

## Automatic Sentiment Analyser Based on Speech Recognition

Dr G Neelima<sup>1</sup>, K Sri Subha<sup>2</sup>, Dr Dhanunjayarao Chigurukota<sup>3</sup>, Dr. B. Santhosh Kumar<sup>4</sup>

<sup>1</sup>Associate Professor, Vignan's Institute of Information Technology, Duvvada, Visakhapatnam

<sup>2</sup>B.Tech Project Scholar, Vignan's Institute of Information Technology, Duvvada, Visakhapatnam

<sup>3</sup>Associate Professor, Dept. of Computer science engineering, Aditya institute of technology and management, Tekkali.

<sup>4</sup>Professor, Department of Computer Science and Engineering, Guru Nanak Institute of Technology, Hyderabad

<sup>1</sup>gullipalli.neelima@gmail.com, <sup>2</sup>srisubha999@gmail.com, <sup>3</sup>rao.chigurukota@adityatekkali.edu.in,

<sup>4</sup>b.santhoshkumar@gmail.com.

**Article History:** Received: 10 January 2021; Revised: 12 February 2021; Accepted: 27 March 2021; Published online: 20 April 2021

**Abstract:** Analysis of the Emotion of a person has been developed over the earlier period decades. The majority of the works in it spun over text emotion analysis with content analyzing strategies. Yet, audio emotion analysis opinion is in the beginning phase of the research network. Our work presents the study of various algorithms of sentiment analysis to identify Emotion by dissecting the acoustic highlights of an individual's voice. The direction of study on datasets and the strategies which are utilized to recognize feeling through voice and actualized the framework to distinguish the best structure for the errand fully expecting and conveying it in a future application.

**Keywords:** Emotion, Sentiment Analysis, Speaker Recognition, Speech Recognition, Mel frequency cepstral coefficient (MFCC), Chroma, Multi-Layer Preceptor (MLP).

### 1. Introduction

Emotions square measure crude states associated with the sensory system welcome on by neuro-physiological changes otherwise connected with concerns, sentiments, conduct reactions, and tier of pleasure or dismay. Emotion is nothing however Associate in Nursing expression is a good/negative event that's connected with a selected illustration of physiological development. Peggy Thoits depicted emotions as including physiological segments, social or emotional labels (anger, stress, and so forth), expressive body exercises, the assessment of conditions and settings.

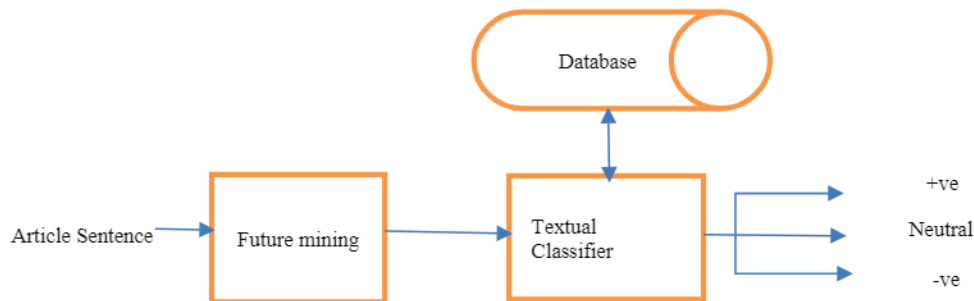
#### 1.1 Sentiment Analysis

Sentiment Analysis is the translation and characterization of feelings (positive or negative or impartial) inside the given information. It can be done through content, sound, video examination systems. Sentiment analysis is that the analysis of individuals' feelings or behaviour towards a circumstance, speech on points, or by and enormous. Thought assessment is moreover employed in numerous applications. Here, the paper comprehends the perspective of people subject to their speech with all others. For associate degree appliance to understand the mentality or perspective of the people through a discussion. It has to acknowledge World Health Organization is interfacing within the speech and what's spoken. To execute a speaker and speak affirmation system, 1st perform the sensation assessment on the information isolated from previous strategies. Understanding the attitude of people is astonishingly helpful in numerous events. as an example, PCs will see and react to human non-lexical correspondence as an example, and sentiments. In such a case, leading to recognizing somebody's emotions, the machine will modify the settings satisfying their desires and tendencies. The investigation organize has gone when dynamic sound materials, as an example, tunes examine, news, political disputes, to content. to boot, the system in like manner worked on sound assessment to think about client bolster phone conversations and numerous conversations including quite one speaker. Since there's quite one speaker within the speech, it gets awkward to look at the sound annals. it's needed to recommend a technique that may think about the presenter character and execute sound examination for solitary speakers and report their inclination.

#### 1.2 Related Background Work

Assessment Analysis has to boot alluded as Sturmarbeitelung, that acknowledges the slant sent in a very book by then assessments it to seek out whether or not report imparts positive or negative opinion. a bigger piece of labor on feeling investigation has centered on procedures, as an example, Innocent Bayesian, call tree, reinforce vector machine, most extraordinary entropy. within the work, the sentences in every record square measure named as crazy and goal, and a short time later customary AI techniques square measure applied for the passionate

elements. With the goal that the furthest purpose classifier ignores the inconsequential or confusing terms. Since social function and naming the information is repetitive at the judgment level, this strategy is nothing but onerous to check. To perform slant examination, we've used the going with procedures – Naive Thomas Bayes, Linear Support Vector Machines, VADER. Besides, a affiliation is created to seek out the winning estimation for our rationalization.



**Figure1.** Arrangement of basic Sentiment Analysis System

### 1.3 Emotion based on the text, sound, and facial expressions

There are different types of emotions, which were expressed based on the reactions of the mind. Based on the mood the mind will react which was shown in the action manner and it is a biological state associated with nerves. If the emotion was serious the person may react in any way. In this context, the recommended system tries to analyze the emotion and try to control the users based on the text, sound, and facial expressions.

Sentiment analyzer detects the emotion of the person based on their speech and helps the person to share his feelings. The main intent of this system is to detect the emotions of a person to find out whether he is feeling sad or happy.

## 2. Proposed System

This paper recommends a model for assessment examination that utilizes feature isolated from the discourse sign to recognize the sentiments of the speakers related to the conversation. The procedure incorporates four phases:

- Phase 1. Pre-processing
- Phase 2. Speech Recognition System
- Phase 3. Speaker Recognition System
- Phase 4. Sentiment Analysis System

To analyse the feeling, we have tendency to come up with some techniques to visualize whether or not the person is feeling unhappy or happy supported the input speech. To avoid or to decrease the speed of this development we have tendency to form a human-specific friend exploitation AI. Here, feeling analysis on speaker discriminated speech transcripts to note the emotions of the individual speakers occupied within the discussion are performed. to check their emotions by their speech input and generates automatic messages to cheer up the person.

The data signal was sent to the Voice Activity Detection Framework, that identifies and disconnects the voices from the sign. The voices are handled as lumps within the record, the items are then conceded to discourse acknowledgment and speaker division structure for seeing the essence and speaker Id. Speaker acknowledgment system names the items with the speaker is, it ought to be seen that the structure works in an exceedingly performance set up, for example, it might notice atmosphere the items are from a similar speaker or one in all a sort and name it as 'Speaker 1' and 'Speaker 2'. The voice acknowledgment system deciphers the items to content. The system any matches the Speaker Id with a decipher matter. it's taken care of as speak within the info. The substance yield from the speak acknowledgment structure unequivocal to solitary speaker fills in needless to say half to live estimation underlined by the individual speaker. the full methodology is representational process pictorially in Figure 2.

## 2.1 Speech Recognition

Discourse acknowledgment was the limit passed to a machine to understand words and articulations within the language spoken by people and alter them to a machine-important setup, which may be what is more used for handling. At present, we've used discourse acknowledgment contraptions. A relationship is created and also the best suit for the planned model is picked.

## 2.2 Speaker Discrimination

Recognizing a personal subject to the assortments and distinctive characteristics within the voice is recommended speaker acknowledgment. it's secured loads of thought from the investigation organize for just about eight decades. Discourse as a proof contains a handful of options that may be removed chronicle, eagerness, speaker unequivocal data, speaker acknowledgment handles the speaker's categorical options from the discourse signal. At this moment, Mel Frequency Cepstral constant (MFCC) is employed for transcription a speaker discriminate structure. The MFCC's for discourse tests from varied speakers square measure removed to boot, appeared otherwise with one another with realize the comparable qualities between the discourse tests.

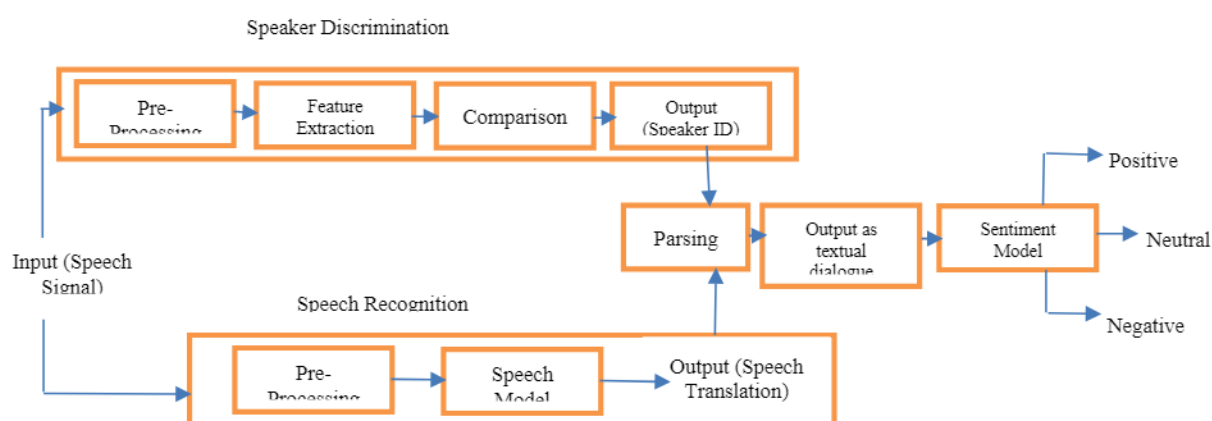


Figure 2. Proposed system

## 2.3 MFCC

Steps to extract data:

1. Take a microphone.
2. Connect it to the preferred device in which you intend to use your speech.
3. To start extraction you need to switch on the microphone.
4. Then the text on the screen appears as 'start speaking'.
5. Then record your voice to the device.
6. In this way the data is collected into the device in which it starts the extraction process

There are various methodologies in detecting emotion from a human. Few main techniques used in detecting these emotions include

- A. Emotion detection based on text.
- B. Emotion detection based on sound.
- C. Emotion detection based on facial expressions.

**A. EMOTION DETECTION BASED ON TEXT:** Determining the emotion of a person based on a text in a particular document, paragraph, or a usual text whether the person is feeling positive, negative, neutral toward the object or more. This in addition entails the classification of opinions, sentiments in keeping with the depth of the speech written. Sentiment analysis may be a procedure study on however opinions, attitudes, emotions, and views square measure expressed in language. This analysis in the main uses a technique referred to as text analysis. Text analysis plays an important role during this sort of detection of emotions supported sentiment analysis. Sentiment analysis may be a procedure study of however feelings, mentalities, stickers, and viewpoints square measure communicated in language. Sentiment Detection, or in its efficient structure – Polarity Classification, may be a slow and troublesome task. discourse changes of extremity demonstrating words, as an example, annulment, mockery even as ineffectual syntactical structures build it inconvenient for the 2 machines and folks to firmly

decide extremity of messages. Opinion Mining aims to work out the polarity and intensity of a given text, i.e., whether it's positive, negative, or neutral, and to what extent. Within the development of machine learning, computing, and tongue process, driven by new innovative prospects, it's potential to modify the investigation of big measures of freely distributed info or information. Text Mining and Social Network Analysis became a desire for investigation information yet because the associations across them. The fundamental target is to recognize the essential data as proficiently as would be prudent, finding the connections between accessible data by applying algorithmic, statistical, and data mining the executive's strategies on the information. To increase the detection of the sentiment based on the text of the person, we can use text analysis with sentiment detection.

**B. EMOTION DETECTION BASED ON SOUND:** The features like linguistic and stress generated for some words of the speech are important tasks to recognition independent speaker emotion and this emotion is based on sound utterance and length of individual speech based on above said features. The acknowledgment of feelings depending on the voice has been read for a considerable length of time. For a machine to grasp the outlook/state of mind of the individuals through a discussion it's to comprehend WHO is associating within the discussion and what's spoken, therefore we tend to execute a speaker and discourse acknowledgment framework initial and perform the emotional examination on the knowledge extracted from earlier procedures. However, the larger a part of the work-concerned is information gathered during a controlled domain during which the knowledge is ideal while not important noise and directly well metameric. What is a lot of, most of such a framework is discourse placed. In reality the procedure is considerably a lot of advanced. There area unit several factors like background and not speech voice sort of a chuckle, a moan, a cry, a sigh, etc., that considerably irritate the outcomes no heritable during a controlled domain. These factors can build the important feeling recognition trained on the information from the verboten setting unsuccessful. Three sorts of speech are watched. Natural speech is just unconstrained speech where all feelings are genuine. Simulated or acted speech is a speech that is communicated in 8 an expertly thought way. At long last, elicited speech is a speech in which the feelings are induced. The inspiring speech is neither neural nor simulated. For instance, depictions of non-experts while impersonating experts produce inspired speeches, which can likewise be a worthy arrangement when an adequatenumber of experts are not available. By using this speech recognition based on the voice of the input or the person. This emotion of a person through the voice input can be recognized using many techniques such as MFCC that is Mel Frequency Cepstral Coefficient, feature extraction, Sentiment analysis, Chroma, MLP Classifier, and many other methods. These various methods play different roles in extracting the sound of the input; these techniques divide the speech into two types of input such as speaker discrimination and speech recognition. Based on these methods pre-processing, further extracting, parsing, speech models, and then the sentiment analysis which then depicts the final result in the positive, negative, and neutral sentiments of the speaker. We use speech detection mechanisms to record the audio and a projected speaker differentiation methodology supported a definite hypothesis to acknowledge the speakers concerned in a very discussion. Further, feeling analysis is performed on the speaker's precise speech knowledge that permits the machine to acknowledge what the humans were discussing and the way they suppose.

**C. EMOTION DETECTION BASED ON FACIAL EXPRESSION:** Humans share associate degree across-the-board and principal set of feelings that area unit expressed through steady facial articulations. associate degree algorithmic program that performs recognition, extraction, also, assessment of those facial expressions can take into account programmed acknowledgment of human reaction in footage and recordings. Introduced here could be a mixture embody extraction and outward look acknowledgment approach that uses Viola-Jones cascade object detectors and Harris corner key-focuses to extract look and countenance from footage and uses head section examination, direct discriminate examination, histogram-of-oriented gradients (Hog) embody removal and support vector machines to coach a multi-purpose indicator for characterizing the seven central human facial appearances. the overall face removal from the image is completed initial utilizing a Viola-Jones cascade object face detector. The Viola-Jones detection framework tries to totally differentiate faces or highlights of a face (or different articles) by utilizes basic options called Hear-like options. the method involves passing module boxes over a picture and registering the distinction of additional element values between close locales. The issue that matters is at that time contrasted and a limit that shows whether or not a commentary is viewed as recognized or not. this needs edges that are prepared ahead of your time for varied component boxes and highlights. express component boxes for facial highlights area unit utilised, with the requirement that the majority faces and therefore the places of interest within it'll meet the overall state of affairs. Basically, in an exceedingly issue locus of interest on the face, it'll for the foremost half hold that some territories are lighter or darker than encompassing territory. This automatic face recognition uses varied methodologies to sight the feeling of someone supported the instant of the face regarding the face options as eyes, mouth, and cheek moments. Following extraction of the eyes and therefore the mouth regions, HOG options area unit calculated and extracted. to see the HOG options, a picture is separated into equally sized and spaced grids. a picture process and classification technique are enforced during which face pictures area unit wont to train a twin classifier predictor that predicts the seven basic human emotions given a take a look at

image. The predictor is moderately effective at anticipating take a look at data from the equivalent knowledge set wont to prepare the classifiers. In any case, the indicator is faithfully poor at characteristic the articulation associated with scorn. this can be probably as a result of a combination of lacking getting ready and take a look at footage that show hate, poor pre-preparing marking of knowledge, and therefore the inherent bother at recognizing scorn. The classifier is in addition not fruitful at foreseeing feelings for take a look at data that have articulations that do not have an area solely with one amongst the seven basic articulations, because it has not been ready for various expressions. Future work ought to involve up the strength of the classifiers by as well as all the additional getting ready footage from numerous knowledge sets, researching progressively actual recognition techniques that despite everything sustain procedure effectiveness, and considering the arrangement of additional nuanced and advanced articulations. People see sound in an exceedingly nonlinear scale, MFCC endeavors to breed the human ear as a numerical model. The veritable acoustic frequencies area unit mapped to Mel frequencies that habitually go between 300Hz to 5KHz. The Mel scale is straight beneath 1KHz and exponent higher than 1KHz. MFCC Constants implies the imperativeness associated with every Mel holder, that is noteworthy to each speaker. This individuality allows United States to understand speakers dependent on their voice.

### 2.3.1 Feature Extraction (MFCC)

The extraction of one-of-a-kind speaker discriminate highlight is crucial to accomplish a superior exactness rate. The exactitude of now is crucial to the subsequent stage since it goes regarding because the info for the subsequent stage. The important endeavour to form a predominant affirmation execution. The potential of this stage is important for the incidental to stage since it impacts its direction. MFCC depends upon human hearing acknowledgments that cannot see frequencies over 1KHz. consequently. MFCC depends upon a far-famed assortment of the human ear's essential info transmission with repeat. MFCC has 2 styles of the channel that area unit scattered licitly at low repeat beneath a thousand cycles/second and power uninflected on top of 1000Hz. A theoretical pitch is obtainable on Mel Frequency Scale to induce the important nature of acoustics in discourse. the final strategy of the MFCC is showed up in Figure one.

### 2.3.2 Pre-emphasis

Pre-emphasis insinuates a structured arrangement planned to increase, inside a group of frequencies, the enormity of a couple frequencies for the degree of the remaining frequencies to refine the general SNR. Therefore, this movement dictates the process of sign through a medium which underscores higher frequencies. This strategy is rend to extend the essentialness of signs at a greater repeat.

### 2.3.3 Framing

The way toward portioning the speech tests got from an ADC into a little edge with the speech length inside the scope of 20 to 40 millisecc. The input voice data is partitioned as casings of N tests. Adjoining outlines are being isolated by M ( $M < N$ ). The general values utilized are for M, N are  $M = 100$  and  $N = 256$ .

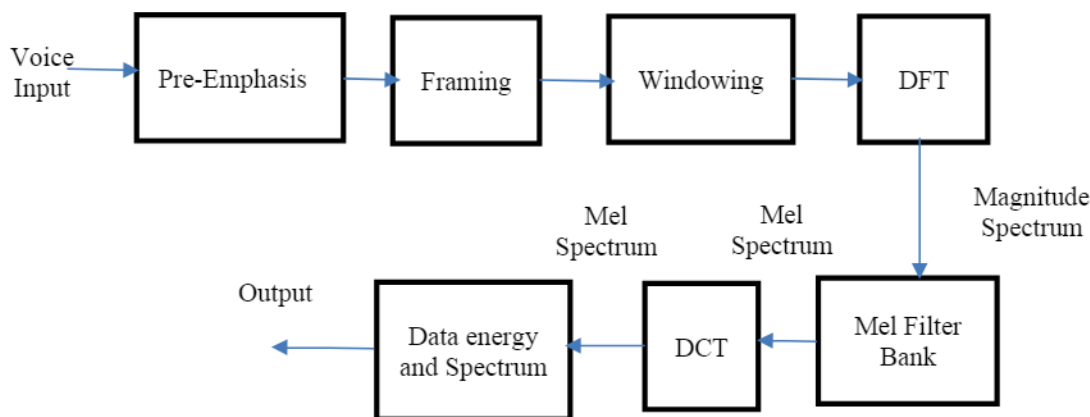


Figure 3. MFCC

### 2.3.4 Hamming windowing

It is utilized as window shape by taking the obtained square in feature extraction handling process and coordinates all the nearest recurrence lines. The Hamming window is spoken to as appeared in "Eq. (1)". If the window is characterized as  $W(n)$ ,  $0 \leq n \leq N-1$  where  $N$  denotes the number of tests in each edge

$Y[n]$  denotes the Output signal

$X(n)$  denotes the input signal

$W(n)$  denotes the Hamming window

The outcome of speech signal that is windowed is demonstrated below:

$$\text{Eq 1: } Y[n] = X(n) * W(n)$$

### 2.3.5 Fast Fourier Transform (FFT)

For converting  $N$  samples of each frame from the time area into the frequency area, FFT is engaged. FFT is utilized to change over the convolution of the glottal pulse  $U[n]$  and the vocal tract motivation reaction  $H[n]$  in the time domain. This statement is described mathematically in "Eq. (2)" as:

$$\text{Eq 2: } Y(w) = \text{FFT}[h(t) * X(t)] = H(w) * X(w)$$

### 2.3.6 Mel Frequency bank processing

The range of frequencies in FFT is very well spread. Also the speech signal doesn't fall under the straight scale. Each channel's greatness recurrence reaction is triangular fit as a fiddle and equal to solidarity at the middle recurrence and reducing directly to zero at focus recurrence of two neighboring filters.

At that point, each channel yield is the aggregate of its separated ghastly segments. Then the accompanying condition as appeared in "Eq. (3)" is utilized to register the Mel for given recurrence  $f$  in HZ:

$$\text{Eq 3: } F(\text{Mel}) = [2595 * \log_{10}(1 + f/700)]$$

### 2.3.7 Discrete cosine Transform (DCT)

In this procedure we change over the logarithmic Mel range into time area utilizing DCT. The consequence of the change can be interpreted as Mel Frequency Cepstrum Coefficient (MFCC). The arrangement of the MFCC is defined as acoustic vectors. Consequently, every information articulation is changed to a succession of an auditory vector.

## 2.4 Delta energy and delta spectrum

The speech signal and the housings difference, for instance, the slope of a formant at its advances. Appropriately, it is mandatory to include highlights associated to the difference in cepstral includes after some time. 13 delta or speed highlights (12 cepstral includes notwithstanding imperativeness), and 39 highlights a twofold delta or speeding up highlight are incorporated. For a individual sign  $x$  in a window, the essentialness in an edge from two time tests,  $t_1$  and  $t_2$  is addressed and showed up underneath in "Eq. (4)".

$$\text{Eq 4: } E = \sum X^2[y] \text{ where } E \text{ denotes Energy}$$

Where R.H.S denotes signals

Every one of the 13 delta highlights addresses the difference between traces identifying with cepstral or imperativeness include, while all of the 39 twofold delta highlights addresses the difference in between plots of the contrasting delta highlights.

The waveform of the discourse signal is as appeared in underneath Figure 4

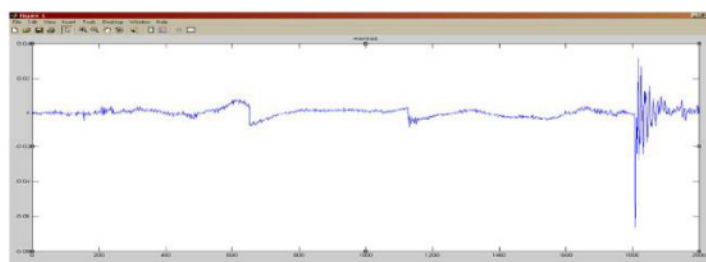


Figure 4. Speech signal

Zero-intersection rate and vitality vector are used to expel the quietness from the sign. Two vitality limits for instance lower and upper edges are resolved. In case, the vitality level of the sign is past or not the greatest or least limit that sign is considered as noise or quiet and in this way evacuated. The fundamental sign got is alluded to as articulation as showed up in the beneath Figure:

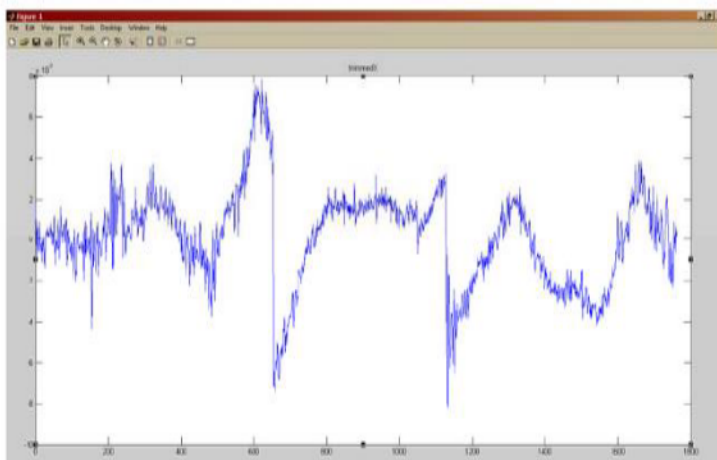


Figure 5. Utterance

The utterance is classified into several small frames as shown in Figure 6

The Utterance is divided into number of frames and then passes through a discrete filter. In the Figure.4 a frame and its output obtained after passing it through discrete filter has been shown. The Utterance is divided into number of frames and then passes through a discrete filter. In the Figure.4 a frame and its output obtained after passing it through discrete filter has been shown.

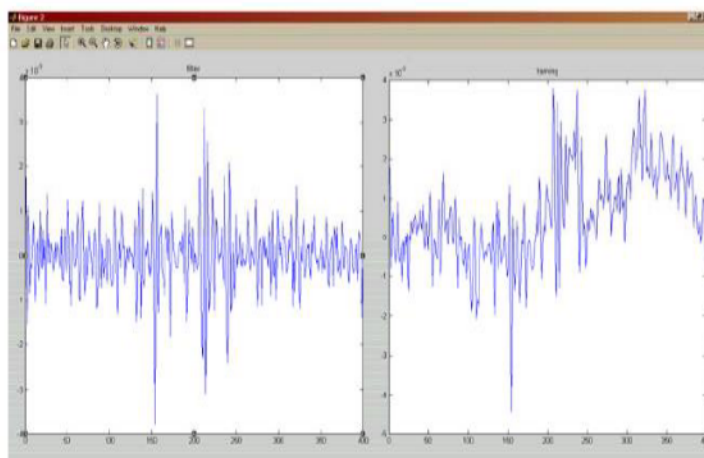


Figure 6. Framing and filtering

Directly this separated sign is experienced the hamming window and a while later to change over this time region signal into recurrence space its 400 point FFT has been found is shown below Figure.7

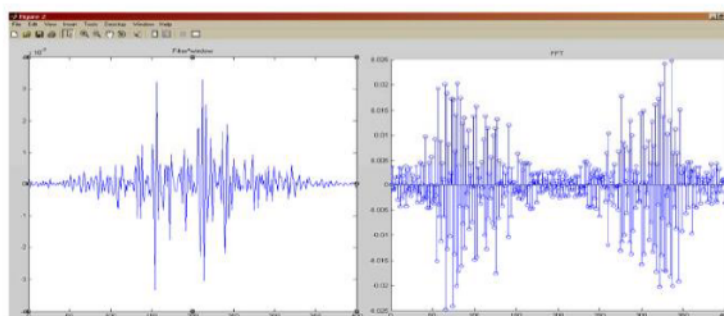


Figure 7. Windowing and its Fast Fourier Transform

This sign can also be experienced in 24 channel Mel bank and 512 length FFT, trying recurrence used is 16kHz and a while later Sparse lattice which contains the channel bank amplitudes is resolved and with its help go as showed up in Figure.6 is gotten which is the most raised and least channels decline towards zero.

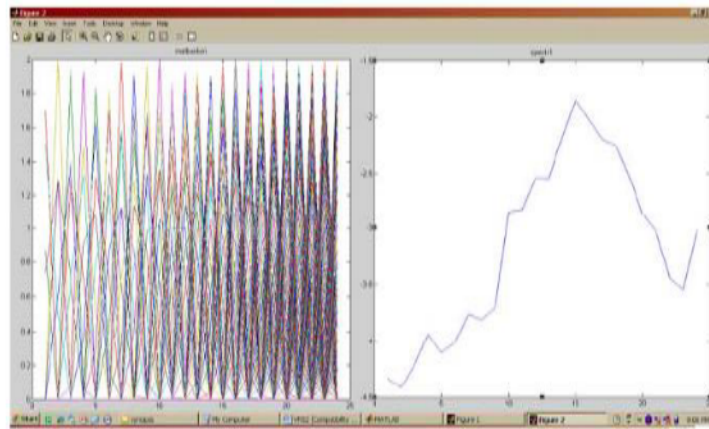


Figure 8. Mel bank processing and spectrum is obtained

### 3. Chroma

It is a notable wonder that human perspective on contributing is irregular as two contributes are viewed as equivalent "concealing" when they differentiate by an octave. Considering the above discernment, a pitch can be classified into 2 sections, which are insinuated as pitch height and chroma. Anticipating the equal-tempered scale, the chromas contrast with the set  $\{C, C\sharp, D, \dots, B\}$  that includes the twelve pitch spelling attributes 1 as used in Western music documentation. As such, a chroma feature is addressed. In the component extraction step, a given sound sign is changed over into a game plan of chroma includes each conveying how the short period of time imperativeness of the sign is spread over the twelve chroma gatherings. Recognizing pitches that differ by an octave, chroma highlights show a significant level of solidarity to assortments in tone and eagerly identify with the melodic piece of concordance. This is the inspiration driving why chroma-based sound highlights, a portion of the time furthermore suggested as pitch class profiles, are a settled instrument for getting ready and separating music information. For example, every agreement affirmation technique relies upon a chroma depiction. In like manner, chroma highlights have gotten the acknowledged standard for assignments, for instance, music synchronization what's more, course of action, similarly as sound structure examination. Finally, chroma highlights have wound up being a mind-boggling mid-level part depiction in content-based sound recuperation, for instance, spread tune ID or sound organizing.

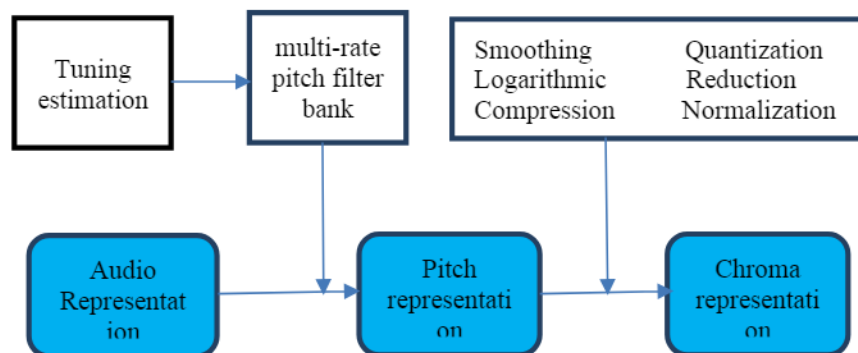


Figure 9. Overview of feature extraction pipeline

#### 3.1 Pitch Representation in chroma:

As an explanation behind the chroma feature extraction, initially, we separate the obtained sound sign into 88 repeat bunches with middle frequencies identifying with the pitches A0 to C8 where MIDI pitches  $p = 21$  to  $p = 108$ . To get sufficient ridiculous objectives for the lower frequencies, either one needs a low testing rate or a gigantic common window. In our instrument compartment, we use a consistent Q multi-rate channel bank using a



looking at the pace of 22050 Hz( high pitches), 4410 Hz (medium pitches), and 882 Hz(low pitches). The used pitch channels have a for the most part wide pass band, while still properly detaching neighboring notes on account of sharp shorts in the advancement gatherings, see Figure 2. Taking everything into account, the pitch channels are solid to deviations of up to  $\pm 25$  pennies 2 from the individual note's center repeat. To avoid tremendous stage mutilations, we use forward-backward isolating with the ultimate objective that the ensuing yield signal has exactly zero phase twisting and an enormity balanced by the square of the channel's size response.

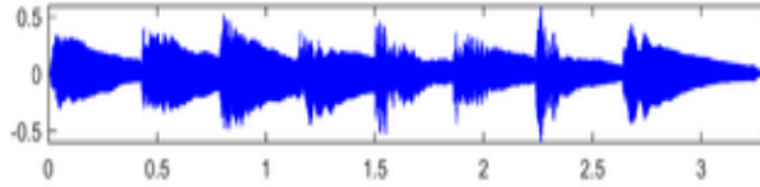


Figure 10. Pitch representation in chroma

In the ensuing stage, for all of the 88 pitch sub-bands, we register the short period (I. e., the instances of each sub-band yield are squared) using a window of a fixed length and a front of 50 %. For example, using a window length of 200 milliseconds prompts a component pace of 10 Hz (10 highlights for consistently). The resulting highlights, which we imply as Pitch, calculate the short period of time essentialness substance of the sound sign inside each pitch sub-band.

### 3.2 Tuning

In order to speak to overall tuning of a narrative, we require to sensibly move the middle frequencies of the sub-band-channels of the multi-rate channel bank. Now, we figure an ordinary spectrogram vector and decide a measure for the tuning deviation by duplicating the filter bank shifts the use of weighted binning techniques. In the available toolbox, we have pre-enrolled 6 particular multi-rate channel banks identifying with a move of  $\sigma \in \{0, 1/4, 1/3, 1/2, 2/3, 3/4\}$  semitones, independently. From these channel banks, the most proper one is picked by the assessed tuning deviation.

Filename	Main parameters	Description
wav_to_audio.m	-	Import of WAV files and conversion to expected audio format.
estimateTuning.m	pitchRange	Estimation of the filterbank shift parameter $\sigma$ .
audio_to_pitch_via_FB.m	winLenSTMSP	Extraction of pitch features from audio data.
pitch_to_chroma.m	applyLogCompr, factorLogCompr $\hat{=}$ $\eta$	Derivation of CP and CLP features from Pitch features.
pitch_to_CENS.m	winLenSmooth $\hat{=}$ $w$ , downsampSmooth $\hat{=}$ $d$	Derivation of CENS features from Pitch features.
pitch_to_CRP.m	coeffsToKeep $\hat{=}$ $n$ , factorLogCompr $\hat{=}$ $\eta$	Derivation of CRP features from Pitch features.
smoothDownsampleFeature.m	winLenSmooth $\hat{=}$ $w$ , downsampSmooth $\hat{=}$ $d$	Post-processing of features: smoothing and downsampling.
normalizeFeature.m	$p$	Post-processing of features: $l^p$ -normalization (default: $p = 2$ ).
visualizePitch.m	featureRate	Visualization of pitch features.
visualizeChroma.m	featureRate	Visualization of chroma features.
visualizeCRP.m	featureRate	Specialized version of visualizeChroma for CRP features.
generateMultiratePitchFilterbank.m	-	Generation of filterbanks (used in audio_to_pitch_via_FB.m).

Table 1: MATLAB functions in CHROMA toolbox

### 3.3 Chroma in Python

Chroma is a Python module for dealing with hues easily. Controlling hues can rapidly grow into a monotonous and muddled undertaking, especially when you become worried about shading frameworks past RGB. Chroma is here to give a basic API to do the truly difficult work, with the goal that you can remain concentrated on the significant pieces of your undertakings. Chroma gives properties to RGB in both buoy and 256 tuple designs. Color RGB yields glide facilitates, extending from 0 to 1, where 1 is white. Color.rgb256 yields number directions running from 0 to 255, where 255 is white.

### 3.4 Audio Matching

As second application circumstance, we tend to think about the trip of sound coordinating with the target to thus recoup all areas from all chronicles within a big sound selection that musically identifies with a given request sound cut. At this moment, the challenge is to regulate to assortments in tone and instrumentation as they seem in numerous understandings, unfold songs, and techniques of slightly of music. during a traditional approach for sound coordinating , the request Q even as every file recording D is initial modified over into color property feature progressions X(Q) and X(D), severally. By then, an in depth by sort of dynamic time

traveling is employed to regionally think about the request course of action  $X(Q)$  with the information progression  $X(D)$  yielding a partition work  $\Delta$ . every close-by least of  $\Delta$  virtually zero shows a district within the information recording that's about to the given inquiry. Considering this coordinating application, going with 2 properties of  $\Delta$  is of basic criticality. From one perspective, the semantically right matches ought to distinction with handy minima of  $\Delta$  virtually zero thus dodging pretend negatives. Then again,  $\Delta$  ought to be over zero outside an area of the right neighborhood minima thus sidestepping sham positives. Considering these needs, the used color property selection accepts a remarkable activity. As associate illustrative model, we tend to think about a chronicle by recording.. The subject of this piece happens on numerous occasions played in four extraordinary instruments (clarinet, strings, trombone, tutti). Demonstrating the four occasions by E1, E2, E3, and E4 and using E3 as the request, the Figure beneath shows a couple of division limits reliant on various chroma variations.

#### 4. Dataset

Here the dataset which is utilized is the RAVDESS dataset; this is the Ryerson Audio-Visual Database of Emotional Speech and Song dataset and is allowed to download. This dataset has 7356 records evaluated by 247 people multiple times on emotional correctness, power, and validity. The whole dataset is 24.8GB from 24 entertainers. The sample rate is brought down on all the files

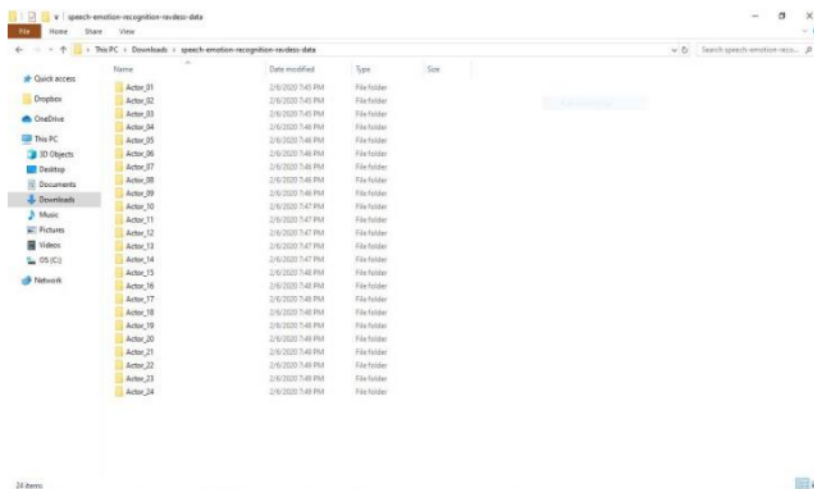


Figure 12. Screenshot of the dataset used

#### 5. Multi-layer Perceptron

This is a Multi-layer Perceptron Classifier; it improves the logarithmic-misfortune work utilizing LBFGS or stochastic gradient descent. Not at all like Support Vector Machines or Naive Bayes, the MLP Classifier has an inside neural system with the end goal of grouping. It is a feed-forward Artificial Neural Network model. Here, it is used to classify given data into different clusters. It fits and trains the data and performs classification.

#### 6. Recommender System

The objective of a Recommender System (RS) is to create significant suggestions to clients about things or items that may hold any importance to them. Recommendation systems are significant intelligent systems that assume a fundamental role in giving specific data to clients. Traditional approaches in recommendation systems incorporate collaborative filtering and content-based sifting. Be that as it may, these methodologies have certain confinements like the need for earlier client history and propensities for playing out the assignment of recommendation. Here we use a recommender system to recommend links, images, videos, or some messages which helps the user to change his state of mind based on the result produced by the sentiment analysis.

Here the main process is that when the system asks the user for the response. The user has to speak something and then the recommended system will grasp and analyze the voice and compare it with different voices based on the categories of the voices and then it is going to be concluded by displaying the result through different emojis and also the system will play the corresponding song relates to the mood of the person. By that, the recommender system will give a positive result to cool down the person if the user may be in a bad or unhappy mood. If the user is in a happy mood based on that related song is going to be played. If the user is in an angry mood then to make him cool that type of song will be played. Hope this system will make it useful to the user.

### 7. Experimental Analysis

Following are the experimental results, figure 13 tells that after submitting the audio, it analyses the emotion and displays the emotion through emojis and greets the user, and gets the data from them.

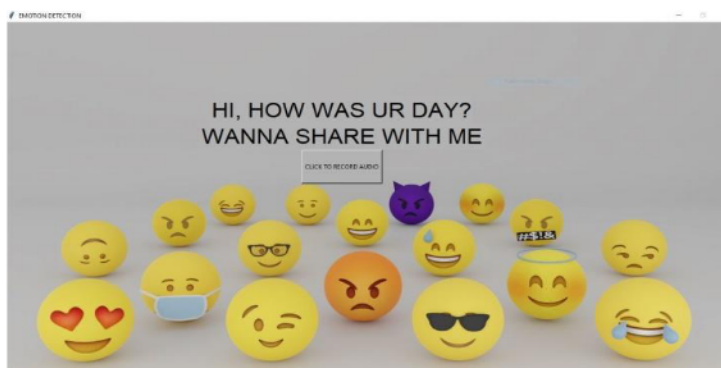


Figure 13: Emotion Analyze

From figure 14, we can see the emotions like smile, angry, love, sleeping, crying, cold, woozy, etc. By using the proposed recommender system, based on the emotion of the user and it recommends the video according to the user's mood. The proposed system will ask or address the user like "HI, HOW WAS UR DAY? WANNA SHARE WITH ME", which was shown in figure 13, if the user responds to the message displayed in the GUI the system can recommend the user to be out of their mood, by seeing the video user may feel relax and they come out of the mood and which was shown in figure 14.

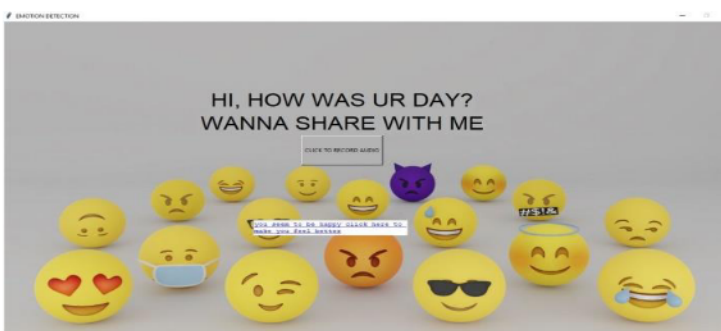


Fig 14: Emotions by using a recommender system

As we have mentioned different emojis in the above paragraph, the mood of the user can be predicted by the related emoji and the related video is going to be played to make the user out of the mood and the user may feel relax and can change the mood and be normal.

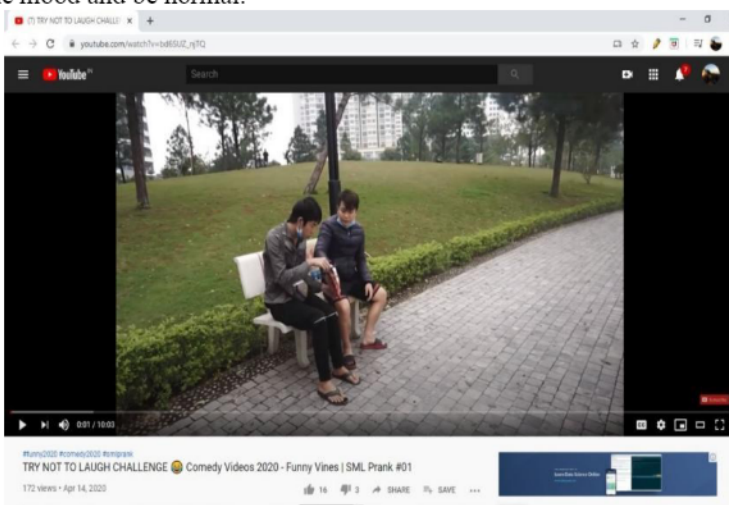


Figure 15: Recommended Video

## 8. Conclusion

The major aim of our work is to perform sentiment analysis of the user's emotion by using speech. We utilized MFCC and Chroma techniques to extract the features of the audio given by the user and later on, we used MLP classifier techniques to classify the emotion and give it to the recommender system which recommends the user some files to view to change the mood of the user and let them feel less lonely.

## References

- A. "Chroma Toolbox: Matlab Implementations for Extracting Variants of Chroma-Based Audio Features" by Meinard Müller - Friedrich-Alexander-University of Erlangen-Nürnberg.
- B. "Sentiment Analysis on Speaker Specific Speech Data" by Maghilnan S, Rajesh Kumar M, Senior IEEE, Member, School of Electronic Engineering, VIT University, Tamil Nadu, India
- C. "VOICE COMMAND RECOGNITION SYSTEM BASED ON MFCC AND DTW" by Anjali Bala et al. / International Journal of Engineering Science and Technology Vol. 2 (12), 2010, 7335-7342
- D. "Detecting Depression from Voice" by Mashrura Tasnim and Elein Stroulia, department of computing science, University of Alberta, Edmonton, AB, Canada.
- E. "Hubert Wassner and Gerard Chollet, "New Time Frequency Derived Cepstral Coefficients For Automatic Speech Recognition", 8th European Signal Processing Conference (Eusipco'96).
- F. Marco Grimaldi and Fred Cummins, "Speaker Identification Using Instantaneous Frequencies, IEEE Transactions On Audio, Speech, And Language Processing", VOL. 16, NO. 6, pp 1097-1111, ISBN: 1558-7916, August 2008.
- G. Mahdi Shaneh and Azizollah Taheri, "Voice Command Recognition System Based on MFCC and VQ Algorithms", World Academy of Science, Engineering and Technology 57 2009, pp 534-538.
- H. Mark A. Bartsch and Gregory H. Wakefield. "Audio thumb nailing of popular music using chroma-based representations". IEEE Transactions on Multimedia, 7(1):96–104, February 2005.
- I. Herbig, T., Gerl, F., & Minker, W. (2010, July). "Fast adaptation of speech and speaker characteristics for enhanced speech recognition in adverse intelligent environment"s. In Intelligent Environments (IE), 2010 Sixth International Conference on (pp. 100-105). IEEE.
- J. Ezzat, S., El Gayar, N., & Ghanem, M. (2012). Sentiment analysis of call centre audio conversations using text classification. Int. J. Comput. Inf. Syst. Ind. Manag. Appl, 4(1), 619-627