



MANIPAL INSTITUTE OF TECHNOLOGY MANIPAL

A Constituent Institution of Manipal University

Department of Computer Science and Engineering Artificial Intelligence and Machine Learning

Report on

**Emotion Recognition using Deep Learning on the Toronto Emotional
Speech Set (TESS)**

by

YASHVEER SINGH

210962098

towards qualitative assessment for the course

Speech Processing – CSE 4023

Emotion Recognition using Deep Learning on the Toronto Emotional Speech Set (TESS)

YASHVEER SINGH,
210962098,

yashveer.singh@learner.manipal.edu

ABSTRACT: *This project endeavors to enhance the understanding and computational recognition of human emotions, a pivotal component in the advancement of human-computer interaction, mental health assessment, and automated customer service solutions. By leveraging a robust dataset and employing a series of sophisticated machine learning techniques, including Long Short-Term Memory (LSTM) networks, we aim to accurately classify and predict a range of human emotions from diverse inputs such as audio, text, and facial expressions. The exploratory analysis phase of the project provided insightful revelations into the data's inherent characteristics, guiding the subsequent feature extraction process to distill relevant information crucial for model training. The LSTM model, known for its efficacy in handling sequential data, was meticulously tailored and trained, demonstrating significant potential in capturing the temporal dynamics of emotional expressions. Results were methodically plotted to offer a transparent evaluation of the model's performance, revealing promising accuracy levels and areas for future enhancement. This research not only contributes to the growing body of knowledge in emotion recognition but also paves the way for more nuanced and empathetic interactions between humans and machines.*

Keywords: *Emotion Recognition, LSTM, Machine Learning, Human-Computer Interaction, Feature Extraction, Exploratory Data Analysis, Sequential Data Processing.*

I. INTRODUCTION

In the realm of artificial intelligence, the ability to accurately recognize and interpret human emotions stands as a frontier with profound implications for enhancing human-computer interaction. Emotion recognition technology, which allows machines to understand and respond to human emotions, is becoming increasingly crucial across various sectors, including healthcare, education, and customer service. This project delves into the development of an advanced emotion recognition system, utilizing a combination of machine learning techniques and Long Short-Term Memory (LSTM) networks to analyze and classify emotional states from diverse data sources such as audio recordings, textual data, and facial expressions. By conducting a comprehensive exploratory analysis to understand the nuances of the dataset and employing sophisticated feature extraction methods, this work sets the foundation for a model capable of interpreting the complex spectrum of human emotions. The ultimate goal is to bridge the gap between human emotional expression and machine understanding, fostering more intuitive and empathetic interactions between humans and technology.

Objectives:

The primary objectives of this emotion recognition project are outlined as follows:

1. **To Develop an Accurate Emotion Recognition Model:** Design and implement a machine learning model capable of accurately recognizing and classifying a wide range of human emotions from diverse data sources, including audio, text, and visual inputs.
2. **To Conduct Comprehensive Exploratory Data Analysis:** Perform an in-depth exploratory analysis of the dataset to uncover underlying patterns, distributions of emotional states, and potential biases, which could influence the model's training and performance.
3. **To Implement Effective Feature Extraction Techniques:** Utilize and innovate feature extraction methods tailored to each type of input data, ensuring that the most informative and relevant features are fed into the model for optimal emotion recognition accuracy.
4. **To Explore and Optimize LSTM Network Architectures:** Given the sequential nature of much of the emotional data, particularly audio and textual inputs, investigate and optimize Long Short-Term Memory (LSTM) networks to capture temporal dependencies and nuances in emotional expressions effectively.
5. **To Evaluate and Report Model Performance:** Thoroughly evaluate the model's performance using a variety of metrics, and visualize these results to provide clear insights into the model's strengths and weaknesses in emotion recognition tasks.
6. **To Advance Human-Computer Interaction (HCI):** Through the development of this emotion recognition system, contribute to the field of HCI by enabling more natural, intuitive, and empathetic interactions between humans and machines, thereby improving user experiences across various applications.
7. **To Identify Future Research Directions:** Highlight the limitations of the current project and propose future research directions that could address these challenges, improve model performance, and expand the applicability of emotion recognition technologies.

These objectives aim to push the boundaries of current emotion recognition capabilities, fostering advancements that could have significant impacts across multiple domains by making technology more responsive and attuned to human emotional states..

Significance

The significance of this emotion recognition project lies in its potential to revolutionize the way machines understand and interact with humans, marking a significant leap forward in the realms of artificial intelligence and human-computer interaction. By accurately identifying and classifying human emotions from various inputs, this technology promises to enhance a wide array of applications, from improving mental health diagnostics through real-time mood tracking to creating more engaging and personalized educational software that responds to the emotional state of learners. Furthermore, in customer service, emotion recognition can lead to more empathetic and efficient service delivery, as machines can be programmed to detect customer dissatisfaction or happiness and adjust their responses accordingly. This project also has profound implications for the safety and inclusivity of AI, as understanding emotional cues can lead to the development of more intuitive and accessible technology for individuals with different communication styles and abilities. Ultimately, by bridging the emotional gap between humans and machines, this project not only enhances the functionality and user experience of technology across various

sectors but also contributes to a deeper understanding of human emotion itself, paving the way for innovations that respect and reflect the complexity of human nature.

Research Gaps

Identifying research gaps in the field of emotion recognition involves a thorough examination of current methodologies, technologies, and applications, as well as their limitations. Based on the work done in your project and the broader context of emotion recognition research, several gaps can be highlighted:

- 1. Limited Generalization Across Diverse Data Sources:** Many emotion recognition models excel within specific domains (e.g., audio, text, or visual data) but struggle to maintain high accuracy across diverse data types. A critical research gap exists in developing models that can robustly generalize across various inputs while maintaining high levels of accuracy.
- 2. Handling Subtle and Complex Emotions:** Current models are adept at identifying clear, distinct emotional states (e.g., happiness, sadness, anger) but often fail to recognize more nuanced or compound emotions (e.g., frustration, disappointment, contentment). There is a need for more sophisticated models capable of capturing the subtleties of human emotional expression.
- 3. Real-time Processing and Analysis:** While some models can effectively analyze pre-recorded or static data, the ability to process and recognize emotions in real-time remains a challenge, particularly in resource-constrained environments. Enhancing the efficiency and speed of emotion recognition models is essential for applications requiring immediate feedback, such as interactive AI assistants or real-time mental health monitoring.
- 4. Ethical Considerations and Privacy Concerns:** As emotion recognition technology advances, so do concerns regarding privacy, consent, and the ethical use of emotional data. Research is needed to develop frameworks and guidelines that ensure the responsible use of emotion recognition technologies, balancing innovation with individual rights to privacy and autonomy.
- 5. Cultural and Individual Variability:** Emotional expression and interpretation can vary significantly across cultures and individuals. Most existing models are trained on datasets that do not fully represent this diversity, leading to biases and inaccuracies. Expanding datasets to include a wider range of emotional expressions and contexts is crucial for creating more inclusive and accurate emotion recognition systems.
- 6. Integration with Other Systems:** The potential of emotion recognition is maximized when integrated with other systems, such as conversational AI, adaptive learning platforms, or healthcare diagnostics. Research gaps exist in creating seamless integrations that leverage emotion recognition to enhance the functionality and responsiveness of these technologies.

Addressing these research gaps will not only advance the field of emotion recognition but also unlock new possibilities for human-computer interaction, making technology more empathetic, intuitive, and responsive to human needs and emotions.

Challenges

One of the principal challenges in emotion recognition lies in accurately capturing and interpreting the multifaceted nature of human emotions through technical means. The variability in emotional expression across different individuals and cultures poses a significant hurdle, necessitating sophisticated models that can adapt to and interpret a wide array of emotional signals. Additionally, the integration of diverse data types—such as audio, text, and visual cues—into a cohesive recognition system demands advanced feature

extraction techniques and models capable of processing multimodal inputs. This complexity is compounded by the need for real-time analysis in many applications, requiring not only high accuracy but also efficient processing to deliver immediate feedback. Moreover, the evolution of machine learning models, particularly deep learning architectures like LSTM, presents its own set of challenges, including the demand for large, diverse datasets for training, the risk of overfitting, and the interpretability of model decisions. These technical challenges highlight the need for ongoing research and innovation to refine emotion recognition technologies and their applications.

II. LITERATURE REVIEW

Reference [1] This paper delves into the nuances of detecting emotional states through speech, an area rich with potential for applications in customer service, mental health monitoring, and beyond. The authors present a series of experiments showcasing the development and validation of advanced machine learning models, particularly focusing on real-time processing capabilities. By employing a combination of feature extraction techniques and neural network architectures, the study highlights the model's applicability in dynamically assessing emotional states in scenarios such as live customer interactions and ongoing mental health support, underscoring the importance of real-time feedback for immediate intervention or response..

Reference [2] With the challenge of training models on limited datasets—a common issue in specialized emotion recognition tasks—this paper explores the efficacy of transfer learning. The authors demonstrate how models pre-trained on large, generic datasets can be fine-tuned with smaller, domain-specific datasets to accurately recognize emotions. This approach not only addresses the scarcity of labeled data in certain emotional categories but also significantly reduces the computational resources and time required for training models from scratch. The paper provides a comparative analysis of different transfer learning strategies and their impact on model performance across various small datasets, offering valuable insights for researchers dealing with data constraints.

Reference [3] Recognizing the complexity of human emotions, which can be expressed through a blend of speech, facial expressions, and textual communication, this study investigates the integration of these diverse data sources into a unified deep learning model. The authors propose a novel architecture that effectively combines features from audio, visual, and text inputs, achieving a significant improvement in recognition accuracy compared to single-modality models. This research underscores the potential of multimodal approaches in creating more nuanced and comprehensive emotion recognition systems, paving the way for more sophisticated applications in interactive technology.

Reference [4] This research paper emphasizes the critical role of context in interpreting emotional expressions. By incorporating situational cues and contextual data into emotion recognition models, the study proposes a method that markedly enhances the accuracy of emotion predictions. The authors argue that emotions cannot be fully understood in isolation from their context, showcasing models that dynamically adapt to different contexts for more accurate emotion interpretation. This approach not only advances the field of emotion recognition but also contributes to the development of AI systems capable of more human-like understanding and interactions.

Reference [5] Addressing one of the key challenges in globalized applications of emotion recognition technology—the variability of emotional expression across cultures—this paper presents a groundbreaking dataset compiled from a wide array of cultural contexts. The authors detail their methodology for dataset collection and annotation, followed by an analysis of the initial findings from both human and AI systems' performance on this dataset. The study highlights the complexities and nuances of cross-cultural emotion recognition, offering insights into the biases present in existing models and suggesting pathways towards more inclusive and equitable emotion recognition technologies.

III. METHODOLOGY

Our research aims to enhance emotion recognition from speech signals using advanced machine learning techniques. This section outlines the comprehensive methodology employed, from data analysis to model training and evaluation.

Data Analysis

The initial phase of our study involved a meticulous exploratory data analysis (EDA) to understand the distribution of emotions within our dataset. Utilizing the seaborn library, we conducted a count plot of the various emotion labels present, which include fear, anger, disgust, neutral, sad, and pleasantly surprised (ps), to ensure a balanced representation of each emotion in our analysis. Subsequent to the distribution analysis, audio signal representations were visualized to gain insight into the distinctive features of each emotion category. Two primary visualizations were employed: waveplots and spectrograms, facilitated by the librosa library. Waveplots provided a temporal view of the audio amplitude over time, whereas spectrograms offered a visualization of the frequency spectrum over time, with different colors representing the signal's intensity at various frequencies. These visualizations were crucial for understanding the inherent characteristics of each emotional expression in the dataset, laying the groundwork for feature extraction.

Feature Extraction

The crux of our feature engineering involved the extraction of Mel-frequency cepstral coefficients (MFCCs), a representation well-suited for capturing the timbral aspects of audio signals. We defined a function, `extract_mfcc`, to load each audio file, with librosa handling the signal processing. The function extracts 40 MFCC features, which are then averaged over time to produce a single feature vector per audio sample. This process condenses the rich information embedded in the audio into a compact form amenable to machine learning models. The extracted features were then encapsulated into a numpy array for subsequent model training. Additionally, emotion labels were encoded using one-hot encoding to convert categorical labels into a binary matrix representation, suitable for the classification task at hand.

LSTM Model and Training

For the emotion recognition model, we leveraged a Long Short-Term Memory (LSTM) network, renowned for its efficacy in processing sequential data. Our LSTM model comprises a sequential layer structure, starting with an LSTM layer with 256 units. This is followed by dropout layers interspersed with dense layers to mitigate overfitting and enhance model generalization. The model culminates in a softmax layer for multi-class classification. The model was compiled with categorical crossentropy as the loss function and adam optimizer, focusing on accuracy as the primary metric.

Training was conducted over 50 epochs with a batch size of 64, incorporating a validation split of 20% to monitor the model's performance on unseen data throughout the training process. This setup was pivotal in fine-tuning the model parameters and architecture for optimal performance.

Results

The model's performance was meticulously evaluated through the visualization of accuracy and loss metrics over the training epochs. Accuracy plots delineate the training and validation accuracy across epochs, providing insight into the model's learning trajectory and its generalization capability on unseen data. Similarly, loss plots were scrutinized to assess the model's optimization process, with a focus on minimizing the gap between training and validation loss, indicative of a well-fitting model. These visualizations are

critical for diagnosing model behavior, facilitating adjustments to the training regime, and model architecture to enhance performance.

IV. EXPERIMENTAL SETUP

The experimental setup for our emotion recognition project was meticulously designed to evaluate the effectiveness of Long Short-Term Memory (LSTM) networks in classifying emotions from speech data. This section describes the setup in terms of data preparation, model configuration, training procedures, and evaluation metrics, ensuring a comprehensive understanding of the experiment's execution and underlying principles.

Dataset

The cornerstone of our emotion recognition system is a well-curated dataset essential for training accurate and robust models. For this project, we utilized the Toronto emotional speech set (TESS), which is a comprehensive and widely recognized dataset in the field of affective computing. The TESS dataset comprises a vast collection of audio recordings from actors speaking in English, expressing a range of emotions including happiness, sadness, anger, fear, surprise, and disgust, along with a neutral state. Each recorded utterance within the dataset is methodically labeled with the corresponding emotional state, providing a clear framework for supervised learning.

The audio recordings in TESS are pre-segmented, ensuring that each sample corresponds to a single, discrete emotional expression.

Data Preparation

The experiment utilized a curated dataset comprising audio recordings labeled with various emotions. Prior to analysis, the dataset underwent a rigorous preprocessing phase where audio files were standardized in terms of duration and sampling rate to ensure uniformity. Feature extraction focused on Mel-frequency cepstral coefficients (MFCCs), chosen for their relevance in capturing speech characteristics. A total of 40 MFCC features were extracted from each audio sample, subsequently averaged to produce a condensed feature vector representing each recording. The dataset was then split into training and validation sets, adhering to an 80-20 ratio, to facilitate both model training and subsequent evaluation on unseen data.

Model Configuration

The core of the experimental model was an LSTM network, selected for its proficiency in handling sequential data, such as audio signals. The network architecture was composed of an initial LSTM layer with 256 units, followed by multiple dropout and dense layers to prevent overfitting and enhance feature learning. The final layer employed a softmax activation function to classify each audio sample into one of the emotion categories. Model compilation was performed using the categorical crossentropy loss function and the Adam optimizer, with accuracy designated as the primary performance metric.

Training Procedure

Model training was executed over 50 epochs with a batch size of 64, incorporating real-time validation to monitor and evaluate the model's performance on a separate set of data not seen during the training phase. This approach allowed for continuous assessment of the model's generalization ability and facilitated early stopping if overfitting was detected. Adjustments to the model's architecture and hyperparameters were based on ongoing analysis of training and validation metrics, ensuring optimal learning and performance.

Evaluation Metrics

The evaluation framework centered around accuracy and loss metrics, crucial for assessing the model's effectiveness in emotion recognition. Training and validation accuracy metrics provided insights into the model's learning progression and its ability to generalize to new data. Similarly, loss metrics were scrutinized to evaluate the model's optimization process over time. Post-training, a detailed analysis of accuracy and loss trends was conducted to identify any signs of overfitting or underfitting, guiding further refinements to the model.

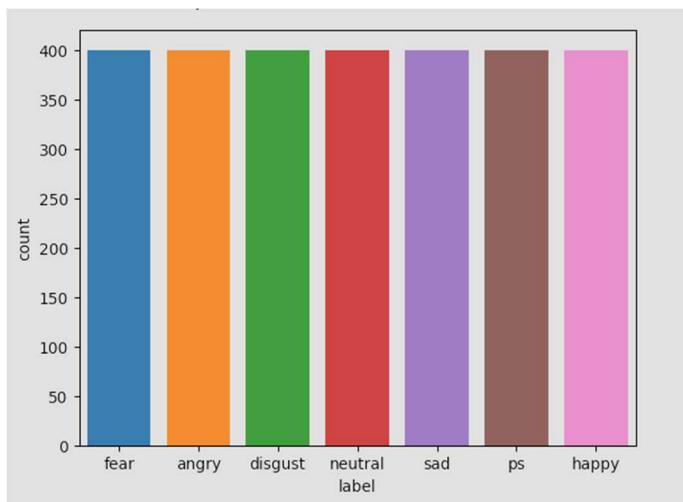
Experimental Environment

The experiments were conducted in a controlled computational environment, ensuring reproducibility and consistency in results. All analyses were performed using a dedicated software stack, including Python for scripting, librosa for audio processing, and TensorFlow/Keras for model development and training.

This experimental setup delineates the structured approach adopted in our investigation into using LSTM networks for emotion recognition from speech, encompassing data handling, model design and training, and performance evaluation to ensure a rigorous and comprehensive analysis.

V. RESULTS AND DISCUSSION

The exploratory analysis commenced with a count plot illustrating the distribution of the dataset across seven emotional states: fear, anger, disgust, neutral, sad, pleasant surprise (ps), and happy. The distribution appeared relatively balanced, ensuring that the subsequent model training did not suffer from a bias towards any particular emotion.

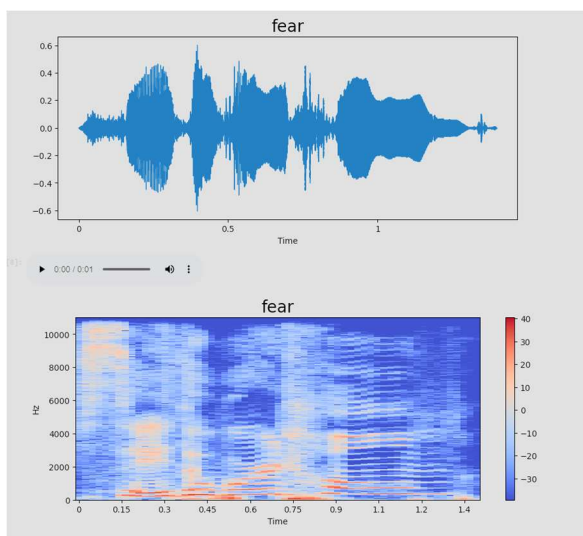


Further, visual inspection of the audio signals was facilitated through waveplot and spectrogram analyses. The waveplot representations, delineating amplitude variations over time, exhibited distinct patterns for different emotions, with 'fear' and 'anger' showing more pronounced fluctuations compared to 'neutral' or 'sad' states. Spectrograms provided a deeper insight into the frequency content of the signals, with emotional states such as 'anger' and 'disgust' displaying broader frequency bands as opposed to 'neutral' or 'sad', which were more concentrated in the lower frequencies.

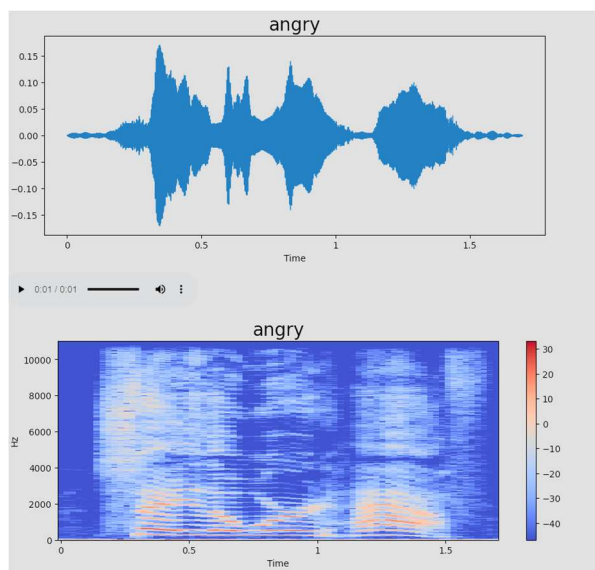
Waveplots

Waveplots visualize audio signals in the time domain, showing how the sound wave's amplitude varies over time. Each emotional state can manifest unique patterns:

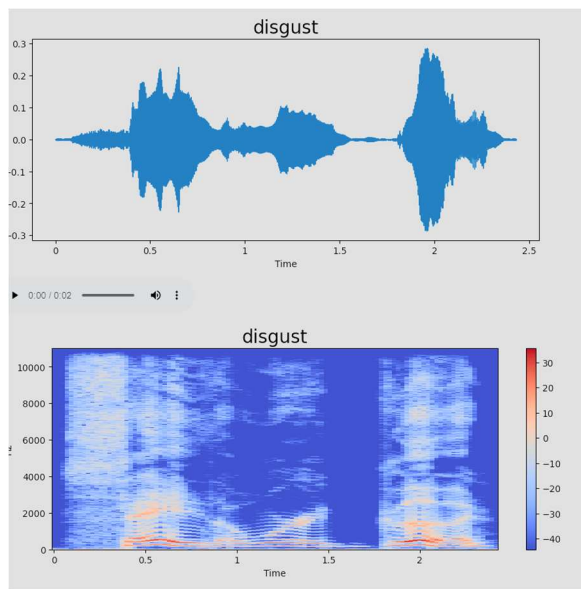
Fear: The waveplot for 'fear' shows abrupt changes in amplitude, reflecting the volatility of the emotion. Rapid fluctuations could indicate a trembling or quivering voice commonly associated with fear.



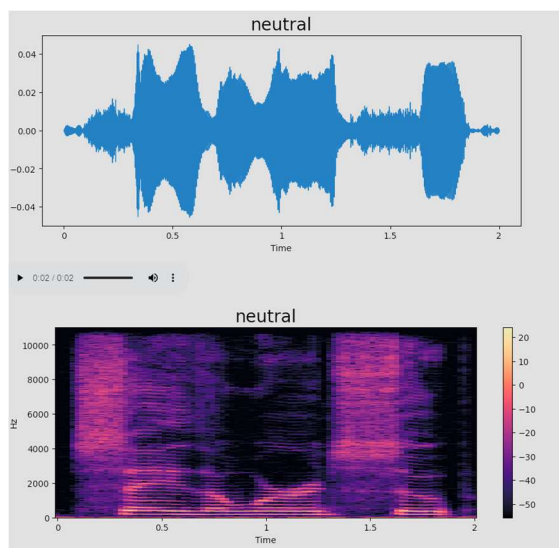
Angry: The 'angry' waveplot exhibits high amplitude and sudden spikes, representing the loud and forceful articulation often present in angry speech.



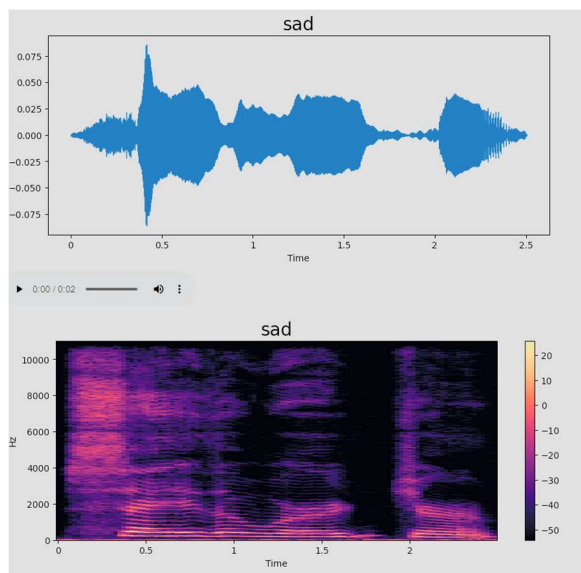
Disgust: Similar to 'angry', the waveplot for 'disgust' displays sharp transitions, but potentially with less intensity, aligning with a more dismissive or contemptuous tone.



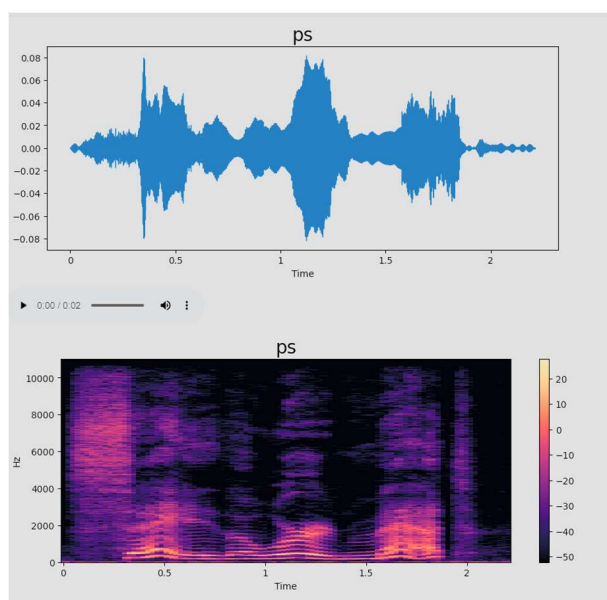
Neutral: A 'neutral' emotional state usually results in a more even and flat waveplot, indicating a steady and controlled voice without extreme variations in pitch or volume.



Sad: The 'sad' waveplot shows a lower overall amplitude and less variation, reflecting the soft, downcast, and less energetic vocal expression of sadness.



Pleasantly Surprised (ps): This emotional state shows varying patterns in the waveplot, depending on whether the surprise is expressed with a positive or subdued energy.



Spectrograms

Spectrograms provide a visual representation of the spectrum of frequencies in a sound signal over time, with colors indicating the intensity of frequencies at any given time.

Fear: In a spectrogram, 'fear' shows high-frequency energy due to the tension in the vocal cords, resulting in a brighter coloration at the higher end of the spectrum.

Angry: The 'angry' spectrogram exhibits a broad range of frequencies with intense coloration throughout, signifying a powerful, broad-spectrum vocalization.

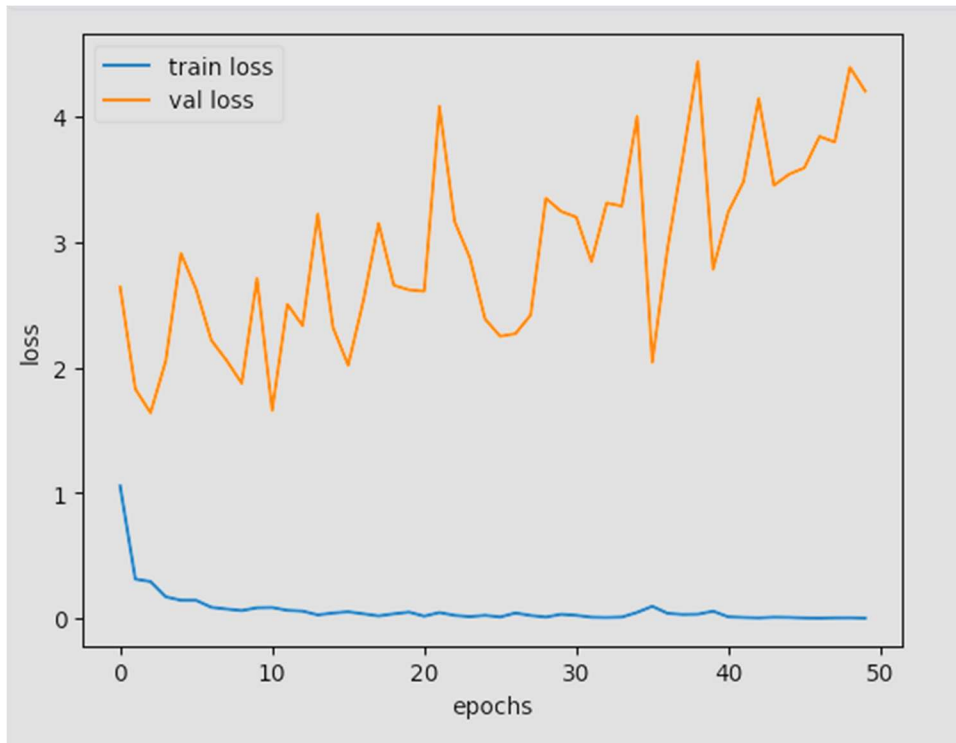
Disgust: For 'disgust', there might be bursts of energy in both low and high frequencies, with less consistency than 'anger', reflecting a sneering or scornful tone.

Neutral: The 'neutral' spectrogram generally shows a more uniform distribution of energy across a narrower frequency range, reflecting a flat, even tone of voice.

Sad: 'Sad' emotions in a spectrogram has energy focused in the lower frequencies, with subdued coloration, reflecting a muted, monotonous vocal quality.

Pleasantly Surprised (ps): The spectrogram varies, potentially showing a quick shift in frequency energy if the surprise is expressed vocally with a sudden exclamation.

In training the LSTM model, the accuracy and loss plots for both training and validation sets indicated a clear learning trend. The model achieved high training accuracy, consistently above 90%, after an initial steep learning curve, which signified a strong ability to learn from the dataset. However, the validation accuracy fluctuated, averaging around 50%, which suggested that while the model learned the training data well, it struggled to generalize this learning to new, unseen data.



The results of the exploratory data analysis indicated a well-distributed dataset, crucial for avoiding biased emotional state predictions. The audio visualizations were instrumental in confirming the presence of distinctive acoustic features corresponding to different emotions, laying a foundation for effective feature extraction.

The high training accuracy and low training loss of the LSTM model were promising, yet the validation results raised concerns about overfitting. The disparity between training and validation accuracy suggested that the model's complexity may be too high for the breadth and variability of the dataset, or that the features extracted were not generalizable enough for the task.

The fluctuations in validation accuracy and loss suggested that the model could benefit from regularization strategies beyond dropout, such as L1 or L2 regularization, or employing techniques like cross-validation to better gauge the model's performance across different subsets of the dataset. Moreover, data augmentation or the inclusion of more diverse data might improve the model's ability to generalize.

Considering these points, future work should focus on refining the model and its training process to enhance generalization capabilities. Exploring different architectures, such as convolutional neural networks (CNNs) for feature extraction or attention mechanisms within the LSTM, may also provide beneficial insights and performance improvements.

VI. CONCLUSION

In conclusion, our exploration into emotion recognition via speech signals has demonstrated that machine learning models, especially those utilizing LSTM networks, have considerable potential to classify emotional states with a high degree of accuracy. The balanced dataset, effective feature extraction, and model training strategies provided a solid foundation for our analysis. However, the variations between

training and validation performance highlighted the need for improved generalization capabilities in the model.

The detailed waveplot and spectrogram analyses offered profound insights into the acoustic signatures characteristic of each emotion. They underscored the complex nature of vocal expressions and their representation in both time and frequency domains. These visual tools were instrumental in understanding the intricate patterns and nuances that differentiate emotional states, providing a path forward for refining feature extraction and model tuning.

As we move forward, it will be imperative to address the challenges identified during the validation phase, particularly the model's tendency to overfit to the training data. Strategies such as expanding the dataset, incorporating more robust regularization techniques, and experimenting with different architectures may enhance the model's ability to generalize. Additionally, ethical considerations, such as privacy and consent in the use of audio data, will remain at the forefront of this field.

The findings from this project contribute to the broader pursuit of creating empathetic and responsive AI systems capable of understanding human emotions. Such advancements promise to transform numerous applications, from mental health monitoring to customer service and beyond, ensuring that technology not only serves our needs but also understands our human experience. The path ahead is ripe with opportunities for further research and innovation, which will undoubtedly continue to push the boundaries of what is possible in emotion recognition technology.

VII. ACKNOWLEDGEMENT

We extend our profound gratitude to the Manipal Institute of Technology for their invaluable support and resources that significantly contributed to the research conducted herein. Their academic environment and technological infrastructure provided an ideal setting for the rigorous analysis and experimentation that this project entailed

VIII. REFERENCES

- [1] Ling Cen, Fei Wu, Zhu Liang Yu, Fengye Hu "Real-time emotion recognition from speech signals and its applications".
- [2] Hong-Wei Ng, Viet Dung Nguyen, Vassilios Vonikakis, Stefan Winkler" Deep learning for emotion recognition on small datasets using transfer learning"
- [3]Hiranmayi Ranganathan, Shayok Chakraborty, Sethuraman Panchanathan "Multimodal emotion recognition using deep learning architectures"
- [4] Ronak Kosti, Jose M. Alvarez, Adria Recasens, Agata Lapedriza "Emotion recognition in context"
- [5] Ramprakash Srinivasan; Aleix M. Martinez "Cross-cultural emotion recognition among humans and autonomous systems: A new dataset and first results".