

1. Dataset Insights: Describe your dataset, including how it was cleaned and labeled.

For this project, I chose the sentiment analysis on IMDB Reviews. Following the project description, I selected 200 entries from the data, and equally splitting them according to the categories: Positive and Negative. Now that I have 100 entries of each categories, I labelled the sentiment column to numerical labels, 0 for negatives and 1 for positives. After that, I tokenized the data, preparing them to be processed by the machine. To ensure length consistency, I truncate and padded the data during the tokenizing stage, achieving structured form.

2. Training Process: Summarize the steps you took to fine-tune the model.

I used the distilbert-base-uncased as my base model to tackle this classification project. I loaded the pre-trained model and tokenizer to process the input. The dataset is huge, but for the focus of this project, I selected 200 entries with equal split before tokenizing the dataset. After I processed the data and have them ready for training, I used train_test_split from sklearn to split the dataset into training (80%) and test (20%) sets, ensuring the split maintained the balance of labels using stratification on the label column. This step is very important as splitting the dataset allows me to test for unseen data to estimate its generalization process. I also stratified the splits to help preserve label distribution, a prevention to biased training by ensuring that both labels are well-represented in both the train split and test split. After that, I followed the tutorial to set up the Hugging Face Trainer, trained the model, and completed the training.

3. Evaluation Results: Present your evaluation metrics and discuss the model's strengths and weaknesses.

	precision	recall	f1-score	support
0	0.61	0.95	0.75	20
1	0.89	0.40	0.55	20
accuracy			0.68	40
macro avg	0.75	0.68	0.65	40
weighted avg	0.75	0.68	0.65	40

This is the result, and the result shows that the accuracy of the model is 68%, which is considerably good given the small dataset. This model might perform better on a huge set of dataset that is also well-represented, but for a 200-entries dataset, the accuracy is sufficient. To break it down deeper, here's the brief summary:

- Precision:
 - Precision for Class 0 (Negative sentiment): 0.61 means that 61 % of the reviews predicted as "negative" were actually negative.
 - Precision for Class 1 (Positive sentiment): 0.89 means that 89% of the reviews predicted as "positive" were actually positive.

- Recall:
 - Recall for Class 0 (Negative sentiment): 0.95 means that 95% of the actual negative reviews were correctly identified by your model.
 - Recall for Class 1 (Positive sentiment): 0.40 means that only 40% of the actual positive reviews were correctly identified by your model.
- F1-Score:
 - F1-Score for Class 0: 0.75 is the harmonic mean of precision and recall for negative sentiment, showing a good balance.
 - F1-Score for Class 1: 0.55 shows a moderate balance between precision and recall for positive sentiment.

Some further improvements to be made:

- Using a bigger dataset can help the machine to learn better
- Further explore parameters and try experimenting with it
- Target the features, further feature engineering such as word-embedding

4. Application and Impact: Explain how this fine-tuned model could be used in a real-world application. Include at least one potential improvement for future iterations.

The model is trained to perform sentiment analysis, classifying positive and negative reviews. This can be used to many businesses, mainly those that sell goods/services to its clients. By being able to easily classify reviews as positive and negative, they can better their services, and if applied to further extent, they can even build an automated reply to express gratitude for positive reviews, or apologize for negative experiences. This type of automation can greatly help the business to improve their overall service, improve customer experience, all done with minimal effort.