



# **iJRASET**

International Journal For Research in  
Applied Science and Engineering Technology



# **INTERNATIONAL JOURNAL FOR RESEARCH**

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

**Volume: 10    Issue: III    Month of publication: March 2022**

**DOI: <https://doi.org/10.22214/ijraset.2022.41099>**

**[www.ijraset.com](http://www.ijraset.com)**

**Call:  08813907089**

**E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)**

# Speech Based Emotion Detection Using Deep Learning

Prof. Martina D'souza<sup>1</sup>, Rohan Adhav<sup>2</sup>, Shivam Dubey<sup>3</sup>, Sachin Dwivedi<sup>4</sup>

<sup>1, 2, 3, 4</sup>Department of Information Technology, Xavier Institute of Engineering, Mumbai University, Maharashtra

**Abstract:** *The Mechanized Discourse Feeling Acknowledgment may be an intense handle because of the hole among acoustic characteristics and human feelings, which depends emphatically on the discriminative acoustic characteristics extricated for a given acknowledgment assignment. Distinctive people have different emotions and through and through a distinctive way to precise it. Discourse feeling do have distinctive energies, pitch variations are emphasized in case considering distinctive subjects. Subsequently, the discourse feeling location may be a requesting assignment in computing vision. Here, the discourse feeling acknowledgment is based on the Convolutional Neural Network (CNN) calculation which employments distinctive modules for the feeling acknowledgment and the classifiers are utilized to distinguish feelings such as joy, astonish, outrage, impartial state, pity, etc. The dataset for the discourse feeling acknowledgment framework is the discourse tests and the characteristics are extricated from these discourse tests utilizing LIBROSA bundle. The classification execution is based on extricated characteristics. At long last able to decide the feeling of discourse flag.*

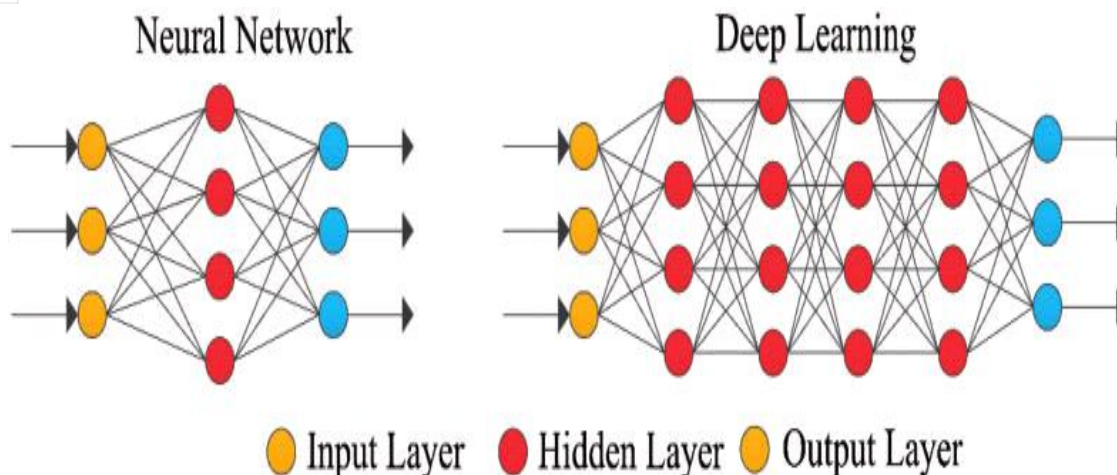
**Keywords:** Deep Learning, Speech emotion, Tensor Flow, CNN

## I. INTRODUCTION

Profound Learning in a single term we are able get it as Human Nervous System. Machine Vision Profound learning sets are made to memorize over a collection of audio/image too known as preparing information, in order to amend a issue. The different profound learning models trains a computer to imagine like a human. Deep learning models based on the inputs to the hubs can visualize. Consequently, arrange sort is like that of a Human Nervous Framework, with each hub performing beneath a bigger arrange as a neuron. So, profound learning models are fundamentally a part of Counterfit Neural Systems. Calculations of Profound learning learns in profundity around the input audio/image because it passes over each Neural Organize Layer. Low-level Characteristics like edges are identified by learning given to the beginning layers, and progressive layers collaborate characteristics from earlier layers in a more philosophical representation. Pictures, sounds, sensor information and other information are those computerized shapes designs which Profound Learning recognizes. For forecast we are pre-training the information and developing a preparing set and testing set (comes about are known). As our forecast gets an ideal hub such that the anticipated hub gives the palatable output. Basis of the neurons are in numerous levels and made to anticipate at each level and the most-optimum forecasts, and from that point for the best-fit result we utilize the information. It is treated as genuine machine intelligence. A Convolutional Neural Organize (CNN) may be a sort of feed-ahead fake arrange in which the joining grouping among its hubs is spurred by showing a creature visual-cortex. Single cortical neurons donate reaction to the boosts at a disallowed region of locale known as the open zones. The responsive zones of different hubs semi-overlap so that they can coordinate the visual range. The answer of a single hub for boosts among its responsive range might be numerically through the convolution operations. Convolutional organize was persuaded by normal methods and are assortments of multi-layer perceptron defined to utilize slightest amount of pre-processing. They have wide utilization in picture and video acknowledgment, suggestion frameworks and NLP.

The measurements of the Characteristics Outline (Convolved Highlights) is controlled by taking after parameters:

- 1) Profundity: Speaking to the channel tally we utilized within the convolution operation.
- 2) Walk alludes to estimate of the channel, on the off chance that the measure of the channel is 5x5 at that point walk is rise to to 5.
- 3) Zero-padding: Cushioning the input framework with 0swas frequently helpful around the border, in arrange to apply channel to 'Input Audio' matrix's bordering components. Utilizing zero cushioning measure of the characteristics outline can be represented.



## II. LITERATURE REVIEW

Paper-1: Speech Enhancement Using Pitch Detection Approach for Noisy Environment, IJEST

Publication Year: 2011

Author(s): Rashmi M, Urmila S, Dr. V.M. Thakare

Summary: Acoustical disorder among planning and testing stages debases exceptionally talk affirmation comes about. This issue has obliged the headway of real-world nonspecific applications, as testing conditions are exceedingly variety or undoubtedly unordinary in the midst of the planning get ready. Along these lines the establishment commotion must be expelled from the boisterous talk hail to amplify the hail comprehensible and to diminish the group of onlookers exhaustion. Update strategies associated, as pre-processing stages; to the systems shockingly make strides affirmation comes approximately. In this paper, a novel approach is utilized to update the seen quality of the discourse flag when the included substance clamor cannot be direct controlled. Instead of controlling the establishment commotion, we propose to brace the talk hail so that it can be tuned in more clearly in disorderly situations. The subjective evaluation shows up that the proposed procedure moves forward perceptual quality of talk in several disorderly circumstances. The subjective evaluation shows up that the proposed strategy makes strides perceptual quality of discourse totally different disorderly circumstances. As in a couple of cases talking may be more accommodating than composing, undoubtedly for fast typists: various numerical pictures are misplaced from the comfort but can be viably talked and recognized. In this way, the proposed system can be utilized in an application arranged for numerical picture affirmation (especially pictures not open on the comfort) in schools.

Paper-2: Emotion Recognition of Affective Speech Based on Multiple Classifiers Using Acoustic-Prosodic Information and Semantic Labels (Extended Abstract), ACII

Publication Year: 2015

Author(s): Chung-Hsien Wu, Wei-Bin Liang

Summary: This work presents an approach to feeling acknowledgment of full of feeling talk based on distinctive classifiers utilizing acoustic-prosodic information (AP) and semantic names (SLs). For AP-based affirmation, acoustic and prosodic highlights are extricated from the distinguished excited eminent segments of the input talk. Three sorts of models GMMs, SVMs, and MLPs are grasped as the base-level classifiers. A Meta Choice Tree (MDT) is at that point utilized for classifier combination to induce the APbased feeling affirmation certainty. For SL-based affirmation, semantic names are utilized to thus remove Feeling Affiliation Rules (EARs) from the recognized word arrangement of the passionate talk. The foremost extraordinary entropy appear (MaxEnt) is from that point utilized to characterize the relationship between enthusiastic states and EARs for feeling acknowledgment. At long final, a weighted thing combination procedure is utilized to facilitated the AP-based and SL based affirmation comes almost for final feeling choice. For evaluation, 2,033 expressions for four eager states were collected. The test comes around reveal that the feeling affirmation execution for AP-based affirmation utilizing MDT finished 80.00%. On the other hand, an ordinary acknowledgment precision of 80.92% was gotten for SL-based acknowledgment. At final, combining AP information and SLs accomplished 83.55% accuracy for feeling affirmation.

**Paper-3: Speech based Emotion Recognition using Machine Learning, ICCMC**

Publication Year: 2019

Author(s): Girija D, Apurva G, Gauri G, Sukanya K

Summary: Feeling affirmation from sound hail requires highlight extraction and classifier planning. The highlight vector comprises of components of the sound hail which characterise speaker specific highlights such as tone, pitch, essentialness, which is significant to plan the classifier appear to see a particular feeling absolutely. The opensource dataset for North American English tongues was physically isolated into planning and testing. Speaker vocal tract information, talked to by Mel-frequency cepstral coefficients (MFCC), was removed from the sound tests in planning dataset. Pitch, Brief Term Vitality (STE), and MFCC coefficients of sound tests in sentiments shock, bliss, and feel sorry for were gotten. These removed highlight vectors were sent to the classifier appear. The test dataset will experience the extraction strategy taking after which the classifier would make a choice with regard to the basic feeling inside the test sound. The planning and test databases utilized were North American English acted and characteristic talk corpus, real-time input English talk, regional tongue databases in Hindi and Marathi. The paper focuses of intrigued the two methodologies associated on incorporate vectors and the affect of extending the number of incorporate vectors bolstered to the classifier. It gives an examination of the exactness of classification for Indian English talk and talk in Hindi and Marathi. The fulfilled precision for Indian English talk was 80 percent.

**Paper-4: Speech Emotion Recognition Using Deep Neural Network considering Verbal and Nonverbal Speech Sounds, ICASSP**

Publication Year: 2019

Author(s): Kun-Yi Huang, Chung-Hsien Wu, Qian-Bei Hong, Ming-Hsiang Su and Yi-Hsuan Chen

Summary: Talk feeling affirmation is getting to be continuously critical for various applications. In genuine life communication, non-verbal sounds interior an verbalization as well play an basic part for people to recognize feeling. In current considers, because it were few feeling affirmation systems considered nonverbal sounds, such as snickering, cries or other feeling interjection, which really exists in our ordinary talk. In this work, both verbal and nonverbal sounds interior an expression were in this way considered for feeling affirmation of real-life discussions. Firstly, an SVM-based verbal/nonverbal sound locator was made. A Prosodic Express (PPh) auto-tagger was progress utilized to remove the verbal/nonverbal parts. For each part, the feeling and sound highlights were independently removed based on convolutional neural frameworks (CNNs) and after that concatenated to make a CNN-based dull incorporate vector. At final, a course of action of CNN-based highlight vectors for a whole trade turn was supported to an careful long short-term memory (LSTM)- based sequence-to-sequence appear to yield an energetic gathering as affirmation result. The test is based on the acknowledgment of seven delighted stages interior the NNIME (The NTHU-NTUA Chinese intelligently multimodal feeling corpus) showed up that the proposed technique finished a area precision of 52.00% beating the customary methodologies.

**Paper-5: Emotion Recognition from Speech based on Relevant Feature and Majority Voting, IEV**

Publication Year: 2014

Author(s): Md. Kamruzzaman S, Kazi Md. Rokibul Alam, Md. Arifuzzaman

Summary: This paper proposes an approach to recognize feeling from human talk utilizing lion's share voting strategy over a few machine learning methods. The commitment of this work is in two folds: firstly it chooses those highlights of talk which is most promising for classification and other than it livelihoods the lion's share voting strategy that chooses the exact course of feeling. Here, lion's share voting methodology has been associated over Neural Organize (NN), Choice Tree (DT), Support Vector Machine (SVM) and K-Nearest Neighbor (KNN). Input vector of NN, DT, SVM and KNN comprises of diverse acoustic and prosodic highlights like Pitch, Mel-Frequency Cepstral coefficients etc. From talk hail various highlight have been removed and because it were promising highlights have been chosen.

To consider a incorporate as promising, Speedy Relationship based highlight choice (FCBF) and Fisher score calculations have been utilized and because it were those highlights are chosen which are exceedingly situated by both of them. The proposed approach has been attempted on Berlin dataset of enthusiastic talk and Electromagnetic Articulography (EMA) dataset. The exploratory result shows up that lion's share voting strategy achieves predominant precision over individual machine learning methods. The work of the proposed approach can viably recognize the feeling of human animals in case of social robot, brilliantly chat client, call-center of a company etc.

**Paper-6: Speech Based Human Emotion Recognition Using MFCC, WiSPNET**

Publication Year: 2017

Author(s): M.S. Likitha, Sri Raksha R. Gupta, K. Hasitha and A. Upendra Raju

Summary: Talk may well be a complex hail comprising of different data, such as information roughly the message to be communicated, speaker, lingo, region, sentiments etc. Discourse Handling is one of the basic branches of computerized flag preparing and finds applications in Human computer interfacing, Media transmission, Assistive progresses, Sound mining, Security and so on. Talk feeling affirmation is crucial to have a ordinary interaction between human being and machine. In talk feeling affirmation, enthusiastic state of a speaker is extricated from his or her talk. The acoustic characteristic of the discourse hail is Highlight. Incorporate extraction is the method that extricates a small entirety of data from the talk flag that can a short time later be utilized to talk to each speaker. Various highlight extraction techniques are open and Mel Repeat Cepstral Coefficient (MFCC) is the commonly utilized procedure. In this paper, speaker sentiments are recognized utilizing the data removed from the speaker voice hail. Mel Repeat Cepstral Coefficient (MFCC) strategy is utilized to recognize feeling of a speaker from their voice. The arranged system was affirmed for cheerful, pitiful and shock sentiments and the capability was found to be around 80%.

**Paper-7: Speech Emotion Recognition Using Support Vector Machine, IJCA**

Publication Year: 2010

Author(s): Yashpalsing Chavan, M.L. Dhore, Pallavi Yesaware

Summary: Modified Talk Feeling Affirmation (SER) may be a current inquire about subject inside the field of Human Computer Interaction (HCI) with wide run of applications. The talk highlights such as, Mel Repeat cepstrum coefficients (MFCC) and Mel Imperativeness Range Enthusiastic Coefficients (MEDC) are removed from discourse expression. The Reinforce Vector Machine (SVM) is utilized as classifier to classify differing excited states such as shock, joy, feel sorry for, impartial, fear, from Berlin eager database. The LIBSVM is utilized for classification of sentiments. It gives 93.75% classification exactness for Sex independent case 94.73% for male and 100% for female discourse.

**Paper-8: Emotion Recognition from Speech Using MFCC and DWT for Security System, ICECA**

Publication Year: 2017

Author(s): Sonali T. Saste, Prof. S.M. Jagdale

Summary: In afterward a long time the feeling affirmation from discourse is zone of more captivated in human computer interaction. There are various different investigators which worked on feeling affirmation from talk with particular systems. This paper endeavors feeling affirmation from talk which is dialect free.

The excited talk tests database is utilized for incorporate extraction. For incorporate extraction MFCC and DWT these two particular calculations are utilized. For classification of differing sentiments like angry, cheery, panicked and fair-minded state SVM classifier is utilized. The classification is based on the highlight vector molded by combination of two calculations. This classified feeling is utilized for ATM security framework.

**Paper-9: Speech Emotion Recognition Using Convolutional Neural Network (CNN), IJPR**

Publication Year: 2020

Author(s): Apoorv Singh, Kshitij Kumar Srivastava, Harini Murugan

Summary: The Mechanized Talk Feeling Affirmation can be a extraordinary plan since of the gap among acoustic characteristics and human sentiments, which depends unequivocally on the discriminative acoustic characteristics extricated for a given affirmation task. Particular individuals have unmistakable feelings and through and through a different way to specific it. Talk feeling do have differing energies, pitch assortments are emphasized within the occasion that considering diverse subjects. Hence, the talk feeling area can be a asking task in computing vision.

Here, the talk feeling acknowledgment is based on the Convolutional Neural Organize (CNN) calculation which livelihoods particular modules for the feeling affirmation and the classifiers are utilized to recognize sentiments such as rapture, shock, shock, fair-minded state, feel sorry for, etc. The dataset for the talk feeling affirmation system is the talk tests and the characteristics are removed from these talk tests utilizing LIBROSA bundle. The classification execution is based on removed characteristics. At long final we'll choose the feeling of talk flag.

Paper-10: Deep learning based Affective Model for Speech Emotion Recognition, UIC-ATC

Publication Year: 2016

Author(s): Xi Zhou, Junqi Guo, Rongfang Bie

Summary: Considering the application regard of feeling, expanding thought has been pulled in on feeling affirmation over the ultimate decades. We commit ourselves to doable talk feeling acknowledgment examine. We develop two full of feeling models based on two significant learning procedures (a stacked autoencoder organize and a significant conviction organize) independently for modified feeling include extraction and feeling states classification. The tests are based on a well-known German Berlin Eager Discourse Database, and the affirmation exactness comes to 65% inside the leading case. In development, we favor the affect of unmistakable speakers and particular feeling categories on affirmation precision.

### III. METHODOLOGY

The discourse feeling acknowledgement application is executed utilizing convolutional neural network. Following is the block diagram of the system.

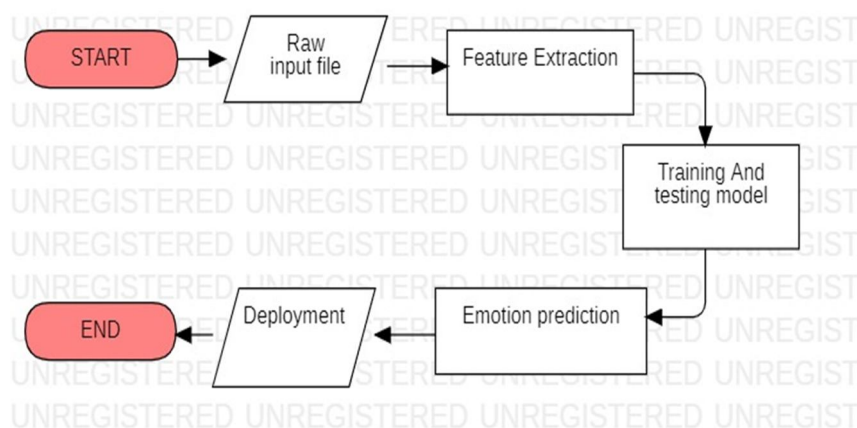


FIGURE 1: Block Diagram of the proposed system.

#### A. Training and Testing Model

A preparing information is gotten to the framework which comprises the expression name and Weight preparing is additionally given for that organize. A sound is taken as an input. From there on, concentrated standardization is connected over the sound. A standardized sound is utilized to train the Convolutional Network, this can be done to guarantee that the effect of introduction arrangement of the cases doesn't influence the preparing execution. The collections of weights come out as a result to this preparing handle and it obtains the most excellent comes about with this learning information. Whereas testing, the dataset gets the framework with pitch and energy, and based on last organize weights prepared it gives the decided feeling. The yield is spoken to in a numerical esteem each compares to either of five expressions.

There are 4 feelings that are being identified based on the person's bpm esteem, those are Anger, Sad, Happy and Neutral. The created art's colors and shapes are parallel to the identified feeling based on the standards of "color psychology" and "shape psychology".

#### B. Algorithm

// Anaconda with Jupyter Notebook in Python Language

Step 1: The test sound is given as input.

Step 2: The Spectrogram and Waveform is plotted from the sound record.

Step 3: Using the LIBROSA, a python library we extricate the MFCC (Mel Recurrence Cepstral Coefficient) more often than not approximately 10–20.

//Processing software

Step 4: Remixing the information, isolating it in prepare and test and there after developing a CNN show and it's taking after layers to prepare the dataset.

Step 5: Predicting the human voice feeling from that prepared information (test no. - anticipated esteem - real esteem)

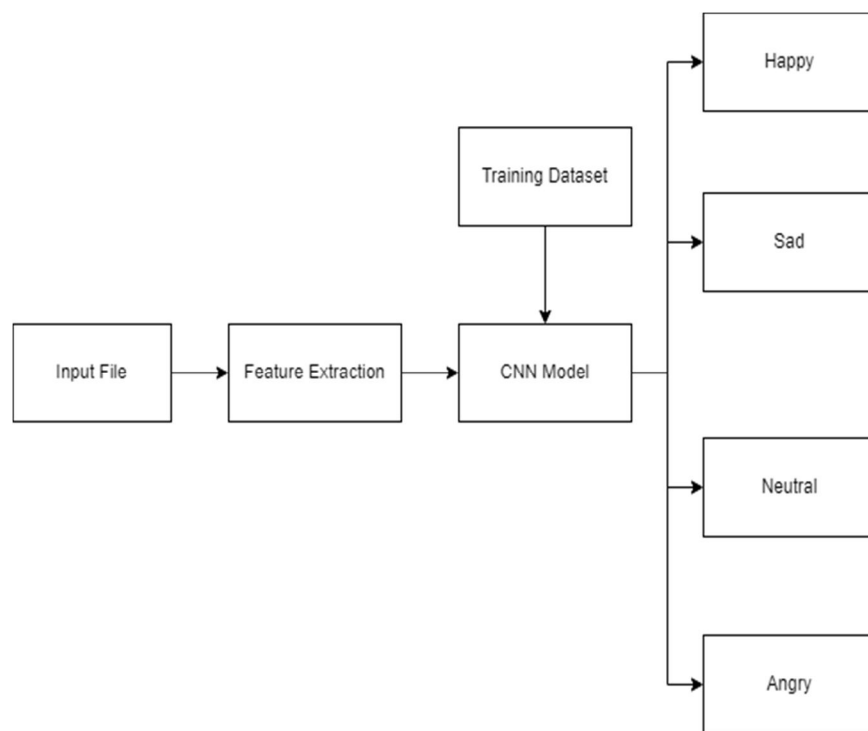


FIGURE 2: Flowchart

### C. Dataset

We are making utilize of RAVDESS dataset. It is downloaded from [kaggle.com](https://www.kaggle.com). It holds “1440 records: 60 trials/actor copied with 24 on-screen characters = 1440 trials”. The RAVDESS comprises of 24 proficient voices (12 ladylike, 12 manly), speaking 2 lexically-matched sentences within the indeed North-American emphasize. Cheerful, pitiful, irate, frightful, calm, nauseate and shock are the different discourse feeling expressions utilized. Each expression is generated in 2 levels of passionate concentrated (light, strong), with a unbiased expression. Each record out of 1440 records has a unique filename. The filename holds a 7-part numerical identifier (e.g., 03-02-05-01-02-02-11.wav). They constitute the inspiring highlights.

“Modality=>01 = full-AV, 02 = video-only, 03 = audio-only.

Vocal channel =>01 = speech, 02 = song

Emotion=>01 = neutral, 02 = calm, 03 = happy, 04 = sad, 05 = angry, 06 = fearful, 07 = disgust, 08 = surprised. Emotional intensity=>01 = normal, 02 = strong.”

NOTE: For the ‘neutral’ expressions there are no strong escalated accessible.

Statement=>

01 = "Kids are talking by the door",

02 = "Dogs are sitting by the door”.

Repetition=>

01 = 1st repetition,

02 = 2nd repetition

Actors => 01 to 24.

01,03,05,07.....23 are Male actors.

02,04,06,08.....24 are female actors.

#### D. Modules

The CNN Model here consists of four crucial layers:

- 1) Convolutional Layer: Distinguishes striking districts at interims, length expressions that are variable and portrays the feature outline sequence.
- 2) Activation Layer: A non-linear Actuation layer work is utilized as standard to the convolutional layer yields.
- 3) Max Pooling Layer: This layer empowers choices with the most extreme esteem to the thick layers. It makes a difference to keep the variable length inputs to a settled measured include array
- 4) Dense Layer

- Sound Highlight Extraction and Visualizations. (module01)

Characteristics extraction is required for classification and delineation. The sound flag could be a 3D flag in which 3 axes indicate time, sufficiency and recurrence. We are going utilize librosa to analyze and extricate characteristics of any sound flag. (.stack) work pulls an sound record and unscrambles it into a 1D cluster which is of time arrangement x, and SR is really examining rate of x. By default SR is 22 kHz. Here I will appear one sound record show with the utilize of (IPython.display) work. Librosa.display is vital to speak to the sound records in different shapes i.e. wave plot, spectrogram and colormap.

Wave plots utilize blaring of the sound at a specific time. Spectrogram shows different frequencies for a particular time with its plentifulness.

- To prepare the demonstrate for exactness calculation. (module02)

Within this module we prepare the show for precision estimations. First, moment essential modules. At that point drag the dataset. We are going get the inspecting rate esteem with librosa bundles and mfcc work. From that point this esteem holds other factors. Presently sound records and mfcc value hold a variable thus it'll include a list. At that point zip the list and hold two factors x & y. At that point we have spoken to (x, y) shape values with the utilize of numpy package.

- Execution handle of CNN show. (module03)

Speech spoken to within the shape of picture with 3 layers. Whereas utilizing CNN, do consider, 1st and 2nd derivatives of discourse picture with time and recurrence. CNN can foresee, analyze the discourse information, CNN can learn from addresses and distinguish words or articulations.

- Classification of discourse feelings. (module04)

When testing we offer the sound input. Following, we run the sound in arrange to listen with ipython.disply packages. From that point plot the sound highlights with librosa.display.waveplot bundles. Extricate the Characteristics using librosa.load. It changes over one information outline and show organized shape. Assist it compares stacked show by predict work clump estimate 32. Eventually it shows the yield from the sound record what sort of expression/emotion that sound record has.

## IV. RESULTS

We tested a sound record to induce its characteristic by plotting the waveform and spectrogram.

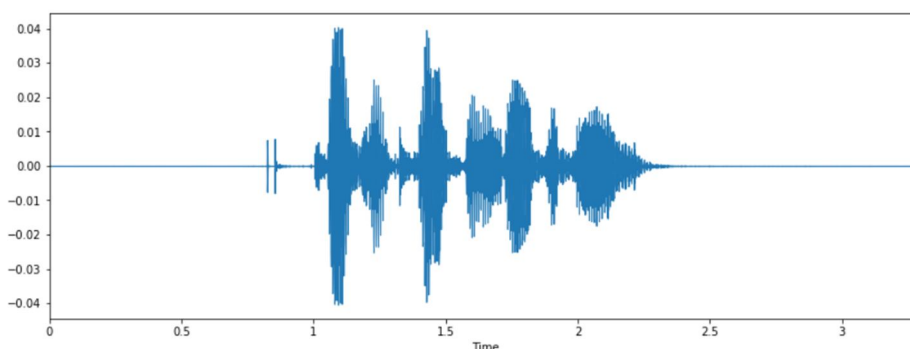


FIGURE 3: Time Domain Plot of Speech Signal

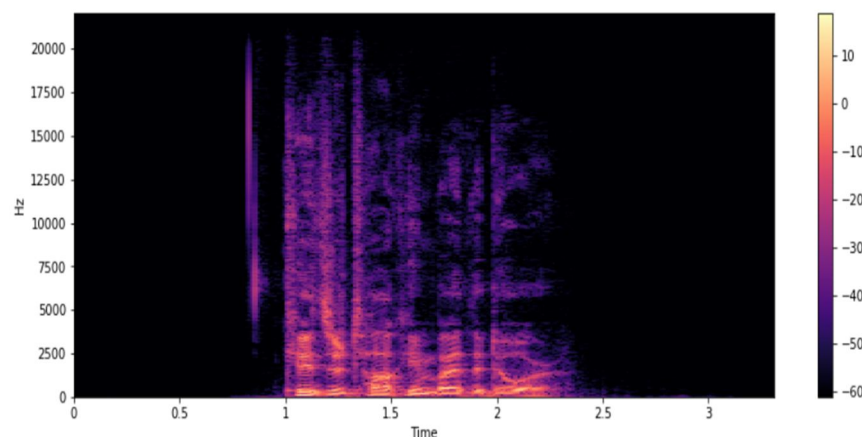


FIGURE 4: Frequency Domain Plot of the Speech Signal

After preparing different models it came out with the foremost ideal accuracy of 78% with SoftMax activation layer, “rmsprop” activation layer, 14 layers, Batch-Size = 30 and with 400 epochs.

```
um.summary()
```

Model: "sequential\_3"

Layer (type)	Output Shape	Param #
conv1d_7 (Conv1D)	(None, 180, 128)	768
activation_9 (Activation)	(None, 180, 128)	0
dropout_7 (Dropout)	(None, 180, 128)	0
max_pooling1d_5 (MaxPooling1D)	(None, 22, 128)	0
conv1d_8 (Conv1D)	(None, 22, 128)	82048
activation_10 (Activation)	(None, 22, 128)	0
max_pooling1d_6 (MaxPooling1D)	(None, 2, 128)	0
dropout_8 (Dropout)	(None, 2, 128)	0
conv1d_9 (Conv1D)	(None, 2, 128)	82048
activation_11 (Activation)	(None, 2, 128)	0
dropout_9 (Dropout)	(None, 2, 128)	0
flatten_3 (Flatten)	(None, 256)	0
dense_3 (Dense)	(None, 8)	2056
activation_12 (Activation)	(None, 8)	0

Total params: 166,920  
 Trainable params: 166,920  
 Non-trainable params: 0

FIGURE 5: Model Summary

The underneath figure appears the preparing and testing loss on our dataset. As the chart says that both “training and testing” mistakes decreases as number of ages to the preparing demonstrate increments.

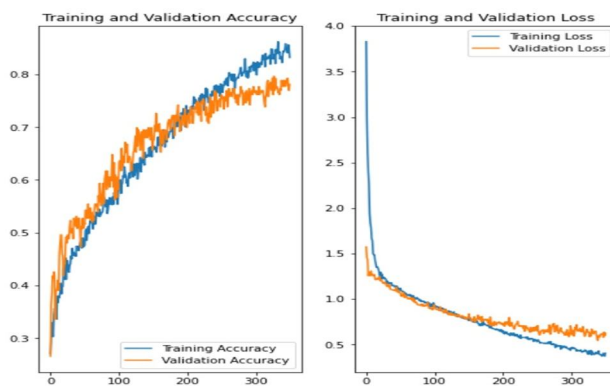


FIGURE 6: Model loss plot

From over plot we are able too induce that the number of reasonable epochs is around 300 as the exactness of test information remains steady after 300 epochs. Post demonstrate preparing we must portray out test information feelings with 75% avg. accuracy and 82.08% at most precision. The table that shows the accuracy of model as well as the confusion matrix and classification report is presented below:

```
loss, acc = um.evaluate(x_testcnn, y_test)
print("Restored model, accuracy: {:.5.2f}%".format(100*acc))
```

```
256/256 [=====] - 0s 448us/step
Restored model, accuracy: 78.12%
```

```
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
prediction = um.predict_classes(x_testcnn)
print(classification_report(y_test,prediction))
print(confusion_matrix(y_test,prediction))
```

	precision	recall	f1-score	support
0.0	0.90	0.63	0.74	82
1.0	0.70	0.82	0.76	68
3.0	0.73	0.75	0.74	32
4.0	0.80	0.92	0.86	74
accuracy			0.78	256
macro avg	0.78	0.78	0.77	256
weighted avg	0.80	0.78	0.78	256

```
[[52 13  2 15]
 [ 5 56  5  2]
 [ 0  8 24  0]
 [ 1  3  2 68]]
```

FIGURE 7: Confusion Matrix of the Model

## V. CONCLUSION

After building different models, we got the superior CNN demonstrate for the feeling refinement assignment. We come to 78% accuracy from the already accessible model. Our model would've performed way better with more information. Moreover, our model performed exceptionally well when recognizing among a manly and ladylike voice.

Our extend can be amplified to coordinated with the robot to assist it to have distant understanding, a much better; a much higher; a much stronger and an improved understanding of the mood the comparing human is in, which is able offer assistance to an improved discussion as well, because it can be coordinate with various music applications to recommend songs to its clients agreeing to his/her emotions, it can be utilized in different online shopping applications such as Amazon to make strides the item suggestion for its clients. Moreover, within the up and coming a long time we will develop a grouping to arrangement show to make voice having distinctive feelings. E.g. a sad voice, an energized one etc. It could be used in various fields of real-time applications like safety risks, cultural & societal and health. Emotion Recognition could be used in call center for classifying calls according to emotions. Feeling Acknowledgment serves as the execution parameter for conversational investigation in this way recognizing the unsatisfied client, client fulfillment so on. It could be used in-car board system based on information of the mental state of the driver can be provided to the system to initiate his/her safety preventing accidents to happen.



## REFERENCES

- [1] Speech Enhancement Using Pitch Detection Approach for Noisy Environment, Rashmi M, Urmila S, Dr. V.M. Thakare, IJEST 2011, ISSN: 0975-5462
- [2] Emotion Recognition of Affective Speech Based on Multiple Classifiers Using Acoustic-Prosodic Information and Semantic Labels (Extended Abstract), Chung-Hsien Wu, Wei-Bin Liang, ACII 2015, 978-1-4799-9953-8/15/\$31.00 ©2015 IEEE
- [3] Speech based Emotion Recognition using Machine Learning, Girija D, Apurva G, Gauri G, Sukanya K, ICCMC, 978-1-5386-7808-4/19/\$31.00 ©2019 IEEE
- [4] Speech Emotion Recognition Using Deep Neural Network considering Verbal and Nonverbal Speech Sounds, Kun-Yi Huang, Chung-Hsien Wu, Qian-Bei Hong, Ming-Hsiang Su and Yi-Hsuan Chen, ICASSP, 978-1-5386-4658-8/18/\$31.00 ©2019 IEEE
- [5] Emotion Recognition from Speech based on Relevant Feature and Majority Voting, Md. Kamruzzaman S, Kazi Md. Rokibul Alam, Md. Arifuzzaman, ELECTRONICS & VISION 2014 978-1-4799-5180-2/14/\$31.00 ©2014 IEEE
- [6] Speech Based Human Emotion Recognition Using MFCC, M.S. Likitha, Sri Raksha R. Gupta, K. Hasitha and A. Upendra Raju, WiSPNET, 978-1-5090-4442-9/17/\$31.00 ©2017 IEEE
- [7] Speech Emotion Recognition Using Support Vector Machine, Yashpalsing Chavan, M.L. Dhore, Pallavi Yesaware, IJCA, ©2010 International Journal of Computer Applications (0975 - 8887) Volume 1 –No. 20
- [8] Emotion Recognition from Speech Using MFCC and DWT for Security System, Sonali T. Saste, Prof. S.M. Jagdale, ICECA, 978-1-5090-5686-6/17/\$31.00 ©2017 IEEE
- [9] Speech Emotion Recognition Using Convolutional Neural Network (CNN), Apoorv Singh, Kshitij Kumar Srivastava, Harini Murugan, International Journal of Psychosocial Rehabilitation, Vol. 24, Issue 08, 2020 ISSN: 1475-7192
- [10] Deep learning based Affective Model for Speech Emotion Recognition, Xi Zhou, Junqi Guo, Rongfang Bie, 978-1-5090-2771-2/16 \$31.00 © 2016 IEEE DOI 10.1109/UICATC-ScalCom-CBDCom-IoP-SmartWorld.2016.42



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)