

SPEECH EMOTION RECOGNITION USING SUPPORT VECTOR MACHINE

Yashpalsing Chavhan
Student
VIT, Pune
India

M. L. Dhore
Professor
VIT, Pune
India

Pallavi Yesaware
Student
VIT, Pune
India

ABSTRACT

Automatic Speech Emotion Recognition (SER) is a current research topic in the field of Human Computer Interaction (HCI) with wide range of applications. The speech features such as, Mel Frequency cepstrum coefficients (MFCC) and Mel Energy Spectrum Dynamic Coefficients (MEDC) are extracted from speech utterance. The Support Vector Machine (SVM) is used as classifier to classify different emotional states such as anger, happiness, sadness, neutral, fear, from Berlin emotional database. The LIBSVM is used for classification of emotions. It gives 93.75% classification accuracy for Gender independent case 94.73% for male and 100% for female speech.

Categories and Subject Descriptors

I.5.0 [Pattern Recognition]: General.

General Terms

Performance, Experimentation, Human Factors.

Keywords

Speech emotion, Emotion Recognition, SVM, MFCC and MEDC.

1. INTRODUCTION

Automatic Speech Emotion Recognition is a very recent research topic in the Human Computer Interaction (HCI) field. As computers have become an integral part of our lives, the need has risen for a more natural communication interface between humans and computers. To achieve this goal, a computer would have to be able to perceive its present situation and respond differently depending on that perception. Part of this process involves understanding a user's emotional state. To make the human-computer interaction more natural, it would be beneficial to give computers the ability to recognize emotional situations the same way as human does.

Automatic Emotion Recognition (AER) can be done in two ways, either by speech or by facial expressions. In the field of HCI,

speech is primary to the objectives of an emotion recognition system, as are facial expressions and gestures. Speech is considered as a powerful mode to communicate with intentions and emotions.

In the recent years, a great deal of research has been done to recognize human emotion using speech information [1], [2]. Many researcher explored several classification methods including the Neural Network (NN), Gaussian Mixture Model (GMM), Hidden Markov Model (HMM), Maximum Likelihood Bayes classifier (MLC), Kernel Regression and K-nearest Neighbors (KNN), Support Vector Machine (SVM) [3], [4].

The Support Vector Machine is used as a classifier for emotion recognition. The SVM is used for classification and regression purpose. It performs classification by constructing an N-dimensional hyperplanes that optimally separates the data into categories. The classification is achieved by a linear or nonlinear separating surface in the input feature space of the dataset. Its main idea is to transform the original input set to a high-dimensional feature space by using a kernel function, and then achieve optimum classification in this new feature space.

A Berlin Emotional database [5] is used for feature extraction and training SVM. The Berlin database of emotional speech was recorded at the Technical University, Berlin. The database contains speech with acted emotions in German language. It contains 493 utterances of 10 professional actors five males and five females who spoke 10 sentences with emotionally neutral content in 7 different emotions. The emotions were wut (anger), langeweile (boredom), ekel (disgust), angst (fear), freude (happiness), trauer (sadness) and neutral emotional state.

Applications of Speech Emotion Recognition include psychiatric diagnosis, intelligent toys, lie detection, learning environment, educational software, and detection of the emotional state in telephone call center conversations to provide feedback to an operator or a supervisor for monitoring purposes.

2. SYSTEM IMPLEMENTATION

The importance of emotions in human-human interaction provides the basis for researchers in the engineering and

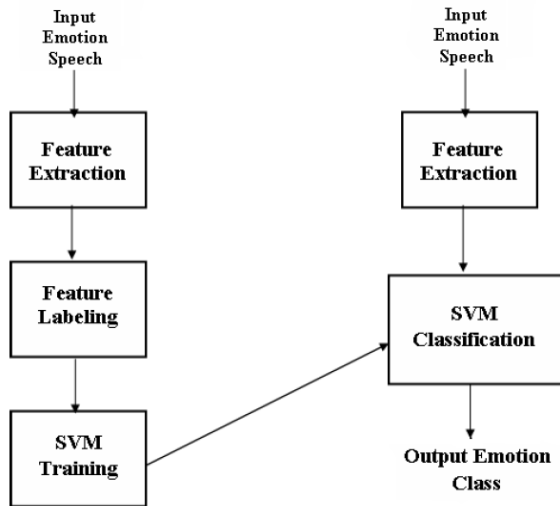


Figure 1. Speech Emotion Recognition System.

computer science communities to develop automatic ways for computers to recognize emotions. As shown in fig. 1 the input to the system is a .wav file from Berlin Emotion Database that contains emotional speech utterance from different emotional classes. After that features extraction process is carried out. In feature extraction process two features are extracted MFCC [6], [7] and MEDC [8]. After that the extracted features and their corresponding class labels are given as input to the LIBSVM classifier. The output of a classifier is a label of a particular emotion class. There are total five classes angry, sad, happy, neutral and fear. Each label represents corresponding emotion class.

2.1 Feature Extraction

In previous works several features are extracted for classifying speech affect such as energy, pitch, formants frequencies, etc. all these are prosodic features. In general prosodic features are primary indicator of speaker's emotional state. Here in feature extraction process two features are extracted Mel Frequency Cepstral Coefficient (MFCC) and Mel Energy spectrum Dynamic coefficients (MEDC). Fig. 2 shows the MFCC feature extraction process. As shown in Fig. 2 feature extraction process contains following steps:

- Preprocessing: The continuous time signal (speech) is sampled at sampling frequency. At the first stage in MFCC feature extraction is to boost the amount of energy in the high frequencies. This preemphasis is done by using a filter.
- Framing: it is a process of segmenting the speech samples obtained from the analog to digital conversion (ADC), into the small frames with the time length within the range of 20-40 ms. Framing enables the non stationary speech signal to be segmented into quasi-stationary frames, and enables Fourier Transformation of the speech signal. It is because, speech signal is known to exhibit quasi-stationary behavior within the short time period of 20-40 ms.

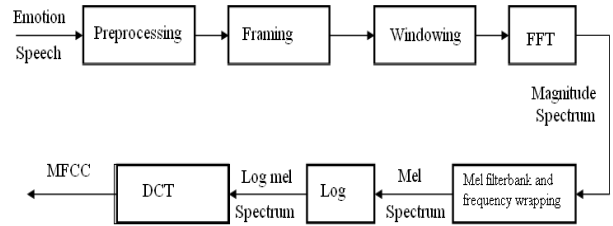


Figure 2. MFCC feature extraction

- Windowing: Windowing step is meant to window each individual frame, in order to minimize the signal discontinuities at the beginning and the end of each frame.
- FFT: Fast Fourier Transform (FFT) algorithm is ideally used for evaluating the frequency spectrum of speech. FFT converts each frame of N samples from the time domain into the frequency domain.
- Mel Filterbank and Frequency wrapping: The mel filter bank [8] consists of overlapping triangular filters with the cutoff frequencies determined by the center frequencies of the two adjacent filters. The filters have linearly spaced centre frequencies and fixed bandwidth on the mel scale.
- Take Logarithm: The logarithm has the effect of changing multiplication into addition. Therefore, this step simply converts the multiplication of the magnitude in the Fourier transform into addition
- Take Discrete Cosine Transform: It is used to orthogonalise the filter energy vectors. Because of this orthogonalization step, the information of the filter energy vector is compacted into the first number of components and shortens the vector to number of components.

Another feature Mel Energy spectrum Dynamic coefficients (MEDC) is also extracted. It is extracted as follows: the magnitude spectrum of each speech utterance is estimated using FFT, then input to a bank of 12 filters equally spaced on the Mel frequency scale. The logarithm mean energies of the filter outputs are calculated $E_n(i)$, $i = 1 \dots N$. Then, the first and second differences of $E_n(i)$ are calculated. MEDC feature extraction process. The MEDC feature extraction process contains following steps shown in figure 3:

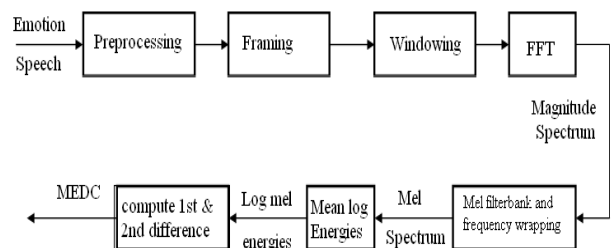


Figure 3. MEDC feature extraction

- Preprocessing, Framing, Windowing, FFT & Mel filterbank and Frequency wrapping processes of MEDC feature extraction are same as MFCC feature extraction.

- Take logarithmic mean of energies: In this process a mean log of every filter energies is calculated. This mean value represent energy of individual filter in a filterbank.

- Compute 1st and 2nd difference: The final Mel energy spectrum dynamics coefficients are then obtained by combining the first and second differences of filter energies.

2.2 Feature Labeling

In Feature labeling each extracted feature is stored in a database along with its class label. Though the SVM is binary classifier it can be also used for classifying multiple classes. Each feature is associated with its class label e.g. angry, happy, sad, neutral, fear.

2.3 SVM Classification

In general SVM is a binary classifier, but it can also be used as a multiclass classifier. LIBSVM [9], [10] is a most widely used tool for SVM classification and regression developed by C. J. Lin. Radial Basis Function (RBF) kernel is used in training phase. Advantage of using RBF kernel is that it restricts training data to lie in specified boundaries. The RBF kernel nonlinearly maps samples into a higher dimensional space, so it, unlike the linear kernel, can handle the case when the relation between class labels and attributes is nonlinear. The RBF kernel has less numerical difficulties than polynomial kernel.

3. EXPERIMENTATION AND RESULTS

Berlin Emotion database contains 406 speech files for five emotion classes. Emotion classes Anger, sad, happy, neutral, fear are having 127, 62, 71, 79 and 67 speech utterance respectively. The LIBSVM is trained on MFCC and MEDC feature vectors using RBF and Polynomial kernel functions. The LIBSVM is used to test these feature vectors. The experimentation is carried out by varying cost values for RBF kernel and degree values for Polynomial kernel. Both gender independent and gender dependent experiments are performed. Using RBF kernel at cost value $c=4$, it gives recognition rate of 93.75% for gender independent case, 94.73% for male and 100% for female speeches. The recognition rate using Polynomial kernel at degree $d=4$ is 96.25% gender independent, 97.36% for male and 100% for female speeches.

The Confusion matrices using RBF kernel gender independent, male and female are shown in Table 1, 2 and 3. Table 4, 5 and 6 shows Confusion matrices using Polynomial kernel gender independent, male and female.

Table 1. Confusion matrix of the RBF LIBSVM classifier (Gender Independent)

Emotion	Emotion Recognition (%)				
	Angry	Sad	Happy	Neutral	Fear
Angry	100	0	0	0	0
Sad	0	100	0	0	0
Happy	0	0	100	0	0
Neutral	0	6.25	0	93.75	0
Fear	0	0	30.76	0	69.24

Table 2. Confusion matrix of the the RBF LIBSVM classifier (Male)

Emotion	Emotion Recognition (%)				
	Angry	Sad	Happy	Neutral	Fear
Angry	100	0	0	0	0
Sad	0	100	0	0	0
Happy	16.66	0	83.34	0	0
Neutral	0	0	0	100	0
Fear	0	0	0	14.85	85.15

Table 3. Confusion matrix of the RBF LIBSVM classifier (Female)

Emotion	Emotion Recognition (%)				
	Angry	Sad	Happy	Neutral	Fear
Angry	100	0	0	0	0
Sad	0	100	0	0	0
Happy	0	0	100	0	0
Neutral	0	0	0	100	0
Fear	0	0	0	0	100

Table 4. Confusion matrix of the Polynomial LIBSVM classifier (Gender Independent)

Emotion	Emotion Recognition (%)				
	Angry	Sad	Happy	Neutral	Fear
Angry	100	0	0	0	0
Sad	0	100	0	0	0
Happy	0	0	100	0	0
Neutral	0	0	0	100	0
Fear	7.69	0	15.18	0	76.92

Table5. Confusion matrix of the Polynomial LIBSVM classifier (Male)

Emotion	Emotion Recognition (%)				
	Angry	Sad	Happy	Neutral	Fear
Angry	100	0	0	0	0
Sad	0	100	0	0	0
Happy	0	0	100	0	0
Neutral	0	0	0	100	0
Fear	0	0	14.28	0	85.72

Table 6. Confusion matrix of Polynomial LIBSVM classifier (Female)

Emotion	Emotion Recognition (%)				
	Angry	Sad	Happy	Neutral	Fear
Angry	100	0	0	0	0
Sad	0	100	0	0	0
Happy	0	0	100	0	0
Neutral	0	0	0	100	0
Fear	0	0	0	0	100

4. CONCLUSION

In this paper Berlin emotion database of German language is used for feature extraction. MFCC and MEDC features are extracted from a speech files in .wav format. From experimentation and result it is proved that system is speaker and text independent. It is also observed that results from LIBSVM by using RBF and Polynomial kernel function are 93.75% and 96.25% respectively. Regarding LIBSVM using RBF and Polynomial kernels it is observed that by changing the parameters of a kernel functions better results can be obtain.

REFERENCES

- [1] Cowie, R., Douglas-Cowie, E., Tsapatsoulis, N., Votsis, G., Kollias, S., Fellenz, W., and Taylor, J. G., Emotion recognition in human-computer interaction, *IEEE Signal Processing magazine*, Vol. 18, No. 1, 32-80, Jan. 2001.
- [2] D. Ververidis, and C. Kotropoulos, Automatic speech classification to five emotional states based on gender information, *Proceedings of the EUSIPCO2004 Conference*, Austria, 341-344, Sept. 2004.
- [3] Christopher. J. C. Burges, A tutorial on support vector machines for pattern recognition, *DataMining and Knowledge Discovery*, 2(2):955-974, Kluwer Academic Publishers, Boston, 1998.
- [4] Tristan Fletcher, *Support Vector Machines Explained*, unpublished.
- [5] Burkhardt, Felix; Paeschke, Astrid; Rolfes, Miriam; Sendlmeier, Walter F.; Weiss, Benjamin A Database of German Emotional Speech. *Proceedings of Interspeech*, Lissabon, Portugal. 2005.
- [6] Fuhai Li, Jinwen Ma, and Dezhi Huang, MFCC and SVM based recognition of Chinese vowels, *Lecture Notes in Artificial Intelligence*, vol.3802, 812-819, 2005
- [7] M. D. Skowronski and J. G. Harris, Increased MFCC Filter Bandwidth for Noise-Robust Phoneme Recognition, *Proc. ICASSP-02*, Florida, May 2002.
- [8] YL. Lin and G. Wei, Speech emotion recognition based on HMM and SVM, proceeding of fourth International conference on Machine Learning and Cybernetics, Guangzhou, 18-21 August 2005.
- [9] Chih-Chung Chang and Chih-Jen Lin, LIBSVM: a library for support vector machines, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [10] C.W Hsu, C.-C. Chang, C.-J. Lin, A Practical Guide to Support Vector Classification, *Technical Report*, Department of Computer Science & Information Engineering, National Taiwan University, Taiwan.