

Choose Files No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please run this cell to enable

```
1 !mkdir ~/.kaggle
2 !cp kaggle.json ~/.kaggle/
3 !cat ~/.kaggle/kaggle.json
4 !chmod 600 ~/.kaggle/kaggle.json

1 !kaggle datasets download -d lakshmi25npathi/imdb-dataset-of-50k-movie-reviews
2 read_reply_from_input(message_id, timeout_sec)
3 !unzip imdb-dataset-of-50k-movie-reviews.zip
4 ---> 97 time.sleep(0.025)

1 import pandas as pd
2
3 df = pd.read_csv('IMDB Dataset.csv')
```

1 df

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production.   The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive
...	...	...
49995	I thought this movie did a down right good job...	positive
49996	Bad plot, bad dialogue, bad acting, idiotic di...	negative
49997	I am a Catholic taught in parochial elementary...	negative
49998	I'm going to have to disagree with the previou...	negative
49999	No one expects the Star Trek movies to be high...	negative

50000 rows x 2 columns

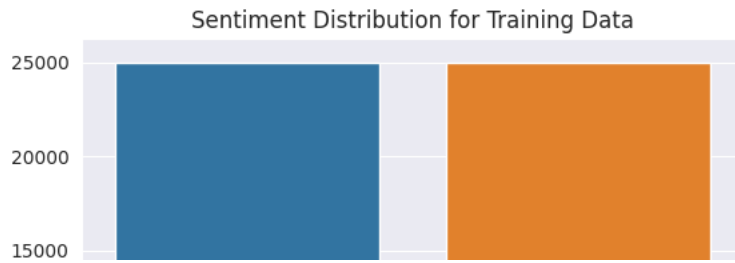
```
1 # get cols of the df
2 df.columns

Index(['review', 'sentiment'], dtype='object')

1 # sentiment is the y value, review is the x value. We split with an 80/20 split
2 from sklearn.model_selection import train_test_split
3
4 X_train, X_test, y_train, y_test = train_test_split(df['review'],
5                                                    df['sentiment'],
6                                                    test_size=0.2,
7                                                    random_state=42)

1 import seaborn as sns
2
3 sns.set_style('darkgrid')
4 sns.countplot(x='sentiment', data=df).set(title='Sentiment Distribution for Training Data')
```

[Text(0.5, 1.0, 'Sentiment Distribution for Training Data')]



## Describe the data set and what the model should be able to predict.

This data set contains IMDB reviews that have a sentiment assigned to them - either positive or negative based on the review's content.

The model should be able to predict accurately the sentiment of a review given the review's content. There is an 80/20 split between training/testing data for the model to predict on.

```

1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.metrics import accuracy_score
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 # vectorize the text data
6 vectorizer = CountVectorizer()
7 X_train_vec = vectorizer.fit_transform(X_train)
8 X_test_vec = vectorizer.transform(X_test)
9
10 # make a Naive Bayes model - no smoothing
11 nb_model = MultinomialNB()
12 nb_model.fit(X_train_vec, y_train)
13
14 # make predictions on the test data
15 y_pred = nb_model.predict(X_test_vec)
16
17 # get accuracy
18 accuracy = accuracy_score(y_pred, y_test)
19 print("Naive Bayes Accuracy without Laplace Smoothing: ", accuracy)
20
21 # make a Naive Bayes model - with Laplace smoothing
22 nb_model = MultinomialNB(alpha=100)
23 nb_model.fit(X_train_vec, y_train)
24
25 # make predictions on the test data
26 y_pred = nb_model.predict(X_test_vec)
27
28 # get accuracy
29 accuracy = accuracy_score(y_pred, y_test)
30 print("Naive Bayes Accuracy with Laplace Smoothing: ", accuracy)

Naive Bayes Accuracy without Laplace Smoothing: 0.8488
Naive Bayes Accuracy with Laplace Smoothing: 0.8343

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import StandardScaler
3
4 # normalize the features
5 scaler = StandardScaler(with_mean=False)
6 X_train_vec_scaled = scaler.fit_transform(X_train_vec)
7 X_test_vec_scaled = scaler.transform(X_test_vec)
8
9 # make a Logistic Regression model - 500 iterations
10 clf = LogisticRegression(random_state=0, max_iter=110).fit(X_train_vec_scaled, y_train)
11
12 # make predictions on the test data
13 y_pred = clf.predict(X_test_vec_scaled)
14
15 # get accuracy
16 accuracy = accuracy_score(y_pred, y_test)
17 print("Logistic Regression Accuracy - 110 iterations: ", accuracy)
18
19 # The following gave the same results; the number of iterations were used to get a local
20 # minimum # of iterations required for convergence:
21 # 100 - too few, no convergence, warning given

```

```
22 # 101 - too few
23 # 110 - converged

Logistic Regression Accuracy - 110 iterations: 0.8662

1 from sklearn.neural_network import MLPClassifier
2
3 # make a Neural Network classifier
4 clf = MLPClassifier(hidden_layer_sizes=(75,), max_iter=1000, alpha=0.0001, random_state=42)
5 clf.fit(X_train_vec, y_train)
6
7 # make predictions on the test data
8 y_pred = clf.predict(X_test_vec)
9
10 # get accuracy
11 accuracy = accuracy_score(y_pred, y_test)
12 print("Neural Network Accuracy: ", accuracy)
13
14 # This took 29 minutes and 6 seconds

Neural Network Accuracy: 0.8877
```

## Analysis

I tried three different approaches when it came to classification of movie review sentiment - Naive Bayes, Logistic Regression, and Neural Networks under the sklearn library.

In order to use the text data for classification, I used a tokenizer to convert the words accordingly for numerical input.

From my analysis, all three models performed fairly well (with >80% accuracy) in increasing order of accuracy. I'll discuss these in detail below.

## Naive Bayes

Naive Bayes was the simplest model and I used two accuracies - with and without laplace smoothing (set at 100). The accuracies are shown below:

- Naive Bayes Accuracy without Laplace Smoothing: 0.8488
- Naive Bayes Accuracy with Laplace Smoothing: 0.8343

As can be seen, the accuracy is lower the higher Laplace smoothing is. This can happen if the dataset is not sparse, i.e., if most of the features occur frequently in the training set (which it does, as there are only two values when it comes to analysis).

## Logistic Regression

Logistic Regression was the next model; unlike Naive Bayes, I only changed the number of max iterations before convergence. The accuracies are shown below:

- Logistic Regression Accuracy - 110 iterations: 0.8662

An interesting issue I noted when I was first testing out this model is that convergence wasn't reached by the default number of max iterations which was 100. Initially, I decided to set it very high at around 500 so I could achieve convergence. But once I got the same accuracy at 250, I decided to keep lowering the number of epochs until the accuracy remained the same, while convergence could still occur. In the end, 110 max iterations was my rough estimate, and it yielded a higher accuracy than Naive Bayes.

## Neural Network - 75 layers, 1000 max iterations, 0.0001 learning rate

Neural Network took by and far the longest time of any of the other models, and for good reason - it's a very complex model that doesn't rely on pure regression alone. But it still yielded the highest accuracy of any of the other models. The accuracy is below:

- Neural Network Accuracy: 0.8877

Even though it took the longest time and had the most complex structure, it performed very well for the time invested, and it performed the best of any other model. Plus, neural networks are generally best for classification, so there is that as well.

## Conclusion

Overall, I had a good time working with the different models to try and figure out what was going on with each during this process.