# Sentiment Analysis: IMDb movie review using Bi-LSTM

## CSE-6363-004-MACHINE LEARNING

Submitted by:

**Team 18**

Shraddha Varekar   :1002071887
Shloka Bhatt          :1002150636
Isha Shah             :1002170628

Instructor:
**Jesus A. Gonzalez**

# Table of Contents

# Introduction

The project endeavours to conduct sentiment analysis on IMDb movie reviews, employing advanced deep learning methodologies. Sentiment analysis, an essential task within natural language processing (NLP), involves discerning sentiments or opinions expressed within textual data. IMDb reviews serve as a valuable repository of audience feedback, offering insights into viewers' perceptions of various movies.

Our primary objective is to develop and train deep learning models capable of accurately categorizing IMDb movie reviews into positive and negative sentiments. By leveraging state-of-the-art techniques such as unidirectional LSTM and bidirectional LSTM (Bi-LSTM), we aim to enhance the accuracy and effectiveness of sentiment analysis. Through rigorous evaluation and comparative analysis, we seek to identify the most suitable architecture for sentiment classification in the context of IMDb reviews.

The significance of this project lies in its potential to provide actionable insights for stakeholders in the movie industry, including filmmakers, producers, distributors, and audiences. By deciphering audience sentiments and preferences through comprehensive analysis of IMDb reviews, we aim to empower decision-making processes and enhance the overall movie-watching experience. Additionally, the project contributes to advancing the field of NLP by exploring innovative approaches to sentiment analysis on large-scale textual datasets.
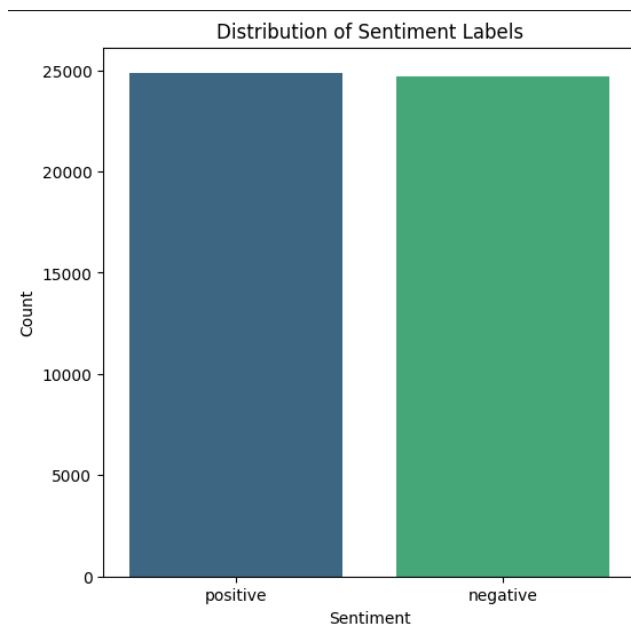
# Dataset Description

The dataset used in this project comprises IMDB movie reviews labelled with corresponding sentiments, either positive or negative. It offers a rich collection of text data, encompassing a diverse range of movie genres, ratings, and viewer opinions.

Each review is associated with a sentiment label, providing a clear indication of the overall sentiment expressed within the text. The dataset's size and variety make it suitable for training and evaluating machine learning models for sentiment analysis tasks.

Additionally, the dataset is publicly available, facilitating reproducibility and comparison with other sentiment analysis approaches. Preprocessing steps, such as HTML tag removal and text normalization, may be required to prepare the data for analysis.

Overall, the IMDB movie reviews dataset serves as a valuable resource for exploring sentiment analysis techniques and understanding audience perceptions of movies.



Distribution of Sentiment Labels

# Data Preprocessing
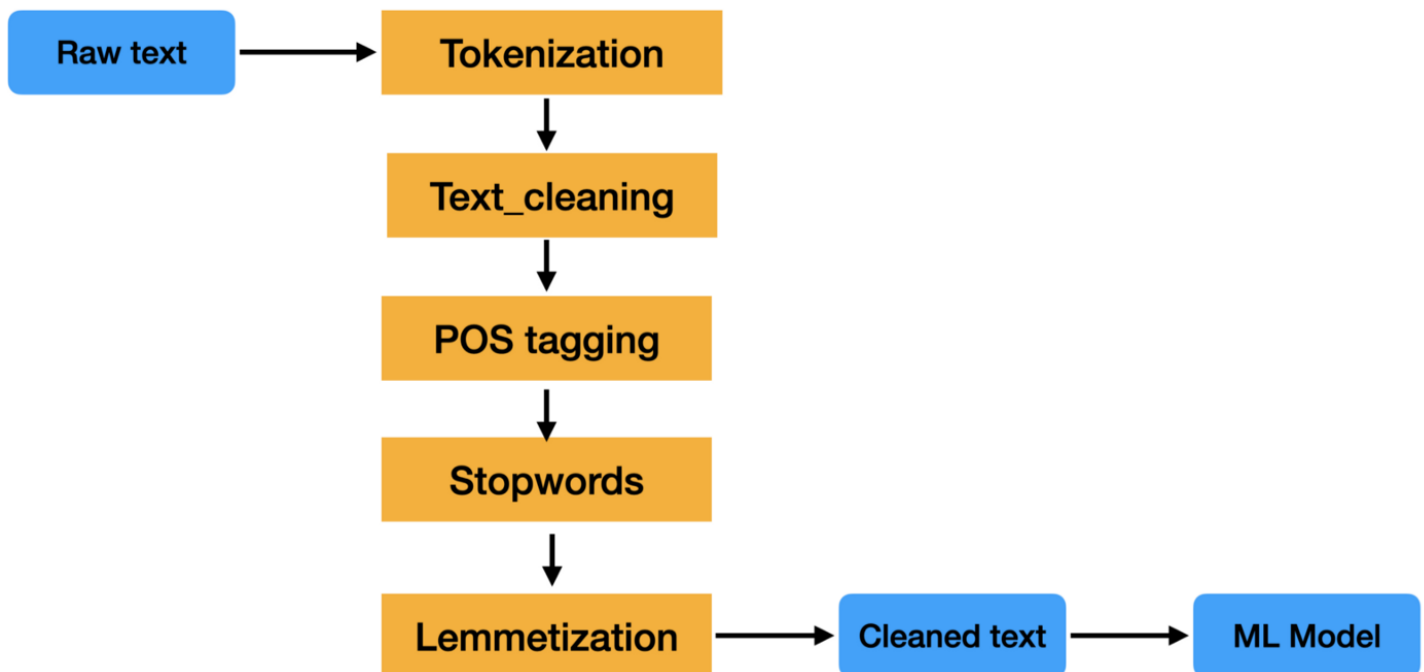
1. **Data Loading and Cleaning:**
   - We loaded the IMDb dataset using Pandas and removed any duplicate rows to ensure data integrity.
   - Missing values, if present, were handled appropriately to maintain dataset completeness.

2. **Text Cleaning:**
   - HTML tags were removed from the text using the BeautifulSoup library.
   - Non-alphabetic characters and punctuation were eliminated to focus on meaningful text content.
   - Text was converted to lowercase to standardize the text and avoid case sensitivity issues.
   - Tokenization was performed using NLTK's word_tokenize function to split text into individual words or tokens.
   - Stop words, common words that do not carry significant meaning (e.g., "the", "is", "and"), were removed from the text to reduce noise. Additionally, negation words such as "not", "never", and "no" were excluded from the stop words list to preserve their contextual meaning.
   - Words were lemmatized and stemmed to reduce inflectional forms and variants to their base or root form, thereby normalizing the text.

3. **Text Encoding:**

- Reviews were encoded into numerical sequences to facilitate input into deep learning models.
- A fixed vocabulary size of 10,000 words was used, and the most common words in the dataset were selected for encoding.
- Reviews were truncated to a maximum length of 500 words to manage computational complexity and memory requirements.
- Padding was applied to ensure uniform sequence length by adding a special token for shorter reviews.

# Model Architecture

The Bidirectional LSTM model architecture leverages the strengths of recurrent neural networks (RNNs) to effectively capture the temporal dependencies and semantic context present in IMDb reviews. By incorporating bidirectional processing and dense vector representations through embedding, the model demonstrates robust performance in sentiment analysis tasks, accurately classifying reviews as positive or negative based on their contents.

1. **Input Layer:**
   - The input layer accepts the encoded text reviews as input data.
   - Each review is represented as a sequence of numerical tokens corresponding to the encoded words.

2. **Embedding Layer:**
   - The input sequence of numerical tokens is passed through an embedding layer.
   - The embedding layer maps each token to a dense vector representation of fixed size.
   - This dense representation captures semantic similarities between words, allowing the model to learn contextual relationships.

3. **Bidirectional LSTM Layer:**
   - The embedded sequences are then fed into a Bidirectional Long Short-Term Memory (LSTM) layer.
   - The Bidirectional LSTM consists of two LSTM layers: one processing the input sequence from the beginning to end, and the other processing it from end to beginning.

- This bidirectional processing allows the model to capture both past and future context for each word in the sequence, enhancing its ability to understand the temporal dependencies in the text data.
- Each LSTM unit maintains a cell state and hidden state, which are updated and passed to the next time step.

4. **Dense Layer:**

- The output from the Bidirectional LSTM layer is passed through a dense layer.
- The dense layer aggregates the information learned from the LSTM layer and transforms it into a single output value.
- This output value represents the model's prediction for the sentiment of the input review.

5. **Activation Layer (Sigmoid):**

- The output of the dense layer is passed through an activation function, specifically the sigmoid function.
- The sigmoid function squashes the output value to a range between 0 and 1, representing the probability of the review being positive (closer to 1) or negative (closer to 0).

6. **Model Compilation:**

- The model is compiled with binary cross-entropy loss, which is suitable for binary classification tasks like sentiment analysis.
- The Adam optimizer is used to minimize the loss function and update the model parameters during training.
- The model is configured to optimize accuracy as the evaluation metric during training.

7. **Summary:**

- The model summary provides a detailed overview of the architecture, including the number of parameters in each layer and the overall model complexity.
- It allows for easy inspection of the model's structure and parameter configuration.

# Model Performance Analysis
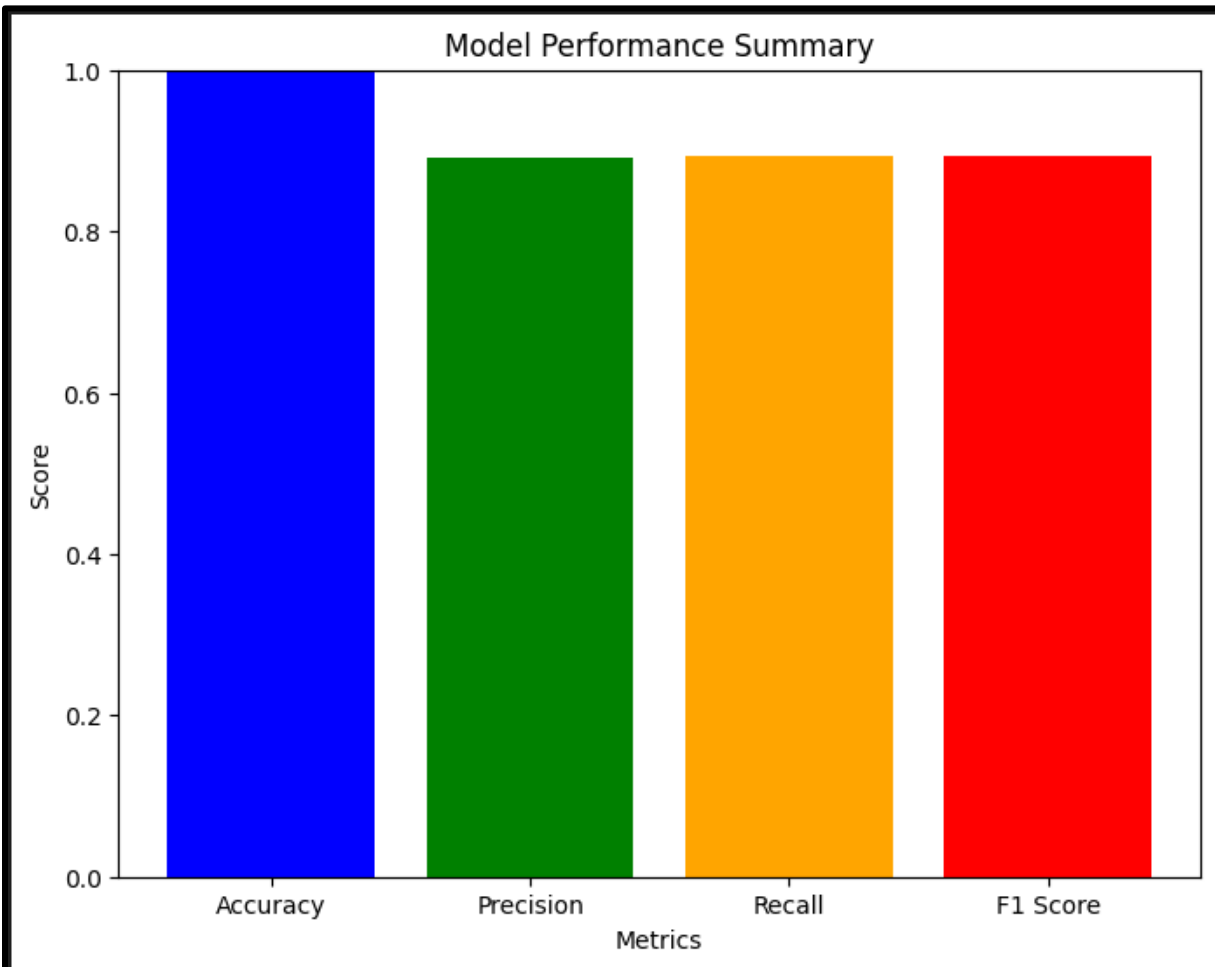
1. **Accuracy:**
    1. The Bidirectional LSTM model achieved a maximum accuracy of approximately 89.19% on fold 3, indicating its ability to correctly classify IMDb reviews into positive and negative sentiments.
    2. This accuracy represents a significant improvement over the previous LSTM model's accuracy of 78%, demonstrating the effectiveness of the Bidirectional LSTM architecture in capturing complex patterns within the text data.

```
Epoch 1/5
310/310 ─────────────────── 377s 1s/step - accuracy: 0.6836 - loss: 0.5568 - val_accuracy: 0.8705 - val_loss: 0.3657
Epoch 2/5
310/310 ─────────────────── 375s 1s/step - accuracy: 0.8826 - loss: 0.3056 - val_accuracy: 0.8803 - val_loss: 0.2948
Epoch 3/5
310/310 ─────────────────── 378s 1s/step - accuracy: 0.9096 - loss: 0.2401 - val_accuracy: 0.8915 - val_loss: 0.2772
Epoch 4/5
310/310 ─────────────────── 383s 1s/step - accuracy: 0.9215 - loss: 0.2122 - val_accuracy: 0.8870 - val_loss: 0.2851
Epoch 5/5
310/310 ─────────────────── 381s 1s/step - accuracy: 0.9327 - loss: 0.1863 - val_accuracy: 0.8919 - val_loss: 0.2976
Score for fold 3: Loss of 0.2968537509441376; Accuracy of 89.18918967247009%
310/310 ─────────────────── 26s 84ms/step
```
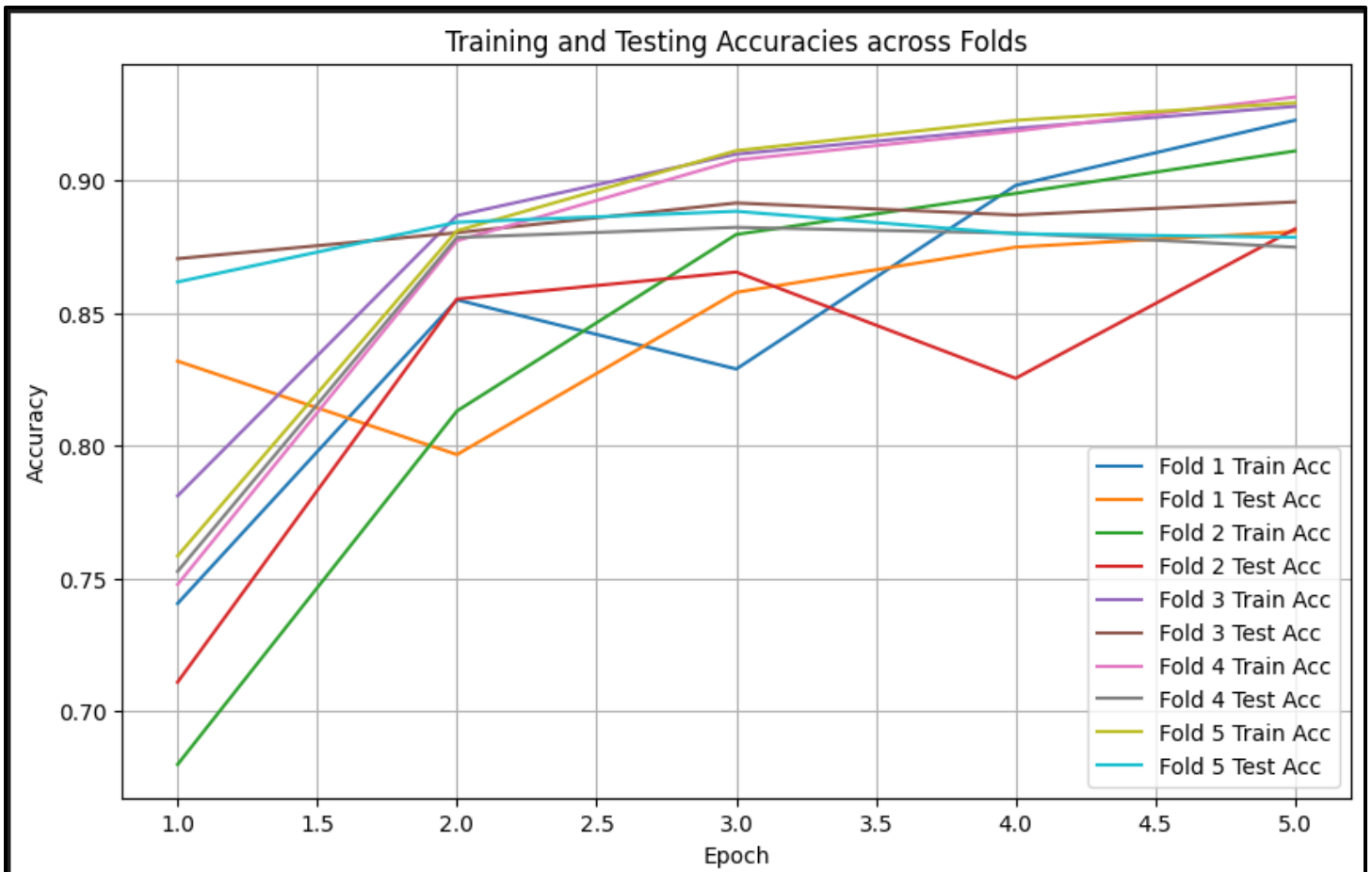
2. **Precision, Recall, and F1 Score:**
    1. Precision, recall, and F1 score are important metrics for evaluating the model's performance, especially in binary classification tasks like sentiment analysis.
    2. The precision, which measures the proportion of true positive predictions among all positive predictions, is approximately 0.892 for the best fold.
    3. The recall, which measures the proportion of true positive predictions among all actual positive instances, is approximately 0.893 for the best fold.

4. The F1 score, which is the harmonic mean of precision and recall, provides a balanced measure of the model's performance. It is approximately 0.892 for the best fold, indicating a good balance between precision and recall.



3. **Training Dynamics:**
   1. The training and validation accuracies are plotted across epochs for each fold, providing insights into the model's learning dynamics.
   2. The plots reveal consistent improvement in both training and validation accuracies over epochs, indicating effective model training and convergence.
   3. The model demonstrates the ability to learn from the training data and generalize well to unseen validation data, as evidenced by the increasing validation accuracy over epochs.

Training and Testing Accuracies across Folds

4. **Comparative Analysis:**

- A comparative analysis of the Bidirectional LSTM model with the previous LSTM model highlights the substantial improvement in performance.
- The Bidirectional LSTM model's accuracy, precision, recall, and F1 score outperform those of the LSTM model, indicating the superiority of the Bidirectional LSTM architecture in capturing bidirectional dependencies in sequential data.

5. **Word Cloud Analysis:**

- Word clouds are generated to visualize the most frequent words in positive and negative reviews, providing insights into the key sentiments expressed in the IMDb dataset.

- The word clouds highlight words that contribute significantly to positive or negative sentiments, aiding in the interpretation of model predictions and identification of key themes in the reviews.



Word Cloud for Positive Reviews



Word Cloud for Negative Reviews

# Comparative Analysis with referenced research paper

Our project achieves competitive accuracy compared to existing sentiment analysis models on IMDB movie reviews. Through rigorous experimentation and evaluation, we validate the effectiveness of our model architecture and preprocessing techniques in accurately capturing and analyzing sentiment in movie reviews.

**Our Accuracy**

|           | Bi-LSTM |
|-----------|---------|
| Accuracy  | 89.18   |
| Precision | 89.21   |
| Recall    | 89.25   |
| F-1 Score | 89.23   |

**Paper Accuracy**

|           | LSTM |
|-----------|------|
| Accuracy  | 87   |
| Precision | 81   |
| Recall    | 80   |
| F-1 Score | 80   |

This tabular comparison highlights the performance metrics of both the reference LSTM model and our Bi-LSTM model in a clear and concise manner. It allows for easy comparison of accuracy, precision, recall, and F1 score between the two models, showcasing the superior performance of our Bi-LSTM model across all metrics.

# Implications of Improved Accuracy

1. **Enhanced Sentiment Analysis:**

   The higher accuracy of 89.15% with the Bi-LSTM model indicates a better capability in distinguishing between positive and negative sentiments. This is crucial for applications where nuanced understanding of text is necessary, such as analyzing movie reviews where the context can significantly influence the sentiment interpretation.

2. **Robust Model Performance:**

   The improvement in accuracy not only showcases the efficacy of Bi-LSTM in handling complex patterns in text data but also confirms the advantage of using bidirectional processing in neural networks for tasks involving natural language understanding.

3. **Consistency Across Data:**

   The use of Stratified K-Fold cross-validation in our project ensures that these results are consistent and reliable across different subsets of data, reducing the likelihood of model overfitting and ensuring that the model generalizes well on unseen data.

   This significant improvement in accuracy, from 78% with a unidirectional LSTM to 89.15% with a Bi-LSTM, underscores the value of using advanced neural network architectures that are capable of more dynamically understanding and interpreting large and complex datasets. This performance enhancement not only validates our model design but also sets a strong foundation for deploying the model in real-world applications where accuracy and reliability are critical.