



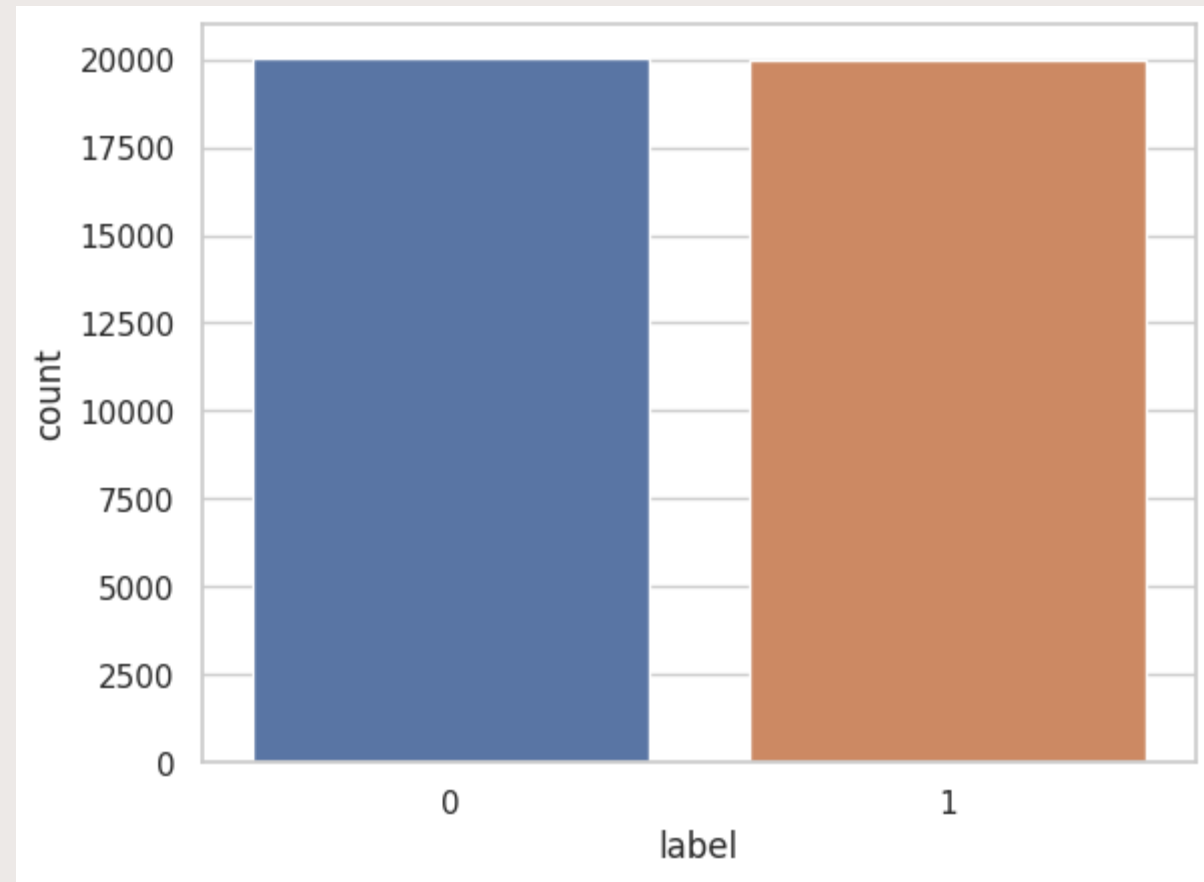
IMDB MOVIE RATINGS SENTIMENT ANALYSIS

OVERVIEW

- Building a machine learning algorithm to predict whether the given review is positive or negative.
- My approach was to try three different model son the dataset and compare the results based on the metrics of F1 score, accuracy, etc.

DATA

- CSV file which included 50,000 data points.
- Features of the CSV file include:
 - *Text*
 - *Label*
- The datapoints were equally distributed and balanced.



PREPROCESSING

The data cleaning process involves several techniques to prepare text data for analysis, including lowering text to lowercase, removing URLs, punctuation and numbers, splitting text into individual words, removing stop words, and reducing words to their root form.

TRAINING

- The data points were split into 80-20 split for better training of the model and testing it upon the remaining data points.
- Three models were chosen to be worked upon for this sentiment analysis.
 - *Logistic Regression Classifier*
 - *RandomForestClassifier*
 - *XGBoost*.

PERFORMANCE

- Comparing the performance of each of the models used.
- Logistic Regression Classifier was the best out of all the three models used.
 - *With the highest accuracy, precision, recall and f1 score.*

	Accuracy	Precision	Recall	F1 Score
Logistic Regression	0.881183	0.871605	0.892767	0.882059
Random Forest	0.845437	0.847527	0.840668	0.844083
XGBoost	0.851982	0.840108	0.867729	0.853695

CONCLUSION

- As we can see the Logistic Regression Classifier was best in predicting the sentiments of the given dataset, it can be best used for the use of movie rating sentiment prediction.

REPOSITORY

