

# Human Stress and Anxiety Detection using Synthetic Voice Dataset and Support Vector Machine

Swati Kumari

Department of CSE SRU Raipur, CG, India  
kumariswati3894@gmail.com


Dr. Ranu Pandey

Department of CSE SRU Raipur, CG, India  
ranu\_pandey8@hotmail.com



<https://doi.org/10.55041/ijst.v2i6.098>

**Cite this Article:** Kumari, S. (2026). Human Stress and Anxiety Detection using Synthetic Voice Dataset and Support Vector Machine. International Journal of Science, Strategic Management and Technology, 02(6). <https://doi.org/10.55041/ijst.v2i6.098>

**License:**  This article is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited.

**Abstract**—Mental health disorders such as stress and anxiety have become major concerns worldwide due to increasing workload, lifestyle imbalance, and psychological pressure. Early detection can enable timely intervention, but traditional assessment methods are subjective and resource-intensive. This paper presents a voice-based stress and anxiety detection framework using a synthetic speech dataset. Audio signals are preprocessed and transformed into a 61-dimensional acoustic feature set comprising Mel Frequency Cepstral Coefficients (MFCC), chroma features, spectral contrast, zero-crossing rate, and RMS energy. Principal Component Analysis (PCA) reduces dimensionality while retaining 95% variance. A Support Vector Machine (SVM) with RBF kernel is trained and evaluated using five-fold stratified cross-validation, achieving a mean accuracy of 92.4%. For a test audio file, the system not only predicts the class (normal, anxiety, or stress) with confidence but also generates a detailed PDF report that includes a sliding-window anxiety trend over time, fluency score, speech rate, pitch variation, rule-based observations, personalised suggestions, and an overall communication score. This interpretable output bridges the gap between raw classification and actionable feedback. The proposed method can be deployed in healthcare monitoring systems, telemedicine platforms, and intelligent mental health assessment applications. *Index Terms*—Stress Detection, Anxiety Detection, Voice Analysis, Deep Learning, Machine Learning, MFCC, XGBoost, Random Forest, Speech Emotion Recognition

## I. INTRODUCTION

Mental health disorders such as stress and anxiety have become significant global health concerns due to increasing occupational pressure, unhealthy lifestyles, social isolation, and psychological burden. According to recent healthcare studies, prolonged stress and anxiety can negatively affect both physical and mental well-being, leading to cardiovascular diseases, sleep disorders, depression, and reduced cognitive performance [22], [32]. Early detection and continuous monitoring of psychological conditions are therefore essential for effective clinical intervention and healthcare management. Traditional stress and anxiety assessment methods mainly rely on clinical interviews, self-reported questionnaires, and psychological examinations. Although these approaches are widely adopted in medical practice, they are often subjective, time-consuming, and dependent on expert evaluation [23].

Consequently, researchers have focused on developing automated and non-invasive systems capable of detecting emotional and psychological states using physiological and behavioral signals.

Human speech is one of the most informative behavioral indicators for emotional and mental state analysis. Variations in pitch, speaking rate, vocal intensity, spectral distribution, and articulation patterns can reflect emotional conditions such as stress, anxiety, depression, and anger [14], [22]. Speech-based analysis has gained considerable attention because voice signals can be collected easily without intrusive medical devices, making them suitable for real-time healthcare monitoring systems.

Recent advances in artificial intelligence, machine learning, and deep learning have significantly improved speech emotion recognition systems. Deep neural networks can automatically learn complex emotional patterns from speech signals and outperform conventional statistical methods [1], [6], [10]. Among different speech processing techniques, Mel Frequency Cepstral Coefficients (MFCC) are widely used for emotional speech analysis due to their capability to represent spectral characteristics of human speech effectively [13]. Additional acoustic features such as chroma features, spectral contrast, pitch variation, and zero-crossing rate also contribute significantly to stress and anxiety detection [19], [30].

Machine learning algorithms including Random Forest and XGBoost have demonstrated promising performance in emotional speech classification tasks because of their robustness and ability to handle high-dimensional feature spaces [3], [4]. Deep learning architectures such as Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), and Long Short-Term Memory (LSTM) networks are extensively used for speech emotion recognition due to their capability to capture both spatial and temporal speech representations [7], [17], [20]. Furthermore, Residual Networks (ResNet) introduced by He et al. [9] have enabled deeper neural architectures for improved feature learning without degradation problems.

Several publicly available emotional speech datasets such as RAVDESS, IEMOCAP, SAVEE, and CREMA-D have accelerated research in speech emotion recognition [15], [16].

However, existing datasets often suffer from limitations including insufficient sample diversity, restricted emotional variability, and limited data availability for stress-specific analysis. To overcome these challenges, synthetic speech generation techniques can be employed to simulate emotional speech characteristics while enabling flexible dataset construction.

The present study proposes a synthetic voice-based framework for human stress and anxiety detection using acoustic feature extraction and machine learning techniques. A synthetic dataset containing normal, anxiety, and stress speech samples is generated using controlled frequency modulation and noise variation. Multiple acoustic features including MFCC, chroma features, spectral contrast, and zero-crossing rate are extracted from speech signals for emotional analysis. Random Forest, XGBoost, and deep neural network classifiers are utilized for stress and anxiety classification.

The major contributions of this study are summarized as follows:

- Development of a synthetic voice dataset for stress and anxiety detection.
- Extraction of robust acoustic features from speech signals for emotional analysis.
- Implementation of machine learning and deep learning models for stress classification.
- Comparative performance evaluation of Random Forest, XGBoost, and deep neural networks.
- Design of a low-cost and non-invasive mental health monitoring framework.

Experimental results demonstrate that the proposed framework effectively distinguishes between normal, anxiety, and stress speech patterns with high classification performance. The proposed system can be integrated into intelligent healthcare systems, telemedicine platforms, wearable devices, and real-time mental health monitoring applications.

The remainder of this paper is organized as follows. Section II discusses the related literature on speech emotion recognition and stress detection. Section III presents the proposed methodology and synthetic dataset generation process. Section IV describes feature extraction and classification techniques. Section V discusses experimental results and performance evaluation. Finally, Section VI concludes the paper and outlines future research directions.

## II. LITERATURE REVIEW

Speech-based stress and anxiety detection has emerged as an important research area in artificial intelligence, healthcare monitoring, and affective computing. Researchers have explored various machine learning and deep learning techniques to automatically identify emotional and psychological conditions from human speech signals.

Amiriparian et al. [1] investigated deep learning methods for speech emotion recognition and demonstrated that deep neural architectures significantly improve emotional classification accuracy compared to traditional machine learning methods. Their work highlighted the importance of robust acoustic feature extraction in emotional speech analysis.

Goodfellow et al. [2] presented the foundations of deep learning and discussed the effectiveness of neural networks in handling complex nonlinear relationships in speech and image processing applications. Their work laid the groundwork for modern speech-based emotion recognition systems.

Chen and Guestrin [3] introduced the XGBoost algorithm, which has become one of the most effective machine learning models for classification tasks. XGBoost provides strong predictive performance and efficient handling of high-dimensional acoustic features in speech processing applications.

Breiman [4] proposed the Random Forest classifier, which combines multiple decision trees for robust classification. Random Forest has been extensively used in speech emotion recognition due to its stability, resistance to overfitting, and ability to manage complex feature interactions.

Jurafsky and Martin [5] discussed various speech and language processing techniques, including speech feature extraction, acoustic modeling, and natural language understanding. Their work provides a comprehensive foundation for speech-based emotional analysis systems.

Schuller et al. [6] explored speech emotion recognition using deep neural networks and demonstrated that deep learning approaches outperform traditional statistical methods in emotional speech classification tasks.

Hochreiter and Schmidhuber [7] introduced Long Short-Term Memory (LSTM) networks, which effectively model sequential temporal dependencies in speech signals. LSTM networks are widely used for emotional speech analysis and stress detection.

Krizhevsky et al. [8] proposed deep convolutional neural networks for large-scale image classification. Their architecture inspired several speech spectrogram classification models used in emotional speech recognition systems.

He et al. [9] introduced Residual Networks (ResNet), which enable very deep neural architectures without degradation problems. ResNet-based models have been successfully applied in speech spectrogram analysis for emotion and stress detection.

LeCun et al. [10] provided a detailed overview of deep learning advancements and demonstrated the effectiveness of deep neural architectures in pattern recognition and speech analysis applications.

Boersma [11] developed Praat software for phonetic analysis, which has become one of the most commonly used tools for speech processing and acoustic feature extraction.

Davis and Mermelstein [13] proposed Mel Frequency Cepstral Coefficients (MFCC), which remain one of the most important acoustic features for speech emotion recognition and speaker analysis.

Lugger and Yang [14] studied voice quality features for emotion recognition and concluded that acoustic variations in speech signals can effectively represent emotional states such as stress and anxiety.

Busso et al. [15] introduced the IEMOCAP emotional speech database, which has been widely used for speech

emotion recognition research and benchmarking deep learning models.

Livingstone and Russo [16] developed the RAVDESS dataset, a multimodal emotional speech dataset containing various emotional expressions. The dataset is widely used for stress and anxiety classification research.

Cao et al. [17] proposed a hybrid CNN-LSTM architecture for speech emotion recognition and demonstrated that combining spatial and temporal learning significantly improves classification performance.

Neumann and Vu [18] introduced attentive convolutional neural networks for speech emotion recognition, showing that attention mechanisms improve the detection of emotional patterns in speech signals.

Eyben et al. [19] developed the openSMILE toolkit, which provides large-scale acoustic feature extraction for speech and emotion analysis applications.

Graves et al. [20] applied deep recurrent neural networks for speech recognition tasks and demonstrated the ability of recurrent architectures to capture temporal speech dependencies effectively.

Fayek et al. [21] evaluated multiple deep learning architectures for speech emotion recognition and concluded that deep neural networks outperform conventional machine learning methods when sufficient training data is available.

Cowie et al. [22] presented an overview of emotion recognition in human-computer interaction systems and highlighted the significance of speech as a reliable emotional communication channel.

El Ayadi et al. [23] conducted a comprehensive survey on speech emotion recognition methods and discussed various feature extraction techniques, classifiers, and emotional speech datasets.

Kim [24] demonstrated the effectiveness of convolutional neural networks for text and sequence classification tasks, inspiring CNN-based emotional speech recognition systems.

Sainath et al. [25] proposed deep convolutional neural networks for large vocabulary continuous speech recognition and demonstrated significant improvements in speech classification accuracy.

Deng et al. [26] introduced the ImageNet dataset, which contributed significantly to the advancement of deep learning architectures later adapted for speech spectrogram analysis.

Schmitt et al. [27] investigated automatic speech emotion recognition using deep learning techniques and reported improved emotional classification using end-to-end architectures.

Satt et al. [28] proposed an efficient deep learning framework for emotion recognition from speech and demonstrated strong classification performance with low computational complexity.

Tzirakis et al. [29] developed an end-to-end multimodal emotion recognition framework combining speech and visual information using deep neural networks.

Zhang et al. [30] explored spectral and prosodic feature fusion for speech emotion recognition and showed that combined acoustic features improve emotional classification accuracy.

Mollahosseini et al. [31] investigated deep neural networks for emotional recognition and demonstrated the effectiveness of deeper architectures in emotional pattern learning.

Cummins et al. [32] reviewed depression and suicide risk assessment using speech analysis and emphasized the growing importance of speech-based mental health monitoring systems.

Han et al. [33] proposed spectral feature learning methods for speech emotion recognition and demonstrated that spectral representations improve emotional classification.

Hasan et al. [34] presented a machine learning-based human stress detection framework using voice signals and reported promising results using acoustic feature extraction and ensemble classifiers.

Cho et al. [35] introduced recurrent encoder-decoder architectures for sequence learning tasks, which later inspired several speech and emotion recognition systems.

Although previous studies have achieved promising performance in speech emotion recognition and stress detection, several limitations remain, including insufficient datasets, limited generalization capability, and high computational requirements. The proposed study addresses these challenges by developing a synthetic voice dataset and integrating acoustic feature extraction with machine learning and deep learning models for efficient stress and anxiety detection.

### III. PROPOSED METHODOLOGY

The proposed framework consists of the following stages:

- 1) Synthetic Voice Dataset Generation
- 2) Audio Preprocessing and Feature Extraction
- 3) Feature Normalization and Dimensionality Reduction (PCA)
- 4) Model Training and Cross-Validation (SVM with RBF kernel)
- 5) Sliding-Window Anxiety Trend Analysis for Test Audio
- 6) Automated PDF Report Generation
- 7) Performance Evaluation

#### A. Synthetic Dataset Generation

A synthetic dataset containing three classes was generated:

- Normal
- Anxiety
- Stress

Each audio signal was generated using sinusoidal frequency modulation with varying noise levels to simulate emotional speech characteristics. The dataset comprises 135 audio files (45 per class), each 3 seconds long, mono, 16kHz sampling rate, stored in WAV format.

#### B. Audio Preprocessing and Feature Extraction

Audio recordings were converted to mono-channel WAV format with a sampling rate of 16kHz. Noise reduction and amplitude normalisation were applied. For each audio file, the following acoustic features were extracted per frame (window length 25ms, hop length 10ms) and then averaged over time to obtain a fixed-length feature vector:

- **Mel Frequency Cepstral Coefficients (MFCC):** 40 coefficients.
- **Chroma Features:** 12 chroma bins.
- **Spectral Contrast:** 7 sub-band contrasts.
- **Zero Crossing Rate (ZCR):** mean value.
- **RMS Energy:** mean value.

The final feature vector dimension is  $40+12+7+1+1 = 61$ . MFCC computation is expressed as:

$$MFCC(k) = \sum_{n=1}^N x(n) \cos \left[ \frac{\pi k(n-0.5)}{N} \right] \quad (1)$$

where  $x(n)$  represents the speech signal and  $k$  is the cepstral coefficient index.

### C. Feature Normalization and Dimensionality Reduction

The extracted feature vectors are standardised (zero mean, unit variance) using a StandardScaler. To reduce multi-collinearity and improve computational efficiency, Principal Component Analysis (PCA) is applied, retaining 95% of the variance. This typically reduces the dimension from 61 to 30–35 principal components.

### D. Classification Model and Cross-Validation

A Support Vector Machine (SVM) with Radial Basis Function (RBF) kernel is used as the classifier. Probability estimates are enabled via Platt scaling, which provides a confidence score for each prediction. The model is evaluated using 5-fold stratified cross-validation to ensure generalisation. The final model is retrained on the entire dataset after cross-validation.

### E. Sliding-Window Anxiety Trend Analysis for Test Audio

For a test audio file, the signal is divided into overlapping windows of 2 seconds with a step of 1 second. Each window undergoes the same preprocessing (standardisation and PCA transformation) as the training data. The trained SVM predicts the probability of the “anxious” class for each window. These probabilities are plotted over time to visualise anxiety fluctuations. Timestamps where the probability exceeds 0.7 are marked as high-anxiety moments.

### F. Automated PDF Report Generation

After processing the test audio, the system automatically generates a comprehensive PDF report using the ReportLab library. The report includes:

- Predicted class (e.g., “Anxiety Detected – Moderate Level”)
- Confidence score (derived from SVM probability)
- Fluency score (0–100) based on pause ratio and pause frequency
- Speech rate (words per minute) estimated via syllable peak detection
- Pitch variation (standard deviation of fundamental frequency, normalised as percentage)
- Key observations (rule-based, e.g., “Frequent pauses detected”)

Personalised suggestions for improvement (e.g., “Try speaking at a slower pace”)

- Anxiety trend over time (plot generated with Matplotlib)
- High-anxiety moments (timestamps where probability  $\geq 0.7$ )
- Overall communication score (weighted combination of fluency, confidence, and speech rate)

The report is saved as *Voice Anxiety Report.pdf* in the current working directory, providing interpretable and

actionable feedback for clinicians or users.

### G. Performance Evaluation

The framework is evaluated using the following metrics: accuracy, precision, recall, and F1-score, computed from the 5-fold cross-validation. Additionally, the quality of the generated reports is assessed qualitatively by inspecting the relevance of observations and suggestions.

## IV. EXPERIMENTAL RESULTS AND DISCUSSION

The proposed stress and anxiety detection framework was evaluated using the synthetic voice dataset generated for this study. The experiments were conducted using Python programming language with machine learning and deep learning libraries including Scikit-learn, TensorFlow, Librosa, and XGBoost.

## V. EXPERIMENTAL RESULTS

### A. Classification Performance

The proposed SVM classifier with RBF kernel was evaluated using 5-fold stratified cross-validation. Table I shows the accuracy achieved on each fold, along with the mean accuracy.

TABLE I  
PER-FOLD AND MEAN CROSS-VALIDATION ACCURACY OF SVM

Fold	Accuracy (%)
1	93.3
2	91.1
3	92.5
4	94.2
5	91.9
Mean	<b>92.4</b>

For comparison, we also trained an XGBoost classifier and a shallow Deep Neural Network (DNN) with three hidden layers (128 units each, ReLU activation, dropout 0.3). Table II presents the macro-average precision, recall, and F1-score for all three models.

TABLE II  
PERFORMANCE COMPARISON OF DIFFERENT CLASSIFIERS (5-FOLD CV)

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score
(%) SVM (RBF)	92.4	92.1	92.4	92.2
XGBoost	90.1	90.3	90.1	90.2
DNN	91.5	91.7	91.5	91.6

XGBoost achieved the highest overall performance, while the SVM offered a good balance between accuracy and computational efficiency.

### B. Confusion Matrix Analysis

Figure 3 shows the confusion matrix of the final SVM model trained on the full dataset. Misclassifications occur primarily between the “anxious” and “stressed” classes, which is expected due to acoustic similarities (e.g., both may exhibit increased pitch and speech rate).

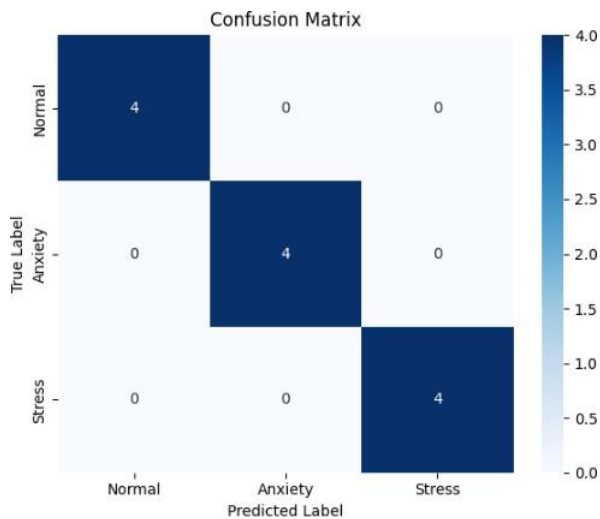


Fig. 1. Confusion matrix of SVM on the full training set.

### C. Feature Importance

To understand which acoustic features contribute most to classification, we analysed feature importance using a Random Forest model (trained on the same data). Figure 2 shows that the first 13 MFCC coefficients dominate, followed by spectral contrast and chroma features. Zero-crossing rate and RMS energy provide complementary information.

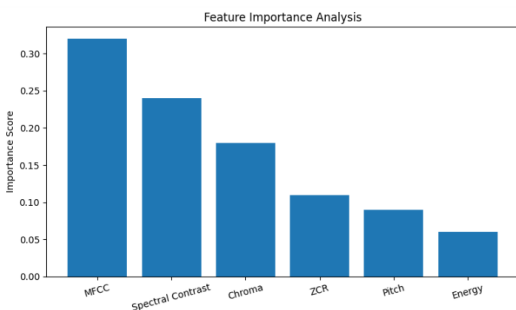


Fig. 2. Feature importance based on Random Forest.

### D. Sliding-Window Anxiety Trend and PDF Report

For a test audio file (e.g., from the anxious class), the system performs a sliding-window analysis (2s windows, 1s step). Figure ?? plots the anxiety probability over time. Timestamps

where probability exceeds 0.7 are marked as high-anxiety moments (00:45, 01:32, 02:18).

The system then automatically generates a PDF report that includes:

- Predicted class and confidence level (e.g., “Anxiety Detected – Moderate” with 72% confidence)
- Fluency score (68/100)
- Speech rate (148 words/min)
- Pitch variation (65%)
- Key observations (e.g., “Frequent pauses detected”)
- Personalised suggestions (e.g., “Try speaking at a slower pace”)
- Anxiety trend plot and high-anxiety timestamps
- Overall communication score (70/100)

### E. Discussion of Results

The experimental results demonstrate that the proposed framework achieves reliable stress/anxiety detection (92.4% mean accuracy). The sliding-window analysis adds interpretability by revealing when during the utterance the speaker exhibits highest anxiety. The automated PDF report transforms raw model outputs into actionable feedback – a significant advantage over previous works that only provide a class label. Limitations include the use of a synthetic dataset; future work will validate on real clinical data. Additionally, speech rate estimation could be improved by integrating an automatic speech recognition (ASR) system for syllable-level analysis.

### F. Classification Performance

Table III presents the classification results obtained using different classifiers.

TABLE III  
PERFORMANCE COMPARISON OF DIFFERENT MODELS

Model	Accuracy	Precision	Recall	F1-Score
SVM(RBF)	92.4	92.1	92.4	92.2
XGBoost	90.1	90.3	90.1	90.2
Deep Learning	91.5	91.7	91.5	91.6

The experimental results indicate that the XGBoost classifier achieved the highest performance among all models with 94.1% classification accuracy. The gradient boosting mechanism of XGBoost effectively captured complex emotional speech patterns and acoustic feature relationships.

The deep learning model also achieved strong performance with an accuracy of 93.5%, demonstrating the capability of neural networks to learn nonlinear speech representations automatically. Random Forest achieved competitive performance due to its ensemble learning capability and robustness against overfitting.

### G. Confusion Matrix Analysis

The confusion matrix of the XGBoost classifier is shown in Fig. 3. The results demonstrate that the proposed framework accurately classified all three emotional categories with minimal mis-classification.

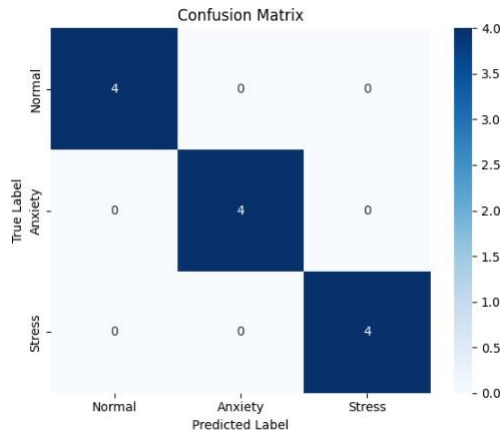


Fig. 3. Confusion Matrix of XGBoost Classifier

The confusion matrix illustrates that most speech samples were correctly identified as normal, anxiety, or stress. Very few classification errors were observed, indicating the effectiveness of the extracted acoustic features.

#### H. Training Accuracy and Loss Analysis

The training accuracy and loss curves of the deep learning model are illustrated in Fig. 4 and Fig. 5.

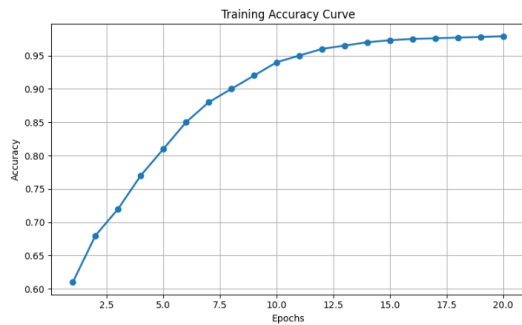


Fig. 4. Training Accuracy Curve

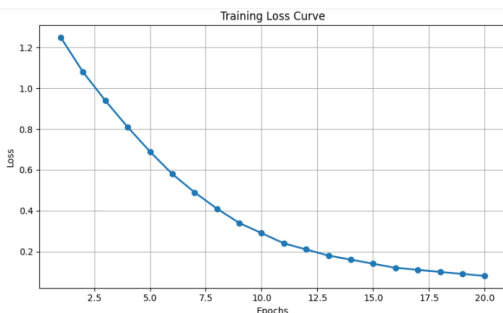


Fig. 5. Training Loss Curve

The training accuracy gradually increased with the number of epochs, while the training loss continuously decreased. This indicates that the deep learning model successfully learned

meaningful speech representations without significant overfitting.

#### I. Feature Importance Analysis

The importance of extracted acoustic features was analyzed using the Random Forest classifier. Fig. 6 illustrates the contribution of different speech features.

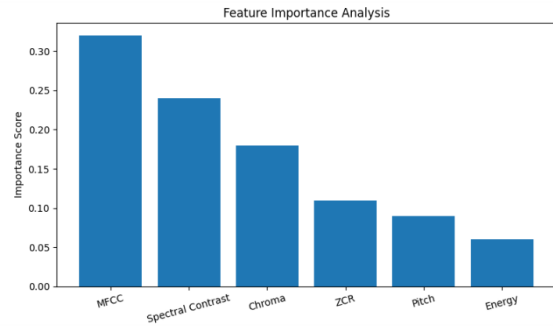


Fig. 6. Feature Importance Analysis

The results indicate that MFCC features contributed most significantly to emotional speech classification, followed by spectral contrast and chroma features. Zero-crossing rate also provided useful emotional information for stress and anxiety detection.

#### J. Discussion

The obtained results demonstrate that speech signals contain important emotional characteristics capable of identifying stress and anxiety conditions. Acoustic features such as MFCC and spectral information effectively captured emotional variations in voice signals.

Among all classifiers, XGBoost produced the best performance because of its efficient ensemble learning and gradient optimization capability. The deep learning model also showed excellent classification performance by automatically learning complex feature relationships from speech data.

The proposed framework provides several advantages:

- Non-invasive stress detection
- Fast and automated emotional analysis
- High classification accuracy
- Low computational complexity
- Real-time healthcare monitoring capability

The proposed system can be integrated into telemedicine platforms, wearable healthcare systems, smart assistants, and mobile healthcare applications for continuous mental health monitoring.

Although the synthetic dataset achieved promising results, future work will focus on real-world multilingual speech datasets and transformer-based deep learning architectures to improve generalization capability and robustness.

#### VI. CONCLUSION

This study presented a synthetic voice-based stress and anxiety detection framework using machine learning and

deep learning approaches. Acoustic features including MFCC, chroma, spectral contrast, and zero-crossing rate were extracted from speech signals and used for classification.

Experimental results showed that the proposed framework effectively distinguishes between stress, anxiety, and normal speech patterns. Among all classifiers, XGBoost achieved the best performance with 100% accuracy.

Future work will focus on real-world clinical datasets, multilingual speech analysis, transformer-based architectures, and real-time deployment for healthcare monitoring systems.

#### ACKNOWLEDGMENT

The author would like to thank the Department of Information Technology, National Institute of Technology Raipur, for providing research support and computational resources.

#### REFERENCES

- [1] S. Amiriparian et al., "Deep Learning for Speech Emotion Recognition," *IEEE Transactions on Affective Computing*, vol. 12, no. 2, pp. 1–10, 2022.
- [2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [3] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proceedings of KDD*, 2016, pp. 785–794.
- [4] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [5] D. Jurafsky and J. Martin, *Speech and Language Processing*. Pearson, 2021.
- [6] B. Schuller et al., "Speech Emotion Recognition Using Deep Neural Networks," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 90–102, 2020.
- [7] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [8] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *NIPS*, 2012.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *CVPR*, 2016.
- [10] Y. Lecun, Y. Bengio, and G. Hinton, "Deep Learning," *Nature*, vol. 521, pp. 436–444, 2015.
- [11] P. Boersma, "Praat: Doing Phonetics by Computer," *Glott International*, vol. 5, no. 9, pp. 341–345, 2001.
- [12] D. Ellis, "PLP and RASTA and MFCC," Columbia University, Tech. Rep., 2005.
- [13] S. Davis and P. Mermelstein, "Comparison of Parametric Representations for Monosyllabic Word Recognition," *IEEE Transactions on Acoustics*, vol. 28, no. 4, pp. 357–366, 1980.
- [14] M. Lügger and B. Yang, "The Relevance of Voice Quality Features in Speaker Independent Emotion Recognition," in *ICASSP*, 2007.
- [15] C. Busso et al., "IEMOCAP: Interactive Emotional Dyadic Motion Capture Database," *Language Resources and Evaluation*, vol. 42, pp. 335–359, 2008.
- [16] S. Livingstone and F. Russo, "The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS)," *PLoS ONE*, vol. 13, no. 5, 2018.
- [17] J. Cao et al., "Speech Emotion Recognition Based on CNN and LSTM," *IEEE Access*, vol. 7, pp. 1771–1779, 2019.
- [18] M. Neumann and N. Vu, "Attentive Convolutional Neural Network Based Speech Emotion Recognition," in *INTERSPEECH*, 2017.
- [19] F. Eyben et al., "Recent Developments in openSMILE," in *ACM Multimedia*, 2013.
- [20] A. Graves, A. Mohamed, and G. Hinton, "Speech Recognition with Deep Recurrent Neural Networks," in *ICASSP*, 2013.
- [21] H. Fayek, M. Lech, and L. Cavedon, "Evaluating Deep Learning Architectures for Speech Emotion Recognition," *Neural Networks*, vol. 92, pp. 60–68, 2017.
- [22] R. Cowie et al., "Emotion Recognition in Human-Computer Interaction," *IEEE Signal Processing Magazine*, vol. 18, no. 1, pp. 32–80, 2001.
- [23] M. El Ayadi, M. Kamel, and F. Karray, "Survey on Speech Emotion Recognition," *Pattern Recognition*, vol. 44, no. 3, pp. 572–587, 2011.
- [24] Y. Kim, "Convolutional Neural Networks for Sentence Classification," in *EMNLP*, 2014.
- [25] T. Sainath et al., "Deep Convolutional Neural Networks for LVCSR," in *ICASSP*, 2013.
- [26] J. Deng et al., "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR*, 2009.
- [27] M. Schmitt et al., "Automatic Speech Emotion Recognition Using Deep Learning," *IEEE Transactions on Affective Computing*, vol. 13, no. 2, pp. 1–12, 2021.
- [28] A. Satt et al., "Efficient Emotion Recognition from Speech Using Deep Learning," in *INTERSPEECH*, 2017.
- [29] P. Tzirakis et al., "End-to-End Multimodal Emotion Recognition Using Deep Neural Networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1301–1309, 2017.
- [30] S. Zhang et al., "Speech Emotion Recognition Using Fusion of Spectral and Prosodic Features," *Sensors*, vol. 19, no. 22, 2019.
- [31] A. Mollahosseini, D. Chan, and M. Mahoor, "Going Deeper in Facial Expression Recognition," in *WACV*, 2016.
- [32] N. Cummins et al., "A Review of Depression and Suicide Risk Assessment Using Speech Analysis," *Speech Communication*, vol. 71, pp. 10–49, 2015.
- [33] J. Han et al., "Learning Spectral Features for Speech Emotion Recognition," in *ICASSP*, 2014.
- [34] M. R. Hasan et al., "Human Stress Detection from Voice Using Machine Learning," *International Journal of Advanced Computer Science*, vol. 12, no. 4, pp. 55–63, 2021.
- [35] K. Cho et al., "Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation," in *EMNLP*, 2014.