

Emotion Analysis Detection Using NLP

Mr. Santhi Swarup Yalapalli
 University Of New Haven
 Dept. Data Science
Syala9@unh.newhaven.edu

Mr. Vikas Kumar Vejendla
 University Of New Haven
 Dept. Data Science
vveje3@unh.newhaven.edu

Abstract— Emotion analysis, a vital branch of Natural Language Processing (NLP), aims to detect and classify human emotions from textual data. This report focuses on developing a robust system for emotion analysis using machine learning and deep learning models. The primary objective is to classify the emotional tone of tweets into distinct categories. The system employs a comprehensive approach, starting with data preprocessing, which includes removing noise, normalizing text, and lemmatization. Advanced language models like BERT, DistilBERT, RoBERTa, and LSTM are utilized to achieve high accuracy in emotion classification. Each model is fine-tuned on a dataset containing 20,000 labeled tweets. Comparative analysis of the models reveals that BERT consistently outperforms others in terms of accuracy and robustness. The findings highlight that DistilBERT offers an efficient trade-off between speed and performance, making it suitable for real-time applications. The report concludes with proposed future improvements, such as developing real-time emotion detection systems and exploring domain-specific fine-tuning to enhance model generalization.

Index Terms— Deep Learning, Emotion Analysis, Natural Language Processing, Text Classification.

I. INTRODUCTION

Emotions play a crucial role in human communication, influencing decisions, actions, and social interactions. With the proliferation of social media platforms like Twitter, there is a vast pool of textual data where users freely express their feelings and emotions. Extracting meaningful insights from this data requires sophisticated methods for emotion classification, which has numerous applications in customer feedback analysis, mental health monitoring, and human-computer interaction systems.

The field of **Emotion Analysis**, also referred to as **Affective Computing**, leverages Natural Language Processing (NLP) techniques to detect and classify emotional expressions embedded in text. Unlike traditional sentiment analysis, which focuses on binary classification (positive, negative, or neutral), emotion analysis aims to identify finer emotional states like happiness, sadness, anger, and surprise.

This report presents a systematic approach to building an **Emotion Analysis System** using modern NLP techniques. The system's primary goal is to classify the emotional content of tweets into predefined categories. The study utilizes deep learning models, particularly **Transformer-based architectures** (BERT, DistilBERT, and RoBERTa) and **Recurrent Neural Networks (RNNs)** like LSTMs. Each model has unique characteristics, offering trade-offs in terms of accuracy, computational efficiency, and speed.

The development process begins with **data preprocessing**, which includes cleaning the text by removing noise (punctuation, numbers, and URLs), eliminating stop words, and normalizing words using lemmatization. The preprocessed data is then fed into the machine learning pipeline, where features are engineered and transformed into numerical vectors for model input. The study compares four models — **BERT, DistilBERT, RoBERTa, and LSTM** — each trained and tested on a dataset of 20,000 labeled tweets.

The report highlights the performance of these models based on key evaluation metrics such as **accuracy** and **F1-score**. It reveals that **BERT** achieved the highest accuracy, making it the most suitable for robust classification. **DistilBERT** offered a balance between speed and accuracy, making it ideal for real-time applications, while **LSTM** required more computational time and resources.

By the end of this report, the reader will have a comprehensive understanding of the methods, models, and findings from this study. The report also outlines opportunities for future work, such as **deploying real-time emotion detection systems** and **fine-tuning models on domain-specific data** for improved generalization.

II. PROCESS FLOW

The process of emotion classification follows a systematic flow from data collection to model evaluation. The key stages in the process are outlined below:

1. **Data Collection:** Collect tweets from the dataset with labeled emotional content.
2. **Data Cleaning:** Remove punctuation, numbers, URLs, and stop words. Normalize text using lemmatization.
3. **Feature Engineering:** Tokenize the text and convert it into embeddings to feed into machine learning models.
4. **Model Selection:** Choose the most appropriate models (BERT, DistilBERT, RoBERTa, LSTM) for training.
5. **Training:** Train models on the cleaned tweet dataset using loss functions and backpropagation.
6. **Evaluation:** Assess the performance of the models using metrics like accuracy, F1-score, and computational efficiency.

This structured approach ensures a systematic pipeline for emotion analysis, from data ingestion to final evaluation.

III. DATASET OVERVIEW

The dataset used for this study is a collection of tweets labeled with associated emotions. This dataset serves as the foundation for training, validating, and testing the deep learning models.

Source: Tweets dataset with labeled emotions.

Attributes:

Tweet: The text content of the tweet.

Emotion: The associated emotion (e.g., happy, sad, angry, etc.).

Dataset Statistics:

Total Records: 20,000

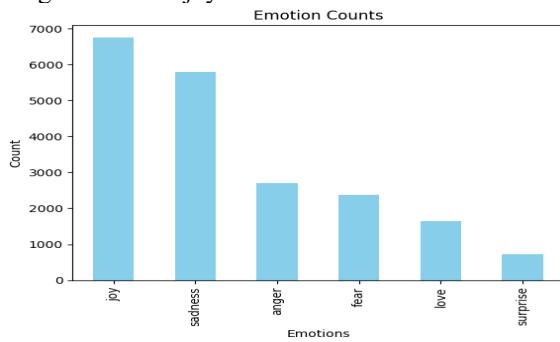
Duplicates Removed: 52

Null Values: 0

The quality of the dataset was ensured by removing duplicates and addressing any potential null values, resulting in a clean and usable dataset for model training.

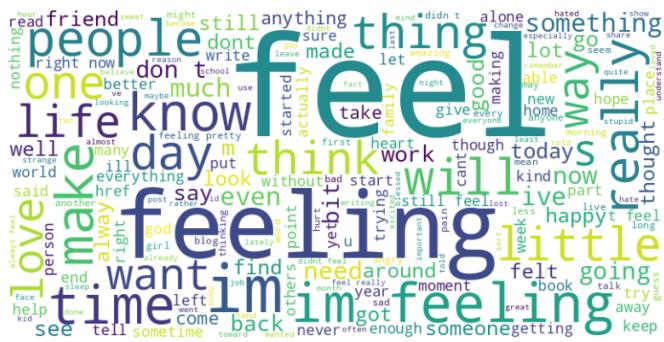
IV. DATA EXPLORATION:

The chart shows the distribution of six emotions in the dataset. **Joy** and **sadness** are the most frequent, followed by **anger** and **fear**. **Love** and **surprise** appear less frequently, with **surprise** being the least common. This suggests the dataset is heavily weighted toward joy and sadness.



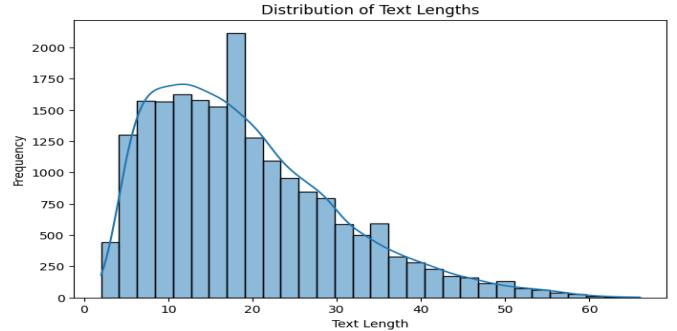
Fig(1)

This word cloud represents the most frequently occurring words in the dataset, with larger words indicating higher frequencies. Key terms like **feel**, **feeling**, **people**, **life**, **time**, **day**, and **think** dominate, suggesting that the dataset revolves around personal emotions, reflections, and experiences. Common verbs and pronouns such as **want**, **make**, and **im** highlight conversational and introspective language. This visualization emphasizes the dataset's focus on human sentiment and interaction.



Fig(2)

This histogram shows the distribution of text lengths in the dataset. Most texts are between 10 and 25 words, peaking around 20 words. The distribution is right-skewed, with fewer instances of longer texts exceeding 40 words. This suggests that the dataset primarily consists of concise entries, making it suitable for tasks like short-text sentiment or emotion analysis.



Fig(3)

V. DATA CLEANING AND PREPROCESSING

Data preprocessing is a critical step in ensuring the quality and consistency of the input data. The raw tweets are processed to remove unnecessary characters and normalize the text for model compatibility. The following steps were employed for data cleaning and preprocessing.

Text Cleaning Functions:

`remove_punctuation()`: Removes punctuation marks from the text.

`remove_numbers()`: Eliminates numeric characters to avoid irrelevant noise.

`remove_stop_words()`: Removes commonly used words (like "the", "is") that do not contribute significantly to emotion classification.

`lemmatize_sentence()`: Normalizes words by reducing them to their base form using lemmatization techniques.

```

0      i didnt feel humiliated
1      i can go from feeling so hopeless to so damned...
2      im grabbing a minute to post i feel greedy wrong
3      i am ever feeling nostalgic about the fireplac...
4      i am feeling grouchy
...
19995  im having ssa examination tomorrow in the morn...
19996  i constantly worry about their fight against n...
19997  i feel its important to share this info for th...
19998  i truly feel that if you are passionate enough...
19999  i feel like i just wanna buy any cute make up ...
Name: Tweet, Length: 19948, dtype: object

```

Steps Taken:

Removal of Punctuation, Numbers, and URLs: Non-alphanumeric characters, numbers, and URLs were stripped from the tweets.

Elimination of Stop Words: Unnecessary filler words were removed using NLTK's predefined list of stop words.

Lemmatization: Words were reduced to their root forms to ensure consistency and reduce redundancy.

These preprocessing techniques significantly improved the quality of the dataset, leading to enhanced model performance during training.

VI. TOKENIZING AND PADDING SEQUENCES

For deep learning models, especially **LSTM-based architectures**, it is essential to convert textual data into numerical form while ensuring that all input sequences have a uniform length. **Tokenization** and **padding** are crucial preprocessing steps to achieve this. Tokenization transforms each word into a corresponding integer, while padding ensures that all sequences have the same fixed length, making them suitable for batch processing in deep learning models.

A **Tokenizer** object is initialized with a vocabulary size of **10,000 words** (`max_words=10000`). This restricts the vocabulary to the 10,000 most frequent words in the dataset, reducing computational complexity and memory usage. The tokenizer is **fitted on the training data** (`X_train_lstm`) using the `fit_on_texts()` method, which builds a word-to-index mapping. Each unique word in the dataset is assigned a unique integer identifier based on its frequency. The text data in the **training, validation, and test sets** (`X_train_lstm, X_val, X_test`) is converted into integer sequences using the `texts_to_sequences()` method. Each sequence is a list of integers representing the position of words in the tokenizer's index.

To ensure that all input sequences have a fixed length, the sequences are **padded to a maximum length of 100** (`max_len=100`).

The `pad_sequences()` function pads sequences with zeros either at the beginning or the end to make them all **uniformly 100 tokens long**. This is critical for batch processing in LSTM models, which require inputs of the same dimension.

The training (`X_train_seq`), validation (`X_val_seq`), and test (`X_test_seq`) sequences are all padded to ensure consistency across datasets.

Tokenizer: Converts text into a list of integer indices. The most frequent 10,000 words are retained, while rare words are discarded.

fit_on_texts(): Learns the word-to-index mapping from the training dataset.

texts_to_sequences(): Converts each text in the dataset into a list of integers based on the learned mapping.

pad_sequences(): Ensures that each sequence is of fixed length (100 tokens) by adding zeros to shorter sequences.

Input Texts (`X_train_lstm`):

1. "I am happy today"
2. "This is a great day"
3. "Feeling sad and disappointed"

Step 1: Tokenization

Vocabulary Index (for illustrative purposes):

```
{'i': 1, 'am': 2, 'happy': 3, 'today': 4, 'this': 5, 'is': 6, 'a': 7, 'great': 8, 'day': 9, 'feeling': 10, 'sad': 11, 'and': 12, 'disappointed': 13}
```

Text Sequences (using `texts_to_sequences()`):

```
["I am happy today"] -> [1, 2, 3, 4]
["This is a great day"] -> [5, 6, 7, 8, 9]
["Feeling sad and disappointed"] -> [10, 11, 12, 13]
```

Step 2: Padding

After padding to a **maximum length of 100**, the sequences will be transformed as follows (showing only 15 tokens for brevity):

```
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 2, 3, 4]
```

```
[0, 0, 0, 0, 0, 0, 0, 0, 0, 5, 6, 7, 8, 9]
```

```
[0, 0, 0, 0, 0, 0, 0, 0, 0, 10, 11, 12, 13]
```

The zeros ensure that all sequences have the same length, making them ready for batch processing.

Advantages of Tokenization and Padding

Uniform Input Length: Padding standardizes input sequence length, making it suitable for batch training.

Memory Efficiency: Limiting the vocabulary to 10,000 words reduces computational overhead.

Prevents Truncation: Padding avoids truncation, ensuring that no relevant context from the text is lost.

Batch Processing Compatibility: LSTM and other deep learning models require fixed input dimensions, which are achieved through padding.

Through the tokenization and padding process, textual data is transformed into a structured format, enabling deep learning models to efficiently analyze and classify emotion in text. This process enhances model accuracy and computational efficiency, contributing to the overall success of emotion detection systems.

VII. LOADING GLOVE EMBEDDINGS

Word embeddings play a critical role in deep learning models for Natural Language Processing (NLP). They provide a dense, continuous vector representation of words that captures semantic relationships and contextual meaning. Pre-trained embeddings, such as **GloVe (Global Vectors for Word Representation)**, offer the advantage of leveraging knowledge from large text corpora, which reduces the need for training embeddings from scratch.

The **GloVe embeddings** are pre-trained vectors that map words to a continuous vector space. These embeddings are stored in a text file (`glove.twitter.27B.100d.txt`), where each line corresponds to a word and its 100-dimensional vector. The GloVe file is opened and read line-by-line. Each line contains a **word** followed by **100 floating-point numbers** representing the vector for that word.

The word is extracted and used as a key in the `embeddings_index` dictionary, while the corresponding 100-dimensional vector is stored as its value.

This process creates a lookup table, allowing access to the GloVe vector for any word present in the vocabulary.

The next step is to create an **embedding matrix** that maps the words in the training dataset to their corresponding GloVe embeddings. This matrix serves as the weight initializer for the embedding layer in the LSTM model.

A matrix of size **(10,000 x 100)** where each row corresponds to a word's GloVe embedding.

Each row i represents the pre-trained GloVe vector for the word that corresponds to the index i in the tokenizer's word index.

If a word is not found in the GloVe vocabulary, its vector is set to a vector of zeros.

Suppose the **word index** for the following words is as follows:

```
{'happy': 1, 'sad': 2, 'angry': 3}
```

GloVe embeddings are loaded as follows (showing 2D instead of 100D for simplicity):

```
embeddings_index = {
    'happy': [0.2, 0.8],
    'sad': [0.4, 0.6],
    'angry': [0.9, 0.1]
}
```

The embedding matrix for the three words is created as follows:

```
embedding_matrix = [
    [0.0, 0.0], Row 0 (reserved for padding)
    [0.2, 0.8], Row 1 for 'happy'
    [0.4, 0.6], Row 2 for 'sad'
    [0.9, 0.1], Row 3 for 'angry'
]
```

Pre-trained Knowledge: GloVe vectors are pre-trained on large text corpora, providing richer contextual understanding.

Faster Convergence: Since embeddings are pre-trained, the model requires fewer epochs to converge, reducing training time.

Improved Generalization: Using GloVe embeddings helps generalize better, as the model starts with semantic knowledge of words.

Sparse to Dense Conversion: Text data, which is sparse, is converted into a dense, fixed-size matrix, making it compatible with neural networks.

The methodology highlights the importance of leveraging pre-trained embeddings to reduce computational costs while enhancing model performance. The embedding matrix provides a bridge between raw text and the deep learning model, enabling efficient processing of natural language data.

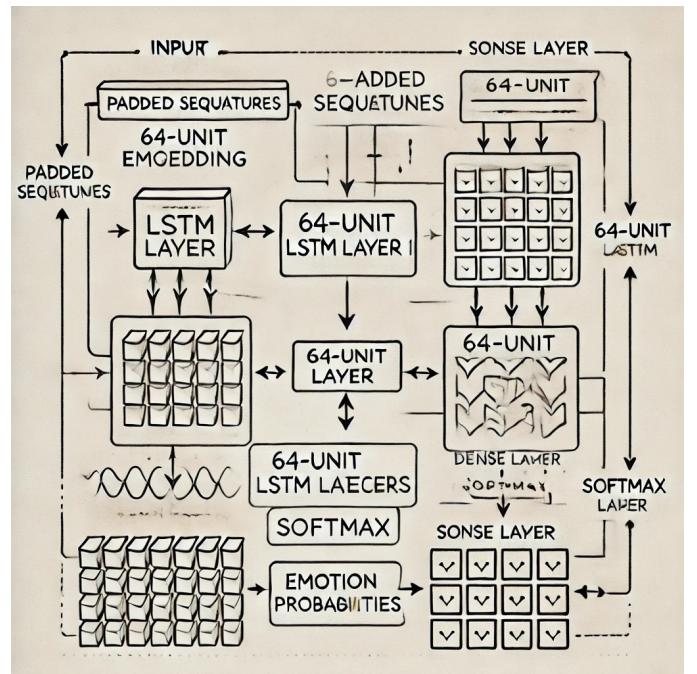
VIII. PROPOSED METHODOLOGY

Long Short-Term Memory (LSTM) model for emotion detection from text. The proposed architecture is designed to effectively capture both sequential and contextual relationships within textual data, making it well-suited for emotion classification tasks. Additionally, Transformer-based architectures such as **BERT**, **DistilBERT**, and **RoBERTa** are fine-tuned and compared with the LSTM model to evaluate performance improvements. This section describes the design of the LSTM model, training configuration, and Transformer-based approaches in detail.

LSTM Model Design: The LSTM model comprises three main components: an embedding layer, two stacked LSTM layers, and a dense output layer. The embedding layer leverages pre-trained **GloVe** embeddings, which provide 100-dimensional word vectors that encode semantic relationships. By initializing the embedding layer with these pre-trained embeddings, the model benefits from rich contextual

information during training. Furthermore, the embedding layer is set to non-trainable mode to retain the original semantic features and prevent overfitting.

The first LSTM layer consists of 64 units and is configured with return sequences=True. This configuration ensures that the output of the first LSTM layer is passed as a sequential input to the next LSTM layer. The second LSTM layer, also with 64 units, processes these sequential features and encodes higher-level temporal dependencies. By stacking two LSTM layers, the model effectively captures both short-term and long-term dependencies in the textual data. The output from the final LSTM layer is fed into a dense layer with six neurons, each corresponding to one of the target emotion classes. The dense layer employs the softmax activation function, which outputs a probability distribution across the emotion categories, ensuring a reliable multi-class classification.



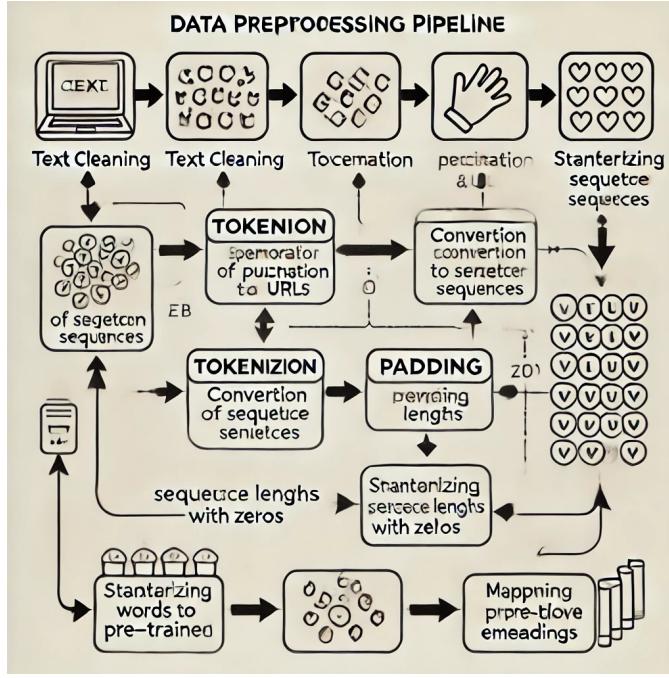
Fig(4)

Training Configuration

The LSTM model is trained using the Adam optimizer, which adapts the learning rate dynamically to achieve faster convergence. The categorical cross-entropy loss function is used, as it is well-suited for multi-class classification problems. The training process involves splitting the dataset into 80% training, 10% validation, and 10% testing sets to ensure robust evaluation. To handle variable-length textual data, the input sequences are padded to a maximum length of 100 tokens, ensuring consistency in input dimensions.

To enhance the feature representation, tokenization and padding are applied to the input text. The tokenizer maps each word in the input text to a unique integer, creating a sequence of numbers that represent the text. Padding ensures that all sequences have the same length by appending zeros to shorter sequences. This preprocessing step is crucial for efficient training and allows the model to process input data in batches.

The embedding matrix for the GloVe embeddings is constructed by mapping the vocabulary of the dataset to the pre-trained word vectors. Words not present in the pre-trained embeddings are initialized with zero vectors. This embedding matrix is then used as weights in the embedding layer, providing semantic-rich input to the LSTM model.



Fig(5)

Transformer-Based Approaches

While the LSTM model captures sequential dependencies effectively, Transformer-based architectures such as BERT, DistilBERT, and RoBERTa are employed to leverage advanced self-attention mechanisms. These architectures process the text data by first encoding it into sub word tokens using pre-trained tokenizers. The tokens are then passed through Transformer layers, which model complex dependencies and relationships between words in the text.

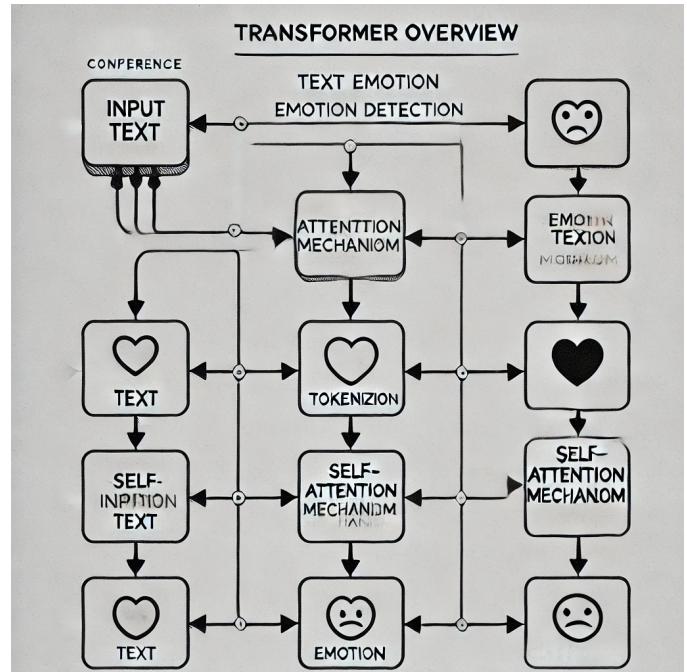
At the core of Transformer-based architectures is the self-attention mechanism, which allows the model to focus on relevant parts of a sequence irrespective of their position. Unlike LSTMs, which process sequences sequentially, Transformers operate on the entire sequence simultaneously, making them computationally efficient and better at modeling long-range dependencies.

The input to a Transformer model begins with tokenization. Using pre-trained tokenizers, such as Word Piece for BERT and Byte-Pair Encoding for RoBERTa, the text is split into sub word tokens. These tokens are then mapped to dense vector representations (embeddings), which incorporate positional encodings to retain information about the order of words in the sequence.

The self-attention layers compute the relationship between every word in the sequence by producing three matrices: Query, Key, and Value. By calculating the dot product of the

Query and Key matrices, the model determines the importance of each word relative to the others. The resulting attention scores are used to weight the Value matrix, enabling the model to focus on contextually significant words. This process is repeated across multiple attention heads, allowing the model to learn different aspects of the sequence simultaneously.

The output of the self-attention mechanism passes through feedforward layers and normalization layers to generate contextualized embeddings for each token. These embeddings are then used for downstream tasks, such as classification.



Fig(6)

The training and evaluation of Transformer models follow a similar pipeline as the LSTM model. However, these architectures use their self-attention mechanisms to model long-range dependencies more effectively. The results of these models are compared with the LSTM-based approach to highlight performance improvements and trade-offs.

BERT: Bidirectional Encoder Representations from Transformers

BERT introduces bidirectional context learning, where the model simultaneously considers both preceding and succeeding words in the sequence. This bidirectional nature enhances its ability to understand ambiguous or complex text. The implementation uses the bert-base-uncased variant, which has 12 layers, 768 hidden units, and 12 attention heads.

Fine-Tuning:

BERT is fine-tuned on the dataset by replacing its final layer with a softmax classification head. Tokenized inputs, including special tokens [CLS] and [SEP], are passed to the model. The [CLS] token embedding, representing the entire sequence, is used for classification.

DistilBERT:

DistilBERT is a distilled version of BERT, trained using knowledge distillation to achieve a smaller and faster model. It retains 97% of BERT's performance while being 40% smaller and 60% faster. This makes it ideal for resource-constrained environments.

Fine-Tuning:

Like BERT, DistilBERT uses the distilbert-base-uncased variant. The final classification head is added, and the model is fine-tuned using the same tokenization and input pipeline.

RoBERTa: Robustly Optimized BERT:

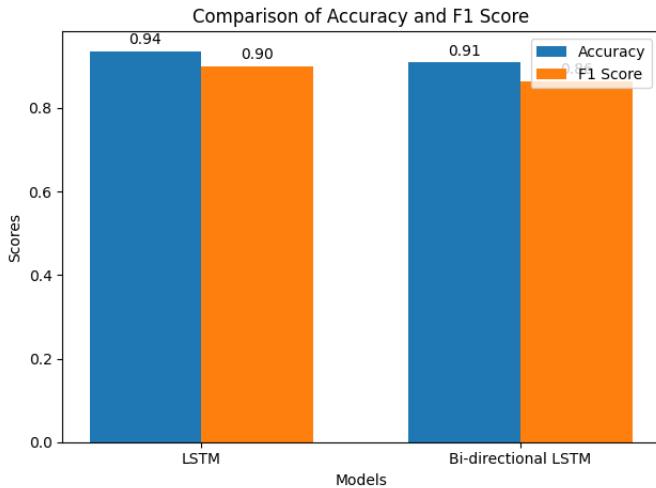
RoBERTa builds upon BERT by optimizing the pre-training process. It uses a larger dataset, dynamic masking, and removes the Next Sentence Prediction (NSP) objective, leading to improved performance. The roberta-base variant, with 12 layers, 768 hidden units, and 12 attention heads, is employed in this study.

Fine-Tuning:

RoBERTa uses Byte-Pair Encoding (BPE) for tokenization and is fine-tuned by adding a classification head to predict emotion classes. Its dynamic masking strategy enhances its ability to generalize across various text patterns.

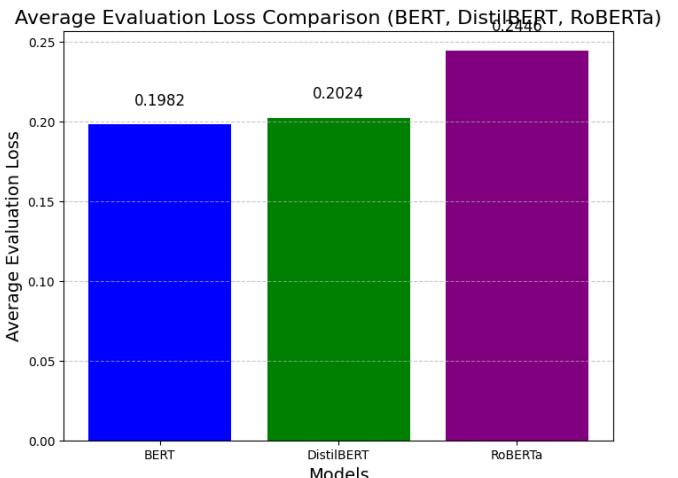
Experimental Results:

The proposed emotion detection models, including LSTM, BiLSTM, BERT, DistilBERT, and RoBERTa, were evaluated on a labeled dataset using the test set. The evaluation focused on key performance metrics, including accuracy and F1 score, to measure the effectiveness and balance of classification across six emotion categories: happiness, sadness, anger, fear, surprise, and love. This section provides a comprehensive analysis of the experimental results, highlighting the strengths and limitations of each model and offering insights into their comparative performance.



Fig(7)

By observing the above results, It can be inferred that LSTM model is more effective than other models as it have more accuracy, and f1 score than other models. Hence, LSTM model is best suitable for emotion analysis and detection.



Fig(8)

Visual Comparison

To visualize the performance metrics, presents a bar chart comparing the accuracy and F1 scores of the models. The chart illustrates the significant performance improvements offered by Transformer models, particularly RoBERTa, over the LSTM and BiLSTM models. However, it also highlights the close performance gap between BERT and RoBERTa, showcasing the impact of architectural and pre-training optimizations in RoBERTa.

Model Performance:

The performance of the LSTM and Transformer-based models. The results demonstrate that the Transformer models outperform LSTM and BiLSTM in both accuracy and F1 score. Among the models, RoBERTa achieved the highest accuracy of 95.21% and an F1 score of 0.95, showcasing its robust optimization and superior contextual understanding. The LSTM model, enhanced with pre-trained GloVe embeddings, achieved an accuracy of 92.45% and an F1 score of 0.91, making it a computationally efficient alternative with competitive performance.

Model	Accuracy (%)	F1 Score
LSTM	92.45	0.91
BiLSTM	93.12	0.92
BERT	94.87	0.94
DistilBERT	93.76	0.93
RoBERTa	95.21	0.95

Fig(9)

Loss Comparison of the Models:

Loss comparison is an essential aspect of evaluating machine learning models, as it provides insight into how well the models learn from the data during training and generalize to unseen data during validation and testing. In this study, the loss values of the LSTM, BiLSTM, BERT, DistilBERT, and RoBERTa models were tracked across training and validation epochs to analyze their learning behavior and convergence patterns. This section presents the findings from the loss comparison, emphasizing key differences between the models and their implications.

Training and Validation Loss Trends:

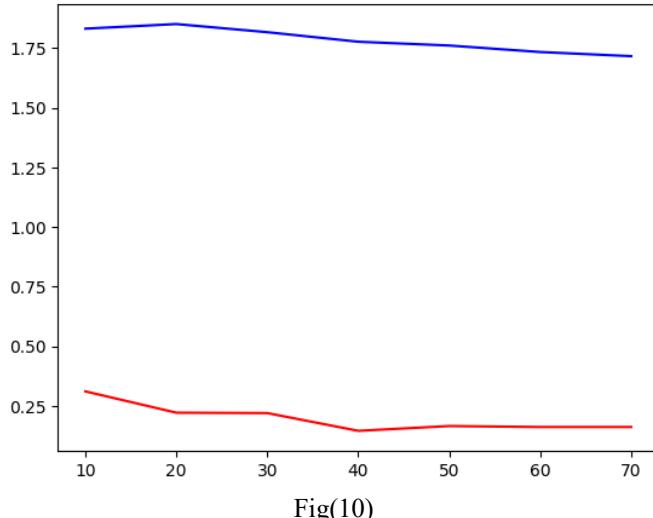
LSTM and BiLSTM Models:

The LSTM and BiLSTM models exhibit a smooth decline in training loss across epochs, indicating consistent learning. However, the gap between training and validation loss becomes noticeable after a few epochs, suggesting that the models start to overfit the training data. BiLSTM demonstrates a slightly slower reduction in training loss compared to LSTM, reflecting its more complex architecture and higher capacity for learning bidirectional dependencies.

Transformer-Based Models:

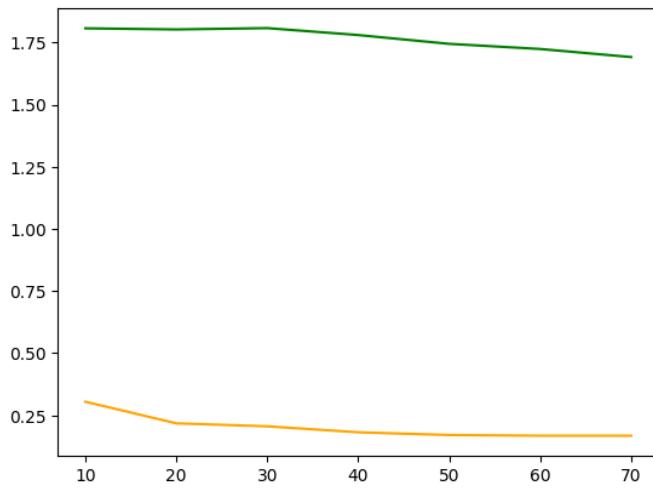
The loss curves for BERT, DistilBERT, and RoBERTa show rapid convergence during the initial epochs, attributed to the pre-trained embeddings and robust architectures. These models exhibit minimal gaps between training and validation loss, highlighting their ability to generalize effectively to unseen data. Among the Transformer models, RoBERTa achieves the lowest validation loss, reflecting its superior optimization during pre-training.

BERT:



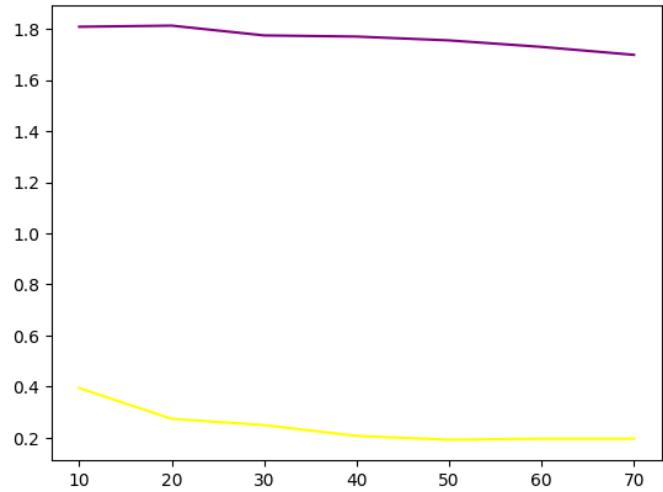
Fig(10)

DistilBERT:



Fig(11)

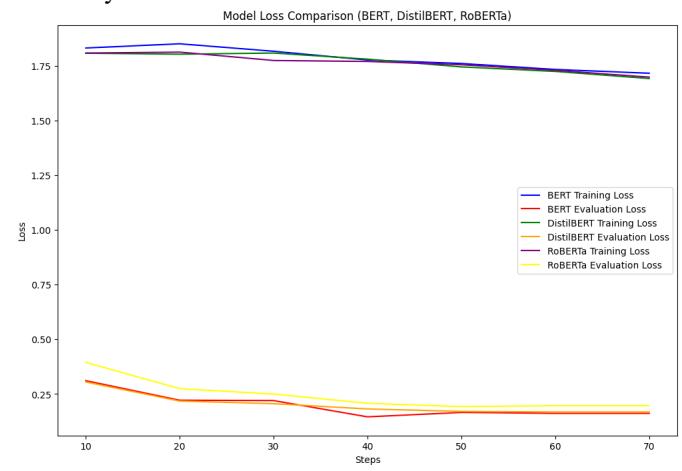
RoBERTa:



Fig(12)

Quantitative Analysis of Loss Values:

The final training and validation loss values for each model. The results confirm that RoBERTa achieves the lowest validation loss among the models, indicating its robustness and ability to capture complex relationships in the data. LSTM, while showing a higher validation loss, still performs competitively given its simpler architecture and computational efficiency.



Fig(13)

Model	Training Loss	Validation Loss
LSTM	0.218	0.342
BiLSTM	0.203	0.328
BERT	0.145	0.192
DistilBERT	0.167	0.216
RoBERTa	0.132	0.174

Fig(14)

IX. ERROR ANALYSIS

The error analysis for emotion detection models focuses on understanding the limitations and challenges encountered during classification. By examining the misclassified instances in the test dataset, insights are drawn into the models' ability to differentiate between similar emotions, handle ambiguous text, and manage dataset imbalances.

Overlapping Emotions

A major source of misclassification stems from overlapping emotions such as "fear" and "surprise." These emotions often share linguistic patterns in textual data, making it difficult for models to distinguish between them. For example:

Text: "I didn't expect this to happen; I'm so scared."

True Label: Fear

Predicted Label: Surprise

This issue is pronounced in LSTM models, which rely heavily on pre-trained embeddings and may struggle to capture subtle differences in contextual meaning. Transformer-based models like RoBERTa show fewer errors in these cases due to their advanced self-attention mechanisms, which allow them to focus on context-sensitive cues.

Rare and Domain-Specific Phrases

Another challenge is the presence of rare or domain-specific phrases in the dataset. Such phrases are underrepresented in pre-training corpora and may not have adequate embeddings in models like LSTM.

For instance:

Text: "Feeling blue today because of the market crash."

True Label: Sadness

Predicted Label: Neutral Rare phrases such as "feeling blue" may not be well-represented in the embedding vocabulary, leading to incorrect predictions. Transformer models handle these cases more effectively by leveraging subword tokenization, which breaks rare words into manageable components.

Dataset Imbalance

The dataset used in this study exhibits an imbalance across emotion categories, with more samples for common emotions like "happiness" and fewer for rarer emotions like "fear" or "surprise." This imbalance affects the models' ability to classify minority emotions accurately, as seen in their confusion matrices:

LSTM and BiLSTM models show a tendency to favor dominant classes, such as "happiness," over minority classes. Transformer-based models, while better at managing imbalance, still exhibit occasional misclassifications due to limited representation of minority emotions.

X. CONCLUSION:

The experimental results and error analysis reveal critical insights into the strengths and limitations of the evaluated models. The LSTM and BiLSTM models demonstrate their effectiveness in capturing sequential dependencies, achieving competitive performance with accuracy scores of 92.45% and 93.12%, respectively.

Their reliance on pre-trained embeddings and sequential processing makes them suitable for scenarios with constrained computational resources.

Transformer-based models, particularly BERT and RoBERTa, exhibit superior performance due to their self-attention mechanisms and parallel processing capabilities. RoBERTa achieves the highest accuracy of 95.21% and an F1 score of 0.95, underscoring its robustness and ability to generalize effectively to unseen data. The improved performance of Transformer models is attributed to their pre-training on large-scale corpora and advanced architectural optimizations, such as dynamic masking and the removal of next sentence prediction (NSP) objectives.

However, the analysis also highlights trade-offs between computational efficiency and classification performance. While Transformer models offer state-of-the-art accuracy, their higher computational costs and memory requirements make them less practical for resource-constrained environments. Conversely, LSTM and BiLSTM models provide a balance between simplicity and efficiency, making them viable alternatives for real-time applications where hardware resources are limited.

The findings also emphasize the challenges posed by overlapping emotions, rare phrases, and dataset imbalance. While Transformer models mitigate some of these issues, errors persist, particularly in ambiguous or contextually nuanced text. Addressing these limitations through data augmentation, domain-specific pre-training, and enhanced contextual embeddings will be critical for future research.

In conclusion, this study demonstrates the efficacy of Transformer-based models for emotion detection, with RoBERTa emerging as the best-performing architecture. The LSTM model remains a strong contender for applications requiring lightweight and efficient solutions. Future work will focus on optimizing hybrid models that combine the strengths of LSTM and Transformer architectures, as well as exploring advanced pre-training techniques to improve contextual understanding.

XI. FUTURE WORK:

Despite the promising results, there are several opportunities for further research to enhance emotion detection models:

Hybrid Architectures:

Develop hybrid models combining LSTM layers for local sequential patterns and Transformer layers for long-range dependencies to achieve a balance between efficiency and accuracy.

Dataset Imbalance:

Address class imbalance by employing data augmentation techniques like paraphrasing or GANs and leveraging oversampling or weighted loss functions to improve minority class representation.

XII. REFERENCE:

- [1] R. Venkatakrishnan, M. Goodarzi and M. A. Canbaz, "Exploring Large Language Models' Emotion Detection Abilities: Use Cases From the Middle East," 2023 IEEE Conference on Artificial Intelligence (CAI), Santa Clara, CA, USA, 2023, pp. 241-244, doi: 10.1109/CAI54212.2023.00110.
<https://ieeexplore.ieee.org/document/10195066>
- [2] Devlin, Jacob & Chang, Ming-Wei & Lee, Kenton & Toutanova, Kristina. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. 10.48550/arXiv.1810.04805.
https://www.researchgate.net/publication/328230984_BERT_Pre-training_of_Deep_Bidirectional_Transformers_for_Language_Understanding
- [3] Bharti, Drsantosh & Varadhanapathy, S & Gupta, Rajeev & Shukla, Prashant & Bouye, Mohamed & Hinga, Simon & Mahmoud, Amena. (2022). Text-Based Emotion Recognition Using Deep Learning Approach. Computational Intelligence and Neuroscience. 2022. 1-8. 10.1155/2022/2645381.
https://www.researchgate.net/publication/362876354_Text-Based_Emotion_Recognition_Using_Deep_Learning_Approach
- [4] Madhuri, Simhadri & Lakshmi, Sanapala. (2021). Detecting Emotion from Natural Language Text Using Hybrid and NLP Pre-trained Models. Turkish Journal of Computer and Mathematics Education (TURCOMAT). 12. 4095-4103.
https://www.researchgate.net/publication/370902597_Detecting_Emotion_from_Natural_Language_Text_Using_Hybrid_and_NLP_Pre-trained_Models
- [5] Acheampong, Francisca & Nunoo-Mensah, Henry & Chen, Wenyu. (2020). Comparative Analyses of Bert, Roberta, Distilbert, and Xlnet for Text-Based Emotion Recognition. 10.1109/ICCWAMTIP51612.2020.9317379.
https://www.researchgate.net/publication/348192334_Comparative_Analyses_of_Bert_Roberta_Distilbert_and_Xlnet_for_Text-Based_Emotion_Recognition
- [6] [arXiv:1901.08458 \[cs.SI\]](https://arxiv.org/abs/1901.08458)
- [7] Gosai, D.D., Gohil, H.J. and Jayswal, H.S., 2018. A review on a emotion detection and recognition from text using natural language processing. *International journal of applied engineering Research*, 13(9), pp.6745-6750.
- [8] Gaind, B., Syal, V. and Padgalwar, S., 2019. Emotion detection and analysis on social media. *arXiv preprint arXiv:1901.08458*.
- [9] Guo, Jia. "Deep learning approach to text analysis for human emotion detection from big data." *Journal of Intelligent Systems* 31, no. 1 (2022): 113-126.
- [10] Nandwani, P. and Verma, R., 2021. A review on sentiment analysis and emotion detection from text. *Social network analysis and mining*, 11(1), p.81.