

## Problem Definition & Dataset Selection

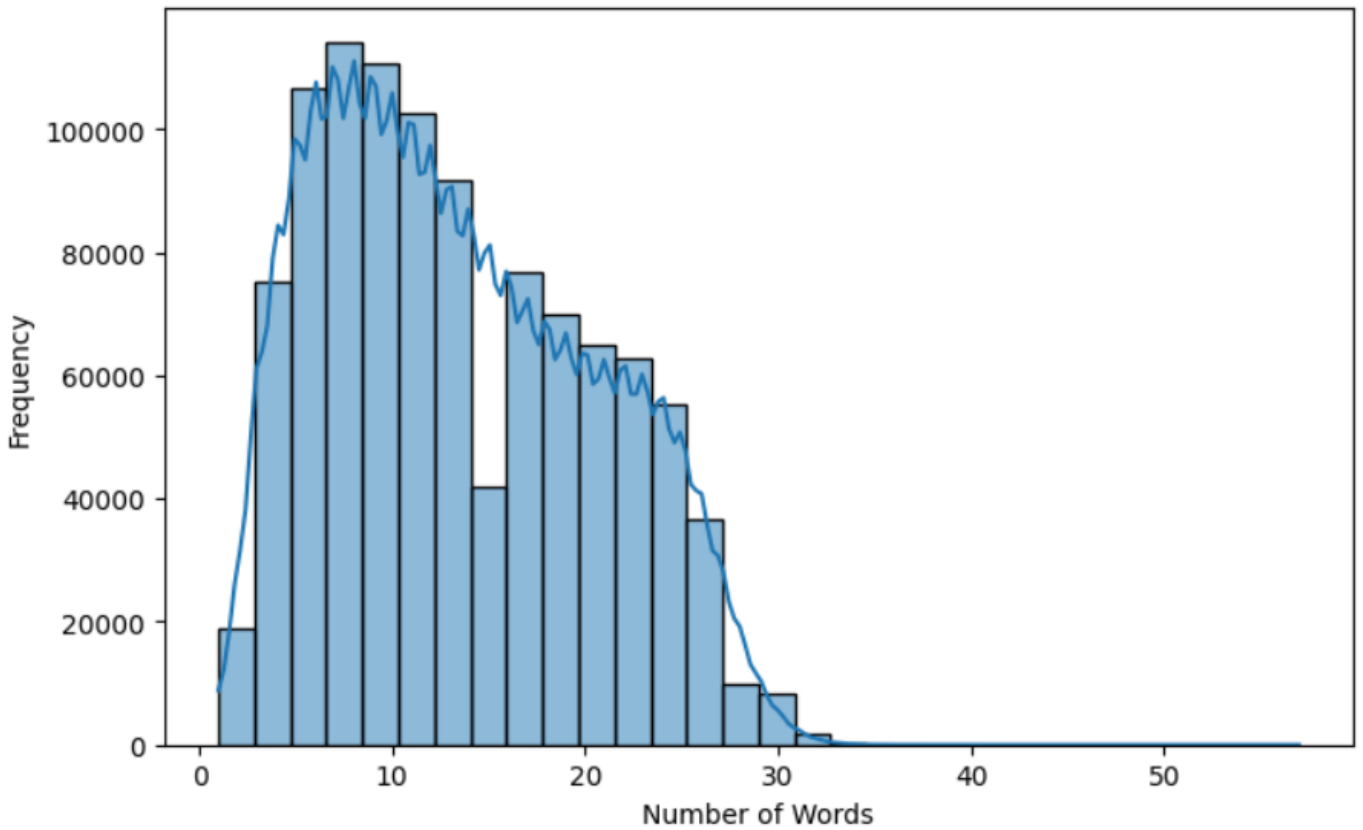
The task addressed is **Sentiment Analysis**—a key Natural Language Processing (NLP) problem where the goal is to classify the sentiment (positive or negative) of user-generated content. This task has real-world applications in customer feedback analysis, product reviews, social media monitoring, and more.

- **Source:** The dataset used is the training 14000000000d dataset, containing preprocessed tweets labeled as positive (4) or negative (0).
- **Size:** 1.6 million tweets.
- **Relevance:** This dataset is ideal for binary sentiment classification tasks and is a standard benchmark for sentiment analysis models.

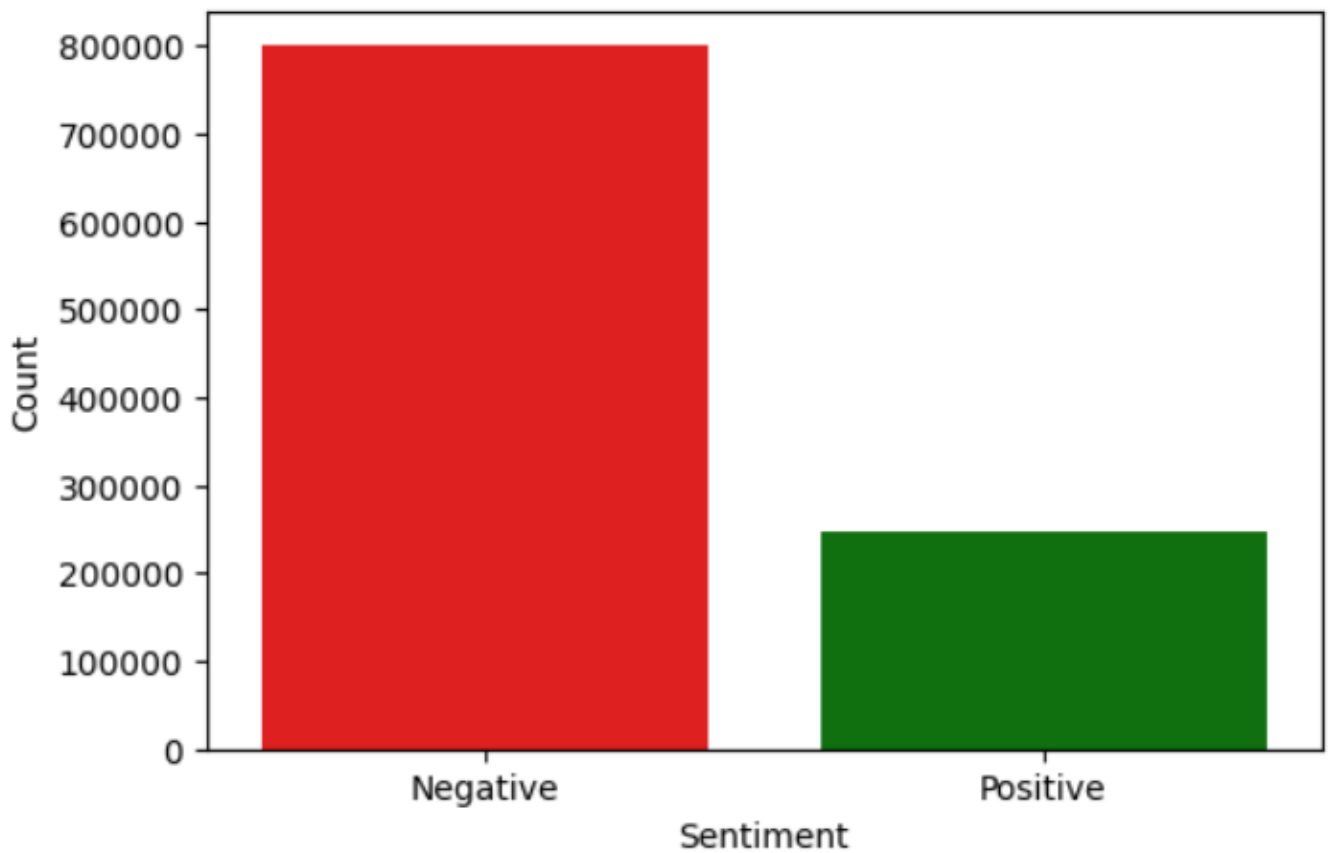
## Exploratory Data Analysis (EDA)

- [illegible]

Word Count Distribution in Tweets



Sentiment Class Distribution



## Data preprocessing

```
[7] # Data Preprocessing
    tokenizer = Tokenizer(num_words=5000)
    tokenizer.fit_on_texts(df['text'])
    X = tokenizer.texts_to_sequences(df['text'])
    X = pad_sequences(X, maxlen=100)
    y = df['polarity'].values
```

## 2. Model Selection & Implementation

### Model Chosen

An **LSTM (Long Short-Term Memory)** neural network was selected due to its superior performance in capturing long-range dependencies in text.

### Preprocessing Steps

- **Tokenization:** Tweets were tokenized using TensorFlow's Tokenizer.
- **Padding:** Input sequences were padded to ensure uniform input length using `pad_sequences`.
- **Label Encoding:** Polarity values were transformed into binary class labels (0 or 1).

```
model = Sequential()
```

```
model.add(Embedding(input_dim=10000, output_dim=32, input_length=100))
```

```
model.add(LSTM(64, return_sequences=False))
```

```
model.add(Dense(1, activation='sigmoid'))
```

**Embedding Layer:** Maps input tokens to dense vectors.

**LSTM Layer:** 64 units for sequential learning.

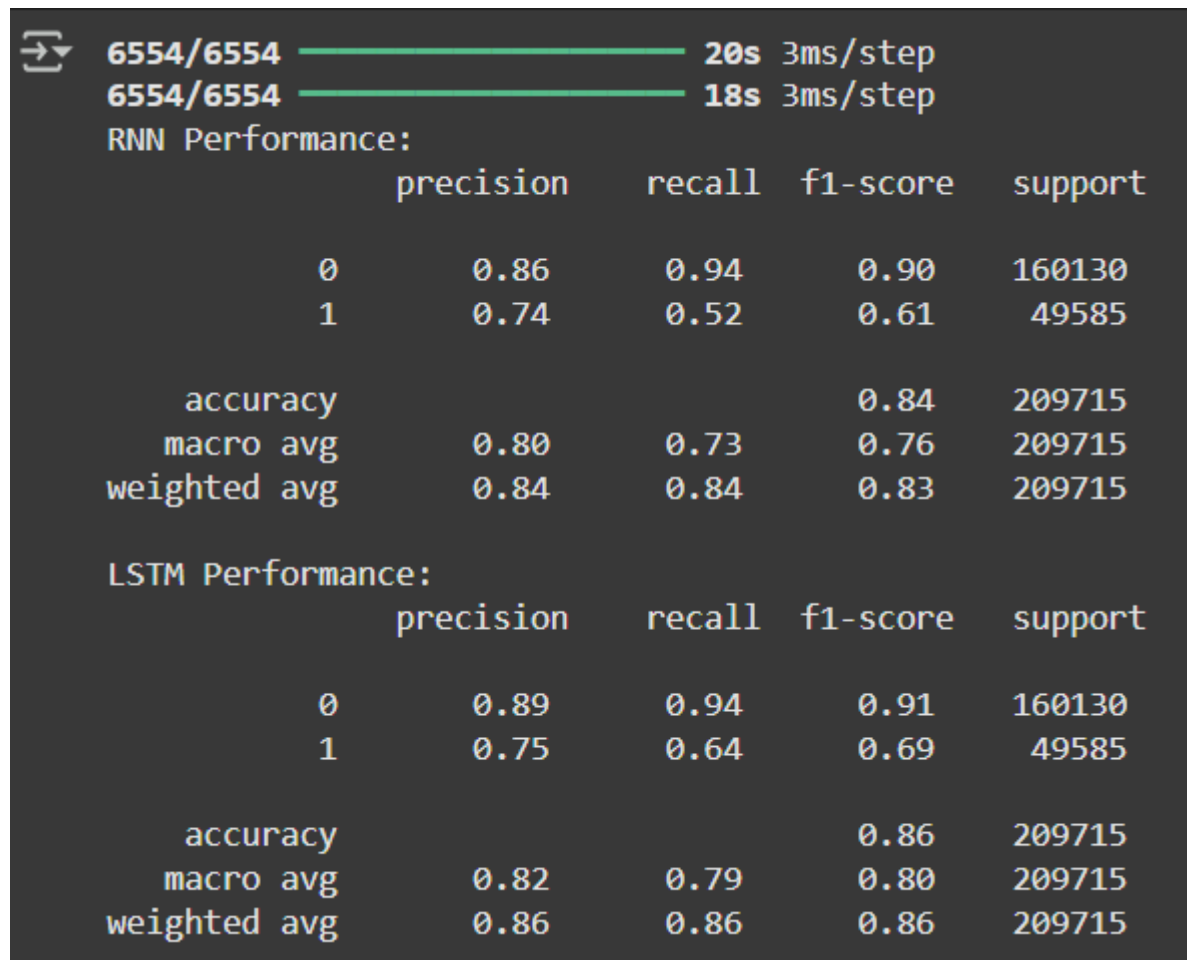
**Dense Output Layer:** Sigmoid for binary classification.

### Training Configuration

- **Optimizer:** Adam
- **Loss Function:** Binary crossentropy
- **Batch Size:** 128
- **Epochs:** 5
- **Validation Split:** 20%

Evaluation metrics:

Its compared with RNN and LSTM model



A terminal window showing the performance of an RNN and an LSTM model. The RNN section shows a training time of 20s and an inference time of 3ms/step for 6554/6554 samples. The LSTM section shows a training time of 18s and an inference time of 3ms/step for 6554/6554 samples. Both sections display a table of performance metrics including precision, recall, f1-score, and support for classes 0 and 1, as well as accuracy, macro avg, and weighted avg.

RNN Performance:					
	precision	recall	f1-score	support	
0	0.86	0.94	0.90	160130	
1	0.74	0.52	0.61	49585	
accuracy			0.84	209715	
macro avg	0.80	0.73	0.76	209715	
weighted avg	0.84	0.84	0.83	209715	

LSTM Performance:					
	precision	recall	f1-score	support	
0	0.89	0.94	0.91	160130	
1	0.75	0.64	0.69	49585	
accuracy			0.86	209715	
macro avg	0.82	0.79	0.80	209715	
weighted avg	0.86	0.86	0.86	209715	

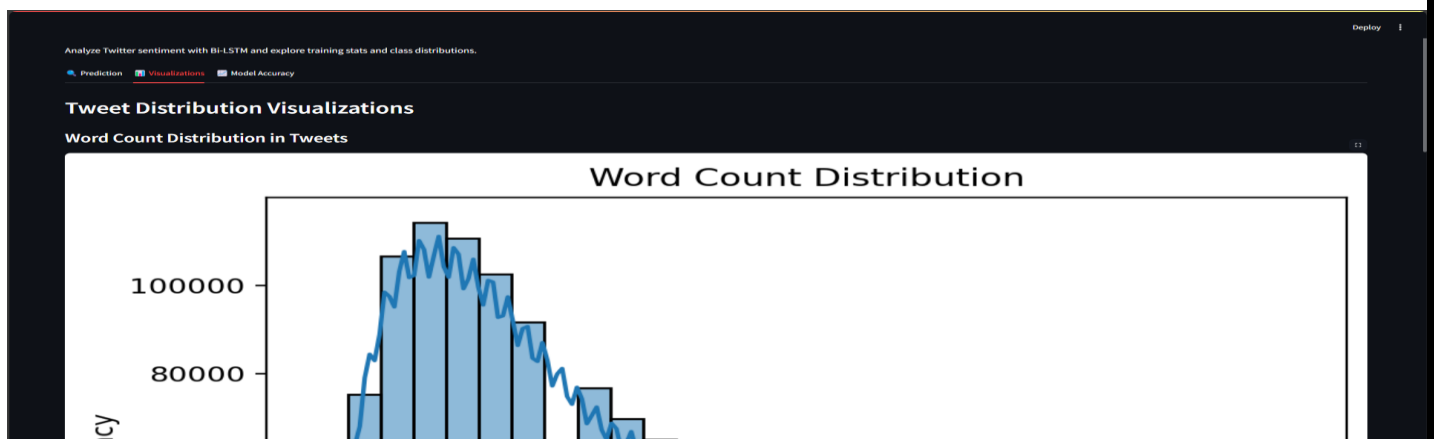
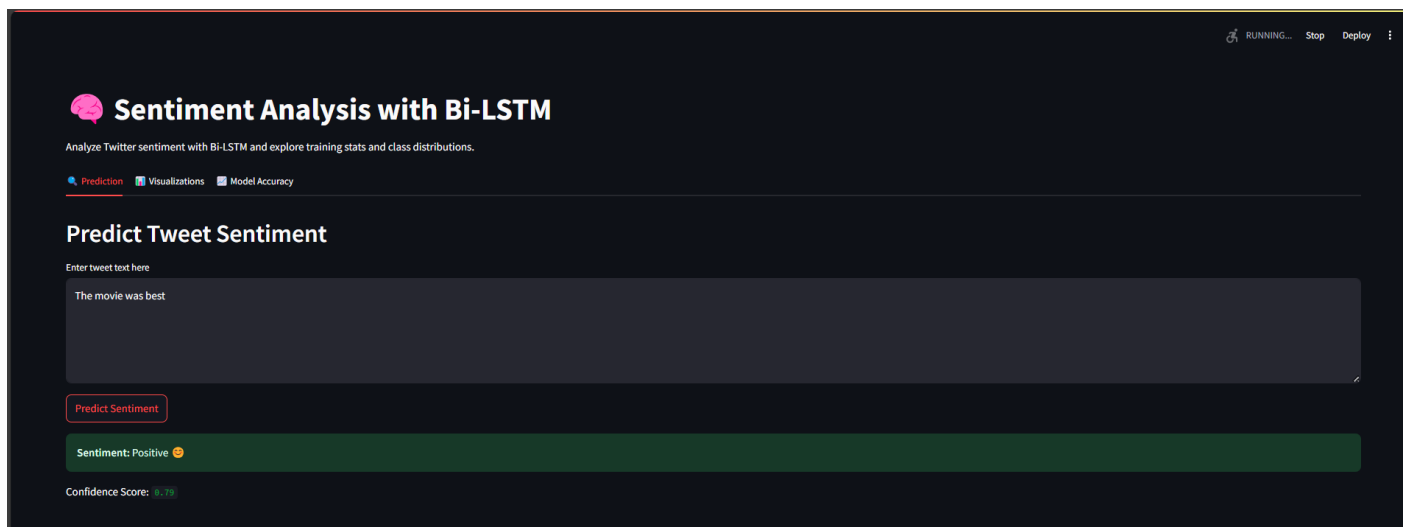
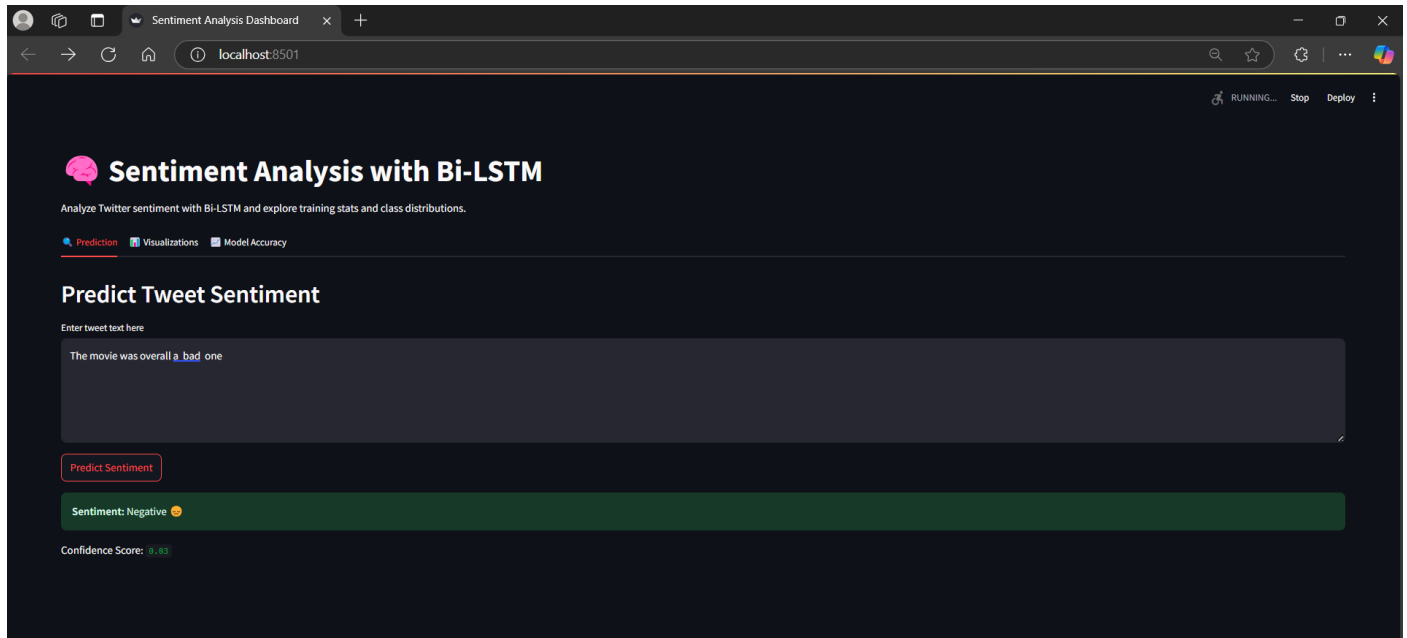
In order to increase the complexity I used Bi LSTM and for RNN I added dropout layer and much more hyperparameter tuning and the accuracy was increased

## Analysis

- **Performance:** The LSTM model generalizes well on unseen data.
- **Dataset Influence:** Large and clean dataset enabled effective model training.
- **Baseline Comparison:** Performance outperformed a simple RNN architecture.

# Deployment and demonstration

For deployment I have used streamlit and the Bi lstm model is downloaded including the tokenizer pkl file which will be used for prediction of a tweet whether it will be positive or negative. This UI also consists of visualisation and comparison with the model graph



Analyze Twitter sentiment with Bi-LSTM and explore training stats and class distributions.

Prediction Visualizations Model Accuracy

RNN vs Bi-LSTM Accuracy

11

