrecovered (2007)	drama

Description
0 Listening in to a conversation between his do...
1 A brother and sister with a past incestuous r...
2 As the bus empties the students for their fie...
3 To help their unemployed father make ends mee...
4 The film's title refers not only to the un-re...

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)

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...

```

1 # Feature Extraction using TF-IDF
2 tfidf = TfidfVectorizer(stop_words='english', max_features=5000)
3
4 # Fit and transform the training data descriptions, and transform the test data descriptions
5 X_train = tfidf.fit_transform(train_df['Description'])
6 X_test = tfidf.transform(test_df['Description'])
7
8 # Labels (Genres)
9 y_train = train_df['Genre']

```

```

1 # Initialize the classifier
2 model = LogisticRegression(max_iter=1000)
3
4 # Train the model
5 model.fit(X_train, y_train)

```



LogisticRegression

LogisticRegression(max\_iter=1000)

```

1 # Predict the genres for the test data
2 y_pred = model.predict(X_test)
3
4 # Add predictions to the test DataFrame
5 test_df['Predicted_Genre'] = y_pred
6
7 # Display the first few rows of the predictions
8 print(test_df.head())

```

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

```

1 import pickle
2
3 # Save the model
4 with open('genre_prediction_model.pkl', 'wb') as model_file:
5     pickle.dump(model, model_file)
6
7 # Save the TF-IDF vectorizer
8 with open('tfidf_vectorizer.pkl', 'wb') as vec_file:
9     pickle.dump(tfidf, vec_file)

```

```
1 uploaded = files.upload()
```

Choose Files test\_data\_solution.txt

- test\_data\_solution.txt(text/plain) - 35436708 bytes, last modified: 6/18/2021 - 100% done

Saving test\_data\_solution.txt to test\_data\_solution.txt

```

1 # Load the test data solution
2 test_solution_df = pd.read_csv('test_data_solution.txt', delimiter=':::', engine='python', names=['ID', 'Title', 'Genre', 'Descr
3
4 # Display the first few rows to ensure it's loaded correctly
5 print(test_solution_df.head())

```

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

```

1 # True genres from the test solution
2 y_true = test_solution_df['Genre']
3
4 # Descriptions for prediction (these should match the descriptions in the test set)
5 X_test_solution = tfidf.transform(test_solution_df['Description'])

```

```

1 # Predict the genres for the test data
2 y_pred = model.predict(X_test_solution)
3
4 # Evaluate the model's accuracy
5 accuracy = accuracy_score(y_true, y_pred)
6 print(f'Accuracy: {accuracy * 100:.2f}%')
7
8 # Detailed classification report
9 print(classification_report(y_true, y_pred))

```

Accuracy: 58.39%

```

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score
_warn_prf(average, modifier, msg_start, len(result))
precision    recall  f1-score   support

```

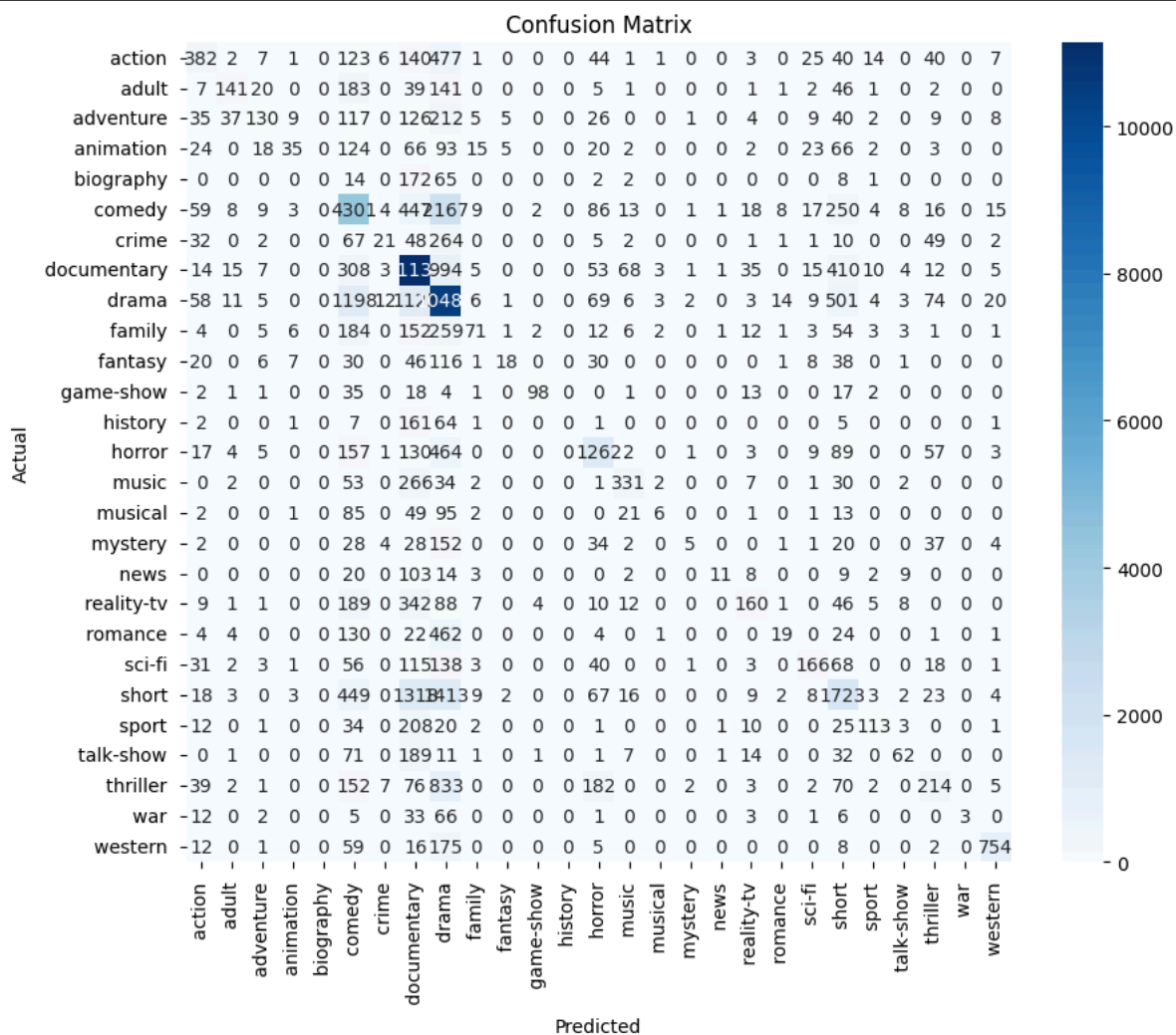
action	0.48	0.29	0.36	1314
adult	0.60	0.24	0.34	590
adventure	0.58	0.17	0.26	775
animation	0.52	0.07	0.12	498
biography	0.00	0.00	0.00	264
comedy	0.53	0.58	0.55	7446
crime	0.36	0.04	0.07	505
documentary	0.67	0.85	0.75	13096
drama	0.54	0.77	0.64	13612
family	0.49	0.09	0.15	783
fantasy	0.56	0.06	0.10	322
game-show	0.92	0.51	0.65	193
history	0.00	0.00	0.00	243
horror	0.64	0.57	0.61	2204
music	0.67	0.45	0.54	731
musical	0.33	0.02	0.04	276
mystery	0.36	0.02	0.03	318
news	0.69	0.06	0.11	181
reality-tv	0.51	0.18	0.27	883
romance	0.39	0.03	0.05	672
sci-fi	0.55	0.26	0.35	646
short	0.47	0.34	0.40	5072
sport	0.67	0.26	0.38	431
talk-show	0.59	0.16	0.25	391
thriller	0.38	0.13	0.20	1590
war	1.00	0.02	0.04	132
western	0.91	0.73	0.81	1032
accuracy			0.58	54200
macro avg	0.53	0.26	0.30	54200
weighted avg	0.57	0.58	0.55	54200

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score
_warn_prf(average, modifier, msg_start, len(result))
```

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 from sklearn.metrics import confusion_matrix
4 import numpy as np
```

```
1 # Predict the genres for the test data
2 y_pred = model.predict(X_test_solution)
3
4 # Compute accuracy
5 accuracy = accuracy_score(y_true, y_pred)
6 print(f'Accuracy: {accuracy * 100:.2f}%')
7
8 # Generate a confusion matrix
9 conf_matrix = confusion_matrix(y_true, y_pred, labels=np.unique(y_true))
10
11 # Plot the confusion matrix
12 plt.figure(figsize=(10, 8))
13 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=np.unique(y_true), yticklabels=np.unique(y_true))
14 plt.title('Confusion Matrix')
15 plt.xlabel('Predicted')
16 plt.ylabel('Actual')
17 plt.show()
18
19 # Detailed classification report
20 report = classification_report(y_true, y_pred, output_dict=True)
21 report_df = pd.DataFrame(report).transpose()
22
23 # Plot accuracy for each genre
24 plt.figure(figsize=(12, 6))
25 sns.barplot(x=report_df.index[:-3], y=report_df['precision'][:-3], palette='viridis')
26 plt.title('Precision for Each Genre')
27 plt.xlabel('Genre')
28 plt.ylabel('Precision')
29 plt.xticks(rotation=45)
30 plt.show()
```

Accuracy: 58.39%



```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score  
_warn_prf(average, modifier, msg_start, len(result))  
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score  
_warn_prf(average, modifier, msg_start, len(result))  
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score  
_warn_prf(average, modifier, msg_start, len(result))  
<ipython-input-19-88cb2754012e>:25: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and  
`sns.barplot(x=report_df.index[:-3], y=report_df['precision'][:-3], palette='viridis')`

