

Automatic Speech Emotion Recognition Using Support Vector Machine

Peipei Shen

Department of Computer
Technology
Shanghai Jiao Tong University
Shanghai, China
shen@sjtu.edu.cn

Zhou Changjun

Department of Computer
Technology
Shanghai Jiao Tong University
Shanghai, China
zchangjun@gmail.com

Xiong Chen

Pudong Branch
China Mobile Group Shanghai
Company Limited
Shanghai, China
chenxiong@sh.chinamobile.com

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 purpose of speech emotion recognition system is to automatically classify speaker's utterances into five emotional states such as disgust, boredom, sadness, neutral, and happiness. The speech samples are from Berlin emotional database and the features extracted from these utterances are energy, pitch, linear prediction cepstrum coefficients (LPCC), Mel Frequency cepstrum coefficients (MFCC), Linear Prediction coefficients and Mel cepstrum coefficients (LPCMCC). The Support Vector Machine (SVM) is used as a classifier to classify different emotional states. The system gives 66.02% classification accuracy for only using energy and pitch features, 70.7% for only using LPCMCC features, and 82.5% for using both of them.

Keywords- *Speech Emotion; Automatic Emotion Recognition; SVM; Energy; Pitch; LPCC; MFCC; LPCMCC*

I. INTRODUCTION

Speech emotion recognition aims to automatically identify the current emotional state of a human being from his or her voice. It is based on in-depth analysis of the generation mechanism of speech signal, extracting some characteristic parameters which contain emotional information from the speaker's voice, using these parameters and taking appropriate pattern recognition methods to identify emotional states from speech. Now, Automatic Speech Emotion Recognition is a very active research topic in the Human Computer Interaction (HCI) field and has a wide range of applications. For distance learning, indentifying students' emotion timely and making appropriate treatment can enhance the quality of teaching. In automatic remote call center, it is used to timely detect customers' dissatisfaction. It is also used to aid clinical diagnosis or to play video games. The research of automatic speech emotion recognition, not only can promote the further development of computer technology, it will also greatly enhance the efficiency of people's work and study, and help people solve their problems more efficiently. It will also further enrich our lives and improve the quality of life.

In recent years, a great deal of research has been done to recognize human emotion using speech information. Many speech databases were built for speech emotion research,

such as BDES (Berlin Database of Emotional Speech) that is German Corpus and established by Department of acoustic technology of Berlin Technical University [1], DES (Danish Emotional Speech) that is Danish Corpus and established by Aalborg University, Denmark [2], SES (Spanish Emotional Speech) that is Spanish Corpus. There are also some Mandarin Affective Speech Databases, such as MASC (Mandarin Affective Speech) and MESC (A Mandarin Emotional Speech Corpora) that is recorded by Tsinghua University, Taiwan. Many researchers have proposed important speech features which contain emotion information, such as energy, pitch frequency [2], formant frequency [3], Linear Prediction Coefficients (LPC), Linear Prediction Cepstrum Coefficients (LPCC), Mel-Frequency Cepstrum Coefficients (MFCC) and its first derivative [4]. Furthermore, many researchers explored several classification methods, such as Neural Networks (NN) [5], Gaussian Mixture Model (GMM), Hidden Markov model (HMM) [6], Maximum Likelihood Bayesian classifier (MLC), Kernel Regression and K-nearest Neighbors (KNN) and Support vector machines (SVM) [7].

In this paper, we use the Berlin Emotional database to train and test our automatic speech emotion recognition system. Prosody and Spectral features have been widely used in speech emotion recognition. In this paper, we compare the recognition rate only using prosody features, with that only using spectral features and using both prosody and spectral features [8]. Furthermore, we compare the recognition rate using energy, pitch, LPCC and MFCC features with that using energy, pitch and LPCMCC features to find which method is better.

The paper is organized as follows. Section II describes the database used in the experiments. Section III introduces the automatic speech emotion recognition system. The speech features are presented in this section. Section IV introduces the Support Vector Machine algorithm. Experiments to assess the proposed system are performed in section V. Section VI concludes this paper.

II. SPEECH DATABASE

A. Berlin Database of Emotional Speech

The database used in this paper is Berlin emotional speech database, which is a simulated speech database. It is

an open source speech database and easy to access, and it is frequently used in fields of speech emotion recognition. This database contains seven basic emotions: anger, boredom, disgust, fear, happiness, sadness and neutral. All of the speech samples are simulated by ten professional native German actors (5 actors and 5 actresses). There are totally about 500 speech samples in this database, in which 286 speech samples are of female voice and 207 samples are of male voice. The length of the speech samples varies from 2 seconds to 8 seconds. These samples are divided into about 1100 segments of 2 seconds in our analyzing [1].

III. SYSTEM IMPLEMENTATION

Like the typical pattern recognition system, our speech emotion recognition system contains four main modules: emotional speech input, feature extraction, SVM based classification, and recognized emotion output.

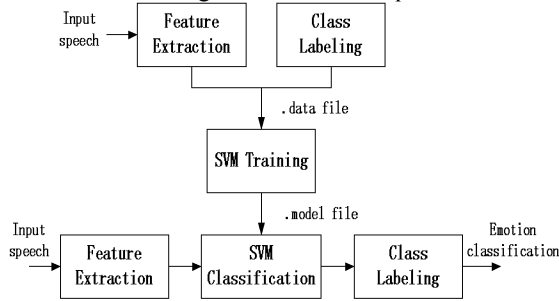


Figure 1. Speech Emotion Recognition System

As Figure 1 shows, the input of this system is the .wav files, which form the different classification of the Berlin emotional speech database. We extract six feature vectors of each speech sample by the feature extraction module: energy, pitch, Linear Prediction Cepstrum Coefficients (LPCC), Mel-Frequency Cepstrum Coefficients (MFCC), Linear Prediction coefficients Mel Cepstrum coefficients (LPCMCC) [9]. After the feature extraction, we give each speech sample with the corresponding emotion class label. After that a “.data” file which contains the class label and feature coefficients is inputted to the LIBSVM classifier, and a “.model” file is got. When an unclassified speech sample come into this system, the system extract the feature coefficients and input to the LIBSVM classifier, the system load the “.model” file and predict the emotion class, the output of a classifier is a label of a particular emotion class.

A. Feature Extraction

The Speech signal contains a large number of parameters that reflect the emotional characteristics, and the different parameters result in changes in emotion. Thus, the most important step in speech emotion recognition is how to extract the feature parameters, which can express mostly the emotion of speech. In recent research of speech emotion recognition, some common features are speech rate, energy, pitch, formant, and some spectrum features, such as Linear Prediction Coefficients (LPC), Linear Prediction Cepstrum

Coefficients (LPCC), Mel-Frequency Cepstrum Coefficients (MFCC) and its first derivative and so on. There have been a handful of research on these features in the past years, and these features have been used in many speech emotion recognition research. In the use of these features parameters while some other features of a good discrimination have been proposed: Linear Frequency Power Coefficients (LFPC), Perspective Linear Prediction (PLP) [4], Mel Energy spectrum Dynamic coefficients (MEDC) [7], Linear Prediction coefficients Mel cepstrum coefficients (LPCMCC) and so on.

In our research, we calculate the statistics of energy, pitch, LPCC, MFCC, LPCMCC, and make use of them to classify the speech emotion.

In the research of speech emotion recognition, seven basic emotions are recognized: anger, happy, sadness, fear, boredom, disgust, and neutral [10]. When people are in different emotional state, their speeches have different changes in speak rate, pitch, energy, and spectrum. Usually, anger has a highest mean value and variance of pitch, and mean value of energy. Disgust and Boredom have a low mean value of pitch and energy. Happy has an improvement of mean value, variation range and variance of pitch, and the mean value of energy. On the contrary, the mean value, variation range and variance of pitch of sadness is decrease, the energy is weak, the speak rate is slow and the decrease of the spectrum in high frequency components. The feature of fear has a high mean value and variation range of pitch and the improvement of spectrum in high frequency components [10]. As a result, we can extract the statistics of pitch, energy and some spectrum features of speech to recognize the emotion in speech.

TABLE I. SUMMARY OF THE EFFECTS OF SEVERAL EMOTION STATES ON SELECTED ACOUSTIC FEATURES

Emotion	Pitch			Energy	Spectrum
	Mean	Variance	Variation range	Mean	High frequency components
Anger	Highest	Highest	Increase	Highest	Most
Disgust	Lowest	-	Increase	Lowest	Decrease
Fear	Highest	-	Increase	Normal	Increase
Boredom	Lowest	-	Decrease	Lowest	-
Happy	Higher	Increase	Increase	Highest	Increase
Sadness	Lower	Decrease	Decrease	Lower	Decrease
Neutral	Normal	Normal	Normal	Normal	Normal

B. Energy and related features

The Energy is an important feature of speech, and the analysis of energy is focused on short-term energy and short-term average amplitude. In order to obtain the statistics of energy feature, we use short-term function to extract the value of energy in each speech frame. Then we can obtain the statistics of energy in the whole speech sample by calculating the energy, such as mean value, max

value, variance, variation range, contour of energy [2]. Each Energy feature vector is 19-dimensional.

- Maximum, Mean, Variance of Energy
- Maximum, Mean, Median duration of rising/falling slopes of Energy
- Maximum, Mean, Median value of rising/falling slopes of Energy
- Interquartile range of rising/falling slopes of Energy
- Interquartile duration of rising/falling slopes of Energy

C. Pitch and related feature

The pitch signal is another important feature in speech emotion recognition. The vibration rate of vocal is called the fundamental frequency F0 or pitch frequency [10]. The pitch signal is also called the glottal wave-form; it has information about emotion, because it depends on the tension of the vocal folds and the sub glottal air pressure, so the mean value of pitch, variance, variation range and the contour is different in seven basic emotional statuses. The method widely used to extract the pitch is based on the short-term autocorrelation function. We calculate the value of pitch frequency in each speech frame, and obtain the statistics of pitch in the whole speech sample. These statistical values reflect the global properties of characteristic parameters. Each Pitch feature vector is 19-dimensional.

- Maximum, Mean, Variance of Pitch
- Maximum, Mean, Median duration of rising/falling slopes of Pitch
- Maximum, Mean, Median value of rising/falling slopes of Pitch
- Interquartile range of rising/falling slopes of Pitch
- Interquartile duration of rising/falling slopes of Pitch

D. Linear Prediction Cepstrum Coefficients (LPCC)

It embodies the characteristics of particular channel of each person, and the same person with different emotional speech will have different channel characteristics, so we can extract these feature coefficients to identify the emotions contained in speech. The advantage is that the computation is small, the algorithm is more efficient and it can better describe the vowels; while the disadvantage is that the description of the consonants are less capable and less noise immunity. The computational method of LPCC is usually a recurrence of computing the linear prediction coefficients (LPC), which is according to the all-pole model.

If we set the exponent number of cepstrum large, more information can be maintained. But as the order increases, the coefficients will become very small and no more effect on it, and will increase the computation. So we extract the first 12-order LPCC coefficients in each frame of speech, and obtain the statistics information by calculating the entire speech signal. That is, 12-order LPCC coefficients of maximum, minimum, mean, variance. Each LPCC feature vector is 48-dimensional.

- Mean, Variance, Maximum, Minimum of each coefficient across all the frames

E. Mel-Frequency Cepstrum Coefficients (MFCC)

There are many researches on MFCC feature parameters at home and abroad, it is widely used in speech recognition and speech emotion recognition studies, and it obtained a good recognition rate. MFCC is based on the characteristics of the human ear's hearing, which uses a nonlinear frequency unit to simulate the human auditory system. Mel frequency scale is the most widely used feature of the speech, with a simple calculation, good ability of the distinction, anti-noise and other advantages [11].

Usually the process of calculating MFCC is shown in Figure 2.

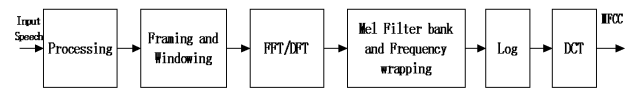


Figure 2. Speech Emotion Recognition System

MFCC in the low frequency region has a good frequency resolution, and the robustness to noise is also very good, but the high frequency coefficient of accuracy is not satisfactory. So we give up the high-level order of the MFCC and use only low-level order as audio feature parameters. In our research, we extract the first 13-order of the MFCC coefficients. For each order coefficients, we compute all the frames of the mean, variance, maximum and minimum in entire speech. Each MFCC feature vector is 52-dimensional.

- Mean, Variance, Maximum, Minimum of each coefficient across all the frames

F. Linear Prediction Coefficients Mel Cepstrum Coefficients (LPCMCC)

The LPCMCC is one of the spectral characteristics; it is based on Mel Frequency and Linear Prediction Cepstrum coefficients (LPCC). The feature combines well the advantage of both Mel Frequency Cepstrum coefficients (MFCC) and Linear Prediction Cepstrum coefficients (LPCC): it not only takes into account channel incentives, but also takes into account the human auditory characteristics; it can be the same as Mel cepstrum features with high-frequency interference shielding, but also has high noise immunity advantages, and can be the same as a linear cepstrum with little computation, the more efficient algorithm [9].

This paper extracts the 14-order LPCMCC coefficients per frame of the voice signal, and calculated for the whole speech signal statistics of mean, variance, maximum and minimum. Each LPCMCC feature vector is 19-dimensional.

- Mean, Variance, Maximum, Minimum of each coefficient across all the frames

IV. SVM CLASSIFICATION ALGORITHM

In recent years in speech emotion recognition, researchers proposed many classification algorithms, such as

Neural Networks (NN), Gaussian Mixture Model (GMM), Hidden Markov Model (HMM), Maximum Likelihood Bayesian classifier (MLC), Kernel Regression and K-nearest Neighbors (KNN) and Support vector machines (SVM). Support vector machine maps the vector to a higher-dimensional vector space, in this space a maximum interval hyper plane to be established. Two parallel hyper planes are established on both sides of the hyper plane. Establish a suitable direction hyper plane to make the distance between the two hyper planes maximization. Its main idea is to use a kernel function to map the original input set to a high dimensional space and then obtain an optimal classification [12] [13]. Since SVM is a simple and efficient computation of machine learning algorithms, and is widely used for pattern recognition and classification problems, and under the conditions of limited training data, it can have a very good classification performance compared to other classifiers [4]. Thus we adopted the support vector machine to classify the speech emotion in this paper.

In our research, we use the system to extract the speech's feature. After the feature extraction, we give each speech sample with the corresponding emotion class label. After that we input them to the LIBSVM classifier and gain a model file by training the data set. When an unclassified speech sample come into this system, the system extract the feature coefficients and use the model file to classify the speech emotion.

V. EXPERIMENTATION AND RESULTS

The Berlin Emotion Database contains 339 speech files for five emotion classes. These five emotions are disgust, border, sad, neutral, and happy, and there are 46, 81, 62, 79, 71 speech utterances for each emotion class respectively. For each emotion, we divide these speech utterances into two subsets as training subset and testing subset. The number of speech utterances for emotion as the training subset of 90%, and 10% as the test subset. Only gender independent is performed.

The LIBSVM is trained only on energy and pitch feature vector using RBF kernel functions. The result is a correct recognition rate of 66.02% when ninety percent of all data are used for training and ten percent of all data are used for testing.

TABLE II. CONFUSION MATRIX OF ONLY USE ENERGY AND PITCH (GENDER INDEPENDENT)

Emotion	Emotion Recognition (%)				
	<i>Disgust</i>	<i>Border</i>	<i>Sad</i>	<i>Neutral</i>	<i>Happy</i>
Disgust	60.87	17.39	8.7	13.04	0
Boredom	11.11	61.73	12.35	14.81	0
Sad	9.68	16.13	59.68	14.52	0
Neutral	10.13	15.19	11.39	63.29	0
Happy	0	8.45	0	7.04	84.51

The LIBSVM is trained only on LPCMCC feature vector using RBF kernel functions. The result is a correct recognition rate of 70.7% is obtained.

TABLE III. CONFUSION MATRIX OF ONLY USE LPCMCC (GENDER INDEPENDENT)

Emotion	Emotion Recognition (%)				
	<i>Disgust</i>	<i>Border</i>	<i>Sad</i>	<i>Neutral</i>	<i>Happy</i>
Disgust	65.22	15.22	8.7	10.87	0
Boredom	7.41	67.9	9.88	12.35	0
Sad	8.06	14.52	66.13	11.29	0
Neutral	7.59	12.66	10.13	68.35	0
Happy	0	7.04	0	7.04	85.92

The LIBSVM is trained both on Energy, pitch and LPCC, MFCC feature vector using RBF kernel functions. The result is a correct recognition rate of 78.17% is obtained.

TABLE IV. CONFUSION MATRIX OF USE ENERGY PITCH LPCC AND MFCC (GENDER INDEPENDENT)

Emotion	Emotion Recognition (%)				
	<i>Disgust</i>	<i>Border</i>	<i>Sad</i>	<i>Neutral</i>	<i>Happy</i>
Disgust	76.09	8.7	6.52	8.7	0
Boredom	6.17	74.07	7.41	11.11	1.23
Sad	4.84	9.68	77.42	8.06	0
Neutral	5.06	11.39	5.06	75.95	2.53
Happy	0	5.63	0	7.04	87.32

The LIBSVM is trained both on Energy, pitch and LPCMCC feature vector using RBF kernel functions. The result is a correct recognition rate of 82.5% is obtained.

TABLE V. CONFUSION MATRIX OF USE ENERGY PITCH AND LPCMCC (GENDER INDEPENDENT)

Emotion	Emotion Recognition (%)				
	<i>Disgust</i>	<i>Border</i>	<i>Sad</i>	<i>Neutral</i>	<i>Happy</i>
Disgust	80.43	8.7	2.17	8.7	0
Boredom	4.94	80.25	4.94	8.64	1.23
Sad	3.23	8.06	80.65	8.06	0
Neutral	3.8	10.13	5.06	81.01	1.27
Happy	0	4.23	0	5.63	90.14

VI. CONCLUSION AND FUTURE WORKS

As can be seen from the experiment, the emotion recognition rate of the system which only uses the spectrum features of speech is slightly higher than that only uses the prosodic features of speech. And the system that uses both spectral and prosodic features is better than that only uses spectrum or prosodic features. The recognition accuracy of the former is higher than the latter. Meanwhile, in the circumstance of combining spectral and prosodic features, the recognition rate of that use energy, pitch, LPCC and

MFCC features is slightly higher than that use energy, pitch and LPCMCC features. Both the spectrum and prosodic features contain speech emotion characteristic and the combination of them can better depict the people's speech emotion.

And the result can be seen from the experiment, happy always has a high recognition rate than other emotions and its misclassification is small and the other four emotions misclassification with each other easily.

To extract the more effective features of speech and enhance the emotion recognition accuracy is our future work. More work is needed to improve the system so that it can be used in real-time speech emotion recognition.

ACKNOWLEDGEMENT

This work was carried out as part of the 'Research on Affective e-Learning Model Based on Multimodal Emotion Recognition' Project Supported by national Natural Science Foundation of China under Grant No. 60873132.

REFERENCES

- [1] <http://www.expressive-speech.net/>, Berlin emotional speech database
- [2] D. Ververidis, C. Kotropoulos, and I. Pitas, "Automatic emotional speech classification", in Proc. 2004 IEEE Int. Conf. Acoustics, Speech and Signal Processing, vol. 1, pp. 593-596, Montreal, May 2004.
- [3] Xiao, Z., E. Dellandrea, Dou W., Chen L., "Features extraction and selection for emotional speech classification", 2005 IEEE Conference on Advanced Video and Signal Based Surveillance (AVSS), pp.411-416, Sept 2005.
- [4] T.-L. Pao, Y.-T. Chen, J.-H. Yeh, P.-J. Li, "Mandarin emotional speech recognition based on SVM and NN", Proceedings of the 18th International Conference on Pattern Recognition (ICPR'06), vol. 1, pp. 1096-1100, September 2006.
- [5] Xia Mao, Lijiang Chen, Liqin Fu, "Multi-level Speech Emotion Recognition Based on HMM and ANN", 2009 WRI World Congress, Computer Science and Information Engineering, pp.225-229, March 2009.
- [6] B. Schuller, G. Rigoll, M. Lang, "Hidden Markov model-based speech emotion recognition", Proceedings of the IEEE ICASSP Conference on Acoustics, Speech and Signal Processing, vol.2, pp. 1-4, April 2003.
- [7] Yashpalsing Chavhan, M. L. Dhore, Pallavi Yesaware, "Speech Emotion Recognition Using Support Vector Machine", International Journal of Computer Applications, vol.1, pp.6-9, February 2010.
- [8] Zhou Y, Sun Y, Zhang J, Yan Y, "Speech Emotion Recognition Using Both Spectral and Prosodic Features", ICIECS 2009. International Conference on Information Engineering and Computer Science, pp.1-4, Dec.2009.
- [9] An X, Zhang X, "Speech Emotion Recognition Based on LPMCC", Sciencepaper Online.2010.
- [10] D. Ververidis and C. Kotropoulos, "Emotional Speech Recognition: Resources, features and methods", *Elsevier Speech communication*, vol. 48, no. 9, pp. 1162-1181, September, 2006.
- [11] Han Y, Wang G, Yang Y, "Speech emotion recognition based on MFCC", Journal of ChongQing University of Posts and Telecommunications(Natural Science Edition),20(5),2008.
- [12] 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>.
- [13] Lin Y, Wei G, "Speech emotion recognition based on HMM and SVM". Proceedings of 2005 International Conference on Machine Learning and Cybernetics, vol.8, pp. 4898-4901. Agu 2005.