

# Speech Emotion Recognition Using Deep Learning Algorithm

Ipsita Roy Barman  
Master of Technology  
Dept of Data Science and Business System  
SRM Institute of Science & Technology  
Chennai, India

Dr. A. Shanthini  
Dept of Data Science and Business System  
SRM Institute of Science & Technology  
Chennai, India

**Abstract**—AI is the emerging technology and is being extensively used now a days in every field and in every organization to gain more insights and hidden patterns in data. Data is the new oil for every organization to run and make huge revenues if one is able to find out the story in the past data and then can recommend based on it in the future. AI not only works effectively in data science related fields but is now taking grip for every work fields as well. AI devices become more interesting when they start to act more like human rather than robotic as the customer interactions increase and so does the revenue of the organization. So, emotion recognition of the users is a very important part in AI development where the device is trained to recognize the correct emotion of the customer behind his / her speech and conversation with the device. We know that random forest uses a bunch of decision trees and tries to correct the errors found in the previous steps to rectify them and then increase the accuracy. In this article, we use CNN based model which is a very popular deep learning neural network-based algorithm to classify every emotion of an individual and then used audios to test the model. In this project, CNN model is used and 79.6% accuracy is achieved.

**Keywords:** Speech emotion, deep learning, CNN, Artificial Intelligence, Decision trees, Random Forest.

## I. INTRODUCTION

Speech and words are the most necessary tool we use to interact with each other. Now a days AI has become an integral part of our life specially the AI based technologies. We love to interact with them the whole day and give them requests or talk some random things. As a result, we as a customer expect some interesting replies from them and that too like a human being with emotions. In this fast-changing world and generation we are left with a few friends and we have no time to meet or interact with them. So, in this new world, AI is our new friend, philosopher, and guide too in many cases. Hence our inputs and commands are considered as data for the above-mentioned problem statement and with the help of CNN we make the audio tuning more perfect and the AI more lifelike. The model is trained in such a way to recognise emotions behind every input of the customers.

### 1.1. Scope

Data is the future of any organisation and thus AI is the emerging technology which uses data as the input and in return gives amazing experience to its customers. The major risk faced in the usage of AI is constant fall of customers interest in using it. This issue can be solved easily by making the output responses of AI more interesting and natural sounding. For this thing to happen, we need to recognise the

emotion behind the customers and then teach the AI to do the same thing. Once it will be able to recognise the emotion behind the speech, it will respond in a more lucid and human like. The neural networks in architecture like CNN works similar to our brain neural structures with the nodes to transfer data from one layer to another. Similarly in neural network, data is passed from input to output layer through a series of hidden layers in between. Hence, this method is proposed in this project so that the AI can be transformed successfully into a complete human like structure. Our model will not only recognise the emotions behind every speech of our customers but can also cheer the customers when they will be sad or depressed. Once this model is built successfully with great precision and accuracy and minimum errors, we can collect data of our customers to classify them into depressed and non-depressed individuals based on the duration of their past data of being in negative emotions such as sadness, anger, fearful, disgust, etc. This entire process will not be possible without the technologies and machine learning algorithms cause data is involved. Data will in fact act as the raw material for this project, without which progress of any kind in this work is just futile. In future this model can be applied on any AI devices or AI applications and collecting customers data from the devices will be very important and can be used to recommend various things for their future needs.

### 1.2 Methods

#### 1.2.1 Deep Learning

This branch of machine learning is very useful as it involves the use of neural networks along with the multiple layers and can process complex data sets. It even is capable for developing artificial intelligence systems that can learn on their own, keep on improving themselves on their own without explicitly being programmed by us for a particular task. The algorithms can recognize patterns and making predictions based on large amount of data and this feature indeed is used in the recognition of emotions behind the speech of our end users.

#### 1.2.2 CNN

It is a type of neural network that is commonly used in image and video analysis tasks and as it has various convolutional layers so it can learn various spatial hierarchies of features and this in turn can scan the input data with a set of filters to extract relevant features. They typically consist of convolution layers, pooling layers, and fully connected layers. They are responsible for detecting low level features and pooling layers downsample the feature maps to reduce the

dimensionality of the data. Fully connected layers combine the extracted features to make a final prediction. As they are responsible for processing structured data, this algorithm is useful in extracting the correct emotion out of the audios we used in this project.

## II. LITERATURE STUDY

Many studies on speech emotion recognition have been conducted and various new algorithms and methods have been proposed by authors. This is a rapidly growing field of research which attracted significant attention from researchers in recent years. Let's look at some of the important significant reviews:

S. Bhattacharyya and S. Poria [1] in 2017 provides a comprehensive review of various techniques and methods and it discusses the challenges associated with the task such as the standardized database for emotions is lacking and the need for robust feature extraction techniques.

V. Balakrishnan and S. Sivakumar [2] in 2019 propose a new method for emotion recognition using speech features such as pitch, energy, and formants. The proposed method was tested on the Berlin Emotional Speech Database and achieves an accuracy of 88.3%.

P. Rao and P. S. S. Avadhani [3] in 2020 proposed recent advances in deep learning based approaches for SER and discussed the use of RNN for feature extraction and classification. The paper also covers the use of transfer learning.

A. Gunavati and R. Bhavani [4] in 2021 presented a comprehensive survey of deep learning techniques and covers recent developments in multimodal emotion recognition which involves combining speech with other modalities such as facial expressions and physiological signals.

R. Mitra [5] in 2019 presented a paper which contains a comprehensive survey of various techniques and methods used for SER., including traditional machine learning approaches, deep learning techniques and hybrid models.

T. A. Rahman [6] in 2020 fine-tuned the pretrained model on the Ryerson Audio Visual database of emotional speech and song dataset using transfer learning technique and achieved 70% accuracy.

M. A. H. Akanda [7] in 2020 investigated the effect of speech enhancement techniques on SER using deep neural networks. The author compared the performance of different enhancement methods such as spectral subtraction, and MMSE-based noise reduction, on the Emo-DB dataset. The results show that speech enhancement can significantly improve the accuracy of emotion recognition.

W. Wang [8] in 2021 proposed a multi task learning approach for SER that integrates both acoustic and lexical features. The authors use a deep neural network with shared and task-specific layers to jointly learn the feature

representations for emotion recognition and sentiment analysis.

G. F. Adewumi [9] in 2021 presented a review of various features, classification techniques, and datasets used for SER. It discusses the challenges associated with SER such as the variability of emotions across speakers and cultures.

M. R. Islam [10] in 2020 compared the performance of various deep learning-based approaches for SER, including LSTM networks and hybrid models. The authors evaluated the models on the IEMOCAP dataset and showed that hybrid model performs the best.

## III. PROPOSED METHODOLOGY

### 3.1 System Proposed

Deep neural networks (DNNs) are a type of machine learning model that has gained significant attention in recent years due to their ability to learn complex patterns and representations from large datasets. DNNs have shown promising results in various applications, including computer vision and speech emotion recognition.

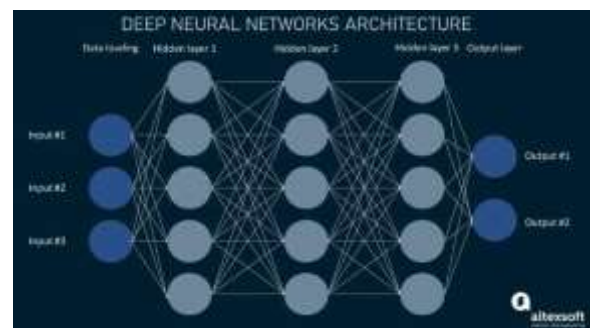


Fig 1. Deep Neural Network

In other words, DNN can be represented as a function  $f(x; a)$  where  $x$  is the input,  $a$  represents the model parameters, and  $f(x; a)$  represents the output of the network for a given input  $x$ .

### 3.2 Dataset-Taken

The dataset contains 1440 files, 7356 recordings, 24 professional actors who acted out different emotional states including, calm, happy, sad, angry, fearful, surprise and disgust.

### 3.3 Dataset Preprocessing

Dataset pre-processing is a crucial step in machine learning and data analysis that involves preparing raw data to be used as input to a model. The goal of pre-processing is to transform the raw data into a format that can be easily analysed and interpreted by a machine learning algorithm. Data pre-processing steps include data cleaning, data normalization, data transformation, data encoding, data splitting and feature selection.

#### 3.3.1 Augmentation Of Data

It is a technique used to artificially increase the size of a dataset by creating new, but similar, versions of existing data. The goal of data augmentation is to improve the performance of machine learning models by providing more diverse examples for the model to learn from. This is particularly used when working with small datasets where overfitting can be a concern. There are various methods of data augmentation and the specific techniques used depend on the type of data being augmented. The purpose of data augmentation is to create new examples that are like the original data but not identical. The augmented data should represent the natural variations in the data that the model is likely to encounter in real-world scenarios. Data augmentation is an important tool in machine learning and can significantly improve the performance of models, especially when working with small datasets. The idea behind data augmentation is that by creating additional training samples that are similar but not identical to the original data, a model can better generalize to new, unseen data. This can help to reduce overfitting and improve the overall performance of the model. Data augmentation is commonly used in computer vision and natural language processing tasks, but it can also be applied to other types of data such as audio or sensor data. The specific techniques used for data augmentation depend on the type of data and the specific task at hand, and it is often necessary to experiment with different augmentation methods to find the best approach for a particular problem statement. Data augmentation is typically done by applying a set of transformation rules to the existing data samples, such as rotating, flipping, or zooming in on images or changing the pitch or speed of audio samples. So this data augmentation is an irreplaceable step in this project.

### 3.4 Deep-Learning

Deep learning for speech emotion recognition refers to the use of deep neural networks to classify and analyze emotional states in speech signals. The goal of SER is to automatically detect the emotional content of speech, such as happiness, sadness, anger, or fear. It typically involves multiple layers of artificial neurons that are trained on large amounts of labelled data to identify patterns and relationships between different features of speech signals, such as pitch, duration, and spectral characteristics. These models can be trained using various architectures such as Convolutional Neural Networks, RNN, or hybrid model that combines both RNN and CNN. The training process for deep learning models involves feeding the network with large amounts of labelled speech data and adjusting the weights of the network's parameters through a process called backpropagation in order to minimize the error between the predicted emotion labels and the true labels. The resulting trained model can be used to classify the emotional content of news speech signals. Deep learning models for SER have shown promising results in recent years and have been used in a variety of applications. However there are still many challenges to be addressed in this field such as dealing with variability

in speech signals, addressing class imbalance in emotional labels, and improving the interpretability of the models.

### 3.5 Convolutional Neural Network (CNN)

The CNN for the speech emotion recognition typically involves several layers that are redesigned to extract relevant features from speech signals and then classify the emotional content.

1. **Input layer:** this layer receives the raw speech signal as an input. This layer can be preprocessed with techniques such as Mel-frequency cepstral coefficients or filter bank energies to extract relevant features.
2. **Convolutional layer:** this layer applies a set of filters to the input signal which extracts local features such as pitch and spectral information.
3. **ReLU activation layer:** this layer applies a non-linear transformation to the output of the convolutional layer, introducing non-linearity into the model.
4. **Pooling layer:** this layer performs a downsampling operation on the output of the ReLU layer, reducing spatial dimensions of the feature map and increasing the model's robustness to small variations in the input signal.
5. **Dropout layer:** it randomly drops out a fraction of the activations in the previous layer during training, preventing overfitting and improving the generalization of the model.
6. **Fully connected layer:** This layer takes the flattened output of the previous layer and applies a linear transformation to it, producing a set of scores that represent the probability of each emotion class.
7. **Softmax activation layer:** applies a softmax function to the output of the fully connected layer, producing a probability distribution over the possible emotion classes.
8. **Output layer:** the output layer of the CNN produces the final predicted emotion class based on the probability distribution from the softmax layer.

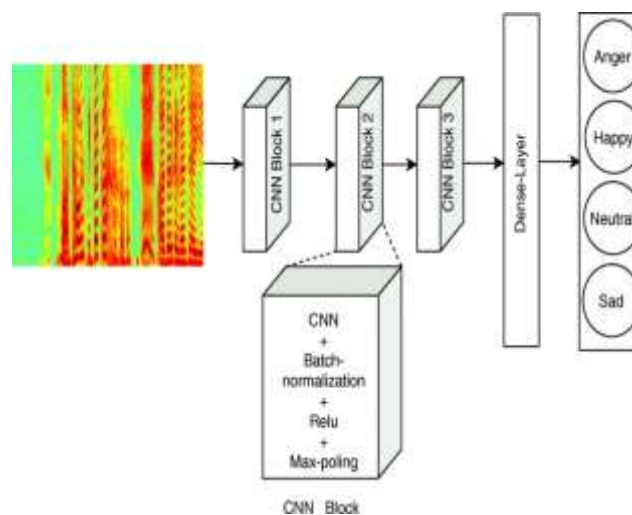


Fig.2. Architecture

IV. MODULES

4.1 CreationOfModels

Creating models for speech emotion recognition using CNN typically involves the following steps:

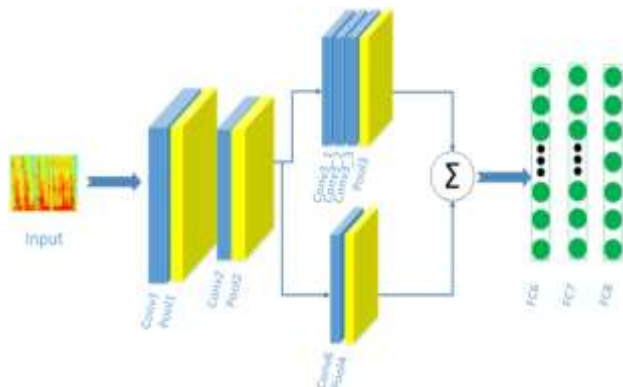


Fig.3. Model Diagram

Our approach consists of five steps:

Step 1 – Data preprocessing:

This involves converting the raw speech signal into a format that can be used by the network, such as Mel-frequency cepstral coefficients (MFCC) or filter bank energies.

Step 2 – Data augmentation:

This is done to improve the performance of the CNN and is often helpful to augment the training data by creating additional synthetic samples.

Step 3 – Model architecture design:

The next step is to design the architecture of the CNN. This is typically done to involve selecting the number and size of convolutional and pooling layers, choosing activation functions, and deciding on the number of fully connected layers.

Step 4 – Training the model:

Once the model architecture is defined, the next step is to train the model on the training data. During training, the weights of the network are adjusted to minimize the error between the predicted emotional labels and the true labels.

Step 5 – Model evaluation:

After training, the performance of the model is evaluated using a separate test set. This allows the accuracy, precision, recall, and F1 score of the model to be calculated and compared with other models.

Step 6 – Fine-tuning:

If the performance of the model is not satisfactory, it may be necessary to fine-tune the model by adjusting the architecture or hyperparameters. This can be done by evaluating the performance of the model on the validation set and adjusting accordingly.

Step 7 – Deployment:

Once the model has been trained and evaluated, it can be deployed in a production environment to perform real-time SER tasks.

To evaluate the optimal efficiency and robustness of the algorithm, metrics such as Precision and Recall rates are evaluated and computed based on the recognition rate. That the proposed system produces the highest recall rate for all types of parameters like speech and then finding the correct emotion behind the speech. The average of all measures for the proposed system.

Overall creating models for SER using CNN involves several steps from data preprocessing to model architecture design, training, evaluation, and deployment. Each step requires careful consideration and experimentation to achieve optimal performance.

4.2 Speech Emotion Recognition

The deployed model takes in new speech signals and predicts the emotional content of the speech in real-time. Overall, SER involves collecting and preprocessing data, augmenting the data, training and evaluating a deep learning model, fine-tuning the model, and deploying the model for real-time SER tasks.

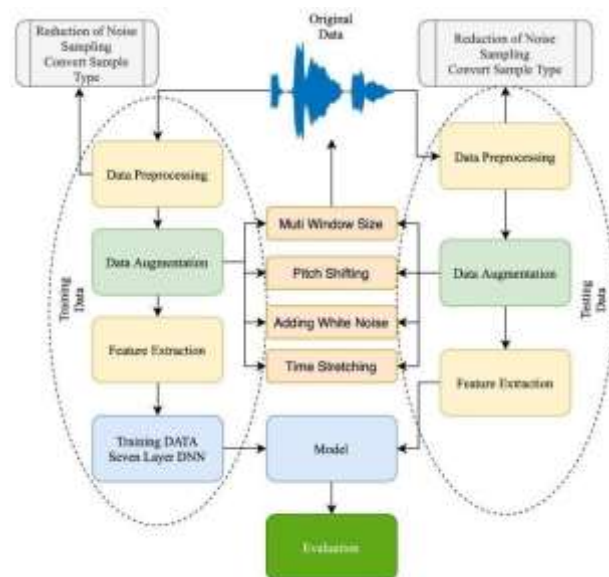


Fig.4 Activity Diagram

4.2.1 Performance of a Speech Emotion Recognition:

The performance can be evaluated using various metrics such as accuracy, precision, recall, and F1 score. Accuracy measures the percentage of correctly predicted emotional labels out of all the labels predicted by the model. Precision measures the percentage of correctly predicted positive emotional labels out of all the positive labels predicted by the model. Recall measures the percentage of correctly predicted positive emotional labels out of all the positive emotional labels in the dataset. F1 Score is the harmonic mean of precision and recall. The performance of the CNN model can be further improved by fine-tuning the model architecture and hyperparameters, augmenting the training data, and using transfer learning techniques to leverage pre-trained models in larger datasets.

TABLE1:PERFORMANCEOFSERSYSTEMUSINGCNN

		Predicted Class			
		Anger	Sad	Neutral	Happy
Actual Class	Anger	43.3	12.6	33.7	10.4
	Sad	9.6	78.3	0	12.1
	Neutral	3.6	0.3	93.3	2.8
	Happy	25.9	0	27	47.1

Architecture

Additionally, the performance of the system can be evaluated on different test sets and compared to other state-of-the-art models to determine its effectiveness in recognizing emotions from speech signals.

Considering a happy track from the dataset with plt.figure explaining the figsize to be (15,5).The model was able to recognisethecorrectemotionbehindtheaudioandwaseven able to figure out the gender of the speech.

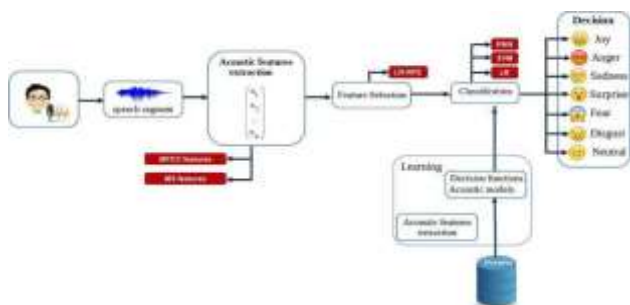


Fig5:Sequence Flowchart

**GENDER DETECTION:** it involves using a machine learning model to predict the gender of the speaker based on the emotional content of their speech. This is achieved by combining deep learning models like CNN and RNN. The CNN is used to extract high-level features from the speech signal, while the RNN is used to model the temporal dynamics of the speech signal. Once the deep learning model has been trained on a dataset of labelled speech samples, it can be used to predict the gender of new speech samples. The model was then fed a pre-processed speech signal, extracted the relevant features using the CNN, and then these features were passed through RNN to model the temporal dynamics. The output of the RNN is then fed into a fully connected layer that predicts the gender of the speaker.

4.2.2 Decision Making

Based on the input data into the model it was decided whether the person was a male or female and the emotion behind the speech of that individual. This step involves making a prediction for the emotional state of the speaker based on the features extracted from the speech signal by the CNN. After the input speech signal has been pre-processed and transformed into a format that can be used by the CNN, the CNN extracts high-level features from the speech signal using a series of convolutional layers. The features learned by the CNN are then fed into

a fully connected layer that makes the final prediction for the emotional state of the speaker. During the decision-making step, the CNN takes the pre-processed speech signal as input, extract the relevant features using the convolutional layers, and passes these features through the fully connected layer to make a prediction for the emotional state of the speaker. The predicted emotional label can then be used for further analysis or to control other systems based on the emotional state of the speaker.

Overall, the decision-making step in SER using CNN involves using a fully connected layer to make a prediction for the emotional state of the speaker based on the features extracted from the speech signal by the CNN.

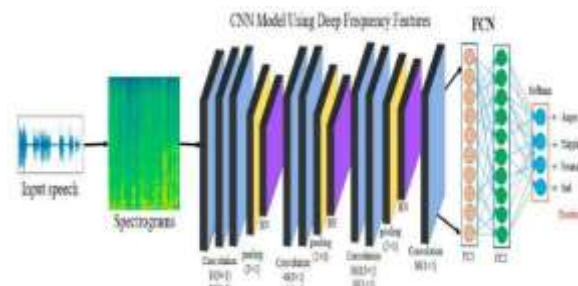


Fig 6:Decision Making Steps In The Speech Emotion Recognition

The speech component of the dataset includes 1440 recordings of 60 sentences spoken in 8 different emotional expressions while the song component includes 1012 recordings of 52 songs sung in 7 different emotional expressions. Each recording is labelled with metadata indicating the actor, gender, expression, and type of recording. The dataset is designed to be used for research and development in areas such as emotion recognition, speech processing and audio analysis. We carried out all these processes in the project.

4.2.3 Fine Tuning The Model

The first step was to load the pre-trained CNN model and remove the last layer, which is typically the output layer that is specific to the original task the CNN was trained on, such as image classification. The remaining layers are frozen, meaning their weights are not updated during training.

The next step was to add a new output layer, which is specific to the sentiment analysis task, with the appropriate number of output nodes for the number of sentiment classes. The weights of the new output layer are randomly initialized.

Finally, the entire network is fine-tuned by continuing the training process with a smaller learning rate. This allows the network to adapt to the domain-specific sentiment dataset and further improve its accuracy.

Fine-tuning a pre-trained CNN can greatly reduce the time and computational resources required to train a new model from scratch while still achieving high accuracy in the domain-specific task of the speech emotion recognition.

V. RESULTS

The obtained results suggested that the dataset is a valuable resource for training and evaluating models for speech emotion recognition, and that there are several effective approaches for achieving high accuracy on this task.

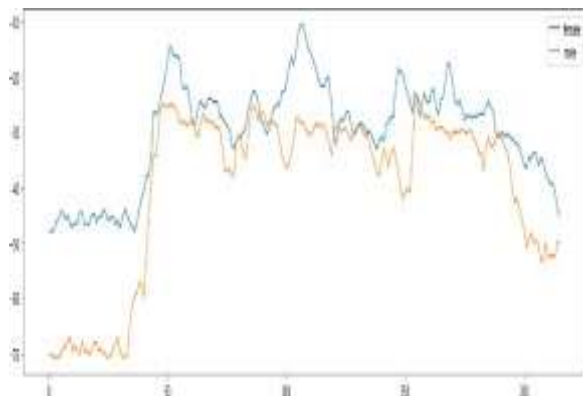


Fig 7: Accuracy Graph

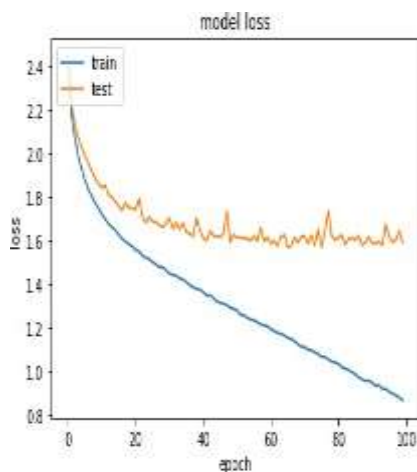


Fig 8: Model Loss vs. no. Of epochs

	actual values	predicted values
170	male_sad	male_sad
171	female_disgust	female_disgust
172	male_angry	male_happy
173	female_disgust	female_disgust
174	male_angry	male_angry
175	female_fear	female_happy
176	male_disgust	male_surprise
177	female_fear	female_happy
178	female_happy	female_happy
179	female_neutral	female_neutral

Fig 9: Actual Vs Predicted Values

Fig.8 and Fig.9 shows that 4 out of 10 random state of the dataset shows incorrect labels where fine tuning of model come into picture.

TABLE 2: TRAINING RESULTS FOR MALE AND FEMALE SAMPLES USING RBF NETWORK

Sample name (Input - Male)	Pitch Observed	Gender Classified Correctly	Emotion Frame Count	Emotion Classified Correctly
Angry	108.0	Yes	255.6	Yes
Angry	188.0	Yes	267.8	Yes
Angry	67.0	Yes	273.9	Yes
Angry	171.0	Yes	250.3	Yes
Angry	72.0	Yes	265.1	Yes
Sad	35.0	Yes	35.9	Yes
Sad	347.0	No	46.8	Yes
Sad	89.0	Yes	78.9	Yes
Sad	171.0	Yes	198.4	No
Sad	188.0	Yes	76.9	Yes
Neutral	209.0	Yes	102.4	Yes
Neutral	134.0	Yes	143.3	Yes
Neutral	190.0	Yes	127.8	Yes
Neutral	455.0	No	109.2	Yes
Neutral	98.0	Yes	145.3	Yes
Happy	189.0	Yes	167.8	Yes
Happy	112.0	Yes	153.4	Yes
Happy	69.0	Yes	155.9	Yes
Happy	41.0	Yes	89.0	Yes
Happy	45.0	Yes	176.3	Yes

Sample name (Input - Female)	Pitch Observed	Gender Classified Correctly	Emotion Frame Count	Emotion Classified Correctly
Angry	388.0	Yes	298.3	Yes
Angry	331.0	Yes	312.5	Yes
Angry	121.0	No	303.2	Yes
Angry	345.0	Yes	278.9	Yes
Angry	72.0	No	289.4	Yes
Sad	335.0	Yes	35.9	Yes
Sad	347.0	Yes	79.0	Yes
Sad	389.0	Yes	59.6	Yes
Sad	341.0	Yes	298.1	No
Sad	338.0	Yes	37.8	Yes
Neutral	393.0	Yes	105.6	Yes
Neutral	334.0	Yes	139.8	Yes
Neutral	397.0	Yes	125.6	Yes
Neutral	355.0	Yes	133.5	Yes
Neutral	348.0	Yes	126.7	Yes
Happy	459.0	Yes	187.2	Yes
Happy	112.0	No	176.3	Yes
Happy	469.0	Yes	199.0	Yes
Happy	451.0	Yes	155.8	Yes
Happy	435.0	Yes	167.8	Yes

ADVANTAGES

The various advantages of the implemented method or system are:

1. Improving Human- compute interaction.
2. Improving mental health.
3. Improving customer service.
4. Enhancing security.
5. Advancing scientific research.

CONCLUSION

We proposed a new method to improve mental health of a person by detecting changes in emotional state over time, which could be helpful in identifying individuals

who may be at risk for mental health issues such as depression or anxiety. We will be collecting all the past data of customers over a good period of time, analyse the data and will be able to figure out who are at risk of calling themselves as depressed or non-depressed person. This is a rapidly developing field with a wide range of potential applications in various industries and fields. As technology continues to advance and more data becomes available, we can expect that SER models will continue to improve, providing more accurate and reliable detection of emotions in speech.

#### FUTUREWORK

There are several directions that future research in speech emotion recognition could take:

1. Incorporating contextual information.
2. Improving model robustness.
3. Handling multilingual speech.
4. Recognizing more complex emotions.
5. Integrating with other technologies.

#### REFERENCES

- [1] M. El Ayadi, M. S. Kamel, and F. Karray, "Survey on speech emotion recognition: Features, classification schemes, and databases," *Pattern Recognit.*, vol. 44, no. 3, pp. 572–587, Mar. 2011. <https://doi.org/10.1016/j.patcog.2010.09.020> [Crossref](#)
- [2] R. Bane and K. R. Scherer, "Acoustic profiles in vocal emotion expression," *Pers. Soc. Psychol.*, vol. 70, no. 3, pp. 572–587, 1996. <https://doi.org/10.1037/0022-3514.70.3.614> [Google Scholar Crossref](#)
- [3] V. Hozjan and Z. Kačič, "Context-independent multilingual emotion recognition from speech signals," *Int. J. Speech Technol.*, vol. 6, no. 3, pp. 311–320, 2003. <https://doi.org/10.1023/A:1023426522496> [Google Scholar Crossref](#)
- [4] F. Burkhardt, A. Paeschke, M. Rolfes, W. F. Sendlmeier, and B. Weiss, "A database of German emotional speech," in *Interspeech*, 2005, vol. 5, pp. 1517–1520. [Google Scholar](#)
- [5] A. Schuller, B. Steidl, S. I. and Batliner, "The interspeech 2009 emotion challenge," *Interspeech*, pp. 312–315, 2009. [Google Scholar](#)
- [6] M. M. H. El Ayadi, M. S. Kamel, and F. Karray, "Speech emotion recognition using Gaussian Mixture Vector Autoregressive Models," in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing - ICASSP '07*, 2007, vol. 4, pp. IV–957–IV–960. [Google Scholar Crossref](#)
- [7] W. Dai, D. Han, Y. Dai, and D. Xu, "Emotion Recognition and Affective Computing on Vocal Social Media," *Inf. Manag.*, Feb. 2015. [Google Scholar](#)
- [8] B. S. Atal, "Effectiveness of linear prediction characteristics of the speech wave for automatic speaker speech wave for automatic speaker identification and verification," *Acoust. Soc. Am.*, vol. 55, no. 6, pp. 1304–1312, 2005. <https://doi.org/10.1121/1.1914702> [Google Scholar Crossref](#)
- [9] S. Wu, T. H. Falk, and W.-Y. Chan, "Automatic speech emotion recognition using modulation spectral features," *Speech Commun.*, vol. 53, no. 5, pp. 768–785, May 2011. <https://doi.org/10.1016/j.specom.2010.08.013> [Google Scholar Crossref](#)
- [10] H. Cao, R. Verma, and A. Nenkova, "Speaker-sensitive emotion recognition via ranking: Studies on acted and spontaneous speech," *Comput. Speech Lang.*, vol. 28, no. 1, pp. 186–202, Jan. 2015. <https://doi.org/10.1016/j.csl.2014.01.003> [Google Scholar Crossref](#) [PubMed](#)
- [11] C. C. Lee, E. Mower, C. Busso, S. Lee, and S. Narayanan, "Emotion recognition using a hierarchical binary decision tree approach," *Speech Commun.*, vol. 53, no. 9–10, pp. 1162–1171, Nov. 2011. <https://doi.org/10.1016/j.specom.2011.06.004> [Google Scholar Crossref](#)
- [12] T. L. Nwe, S. W. Foo, and L. C. De Silva, "Speech emotion recognition using hidden Markov models," *Speech Commun.*, vol. 41, no. 4, pp. 603–623, Nov. 2003. [https://doi.org/10.1016/S0167-6393\(03\)00099-2](https://doi.org/10.1016/S0167-6393(03)00099-2) [Google Scholar Crossref](#)
- [13] L. Chen, X. Mao, Y. Xue, and L. L. Cheng, "Speech emotion recognition: Features and classification models," *Digit. Signal Process.*, vol. 22, no. 6, pp. 1154–1160, Dec. 2012. <https://doi.org/10.1016/j.dsp.2012.05.007> [Google Scholar Crossref](#)
- [14] J. H. Yeh, T.-L. Pao, C.-Y. Lin, Y.-W. Tsai, and Y.-T. Chen, "Segment-based emotion recognition from continuous Mandarin Chinese speech," *Comput. Human Behav.*, vol. 27, no. 5, pp. 1545–1552, Sep. 2011. <https://doi.org/10.1016/j.chb.2010.10.027> [Google Scholar Crossref](#)
- [15] Gomathy, V., Janarthanan, K., Al-Turjman, F., Sitharthan, R., Rajesh, M., Vengatesan, K., & Reshma, T. P. (2021). Investigating the spread of coronavirus disease via edge-AI and air pollution correlation. *ACM Transactions on Internet Technology*, 21(4), 1–10.
- [16] E. M. Albornoz, D. H. Milone, and H. L. Rufiner, "Spoken emotion recognition using hierarchical classifiers," *Comput. Speech Lang.*, vol. 25, no. 3, pp. 556–570, Jul. 2011. <https://doi.org/10.1016/j.csl.2010.10.001> [Google Scholar Crossref](#)
- [17] C.-H. Wu and W.-B. Liang, "Emotion Recognition of Affective Speech Based on Multiple Classifiers Using Acoustic-Prosodic Information and Semantic Labels," *IEEE Trans. Affect. Comput.*, vol. 2, no. 1, pp. 10–21, Jan. 2011. <https://doi.org/10.1109/T-AFFC.2010.16> [Google Scholar Crossref](#)
- [18] S. S. Narayanan, "Toward detecting emotions in spoken dialogs," *IEEE Trans. Speech Audio Process.*, vol. 13, no. 2, pp. 293–303, Mar. 2005. <https://doi.org/10.1109/TSA.2004.838534> [Google Scholar Crossref](#)
- [19] B. Yang and M. Lugger, "Emotion recognition from speech signals using new harmony features," *Signal Processing*, vol. 90, no. 5, pp. 1415–1423, May 2010. <https://doi.org/10.1016/j.sigpro.2009.09.009> [Google Scholar Crossref](#)
- [20] C. C. Lee, E. Mower, C. Busso, S. Lee, and S. Narayanan, "Emotion recognition using a hierarchical binary decision tree approach," *Interspeech*, vol. 53, pp. 320–323, 2009. [Google Scholar](#)
- [21] S. Bjorn, S. Steidl, and A. Batliner, "The INTERSPEECH 2009 Emotion Challenge," 2009. [Google Scholar](#)
- [22] J. P. Arias, C. Busso, and N. B. Yoma, "Shape-based modeling of the fundamental frequency contour for emotion detection in speech," *Comput. Speech Lang.*, vol. 28, no. 1, pp. 278–294, Jan. 2014. <https://doi.org/10.1016/j.csl.2013.07.002> [Google Scholar Crossref](#)
- [23] Rajesh, M., & Sitharthan, R. (2022). Introduction to the special section on cyber-physical system for autonomous process control in industry 5.0. *Computers and Electrical Engineering*, 104, 108481.