# Machine learning project

## MOVIE GENRE CLASSIFICATION

**Dataset Link :**

https://www.kaggle.com/datasets/hijest/genre-classification-dataset-imdb

https://www.kaggle.com/datasets/hijest/genre-classification-dataset-imdb

https://www.kaggle.com/datasets/hijest/genre-classification-dataset-imdb

https://www.kaggle.com/datasets/hijest/genre-classification-dataset-imdb

**Link of the google collab file :**

https://colab.research.google.com/drive/1im6OjSWedMcKzLkkKko2jTm8spSajOf-#scrollTo=dvr2sTW-1hiz

**Github Repository Link :**

https://github.com/AbhishekYadav-01/Encryptix/tree/main/MOVIE%20GENRE%20CLASSIFICATION

**Aim :**

Create a machine learning model that can predict the genre of a movie based on its plot summary or other textual information. You can use technique TF-IDF with classifiers such as Naive Bayes, Logistic Regression, or Support Vector Machines.

**Sol :**

**Data Preprocessing** :

- Text Cleaning : Removal of stopwords, punctuation, and special characters.
- Tokenization : Splitting sentences into words or tokens.
- TF-IDF Vectorization : Convert text data into numerical features using TF-IDF.

**Dataset Overview :**

The dataset consists of :

- Number of Training Samples : 9560
- Number of Test Samples : 9822

**Genre Distribution in Training Set**

The training set includes a variety of movie genres, with the following distribution:

- drama         2365
- documentary   2307
- comedy        1308
- short          897

- horror               379
- thriller           313
- action              236
- western             194
- reality-tv       159
- family              138
- music               130
- adventure        129
- romance             114
- sci-fi              114
- adult               114
- animation         85
- sport                78
- crime                74
- talk-show         71
- fantasy             65
- mystery             56
- musical             52
- biography         47
- ...
- game-show         37
- news                 31
- war                  22

**Model Training**:

Implemented several classifiers:

- ■
  - **Cross-Validation**: Utilized to ensure robustness of models against overfitting.
2. **Model Evaluation:**
   - Metrics used: Accuracy, Precision, Recall, and F1-Score.
   - Confusion matrices were plotted to visualize performance across different genres.

# Now we will train all models and see the results on test data:

```
Evaluating Logistic Regression model:
Accuracy (Logistic Regression): 0.5280058193910423
Evaluating Naive Bayes model:
Accuracy (Naive Bayes): 0.47428036994700196

Evaluating Support Vector Machine with linear kernel:
Accuracy (SVM (linear)): 0.5408916138418373
```

|              | precision | recall | f1-score | support |
|-------------:|----------:|-------:|---------:|--------:|
| action       | 0.49      | 0.18   | 0.26     | 237     |
| adult        | 0.50      | 0.13   | 0.20     | 102     |
| adventure    | 0.50      | 0.18   | 0.26     | 119     |
| animation    | 1.00      | 0.03   | 0.05     | 108     |
| biography    | 0.00      | 0.00   | 0.00     | 49      |
| comedy       | 0.45      | 0.51   | 0.48     | 1288    |
| crime        | 0.00      | 0.00   | 0.00     | 72      |
| documentary  | 0.63      | 0.83   | 0.72     | 2331    |
| drama        | 0.50      | 0.76   | 0.60     | 2473    |
| family       | 0.83      | 0.03   | 0.06     | 148     |
| fantasy      | 0.00      | 0.00   | 0.00     | 51      |
| game-show    | 0.89      | 0.55   | 0.68     | 29      |
| history      | 0.00      | 0.00   | 0.00     | 42      |
| horror       | 0.67      | 0.47   | 0.55     | 390     |
| music        | 0.58      | 0.32   | 0.41     | 133     |
| musical      | 0.00      | 0.00   | 0.00     | 41      |
| mystery      | 0.00      | 0.00   | 0.00     | 54      |
| news         | 0.67      | 0.06   | 0.11     | 34      |
| reality-tv   | 0.71      | 0.13   | 0.22     | 152     |
| romance      | 0.00      | 0.00   | 0.00     | 133     |
| sci-fi       | 0.40      | 0.16   | 0.23     | 115     |
| short        | 0.44      | 0.23   | 0.30     | 933     |
| sport        | 0.79      | 0.15   | 0.26     | 71      |
| ...          |           |        |          |         |
| weighted avg | 0.52      | 0.54   | 0.49     | 9623    |

```
Evaluating Support Vector Machine with poly kernel:
Accuracy (SVM (poly)): 0.42949184246077104
Evaluating Support Vector Machine with rbf kernel:
Accuracy (SVM (rbf)): 0.5128338356022031

Evaluating Support Vector Machine with sigmoid kernel:
Accuracy (SVM (sigmoid)): 0.5370466590460355


Evaluating Random Forest model:
Accuracy (Random Forest): 0.47448820534136965
Model Comparison:
Logistic Regression: 0.5280058193910423
Naive Bayes: 0.47428036994700196
Random Forest: 0.47448820534136965
SVM (linear): 0.5408916138418373
SVM (poly): 0.42949184246077104
SVM (rbf): 0.5128338356022031
SVM (sigmoid): 0.5370466590460355
```
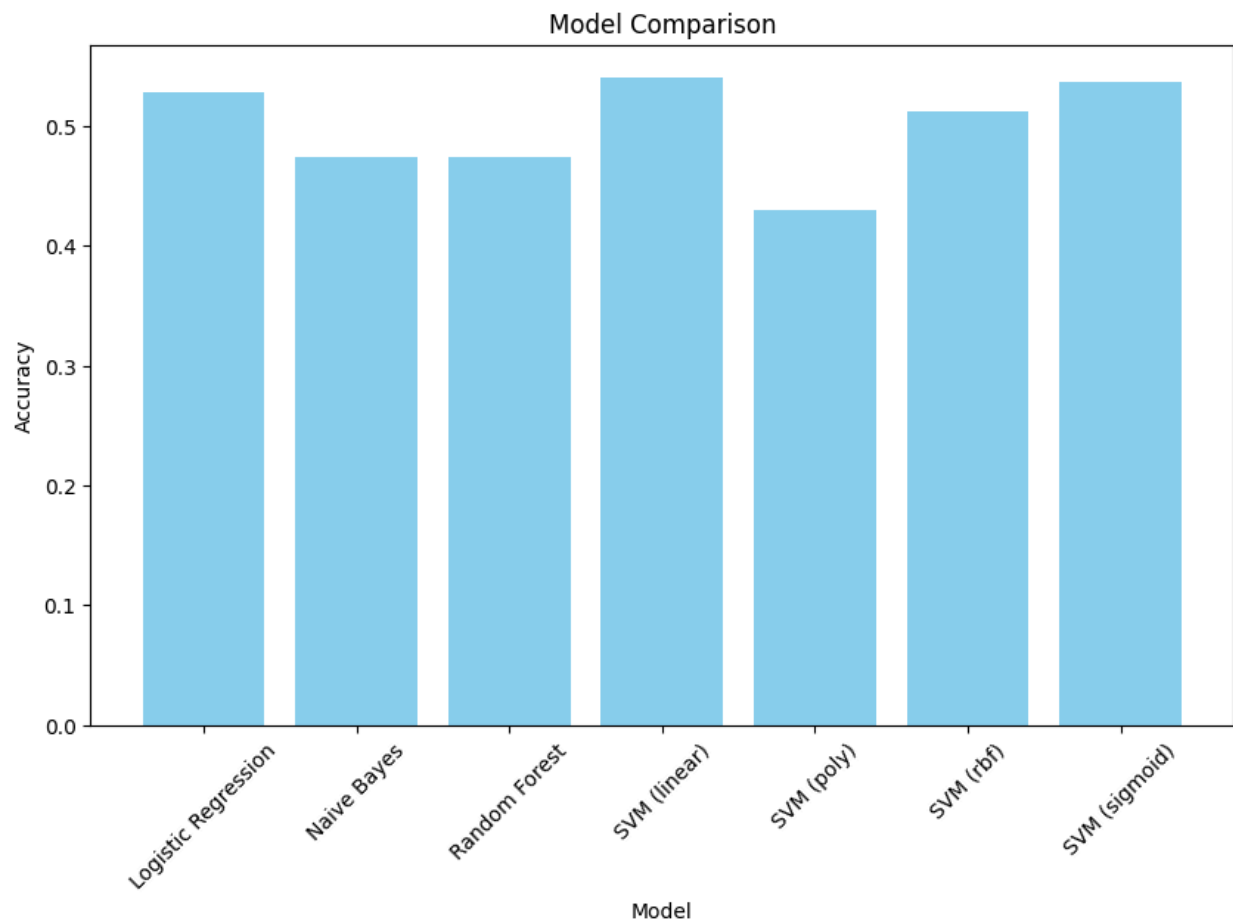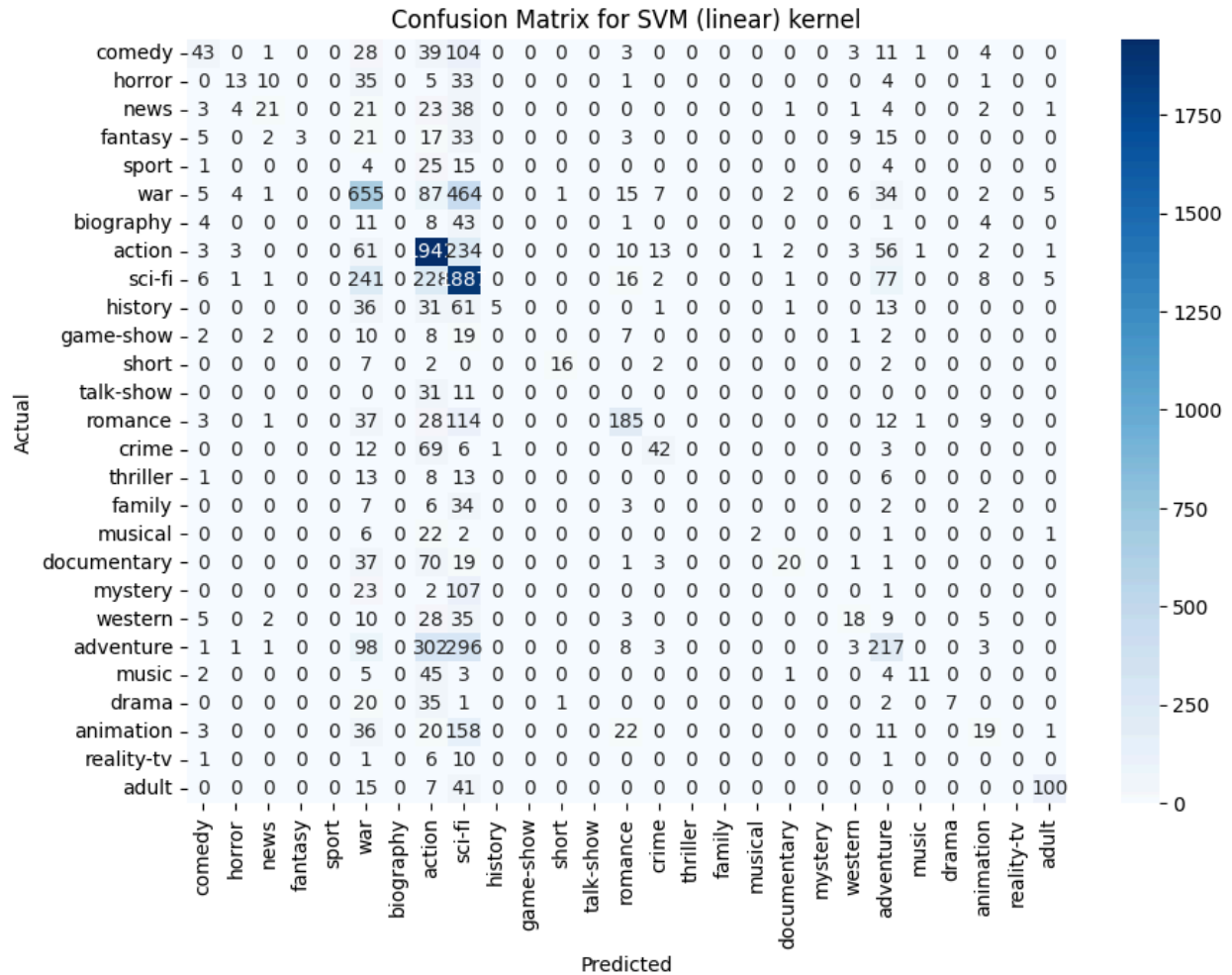
Model Comparison

**This is the confusion matrix for SVM linear kernal :**
**(For other models confusion matrix check on github)**

Confusion Matrix for SVM (linear) kernel

**Conclusion :**

In this project, we implemented and evaluated various machine learning models to classify movie genres based on plot summaries. The Support Vector Machine with a linear kernel achieved the highest accuracy (0.541), outperforming other models. Logistic Regression and SVM with a sigmoid kernel also performed relatively well. However, it is important to note that while accuracy is a useful metric, other factors such as precision, recall, and the confusion matrix should also be considered to understand the models' performance better. Future improvements could include exploring deep learning models, optimizing hyperparameters, and incorporating additional features to enhance the classification accuracy.