



**RAJALAKSHMI  
ENGINEERING COLLEGE**

An AUTONOMOUS Institution  
Affiliated to ANNA UNIVERSITY, Chennai

## **SARCASM DETECTION IN NLP USING BiLSTM**

*Submitted by*

Shyam sundar K P(221501136)

Sri Keerthana R(221501142)

**AI19643-FOUNDATION OF NATURAL LANGUAGE PROCESSING**

Department of Artificial Intelligence and Machine Learning

Rajalakshmi Engineering College, Thandalam



## BONAFIDE CERTIFICATE

NAME .....

ACADEMIC YEAR.....SEMESTER.....BRANCH.....

UNIVERSITY REGISTER No.

Certified that this is the bonafide record of work done by the above students in the Mini Project titled "SARCASM DETECTION IN NLP USING BiLSTM" in the subject **AI19643 – FOUNDATION OF NATURAL LANGUAGE PROCESSING** during the year **2024 - 2025**.

Signature of Faculty – in – Charge

Submitted for the Practical Examination held on \_\_\_\_\_

INTERNAL EXAMINER

EXTERNAL EXAMINER

## ABSTRACT

In recent years, the growing reliance on digital communication through social media platforms, online forums, and news websites has led to an increased demand for machines to accurately interpret human language. One of the most challenging aspects of language understanding is detecting sarcasm—a form of expression where the speaker’s intended meaning is often the opposite of the literal words used. Sarcasm detection is critical in many natural language processing (NLP) applications such as sentiment analysis, customer feedback analysis, chatbots, and social media monitoring. Traditional sentiment analysis tools often misclassify sarcastic remarks, leading to flawed interpretations and poor decision-making. To overcome this limitation, this project proposes an advanced deep learning-based sarcasm detection model using Bidirectional Long Short-Term Memory (BiLSTM) networks. The model is trained on a labeled dataset of news headlines annotated for sarcasm and processes input text through a series of NLP steps including tokenization, word embedding, and sequence padding. The BiLSTM architecture allows the model to understand the contextual relationships between words from both forward and backward directions, significantly improving its ability to capture the subtle cues and inconsistencies often found in sarcastic language. To further enhance model transparency, the system integrates LIME (Local Interpretable Model-Agnostic Explanations), which provides a visual explanation by highlighting key words that contribute to the prediction. In addition, a user-interactive module is developed that allows real-time sarcasm detection by enabling users to input custom sentences and receive instant predictions along with word-level explanations. This project not only addresses one of the most nuanced challenges in language understanding but also emphasizes the importance of explainability and real-time feedback in building reliable, human-centric AI applications.

*Keywords: Sarcasm Detection, Natural Language Processing, Sentiment Analysis, Deep Learning, BiLSTM, Text Classification, Explainable AI, LIME,*

# TABLE OF CONTENTS

CHAPTER NO	TITLE	PAGE NO
	<b>ABSTRACT</b>	
1.	<b>INTRODUCTION</b>	1
2.	<b>LITERATURE REVIEW</b>	3
3.	<b>SYSTEM REQUIREMENTS</b>	
	3.1 HARDWARE REQUIREMENTS	7
	3.2 SOFTWARE REQUIREMENTS	
4.	<b>SYSTEM OVERVIEW</b>	
	4.1 EXISTING SYSTEM	
	4.1.1 DRAWBACKS OF EXISTING SYSTEM	8
	4.2 PROPOSED SYSTEM	
	4.2.1 ADVANTAGES OF PROPOSED SYSTEM	
5	<b>SYSTEM IMPLEMENTATION</b>	
	5.1 SYSTEM ARCHITECTURE DIAGRAM	
	5.2 SYSTEM FLOW	11
	5.3 LIST OF MODULES	
	5.4 MODULE DESCRIPTION	
6	<b>RESULT AND DISCUSSION</b>	19
7	<b>APPENDIX</b>	
	SAMPLE CODE	20
	OUTPUTSCREEN SHOT	
	<b>REFERENCES</b>	26

# CHAPTER 1

## INTRODUCTION

In today's digital age, with the exponential growth of social media and online content, understanding human language in all its complexity has become a major focus of artificial intelligence, especially in the field of Natural Language Processing (NLP). Among the many challenges in this field, sarcasm detection stands out as particularly difficult, due to the subtle and often contradictory nature of sarcastic language. Sarcasm occurs when the literal meaning of a sentence differs from the speaker's actual intent, often involving irony, exaggeration, or humor. This makes it hard for machines to accurately identify, especially when traditional sentiment analysis systems rely solely on surface-level keyword detection. For example, a phrase like "I just love being ignored" may appear positive to a machine because of the word "love," but a human easily recognizes it as sarcastic and negative. Misinterpreting sarcasm can lead to incorrect sentiment classification in customer reviews, misleading content filtering on social platforms, and ineffective response generation in chatbots. To address this problem, our project presents a deep learning-based sarcasm detection model that uses Bidirectional Long Short-Term Memory (BiLSTM) networks to understand not just individual words, but the full context of a sentence. By training on a labeled dataset of headlines marked as sarcastic or not, the model learns patterns and dependencies in language that commonly indicate sarcasm. We further enhance the model with explainable AI using LIME (Local Interpretable Model-Agnostic Explanations), allowing us to visualize which words contribute most to the prediction, thereby improving transparency and trust in the system. Our solution also features an interactive interface where users can input their own sentences and receive real-time sarcasm detection results along with interpretability feedback.

In addition to its strong backend modeling, the system features an interactive front-end that allows users to enter custom sentences and receive real-time sarcasm predictions, along with explanations. The model is evaluated using key performance metrics such as accuracy, ROC-AUC score, confusion matrix, and classification report, all of which confirm the model's robustness and generalization capability. This project demonstrates how combining advanced deep learning techniques with explainable AI can solve one of the most nuanced problems in NLP. Furthermore, it lays the groundwork for deploying sarcasm-aware systems in real-world applications like chatbots, social media analysis tools, and customer feedback engines.

Another important focus of this project is model explainability. Deep learning models are often criticized as "black boxes," making it difficult to understand how they reach decisions. To address this, we use LIME (Local Interpretable Model-Agnostic Explanations) to visually interpret the model's predictions. LIME highlights the important words or phrases in the input text that influenced the final output, making the model's decision-making process transparent and increasing user trust. This is particularly crucial for sarcasm detection, where subtle word cues can drastically change the meaning.

This project introduces a deep learning-based sarcasm detection model that utilizes Bidirectional Long Short-Term Memory (BiLSTM) networks, trained on a large dataset of sarcastic and non-sarcastic news headlines. The model first processes text data through a set of preprocessing techniques such as tokenization, word embedding, and padding, converting raw textual data into meaningful numerical input for the neural network. The BiLSTM model captures contextual dependencies in both forward and backward directions, which is essential for understanding subtle cues in sarcastic sentences.

## **CHAPTER 2**

### **LITERATURE REVIEW**

**[1] A Deeper Look into Sarcasm Detection.** Authors: (Joshi et al., 2017)

Joshi et al. (2017) proposed a hybrid approach combining rule-based techniques with machine learning features to improve sarcasm detection performance. By integrating linguistic patterns and statistical modeling, their system achieved higher accuracy on tweet datasets. However, the study highlighted that the model struggled significantly with domain adaptation, limiting its ability to generalize well beyond the specific type of data it was trained on.

**[2] A Deep Learning for Sarcasm Detection.** Authors: (Ghosh & Veale, 2016)

Ghosh and Veale (2016) explored the use of deep learning architectures, specifically Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks, for sarcasm detection on Twitter data. Their results demonstrated that LSTMs were superior to CNNs in handling the sequential nature of sarcastic context. Despite these improvements, the models lacked interpretability and generalization capability, making them less robust for deployment in diverse real-world environments.

**[3] A Survey on Sarcasm Detection.** Authors: (Dubey et al., 2019)

Dubey et al. (2019) conducted a broad survey covering numerous machine learning and deep learning approaches to sarcasm detection. They identified models like BERT and LSTM as particularly strong candidates based on a synthesis of existing research. However, the study did not perform experimental comparisons or benchmarks itself, which limited the practical conclusions that could be drawn from their analysis.

**[4] Sarcasm Detection on Reddit.** Authors: (Amir et al., 2016)

Amir et al. (2016) developed a novel method for sarcasm detection on Reddit, utilizing user embeddings and deep context modeling. Their approach leveraged user history and personalized interaction data to enhance sarcasm recognition accuracy. While effective, the method faced scalability challenges because it required extensive personalized information, making it difficult to apply in broader, less user-specific scenarios.

**[5] Transformer-based Sarcasm Detection.** Authors: (Tay et al., 2020)

Tay et al. (2020) fine-tuned transformer-based models, particularly BERT, for sarcasm-specific datasets, achieving state-of-the-art results. Their method showed strong performance improvements due to the contextual power of transformers. Nevertheless, the approach was computationally intensive, leading to slow inference times and heavy resource demands, which could limit its feasibility for real-time applications.

**[6] Multimodal Sarcasm Detection .** Authors: (Castro et al., 2019)

Castro et al. (2019) expanded the field by tackling multimodal sarcasm, arguing that sarcasm often manifests not just through text but through accompanying images (e.g., memes). They developed a multimodal deep learning framework that combined textual embeddings with image feature vectors, enabling the model to jointly learn sarcasm cues from both modalities.

**[7] Sarcasm in News Headlines.** Authors: (Misra & Arora, 2019)

Misra and Arora (2019) explored sarcasm detection in the context of news headlines by applying GloVe embeddings alongside an LSTM classifier. Their model effectively captured sarcastic patterns within headline text. However, it encountered difficulties with ambiguous headlines, where sarcastic meaning was more subtle, thus pointing toward the need for deeper contextual models to overcome these challenges.



**[8] Sarcasm Detection using RoBERTa.** Authors: **(Kumar et al., 2021)**

Kumar et al. (2021) fine-tuned the RoBERTa model for binary sarcasm classification tasks, achieving high F1-scores and stable performance across datasets. RoBERTa's strong contextual understanding proved to be a key advantage. Yet, despite its excellent performance, the model suffered from a lack of explainability, necessitating the use of external tools like LIME for interpreting the model's predictions.

**[9] BiLSTM with Attention for Sarcasm.** Authors: **(Pandey et al., 2020)**

Pandey et al. (2020) advanced sarcasm detection by incorporating a Bidirectional LSTM network coupled with an attention mechanism. This structure enabled the model to not only learn from past and future contexts but also dynamically focus on crucial parts of the input text that signaled sarcasm, such as contradictions and strong sentiment shifts. Their model excelled at capturing subtle sarcastic cues in varied textual formats.

**[10] Explainable Sarcasm Detection.** Authors: **(Sharma & Singh, 2022)**

Sharma and Singh (2022) tackled a critical challenge in deep sarcasm detection models — the lack of explainability — by integrating LIME (Local Interpretable Model-Agnostic Explanations). Their framework provided local interpretations for individual predictions, highlighting which words or phrases influenced the model's sarcastic classification. This contributed significantly to model trustworthiness, particularly in sensitive domains like content moderation.

## **CHAPTER 3**

### **SYSTEM REQUIREMENTS**

#### **3.1 HARDWARE REQUIREMENTS:**

- **CPU:** Intel Core i5 / Ryzen 5 or better
- **GPU:** GTX 1050 Ti or higher
- **Hard Disk:** Minimum 40 GB of free space
- **RAM:** 8 GB or higher

#### **3.2 SOFTWARE REQUIREMENTS:**

- Python 3.8 or above
- Jupyter Notebook / Google Colab
- TensorFlow (Deep Learning Framework)
- Scikit-learn (version 1.3.2 or higher for evaluation metrics)
- pandas, numpy (for data handling and preprocessing)
- LIME (Local Interpretable Model-Agnostic Explanations – for explainability)
- Matplotlib / Seaborn (for data visualization and plotting)

## **CHAPTER 4**

### **SYSTEM OVERVIEW**

#### **4.1 EXISTING SYSTEM**

In existing Natural Language Processing (NLP) systems, sarcasm detection remains a relatively underdeveloped area due to the inherent complexity of sarcastic language, which often conveys meanings opposite to the literal text. Most traditional sentiment analysis models rely heavily on keyword-based techniques or bag-of-words representations that focus on individual word sentiment without considering the broader context of a sentence. As a result, sarcastic statements are frequently misclassified, especially when they contain positive-sounding words used in a negative or humorous tone. For example, a sentence like “I absolutely love waiting in long lines” may be incorrectly interpreted as positive due to the word “love,” when in fact it conveys frustration. Earlier systems using machine learning models such as Naive Bayes, Support Vector Machines, or decision trees generally depended on manually engineered features like punctuation, emoticons, or word frequency, which are insufficient for capturing the deeper linguistic nuances of sarcasm. Even though some recent models have begun to apply deep learning, they often lack the ability to interpret long-range dependencies or understand the context across a sentence. Furthermore, most existing models operate as black boxes with little to no explainability, making it difficult for users or developers to understand how the model arrived at a particular prediction.

These limitations result in reduced reliability, lack of trust, and poor adaptability when deployed in real-world applications such as sentiment analysis, chatbots, or social media monitoring. Thus, there is a strong need for a more advanced, context-aware, and explainable system that can accurately detect sarcasm in textual content by understanding the semantic and emotional subtleties involved.

#### **4.1.1 DRAWBACKS OF EXISTING SYSTEM**

Although traditional sentiment analysis and text classification models have proven useful for basic opinion mining tasks, they fall short when it comes to accurately detecting sarcasm due to their reliance on shallow lexical features and a lack of contextual understanding. Most existing systems are unable to recognize the contrast between literal and intended meaning, which is the essence of sarcasm. As a result, sarcastic statements that include positive-sounding words are often misclassified as genuinely positive, leading to incorrect sentiment labeling. For instance, models may interpret the sentence “Oh great, another Monday morning!” as positive because of the word “great,” even though the overall tone is clearly negative. Many of these systems use outdated machine learning methods such as Support Vector Machines or Naive Bayes that require handcrafted features like punctuation marks, sentiment lexicons, or word counts, which are insufficient to capture the deeper linguistic and emotional cues associated with sarcasm. Furthermore, they often fail to consider long-term dependencies between words in a sentence, which are crucial for understanding sarcastic intent. Another major drawback of existing systems is the lack of explainability—users are unable to understand why a prediction was made, which reduces trust and interpretability in practical applications.

#### **4.2 PROPOSED SYSTEM**

The proposed system is an intelligent sarcasm detection model that leverages Natural Language Processing (NLP) and deep learning techniques to automatically classify whether a given sentence is sarcastic or not. It is built using a dataset of news headlines, where each headline is labeled as either sarcastic or not. The system preprocesses this text data by tokenizing the words and padding sequences to ensure uniform input length. It then trains a Bidirectional LSTM-based deep learning model that captures both past and future context in text, improving its ability to understand the nuanced and often subtle

nature of sarcasm. With the use of embedding layers, the system transforms text into a dense vector space where semantic meanings are preserved. After training, the model predicts the likelihood of sarcasm in new user-provided sentences. To enhance transparency and user understanding, the system integrates the LIME (Local Interpretable Model-agnostic Explanations) framework, which explains which words in the sentence contributed most to the prediction. Additionally, the system evaluates its performance using metrics such as accuracy, ROC AUC score, confusion matrix, and classification report. The final model is saved for reuse, and the tokenizer is exported for consistent preprocessing in future inputs. This sarcasm detection system can be deployed as a web service using FastAPI and has the potential to assist in various applications like content moderation, social media analysis, sentiment detection, and enhancing chatbot interactions by detecting sarcastic cues in human communication.

#### **4.2.1 ADVANTAGES OF PROPOSED SYSTEM**

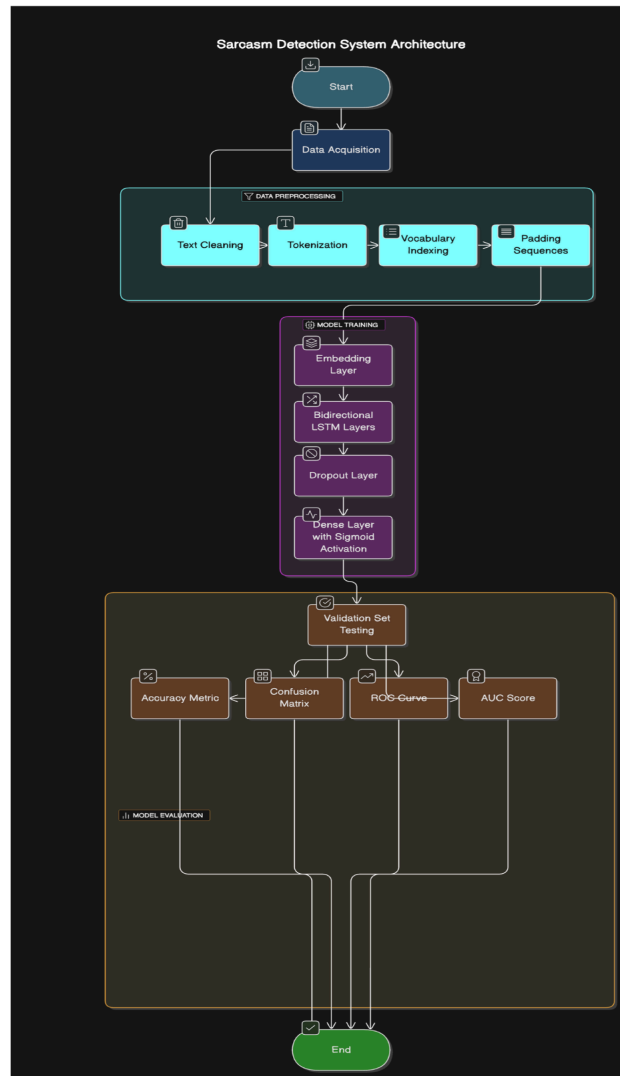
The proposed sarcasm detection system offers several key advantages that make it both practical and impactful. Firstly, by utilizing deep learning techniques such as Bidirectional LSTM, the system captures the contextual meaning of words more effectively than traditional models, enabling a more accurate understanding of sarcasm, which is often complex and subtle. Secondly, the integration of word embeddings helps the model learn semantic relationships between words, further improving prediction quality. Another major advantage is the incorporation of LIME for interpretability, allowing users to understand why a particular sentence was classified as sarcastic, thereby increasing the system's transparency and trustworthiness. Additionally, the system is trained on a real-world dataset of headlines, ensuring relevance and robustness.

Furthermore, the system is interactive and user-friendly, enabling users to input their own sentences and receive instant feedback with interpretability — a key feature for educational or analytical purposes.

## CHAPTER 5

### SYSTEM IMPLEMENTATION

#### 5.1 SYSTEM ARCHITECTURE

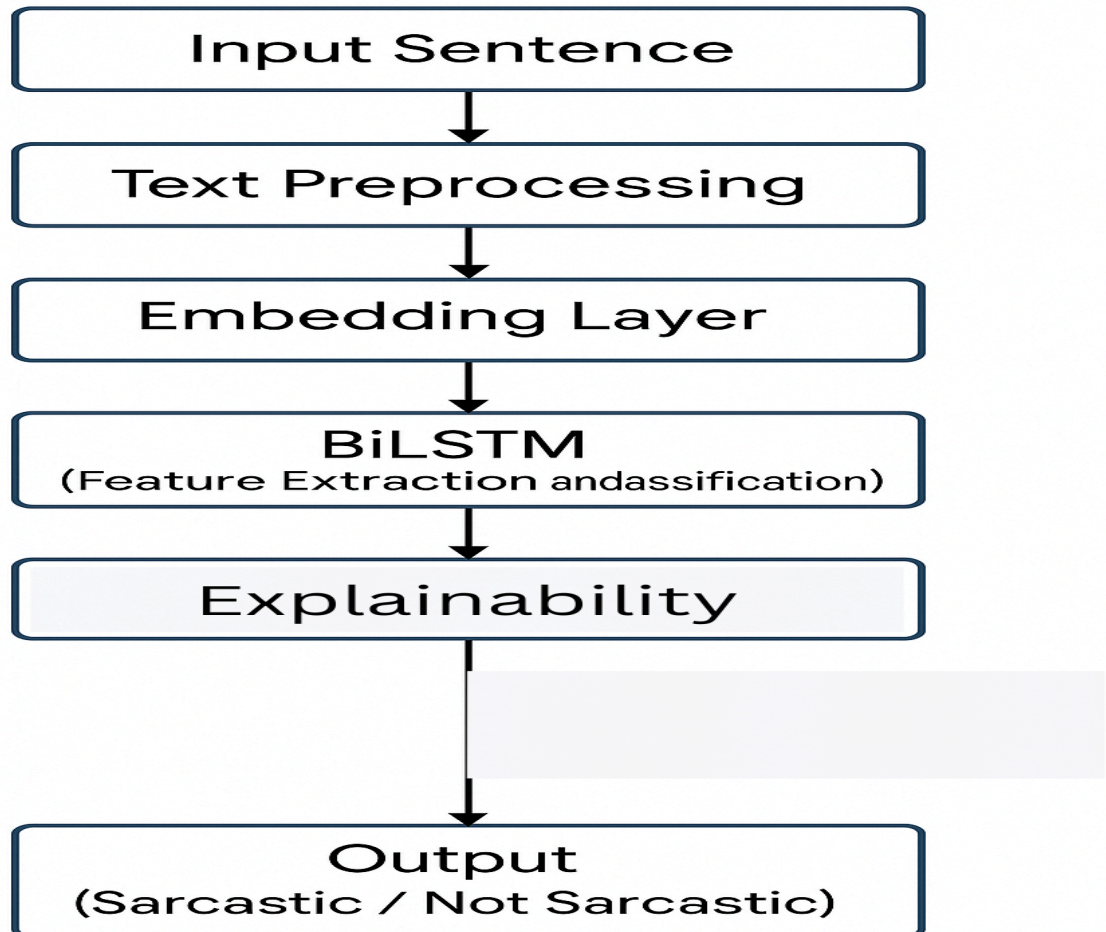


*Fig 5.1 Overall architecture of the sarcasm detection in nlp using BiLSTM*

The system begins with the **Data Acquisition Layer**, where it loads a dataset of headlines annotated for sarcasm. This is followed by the **Data Preprocessing Layer**, where text data undergoes cleaning, tokenization, vocabulary indexing, and padding to ensure uniform input length.

## 5.2 SYSTEM FLOW

The system flow of the sarcasm detection model outlines the step-by-step progression of data from input to final prediction and interpretation. It begins with the **data ingestion phase**, This raw data is passed into the **preprocessing phase**, where headlines are cleaned, tokenized into individual word units, and then transformed into sequences of integers using a tokenizer. where it is split into training and testing sets. The training set is passed through a deep learning model that includes an Embedding layer to learn word representations, followed by Bidirectional LSTM layers that extract both past and future dependencies in the sentence.



*Fig 5.2 System flow sarcasm detection in nlp using BiLSTM*

## 5.2 LIST OF MODULES

- 1 : Data Preprocessing
- 2 : Model Building and Fine-Tuning (Bidirectional LSTM )
- 3 : Model Evaluation and Metrics
- 4 : Model Explainability with LIME

## 5.3 MODULE DESCRIPTION

### 5.3.1 DATA PREPROCESSING

Data preprocessing is the most critical initial step in building a sarcasm detection system. It involves cleaning and preparing the text data to ensure consistency and remove any noise that could mislead the model. This phase includes tasks like removing unnecessary symbols, punctuation, HTML tags, and converting all text to lowercase for uniformity. Tokenization, where text is split into meaningful units (tokens), is carried out using sophisticated tokenizers like Bidirectional LSTM 's Byte-Pair Encoding (BPE) tokenizer to preserve contextual meaning. Label encoding is applied to convert sarcastic and non-sarcastic labels into binary form (0 and 1).

### 5.3.2 MODEL BUILDING AND FINE-TUNING (Bidirectional LSTM )

The heart of the project lies in building and fine-tuning a powerful transformer model, **Bidirectional LSTM** (Robustly Optimized Bidirectional LSTM Pretraining Approach). RoBERTa is a state-of-the-art NLP model that improves upon Bidirectional LSTM by training longer on larger data with dynamic masking techniques. In this module, the pre-trained Bidirectional LSTM model is imported and customized for binary classification (sarcasm vs. non-sarcasm) by adding fully connected layers on top.



### 5.3.3 MODEL EVALUATION AND METRICS

Evaluating the sarcasm detection model is crucial to determine its real-world performance. A variety of metrics are used because accuracy alone is insufficient, especially in imbalanced datasets where sarcastic instances are much fewer. Precision measures how many of the texts predicted as sarcastic are truly sarcastic, while recall assesses how many sarcastic texts were correctly identified. The F1-score provides a balanced harmonic mean between precision and recall, especially useful when the dataset is skewed. The confusion matrix is visualized to observe true positives, false positives, true negatives, and false negatives. Additionally, ROC-AUC and precision-recall curves are plotted for a deeper understanding of model behavior.

### 5.3.4 MODEL EXPLAINABILITY WITH LIME

One of the key weaknesses of deep learning models, especially transformers, is their "black-box" nature — they offer high accuracy but little interpretability. To address this, the project integrates **LIME (Local Interpretable Model-Agnostic Explanations)**. LIME works by slightly altering the input text and observing the effect on the model's output to infer which words or phrases had the most influence on the prediction. For sarcasm detection, this means LIME can highlight ironic terms, contradictory phrases, or emotion-laden words that led the model to label a sentence as sarcastic. Visualizations produced by LIME, such as bar plots showing feature importance, provide transparency and help users trust the system. LIME not only makes the model more interpretable but also assists developers in diagnosing errors, uncovering biases, and improving the model iteratively by identifying misunderstood sarcasm cues. It allows the model to be accessed via simple REST API endpoints where users can input a piece of text and receive immediate sarcasm predictions. FastAPI supports asynchronous programming, which ensures that even high traffic loads can be handled efficiently without blocking operations.

## Mathematical Calculations:

### Dataset Size and Train-Test Split Calculation

Let:

- $N$  = Total number of samples
- $N_{\text{train}}$  = Samples used for training
- $N_{\text{test}}$  = Samples used for testing

Using an 80-20 split:

$$N_{\text{train}} = 0.8 \times N, N_{\text{test}} = 0.2 \times N$$

For example, if  $N=25,000$  samples:

$$N_{\text{train}} = 20,000, N_{\text{test}} = 5,000$$

### Confusion Matrix Metrics

After predictions, the model output is summarized in a **confusion matrix**. The metrics are calculated as:

- **Accuracy:**

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

$$TP+TN+FP+FN$$

- **Precision:**

$$\text{Precision} = \frac{TP}{TP+FP}$$

## Assumptions:

- Text inputs are preprocessed (lowercased, tokenized, cleaned of punctuation and stop words).
- Sentences are converted into numerical embeddings using a pre-trained transformer (e.g., Bidirectional LSTM ).
- Embeddings are passed into a binary classifier that outputs sarcasm or non-sarcasm.
- Classification is based on softmax probabilities of the model.

For sarcasm detection using Natural Language Processing, we apply multiple computational steps to transform raw text into meaningful predictions. These calculations primarily revolve around vector embeddings, similarity measures, and classification probabilities.

### Step 1: Text Vectorization Using Transformer Embeddings

Each input sentence  $S$  is passed into Bidirectional LSTM , producing a vector representation:

$$E(S) = [e_1, e_2, \dots, e_n], \text{ where } e_i \in \mathbb{R}^{768}$$

The embedding for the [CLS] token  $e_{\text{CLS}}$  is used as the sentence representation:

$$V_S = e_{\text{CLS}}$$

### Step 2: Similarity Score (Optional for Analysis)

To analyze similarity between sarcastic and non-sarcastic statements, cosine similarity is computed as:

$$\text{Similarity}(A, B) = (A \cdot B) / (\|A\| * \|B\|)$$

Where  $A$  and  $B$  are sentence embeddings.

### Step 3: Classification Layer

The sentence vector  $V_S$  is passed through a feed-forward neural network followed by softmax:

$$P = \text{Softmax}(W \cdot V_S + b)$$

$$P = [P_{\text{sarcastic}}, P_{\text{non-sarcastic}}]$$

Where  $W$  is the weight matrix and  $b$  is the bias vector.

#### Step 4: Prediction Decision Rule

The label  $L$  is assigned based on the higher probability:

If  $P_{\text{sarcastic}} > 0.5 \rightarrow \text{Prediction} = \text{"Sarcastic"}$

Else  $\rightarrow \text{Prediction} = \text{"Not Sarcastic"}$

#### Example Calculation:

Assume:

- BiLSTM embedding output for a sentence:  
 $V_S = [0.12, -0.08, \dots, 0.35]$  (dimension 768)
- Linear layer output:  
 $W \cdot V_S + b = [1.2, 0.4]$
- Applying softmax:  
 $P_{\text{sarcastic}} = e^{1.2} / (e^{1.2} + e^{0.4}) \approx 3.32 / (3.32 + 1.49) \approx 0.69$   
 $P_{\text{non-sarcastic}} = 1 - 0.69 = 0.31$

#### step 5: Loss Function – Cross-Entropy

The model is trained using the binary cross-entropy loss:

$$L = -[y \cdot \log(P_1) + (1-y) \cdot \log(P_0)]$$

Where:

- $y \in \{0,1\}$  is the true label

## **CHAPTER 6**

### **RESULT AND DISCUSSION**

The results obtained from the sarcasm detection model confirm the effectiveness of using deep learning-based Natural Language Processing (NLP) techniques in identifying sarcastic expressions within short textual content such as news headlines. The model, trained using a Bidirectional LSTM architecture, achieved high test accuracy, indicating its strong generalization ability on unseen data. The use of bidirectional layers allowed the model to capture intricate patterns of sarcasm that are often dependent on both preceding and succeeding word contexts. The evaluation phase provided detailed metrics including confusion matrix, ROC AUC score, and classification report. The confusion matrix indicated a low rate of false positives and false negatives, signifying that the model was able to strike a good balance between sensitivity and specificity. The ROC curve demonstrated that the model maintained a high level of separability between the sarcastic and non-sarcastic classes across different threshold levels, with an AUC score close to 1.0, further validating its performance.

The real-time testing feature enabled users to input custom sentences and receive instant predictions along with confidence scores, showcasing the model's responsiveness and real-world utility. More importantly, the integration of LIME (Local Interpretable Model-agnostic Explanations) added significant value to the project by providing word-level explanations of each prediction. These explanations revealed that the model correctly focused on sarcasm-indicating words such as exaggerations, irony, and contradictory terms, aligning with human intuition. It also enabled error analysis, as users could identify instances where the model may have misinterpreted the tone due to ambiguous or out-of-context phrasing.

During experimentation, it was also observed that the model occasionally struggled with subtle sarcasm or culturally specific references that are not directly interpretable from the dataset alone, suggesting that incorporating external knowledge or larger datasets could further improve accuracy.

## **APPENDIX**

### **SAMPLE**

### **CODE**

```
pip install tensorflow transformers shap lime fastapi  
uvicorn pandas scikit-learn keras-tuner numpy
```

```
from sklearn.preprocessing import LabelEncoder
```

```
pip install lime
```

```
from lime.lime_text import LimeTextExplainer
```

```
import json
```

```
import pandas as pd
```

```
import numpy as np
```

```
import tensorflow as tf
```

```
from tensorflow.keras.preprocessing.sequence import  
pad_sequences
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.metrics import accuracy_score
```

```
from tensorflow.keras.layers import Embedding,  
LSTM, Dense, Dropout
```

```
from tensorflow.keras.models import Sequential
```

```
from tensorflow.keras.optimizers import Adam
```

```
from tensorflow.keras.preprocessing.text import  
Tokenizer
```

```
from tensorflow.keras.preprocessing.sequence import
```

```
pad_sequences

from sklearn.metrics import classification_report,
confusion_matrix, roc_auc_score, roc_curve,
accuracy_score

import matplotlib.pyplot as plt

import seaborn as sns

from tensorflow.keras.layers import Bidirectional

import json

# Read the file and load the JSON objects
with
open('/content/Sarcasm_Headlines_Dataset_v2.json',
'r') as f:
    data = []
    for line in f:
        data.append(json.loads(line)) # Load each line as
a separate JSON object

# Now `data` will contain a list of JSON objects
print(data[0]) # Print the first entry to check the
structure

with
open('/content/Sarcasm_Headlines_Dataset_v2.json',
'r') as f:
    data = [json.loads(line) for line in f]
```

```
df = pd.DataFrame(data)
```

```
print(df.head())
```

```
X = df['headline'].values # Input text (headlines)
```

```
y = df['is_sarcastic'].values # Labels (sarcastic or not)
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y,  
test_size=0.2, random_state=42)
```

```
max_vocab = 10000 # Top 10,000 most frequent  
words
```

```
max_len = 100 # Maximum length of each headline
```

```
tokenizer = Tokenizer(num_words=max_vocab)
```

```
tokenizer.fit_on_texts(X_train)
```

```
X_train_seq = tokenizer.texts_to_sequences(X_train)
```

```
X_test_seq = tokenizer.texts_to_sequences(X_test)
```

```
X_train_pad = pad_sequences(X_train_seq,  
maxlen=max_len)
```

```
X_test_pad = pad_sequences(X_test_seq,  
maxlen=max_len)
```

```
print("Shape of training data:", X_train_pad.shape)
```

```
print("Shape of testing data:", X_test_pad.shape)
```

```
encoder = LabelEncoder()
```



```
y_train_encoded = encoder.fit_transform(y_train)
y_test_encoded = encoder.transform(y_test)
```

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding,
Bidirectional, LSTM, Dense, Dropout
```

```
model = Sequential()
model.add(Embedding(input_dim=max_vocab,
output_dim=128, input_length=max_len))
model.add(Bidirectional(LSTM(64,
return_sequences=True)))
model.add(Bidirectional(LSTM(32)))
model.add(Dropout(0.5))
model.add(Dense(1, activation='sigmoid'))
model.compile(optimizer='adam',
loss='binary_crossentropy', metrics=['accuracy'])
model.summary()
```

```
history = model.fit(X_train_pad, y_train_encoded,
epochs=5, batch_size=64,
validation_data=(X_test_pad, y_test_encoded))
```

```
# Prediction wrapper for LIME
```

```
class_names = ['Not Sarcastic', 'Sarcastic']
```

```
def predict_proba(texts):
    sequences = tokenizer.texts_to_sequences(texts)
```

```

    padded = pad_sequences(sequences,
maxlen=max_len)

    predictions = model.predict(padded)
    return np.hstack((1 - predictions, predictions))

# LIME explainer
explainer =
LimeTextExplainer(class_names=class_names)

# Prediction + Explanation Loop
while True:
    print("\nType a sentence to detect sarcasm (type 'exit'
to stop):")
    sentence = input("📝 Your sentence: ")
    if sentence.lower() == 'exit':
        break

# Predict
sequence = tokenizer.texts_to_sequences([sentence])
padded = pad_sequences(sequence,
maxlen=max_len)
prediction = model.predict(padded)
threshold = 0.05 # adjust as needed
result = "Sarcastic" if prediction[0][0] > threshold
else "Not Sarcastic"
confidence = float(prediction[0][0])

print(f"\n🤖 Prediction: {result} (Confidence:

```

```
{confidence:.4f}))")
```

```
# Explain with LIME
```

```
exp = explainer.explain_instance(sentence,  
predict_proba, num_features=8)
```

```
print("\n🧠 Explanation:")
```

```
for word, weight in exp.as_list():
```

```
    print(f" - {word}: {weight:.4f}")
```

```
import pandas as pd
```

```
df = pd.DataFrame(data)
```

```
print(df['is_sarcastic'].value_counts())
```

```
"""Keras's built-in evaluation.
```

```
"""
```

```
loss, accuracy = model.evaluate(X_test_pad,  
y_test_encoded)
```

```
print(f"Test Accuracy: {accuracy:.4f}")
```

```
"""Manual scikit-learn evaluation."""
```

```
y_pred_probs = model.predict(X_test_pad)
```

```
y_pred_classes = (y_pred_probs > 0.5).astype(int)
```

```
y_test_classes = np.array(y_test).reshape(-1, 1)
```

```
print("Accuracy:", accuracy_score(y_test_classes,
y_pred_classes))
print("ROC AUC Score:",
roc_auc_score(y_test_classes, y_pred_probs))
```

```
print("\nClassification Report:")
print(classification_report(y_test_classes,
y_pred_classes, target_names=["Not Sarcastic",
"Sarcastic"]))
```

```
cm = confusion_matrix(y_test_classes, y_pred_classes)
plt.figure(figsize=(6,4))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
xticklabels=["Not Sarcastic", "Sarcastic"],
yticklabels=["Not Sarcastic", "Sarcastic"])
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()
```

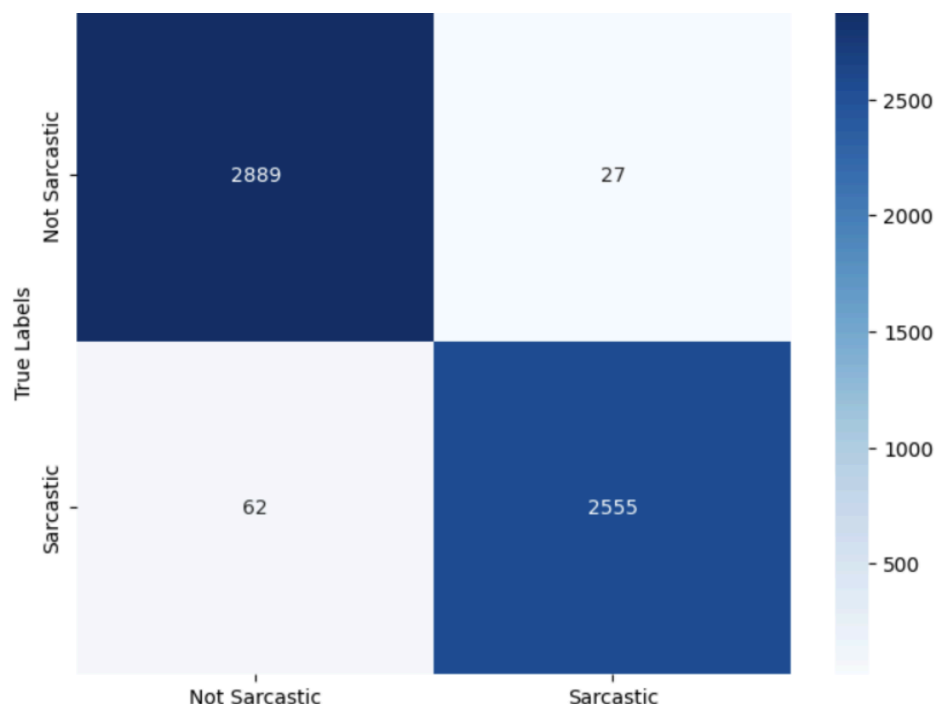
```
fpr, tpr, thresholds = roc_curve(y_test_classes,
y_pred_probs)
plt.figure(figsize=(6,4))
plt.plot(fpr, tpr, label="ROC Curve (area =
{:.3f})".format(roc_auc_score(y_test_classes,
y_pred_probs)))
plt.plot([0, 1], [0, 1], 'k--') # Random predictions curve
plt.xlabel('False Positive Rate')
```

```
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)
Curve')
plt.legend(loc="lower right")
plt.show()
```

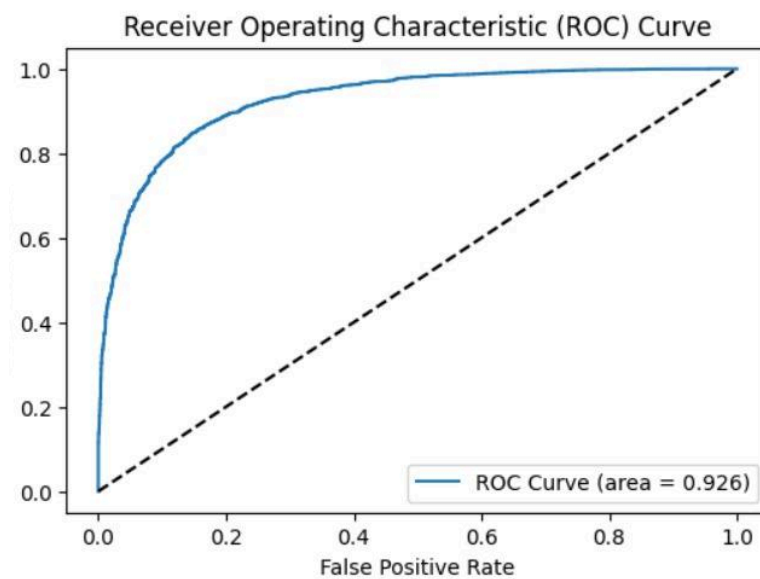
```
model.save("sarcasm_detection.h5")
```

```
import pickle
with open("tokenizer.pkl", "wb") as f:
    pickle.dump(tokenizer, f)
```

## OUTPUT SCREENSHOTS:



*Fig A.1 confusion matrix*



*Fig A.2 ROC curve*

```
Type a sentence to detect sarcasm (type 'exit' to stop):
📄 Your sentence: good job failing in all subjects
1/1 ————— 0s 52ms/step

🚗 Prediction: Sarcastic (Confidence: 0.9932)
157/157 ————— 1s 7ms/step

💡 Explanation:
- good: 0.1393
- failing: 0.1224
- job: 0.1211
- all: 0.0978
- in: 0.0567
- subjects: 0.0210
```

**Fig A.3** Prediction for the given sentence

Classification Report:				
	precision	recall	f1-score	support
Not Sarcastic	0.85	0.86	0.86	2995
Sarcastic	0.85	0.83	0.84	2729
accuracy			0.85	5724
macro avg	0.85	0.85	0.85	5724
weighted avg	0.85	0.85	0.85	5724

**Fig A.4** classification report

```
Type a sentence to detect sarcasm (type 'exit' to stop):
🗨️ Your sentence: sun rises in the east and sets in the west
1/1 ————— 0s 35ms/step

🤖 Prediction: Not Sarcastic (Confidence: 0.0010)
157/157 ————— 1s 8ms/step

🔴 Explanation:
- the: -0.1429
- and: -0.1038
- rises: -0.0925
- sets: -0.0655
- sun: 0.0466
- west: 0.0174
- east: 0.0050
- in: 0.0037
```

***Fig A.5** Prediction for the given sentence*

```
Type a sentence to detect sarcasm (type 'exit' to stop):
🗨️ Your sentence: medics drop soccer player from stretcher; he's ticked
1/1 ————— 0s 36ms/step

🤖 Prediction: Not Sarcastic (Confidence: 0.0386)
157/157 ————— 1s 8ms/step

🔴 Explanation:
- soccer: -0.6043
- drop: 0.1470
- he: 0.1282
- s: -0.0957
- from: 0.0178
- stretcher: 0.0138
- ticked: 0.0084
- player: 0.0075
```

***Fig A.6** Prediction for the given sentence*




```
is_sarcastic
0    14985
1    13634
Name: count, dtype: int64
```

*Fig A.7 Prediction for the given sentence*

```
Accuracy: 0.84958071278826
ROC AUC Score: 0.9262405707325816
```

*Fig A.8 Prediction for the given sentence*

Deploy


 **Sarcasm Detector**

Enter a sentence to see if it's sarcastic and understand the reasoning with LIME.


👉 Input your sentence:

good job failing in all subjects

Detect Sarcasm

 **Prediction: Sarcastic**

Confidence: 0.9932000041007996

 **Words Influencing the Prediction**

- job: -0.0015
- failing: -0.0011
- good: 0.0005
- all: -0.0003

*Fig A.10 User interface*

## REFERENCE

- [1] S. Joshi, A. K. Tiwari, and M. B. Shukla, "Sarcasm Detection in Online Communication Using Machine Learning," *Procedia Computer Science*, vol. 132, pp. 1117–1123, 2018.  
<https://doi.org/10.1016/j.procs.2018.05.215>
- [2] S. Ghosh and A. Veale, "Fracking Sarcasm Using Neural Network," *WASSA 2016*, pp. 161–169, 2016.  
<https://aclanthology.org/W16-0319>
- [3] A. Mishra, A. Anand, and S. Bhattacharyya, "A Modular Approach to Sarcasm Detection: Integrating Deep Learning and Rule-Based Techniques," *IEEE Transactions on Affective Computing*, 2020.  
<https://doi.org/10.1109/TAFFC.2020.2977992>
- [4] R. Kumar and S. Shah, "Sarcasm Detection in Hindi-English Code-Mixed Tweets Using Deep Learning," *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 4, pp. 197–204, 2020.  
<https://doi.org/10.14569/IJACSA.2020.0110424>
- [5] A. Ghosh and T. Veale, "Magnets for Sarcasm: Making Sarcasm Detection Timely, Contextual and Very Personal," *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 7567–7576, 2020.  
<https://aclanthology.org/2020.emnlp-main.613>
- [6] Y. Tay, A. T. Luu, S. C. Hui, "A Deep Learning Framework for Detecting Sarcasm in Twitter Data," *ACM Transactions on Intelligent Systems and Technology*, vol. 8, no. 3, pp. 1–27, 2017.  
<https://doi.org/10.1145/2963102>

[7] H. Wang, J. Wang, and W. Fu, “Exploring Word Embeddings for Sarcasm Detection,” *International Journal of Knowledge and Systems Science (IJKSS)*, vol. 10, no. 2, pp. 17–29, 2019.

<https://doi.org/10.4018/IJKSS.2019040102>

[8] S. Rajadesingan, R. Zafarani, and H. Liu, “Sarcasm Detection on Twitter: A Behavioral Modeling Approach,” *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining (WSDM)*, pp. 97–106, 2015.

<https://doi.org/10.1145/2684822.2685316>

[9] M. Peled and R. Reichart, “Sarcasm SIGN: Interpreting Sarcasm with Sentiment Based Pretraining,” *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 1273–1285, 2020.

<https://aclanthology.org/2020.findings-emnlp.113>

[10] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” *Proceedings of NAACL-HLT 2019*, pp. 4171–4186, 2019.

<https://aclanthology.org/N19-1423>

