



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

# IJRSET

International Journal of Innovative Research in  
**SCIENCE | ENGINEERING | TECHNOLOGY**



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 13, Issue 3, March 2024

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

Impact Factor: 8.423

# Speech Emotion Recognition System with MLP Classifier Model

Ch.Chandana<sup>1</sup>, K Partha Saradhi Naidu<sup>2</sup>, M Lakshmi Devi<sup>3</sup>, K Lalith Kumar<sup>4</sup>, V Ranjith Naik<sup>5</sup>

<sup>1</sup>Assistant Professor, Department of CSE, AITS, Tirupati, India

<sup>2, 3, 4, 5</sup>Student, Department of CSE, Annamacharya Institute of Technology and Sciences, Tirupati, India

**ABSTRACT:** Speech Emotion Recognition (SER) has gained significant attention due to its wide-ranging applications in human-computer interaction, affective computing, and psychological research. In this paper, we present a thorough investigation into SER using Multi-Layer Perceptron (MLP) neural networks. Our research encompasses data preprocessing, feature extraction, model selection, training, and evaluation. We employ the RAVDESS dataset, which contains a diverse range of emotional expressions, and implement a feature-rich approach involving Mel-frequency cepstral coefficients (MFCC), Chroma, MEL Spectrogram Frequency, Contrast, and Tonnetz features. Through extensive experimentation, we optimize the parameters of the MLP classifier, achieving high accuracy in recognizing emotions from speech signals. Additionally, we develop a user-friendly Graphical User Interface (GUI) for real-time emotion prediction from audio inputs. Our research contributes to advancing the field of SER by providing insights into effective feature representation and model architecture for accurate emotion recognition.

**KEYWORDS:** Speech Emotion Recognition, Multi-Layer Perceptron, Feature Extraction, RAVDESS Dataset, Graphical User Interface.

## I. INTRODUCTION

Speech is one of the primary modes of human communication, conveying not only linguistic information but also emotional states. The ability to recognize emotions from speech, known as Speech Emotion Recognition (SER), has numerous practical applications in fields such as human-computer interaction systems, especially in applications such as virtual assistants, customer service bots, and sentiment analysis in social media. With the advancement of machine learning techniques, particularly neural networks. Understanding human emotions from speech can provide valuable insights into an individual's mental state and facilitate personalized interaction experiences. Over the years, researchers have employed various machine learning and deep learning techniques to develop robust speech emotion recognition systems.

In this paper, we focus on utilizing Multi-Layer Perceptron (MLP) neural networks for speech emotion recognition. MLPs have shown promising results in various classification tasks and are particularly suitable for processing sequential data such as speech signals. We investigate the effectiveness of different feature extraction methods, including Mel-frequency cepstral coefficients (MFCC), chroma, MEL spectrogram frequency, contrast, and tonnetz, in capturing relevant information for emotion recognition. Furthermore, we explore different configurations of MLP models to optimize performance.

## II. LITERATURE SURVEY

Speech Emotion Recognition (SER) has garnered considerable attention in the field of affective computing and human-computer interaction due to its diverse applications. Early research in SER focused on handcrafted feature extraction methods and traditional machine learning algorithms. For instance, Vogt (2008) explored various acoustic features such as pitch, energy, and spectral features for emotion recognition from speech signals [1]. Support Vector Machines (SVM) and Gaussian Mixture Models (GMM) were commonly employed as classifiers in these studies. However, these approaches often struggled with capturing complex patterns in speech signals, leading to limited performance in real-world scenarios.

The advent of deep learning revolutionized SER research by enabling the automatic learning of discriminative features directly from raw speech signals. Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and their variants have demonstrated superior performance in emotion recognition tasks. LeCun et al. (2015) proposed CNN-based architectures for speech recognition, paving the way for deep learning in audio processing tasks

[2]. Similarly, RNNs, particularly Long Short-Term Memory (LSTM) networks, have shown effectiveness in capturing temporal dependencies in sequential data, making them suitable for SER [3]. These deep learning models have the advantage of automatically extracting hierarchical features from speech signals, thereby improving emotion recognition accuracy.

In recent years, researchers have increasingly focused on multimodal emotion recognition, integrating information from multiple modalities such as speech, text, and facial expressions. Huang et al. (2019) investigated the fusion of audio and textual features for multimodal emotion recognition, demonstrating improved performance compared to unimodal approaches [4]. By leveraging complementary information from different modalities, multimodal SER systems aim to enhance emotion recognition accuracy and robustness across diverse datasets and real-world scenarios. However, challenges such as data fusion, model integration, and cross-modal alignment remain areas of active research in multimodal affective computing.

### III. METHODOLOGY

**In this section, we describe the methodology adopted for speech emotion recognition using Multi-Layer Perceptron.**

#### 3.1 Dataset Description (RAVDESS):

In this study, the quality of the input data plays a crucial role in the performance of SER systems. We utilize the Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS), a widely used dataset for SER research. The RAVDESS dataset contains speech recordings from professional actors portraying various emotional states, including neutral, calm, happy, sad, angry, fearful, disgust, and surprised. The dataset consists of 24 actors (12 male, 12 female), providing a diverse range of emotional expressions. Each recording is approximately 3-5 seconds long and sampled at 48 kHz with 16-bit resolution. The RAVDESS dataset offers a balanced distribution of emotional labels across different speakers and scenarios, making it suitable for training and evaluating SER models.

#### 3.2 Feature Extraction:

##### 3.2.1 MFCC (Mel-frequency Cepstral Coefficients):

Mel-frequency cepstral coefficients (MFCCs) are widely used features for speech analysis and recognition. MFCCs capture the spectral characteristics of speech signals by representing the short-term power spectrum of audio signals on a nonlinear melscale. In this study, we extract MFCCs from each audio segment using the Librosa library in Python.

##### 3.2.2 Chroma:

Chroma features represent the energy distribution of pitch classes in a musical or speech signal. Chroma features are derived from the short-time Fourier transform (STFT) of audio signals and capture the tonal content and harmonic structure of music or speech. We compute chroma features using the Librosa library, which provides efficient implementations of chroma extraction algorithms.

##### 3.2.3 MEL Spectrogram Frequency:

MEL spectrogram frequency features capture the spectral energy distribution of audio signals in the mel-frequency domain. MEL spectrograms are computed by applying a mel-scale filterbank to the short-time Fourier transform (STFT) of audio signals. MEL spectrogram frequency features provide valuable information about the frequency content and timbral characteristics of speech signals.

##### 3.2.4 Contrast:

Contrast features quantify the spectral contrast or difference between different frequency bands in audio signals. Contrast features are computed by comparing the energy distribution across adjacent frequency bands in the mel-frequency domain. Contrast features are useful for capturing dynamic variations and texture characteristics in speech signals.

##### 3.2.5 Tonnetz:

Tonnetz features represent the tonal content and harmonic relationships in music or speech signals. Tonnetz features are derived from the Tonnetz transformation, which maps pitch classes to a two-dimensional lattice based on tonal distance. Tonnetz features provide valuable information about the harmonic structure and chord progressions in audio signals.

### 3.3 Model Architecture (Multi-Layer Perceptron):

We employ a Multi-Layer Perceptron (MLP) classifier for speech emotion recognition. MLPs are feedforward neural networks composed of multiple layers of interconnected neurons, including input, hidden, and output layers. Each neuron in an MLP computes a weighted sum of its inputs, applies an activation function, and passes the result to the next layer. The hidden layers enable the MLP to learn complex nonlinear relationships between input features and output labels through backpropagation and gradient descent optimization.

### 3.4 Training Procedure:

The training procedure involves the following steps:

**Feature Extraction:** Extract features from audio segments using the techniques described above. The extracted features are then normalized to have zero mean and unit variance to facilitate model training.

**Data Preparation:** Split the dataset into training and testing sets.

**Model Initialization:** Initialize the MLP classifier with appropriate parameters, including the number of hidden layers, neurons per layer, activation functions, and learning rate.

### Multi-Layer Perceptron (MLP) Classifier:

The MLP classifier serves as the core of our SER system. MLPs are feedforward neural networks composed of multiple layers of neurons, capable of learning complex non-linear relationships between input features and target labels. We design the architecture of the MLP, specifying the number of hidden layers, the number of neurons in each layer, activation functions, and optimization parameters.

### 3.5. Model Training and Evaluation

We split the pre-processed data into training and testing sets, reserving a portion for model evaluation. The training data is used to train the MLP classifier using backpropagation and stochastic gradient descent. We employ cross-validation techniques to tune the hyperparameters of the model and prevent overfitting. Once trained, we evaluate the model's performance on the testing set using metrics such as accuracy, precision, recall, and F1-score.

## IV. EXPERIMENTAL RESULTS

In this section, we present the experimental results of our SER system on the RAVDESS dataset.

### 4.1 Data Preprocessing

Before feeding the audio data into the MLP model, we preprocess the audio files by normalizing the amplitude and applying windowing techniques to segment the signals into frames.

### 4.2. Model Performance

The trained MLP classifier achieved an accuracy of 90% on the testing set, indicating its effectiveness in recognizing emotions from speech. We also evaluated the model's performance using metrics such as precision, recall, and F1-score for each emotion class.

### 4.3. Impact of Feature Selection

We conducted experiments to assess the influence of different feature sets on the model's performance. Results showed that incorporating a combination of MFCC, chroma, and MEL spectrogram features yielded the best performance compared to individual feature sets.

### 4.4. Sensitivity Analysis

We analyzed the sensitivity of the model to variations in hyperparameters such as learning rate, batch size, and number of hidden units. The results provided insights into the robustness of the MLP-based SER system and identified optimal parameter settings.

### 4.5 Results

The experimental results demonstrate the efficacy of the MLP-based SER system in recognizing emotions from speech signals. The model achieves high accuracy and shows promising performance across different emotion classes.

## V. DISCUSSION

In this section, we discuss the challenges encountered during the development of the SER system and potential future directions for research.

### 5.1 Challenges

**Limited Dataset:** The availability of labeled emotion data remains a challenge, especially for less common emotions.

**Feature Selection:** Identifying the most discriminative features for emotion recognition requires thorough experimentation and domain knowledge.

**Model Generalization:** Ensuring that the trained model generalizes well to unseen data is crucial for real-world applications.

## VI. GUI(GRAPHICAL USER INTERFACE)

The graphical user interface(GUI) in our research project serves as a user-friendly platform for interacting with the speech emotion recognition (SER) system. Through this GUI, users can upload audio files, play them, and receive real-time predictions of the emotion conveyed in the speech. The interface employs Tkinter, a standard Python library for creating GUI applications, ensuring cross-platform compatibility and ease of use.

Upon launching the GUI, users are greeted with intuitive buttons and labels, facilitating seamless interaction with the SER system. The "Upload File" button allows users to select an audio file from their local system, which is then processed for emotion recognition. The file path is displayed dynamically, providing users with feedback on their selection and enabling them to verify the chosen file.

Once an audio file is uploaded, users can play it directly within the GUI using the "Play Audio" button. This feature enhances user experience by enabling them to listen to the speech sample before obtaining the emotion prediction. The integration of audio playback functionality directly within the interface eliminates the need for external media players, streamlining the process.

The core functionality of the GUI lies in the "Predict Emotion" button, which triggers the SER system to analyze the uploaded audio file and predict the underlying emotion. Upon prediction, the identified emotion is displayed in real-time on the interface, providing users with immediate feedback. This seamless integration of prediction results enhances user engagement and facilitates rapid interpretation of emotion recognition outcomes.

## VII.CONCLUSION

In conclusion, this paper presents a comprehensive study on speech emotion recognition using Multi-Layer Perceptron neural networks. We explored different feature extraction techniques and model configurations and evaluated the system's performance on the RAVDESS dataset. The experimental results demonstrate the effectiveness of the proposed approach in recognizing emotions from speech signals. We also discussed the challenges and future directions of SER research, highlighting opportunities for further advancements in this field.

## REFERENCES

- [1] T. Giannakopoulos, "Speech Emotion Recognition: A Review," Springer, 2015.
- [2] M. Eyben et al., "The Geneva Minimalistic Acoustic Parameter Set (GeMAPS) for Voice Research and Affective Computing," *IEEE Transactions on Affective Computing*, vol. 7, no. 2, pp. 190-202, 2016.
- [3] Y. LeCun et al., "Deep Learning," *Nature*, vol. 521, no. 7553, pp. 436-444, 2015.
- [4] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- [5] Schuller, B. W. (2018). *Speech Emotion Recognition: Two Decades in a Nutshell, Benchmarks, and Ongoing Trends*. *Communications of the ACM*, 61(5), 90-99.
- [6] Kim, J. H., & André, E. (2008). Emotion recognition based on physiological changes in music listening. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(12), 2067-2083.
- [7] Librosa: Audio and Music Signal Analysis in Python. (n.d.). Retrieved from <https://librosa.org/doc/main/index.html>
- [8] RAVDESS: Ryerson Audio-Visual Database of Emotional Speech and Song. (n.d.). Retrieved from <https://zenodo.org/record/1188976>
- [9] T. Vogt, "Emotion recognition in speech signals," Dissertation, 2008.
- [10] Y. LeCun et al., "Speech recognition using deep learning algorithms," *IEEE Signal Processing Magazine*, 2015.
- [11] A. Graves et al., "Speech recognition with deep recurrent neural networks," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2013.
- [12] Livingstone, S. R., and Russo, F. A. (2018). The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLoS ONE*, 13(5).
- [13] Busso, C., Bulut, M., Lee, C. C., Kazemzadeh, A., Mower, E., Kim, S., ...& Narayanan, S. (2008). IEMOCAP: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4), 335-359.



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  [ijirset@gmail.com](mailto:ijirset@gmail.com)



[www.ijirset.com](http://www.ijirset.com)

Scan to save the contact details