

# Speaker Dependent Speech Emotion Recognition using MFCC and Support Vector Machine

Prajakta P. Dahake  
Department of E&TC  
D. Y. Patil College of Engineering  
Akurdi, Pune, India  
prajakta17dahake@gmail.com

Kailash Shaw  
Department of CSE  
D. Y. Patil College of Engineering  
Akurdi, Pune, India

Mrs. P. Malathi  
Department of E&TC  
D. Y. Patil College of Engineering  
Akurdi, Pune, India

**Abstract**— In human computer interaction, speech emotion recognition is playing a pivotal part in the field of research. Human emotions consist of being angry, happy, sad, disgust, neutral. In this paper the features are extracted with hybrid of pitch, formants, zero crossing, MFCC and its statistical parameters. The pitch detection is done by cepstral algorithm after comparing it with autocorrelation and AMDF. The training and testing part of the SVM classifier is compared with different kernel function like linear, polynomial, quadratic and RBF. The polish database is used for the classification. The comparison between the different kernels is obtained for the corresponding feature vector.

**Keywords**— Energy , Pitch, MFCC, Speech Emotion , SVM.

## I. INTRODUCTION

Humans interact with their environment using many different sensing mechanisms. Detecting an unpleasant state during the task and intervening the process is possible with real time affective systems. In human computer interaction, the main task is to keep users' level of satisfaction as high as possible. A computer with affective properties could detect the users' emotion and could develop a counter response to increase user satisfaction. Speech and gesture recognition are the most popular affective computing topics. Speech and gesture recognition are possible with passive sensors.

Human emotions are hard to guess. The foremost way to recognize emotions is by facial recognition but with growing age humans learn to control the expression. Few gestures like disgust, boredom they fail to identify. To overcome this emotion recognition methods depending on speech are in the field of research. On paralinguistic basis SER system differentiate between emotions [2]. The system must identify and act accordingly depending on the emotions. In emotion recognition primary objective of recognizing is through speech.

Speech recognition systems have a framework performing features extraction selection classification and identification of it. Pre-processing consist of features extraction were formants, pitch are extracted. It is important to normalize the extracted features by reducing the formant, pitch frequency and dependency on phonetic and speaker. Intonation, pause, stress, pitch and rhythm are prosodic feature

examples. Spectral features investigate frequency components of speech signal. Used classification algorithm varies from algorithm to algorithm. The second part involves classification of extracted features which includes training of emotional modules. Another aspect is of training the modules with the help of database. The classifier must select appropriate features given by the user. There is uncertainty for using classifier to detect emotions correctly. Invariant to language, speaker and contents performance of audio emotion identification rely on how relevant features are extracted.

The need of SER is as follows. One of the application is human robotic interaction: robots can be trained to interact with humans and differentiate human emotion [5]. They must understand the spoken words as well as other information like health and emotional status of the user and respond accordingly. Another use is in-car board system will be provided for the safety issues of the driver with information about the emotional status, it will provide and resolve the errors with the solution by communicating with the driver.

Thus the objective of the paper can be summarized as follows: 1) extracting features using pitch, formants, MFCC. 2) improve speaker dependent SER by comparing the results with different kernels of SVM classifier. This paper is organized as follows: section II Basic frame work is given, section III SVM Classification is given, in section IV about emotion Database & Acoustic characteristics of emotions in speech, section V results and section VI conclusion are given respectively.

## II. BASIC FRAME WORK

The primary task after taking speech as input is the feature extraction. The extracted features are related to the fundamental frequency, formants and the mel frequency cepstrum coefficients. For the speech the acoustic features like mean, variance can be considered.

There are two categories linguistic and paralinguistic. Paralinguistic consist of prosodic features, spectral features, and voice quality features.

### a) Prosodic features

The pitch is caused due to the vibration of the vocal cord which carries the information of the emotions due to

stress on it. The segments of speech consist of phonemes, consonant, vowels etc.

#### b) Spectral features

These features are produced due to the air flow from the vocal cord, in the state of anger the air flow is faster and in state of calm mood the air flow is slower. So there is dependency on air flow of the vocal chords. The energy based features is used to measure the flow of air.

#### c) Voice quality

It has strong relation between the features and emotions detection. The voice has many types which include harsh, pitch, calm depending on the mood.

We grouped the feature extraction into three groups of frequency features consisting of fundamental frequency and formants, energy features consist of the energy contour, MFCC features consist the mel 24 coefficients. Before starting with the pre-processing the sample is executed with frame level breakdown of size 20ms.

#### A. Frequency-based features

The F0 range is 60Hz to 450 Hz. The size of window is 20ms, so the signal will be overlapped for 10ms to determine F0 and formant[7]. The signal is stationary in time domain hence used for detection of the speech features. The F0 is determined by cepstral method and formants by the LPC method. The consonants are not considered for computation of the F0 and the formants and assumed as zero.

#### B. Energy-based features

The ratio of low frequency energy is used for the energy distribution over the spectrum. The energy values in the energy contour are also calculated over windows of 20 ms with 10 ms overlap. The edge points of the plateaus of the energy contours are defined as the points at three db lower than the peak points. The valley of the energy contours are obtained with similar method as the points at three db higher than the local lowest points. The energy plateaus and the slopes are obtained by approximating the energy contour with straight lines.

#### C. MFCC features

These coefficients are derived from the cepstral coefficients and methods [6]. The cepstral coefficient can be filtered by the auditory filters like bark and mel filters. The cepstral features are derived from the filter bank designed as per the auditory system. It is based on pitch perception and the filters are placed in triangular shaped. In this filters the edge of the adjacent filter coincides with the center of the previous filter.

### III. SVM CLASSIFICATION

SVM is one of the classifier used for the emotion recognition. It recognizes the pattern and analyze the data. The hyper plane separates the data into different classes [5]. It separates linear non-linear features of the input signal. The idea is to transform the output into high dimensional features.

There are many different classifiers up to this date but SVM is definitely one of the most used ones in speech applications. It is a supervised machine learning algorithm and can be used for classification as well as regression. It uses some kernel tricks to transform the data and then it tries to find decision boundaries in the data. Depending on which side of the boundary the test sample lays, the test sample is classified. This is true for binary classifications and in this case the data is linearly separable. However in many cases, the data is not linearly separable and so the SVM maps the data into multi-dimensional spaces and tries to find an optimal separating plane.

The optimization objective function for Support Vector Machines can be described as below:

$$\min_{\theta} C \sum_{i=1}^m [y^i \text{cost}_1(\theta^T f^i) + (1 - y^i) \text{cost}_0(\theta^T f^i)] + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (1)$$

where the idea is to try to minimize this cost function.

And also the hypothesis of SVM is:

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

If the parameter vector  $\theta$  transpose times  $x$  is greater or equal than zero, it will be classified as 1 and otherwise it will be classified as zero. SVM really does a good job in defining the decision boundaries of the training data [4]. Even if there are some outliers in the training set, meaning that there are some positive examples near the negative examples and vice versa, not chosen big regularization parameter  $C$  will keep the original decision boundary without the outliers.

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (4)$$

$$\begin{aligned} \theta^T x^i &\geq 1 & \text{if } y^i &= 1 \\ \theta^T x^i &\leq -1 & \text{if } y^i &= 0 \end{aligned}$$

SVM has different types of kernel functions and each of these kernel functions give different classification results. For SVM classifier, actually it is a good idea to try different kernels on the training and test data set and see the performance of each other. However the two most popular and most widely used kernels are linear kernel and RBF kernel.

There are kernels for SVM such as Radial Basis Function (RBF) kernel, or linear kernel, Gaussian kernel or

polynomial kernel. There is no general rule for selecting a kernel when using SVM for classification problems. Sometimes, RBF might give good results, sometimes a different kernel might give better results. But if number of features  $n$  is large, and number of training examples  $m$  is small, using a linear kernel might be a good option. If  $n$  is small and  $m$  is large, using a RBF kernel can be better.

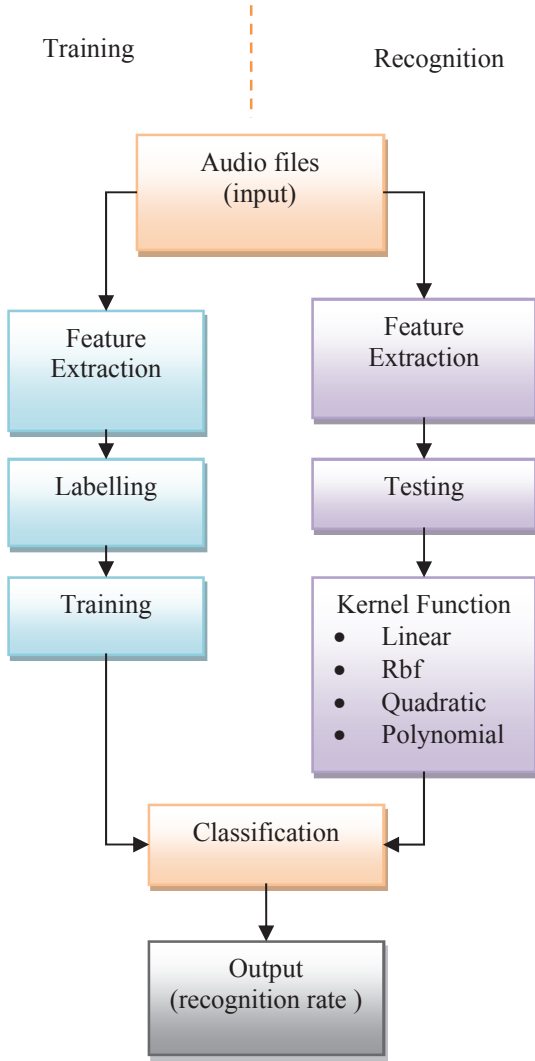


Figure 1: Block Diagram of Emotion Recognition

Given some similarity function between the actual training example and a landmark, the training algorithms of a SVM can be written as below:

$$f_i = \text{similarity}(x, l^i) = e^{\left(-\frac{\|x-l^i\|^2}{2\sigma^2}\right)} \quad (5)$$

$$\min_{\theta} C \sum_{i=1}^m [y^i \text{cost}_1(\theta^T f^i) + (1 - y^i) \text{cost}_0(\theta^T f^i)] + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (6)$$

There are many kernel functions for SVM few of them are quadratic, radial basis function, linear, spline, exponential radial function. The equation of kernel function are given in equation (7) – (10) and accordingly the respective parameters are set for higher accuracy.

Linear kernel function

$$\text{kernel}(a_i, a_j) = (a_i \cdot a_j) \quad (7)$$

There are actually two parameters that should be considered when training a SVM[8]. These parameters are  $C$  (regularization parameter) and  $\gamma$  (parameter of RBF). Choosing a large  $C$  value leads to lower bias and high variance. Choosing a small  $C$  value leads to higher bias and lower variance. In terms of selecting the  $\gamma$  value, a large  $\gamma$  value leads to features changing more smoothly and so higher bias and lower variance and a small  $\gamma$  is the opposite.

Radial basis function

$$\text{kernel}(a_i, a_j) = e^{\left(-\frac{\|x-l^i\|^2}{2\sigma^2}\right)} \quad (8)$$

Exponential Radial basis kernel function

$$\text{kernel}(a_i, a_j) = e^{\left(-\frac{\|x-l^i\|}{2\sigma^2}\right)} \quad (9)$$

Polynomial kernel function

$$\text{kernel}(a_i, a_j) = (a_i \cdot a_j + 1)^d \quad (10)$$

#### IV. EMOTION DATABASE

The criteria is present to analyse how emotional database determine the emotion in real world environment. Now the criteria depends on who is the speaker, how he utters the emotion, whether utterances are simulated, distributed, balanced or unbalanced.

For public use the database availability is scarce and costly and unavailable for free. Thus, benchmark is set for the database according to the criteria discussed above.

##### A. Types of database:

Earlier the database used for automatic speech emotion recognition research was based on acted speech and the researcher tends to have real based data [8]. Database is divided into three types.

1. Acted speech: actors are asked to express deliberately the human emotions which are predefined.
2. Real life speech: the natural response to the conversion of human with spontaneous reaction which are authentic in nature. Example : call center.
3. Elicited emotional speech in which the emotions are induced with self-report instead of labelling, where emotions are provoked and self-report is used for labelling control. The speech which is neither neutral nor simulated.

*B. Acoustic characteristic of emotions in speech:*

The prosodic features consist of pitch, voice quality, intensity, speaking rate important to identify the various types of emotions [12]. In particular pitch and intensity seem to be correlated to the amount of energy required to express a certain emotion.

When the person is in state of anger, the pitch range is enunciated with very high frequency and exactly opposite frequency range in the state of sadness.

The voice quality produced during joy is modal or the tone is tensed while in state of anger its moderately blaring timbre and for sadness its resonant timbre.

V. RESULTS

The emotion recognition system is built using MATLAB R2014a software. The sampling frequency of the audio signal is 44.1 KHz. The frame size is 236 samples with 118 sample frame overlapping. 12 mel coefficients with 300Hz and 3700 Hz as low frequency and high frequency resp. The 20 mel filter bank are used as shown in the Fig.2.

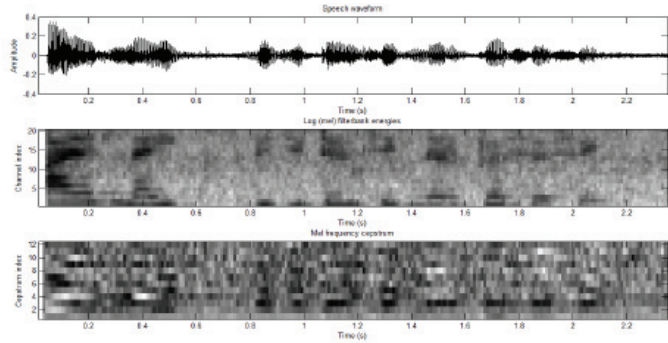


Figure 2: Feature Extraction using MFCC

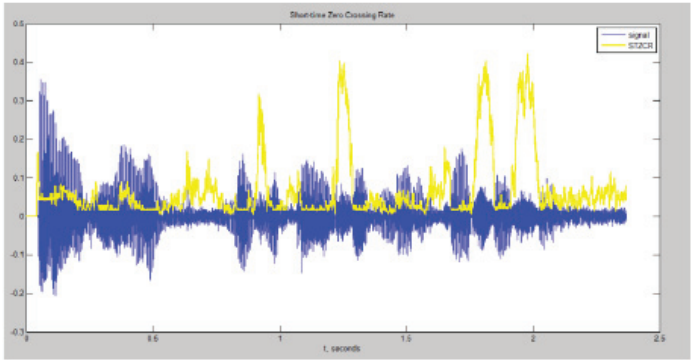


Figure 3: Zero-Crossing Rate Feature Extraction

The feature extraction by taking mean of the zero crossing and energy of the input speech is used. The database is having 12 MFCC coefficient, mean of pitch, zero crossing and energy shown in Fig. 3 ZCR and Fig. 4.

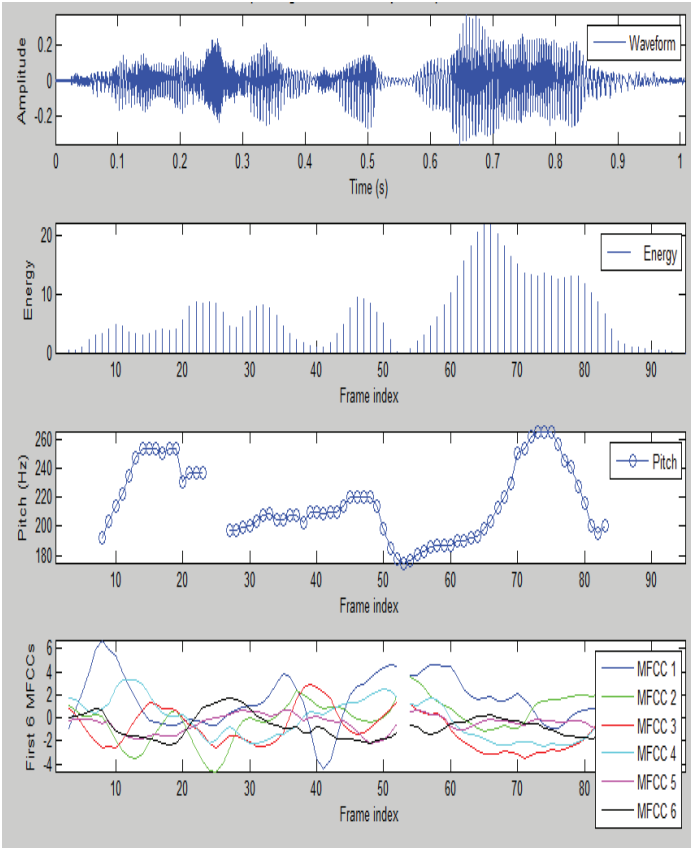


Figure 4: Speech utterance with its Energy, Pitch, MFCC

## REFERENCES

Table 1: Accuracy of Different kernel functions in SVM

Kernel Emotions	Linear	Quadratic	RBF	Polynomial
Anger	75	83.33	83.33	83.33
Fear	83	91.67	84.72	90.28
Joy	84.72	90.28	83.33	86.11
Sadness	90.28	90.28	83.33	75
Boredom	76.39	79.17	84.72	86.11
Neutral	79.17	80.56	83.33	77.78

Table I shows the accuracy of different kernels for six emotions. It is observed the accuracy is highest for sadness in linear kernel function. For quadratic kernel the accuracy for emotions fear, joy and sadness is more than 90%. For radial basis function (rbf) gives overall accuracy for all the emotion. In polynomial kernel only fear is having good accuracy while for other emotion it gives very less accuracy.

## VI. CONCLUSION

Feature extraction is the core part of speech recognition. For pitch detection i.e. fundamental frequency the best performance is obtained by cepstral when compared with autocorrelation and AMDF algorithms. The fusion feature vector of cepstral with MFCC and formants gives good accuracy for emotion recognition. For speech recognition, the traditional classifiers used are SVM, HMM, kNN etc. But the simplest and the most effective with less complexity and great accuracy is SVM. Comparing the different kernel function is good depending on the application. The best recognition for all emotions is done by RBF kernel function with 84% recognition rate. Linear and quadratic function shows highest recognition rate for fear, joy and sadness. And polynomial is worst for speech emotion recognition.

- [1] Jihyuck Jo, Hoyoung Yoo, and In-Cheol Park, "Energy-Efficient Floating-Point MFCC Extraction Architecture for Speech Recognition Systems", IEEE Trans (VLSI), vol. 24, no. 2, march 2016.
- [2] Jia-Ching Wang, Li-Xun Lian, Yan-Yu Lin, and Jia-Hao Zhao, "VLSI Design for SVM-Based Speaker Verification System" IEEE Trans (VLSI) systems, vol. 23, no. 7, Feb 2015
- [3] Kunxia Wang, Ning An, Bing Nan Li, Yanyong Zhang, "Speech Emotion Recognition Using Fourier Parameters", IEEE Trans. On Affective Computing, Vol. 6, No. 1, June 2015
- [4] Jagvir Kaur Abhilash Sharma, "Emotion Detection Independent of User Using MFCC Feature Extraction", IJARCSSE Volume 4, Issue 6, Jan 2014
- [5] M. Li, K. Han, and S. Narayanan, "Automatic speaker age and gender recognition using acoustic and prosodic level information fusion," Computer speech and language, Vol. 27, No. 1, pp. 151-167, Nov 2013
- [6] Bhoomika Panda, Debananda Padhi, Kshamamayee Dash, "Use of SVM Classifier & MFCC in Speech Emotion Recognition System", IJARCSSE, Vol. 2, Issue 3, March 2012.
- [7] Rabiner, L.R, Schafer, R.W, Digital Processing of Speech Signals, Pearson education, 1st Edition, 2004.
- [8] Bhoomika Panda, Debananda Padhi, Kshamamayee Dash, "Use of SVM Classifier & MFCC in Speech Emotion Recognition System", IJARCSSE, Vol. 2, Issue 3, March 2012.
- [9] Iker Luengo, Eva Navas, and Inmaculada Hernez., "Feature Analysis and Evaluation for Automatic Emotion Identification in Speech," IEEE Transactions On Multimedia, Vol. 12, No. 6 Jan 2010
- [10] R. Mller, B. Schuller, and G. Rigoll, "Enhanced robustness in speech emotion recognition combining acoustic and semantic analyses", Proc. From Signals to Signs of Emotion and Vice Versa, Santorino, Greece, 2004
- [11] H Meinedo and I Trancoso, "Age and Gender Classification Using Fusion of Acoustic and Prosodic Features", Proc. INTERSPEECH, pp. 2818-2821
- [12] R. Nisimura, A. Lee, H. Saruwatari, and K. Shikano., "Public speech-oriented guidance system with adult and child discrimination capability," Proc. ICASSP2004, vol. 1, pp. 433-436, june 2004