



Robust Speech Recognizer using multiclass SVM

Inge Gavat, Gabriel Costache, Claudia Iancu

University Politehnica of Bucharest

igavat@alpha.imag.pub.ro

ABSTRACT. In this paper a robust speech recognizer is presented based on features obtained from the speech signal and also from the image of the speaker. The features were combined by simple concatenation, resulting composed feature vectors to train the models corresponding to each class. For recognition, the classification process relies on a very effective algorithm, namely the multiclass SVM. Under additive noise conditions the bimodal system based on combined features acts better than the unimodal system, based only on the speech features, the added information obtained from the image playing an important role in robustness improvement.

Keywords: robust speech, bimodal system, support vector machines, neural networks.

1. Introduction

The main problem of many classification systems is that there are not robust, their performances are not constant especially when the conditions (environment, user, application) are changed. There are two causes for that: first the source of the signals that should be classified can be corrupted with noisy unwanted components and second, the classifiers cannot deal properly with new variants of the same pattern. Concerning the first problem, in image classification systems (especially face recognition or detection) for example, different illuminations and positions of the object to be recognized can be seen as introducing unwanted components. In the audio classification systems such unwanted components are represented by the inherent noise that is captured along with the signal to be classified. The usual solution for this kind of problems is a preprocessing of the signal before classification in order to eliminate the unwanted components, with the drawback to affect also the original signal. Another more viable alternative could be the use of features obtained from more sources, connected with the object to be classified, acting in a multimodal way. Concerning the second problem, Artificial Neural Networks[4] and many statistical methods offer solutions by allowing to form models of one pattern using more variants of the pattern. Furthermore these models can be re-trained using new particular occurrences of the pattern so that the system is able to learn from examples.

In this paper is proposed a robust speech recognition system, based on a bimodal structure using features obtained from two sources: the speech signal and the speaker image. For classification is applied the Support Vector Machines[8] algorithm that combines the advantages of ANNs and statistical approaches by having

good generalization and learning properties. SVMs was successfully used in a multimedia classifier[2].

A bimodal system is a particular case of multimodal system, namely that system that uses features obtained not only from the signals that should be classified but also from other signals related with them.

The bimodal systems act in two major steps like each unimodal classification system. In the first step feature extraction is performed, where are determined only the important characteristics of the signal, in the second, the recognition is realized, where based on a classification algorithm is made a decision. There are two main strategies to build multimodal system[1].

The first method is to apply decision fusion and means taking a decision for each source of information and combine those two to make the final decision. The most common way to implement the decision fusion algorithm is using neural networks or Markov models where the entries of the network are the output of each classifier from each source.

The other method to construct a multimodal system is using the feature fusion. This means that after feature extraction from each source a combined feature set is realized as basis for multimodal models and then applying any classification method we make the decision. The main disadvantage of this second method is that we have to synchronize the signals from the different sources.

For each source we can use different parameterization methods depending on the signals. Depending on the application, the signals can be images, audio signals and others.

2. Architecture

The recognition system we have experimented is given in Figure 1 and is based on fusion of parameters obtained from speech and from image. In order to combine the feature vectors, the two signals have to be synchronized, this being the main weakness of this type of bimodal system architecture. Because the database we used had synchronization between speech and image, we can apply without problems the architecture based on parameters fusion.

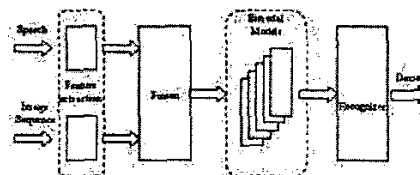


Figure 1 Bimodal speech recognition system

The first step in the system is feature extraction where we extract only the important characteristics of the signals. For speech parameterization we used perceptual approaches of two well known methods: linear prediction and cepstral analysis. From image we extract geometric features of speaker's mouth. For that, first a face tracking algorithm based on Gaussian Mixture Models and then a deformable template was used to model the face. The deformation was calculated so the template would contain as many pixels from the face as possible. The decision for each pixel to be or not in the face class was taken using the Bayes statistical criteria.

Features were combined by simple concatenation of feature vectors for each analyzed window (or frame). After fusion we construct bimodal models for the patterns we want to classify.

For classification we choused to use a statistical approach called Support Vector Machines. SVM is a binary decision method with a good generalization property and is based on finding an optimal hyperplane as a decision boundary between classes. Also SVM is a kernel method meaning that the hyperplane is found in a feature space using a non-linear transformation which transform the input space in a feature space which has a much bigger dimension and we don't have to calculate the transformation for each data sample, we have to calculate only some kernel products in order to find the hyperplane.

In order to extend the binary algorithm to multiclass decision we combined several binaries SVMs using Directly Acyclic Graph SVM (DAGSVM) algorithm.

The first stage in the classification process is to train the support vector network (find the hyperplane) using some of the data samples (bimodal models) from the database and next we test the trained network using the other models from database or the same models used in the training process.

3. Support Vector Machines

The foundations of Support Vector Machines (SVM) have been developed by Vapnik[8]. The formulation is based on Structural Risk Minimization (SRM) principle, which minimizes an upper bound on the generalization error, as opposed to Empirical Risk Minimization (ERM) which minimizes the error on the training data.

Support Vector Machines (SVM) is a statistical algorithm with a great potential to generalize, that can successfully be used in pattern recognition and information retrieval tasks. The main idea in training a SVM system is finding a hyperplane as a decision boundary between two classes. Fundamentally SVM is a binary decision method, but there are several techniques that allow the use in classification tasks with more than 2 classes. In the case of separable patterns, Figure 2 represents in two dimensional orientation what support vectors are.

The equation that is verified for each data sample is :

$$d_i(w^T x_i + b) \geq 1 \quad \text{for } i=1,2,\dots,N \quad (1)$$

where d_i is the label for sample data x_i and it can be +1 or -1 and w_i and b are the weights and the bias which

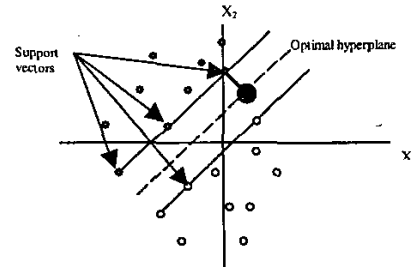


Figure 2 Separable patterns

describe the hyperplane. The support vectors are the data samples for which the eq 1 is verified with the equal sign. After the training process only the support vectors will be kept from all data. In the case of non separable patterns, Figure 3 is representative

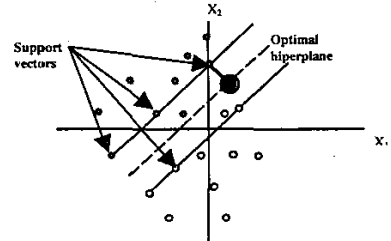


Figure 3 Non separable patterns

In this case, eq 1 becomes:

$$d_i(w^T x_i + b) \geq 1 - \xi_i \quad \text{for } i=1,2,\dots,N \quad (2)$$

where ξ_i represents the number of data samples left inside the decision area, giving the number of training errors. The problem of finding the optimal hyperplane becomes a problem of minimizing the cost function described by the eq 3:

$$\Phi(w, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i \quad (3)$$

where minimizing the first term means maximizing the distance between the two classes and minimizing the second term means reducing the number of training errors. Under those circumstances, parameter C becomes a balance between a smaller training error and a bigger distance between classes. The minimization of eq 3 is done using Lagrange multipliers method.

Another important part of SVM is the use of the inner product kernel functions. Cover's theorem says that giving a input space where the patterns are non separable, there is a transformation that will lead to another space where with high probability the patterns are separable with two conditions: one, the transformation is non linear and two, the dimension of the output space is high enough. We can use this theorem in solving the Lagrange multipliers systems. We will not calculate the transformation for each data sample in the output space,

we will only have to calculate products called inner product kernels, like in the equation 4:

$$K(x, x_i) = \phi^T(x_i)\phi(x) = \sum_{j=0}^{m_i} \phi_j(x)\phi_j(x_i) \text{ for } i=1,2,\dots,N \quad (4)$$

We can use any type of kernels: polynomial, radial basis function, two layers perceptron and so on. Figure 4 gives an example of how a polynomial kernel works:

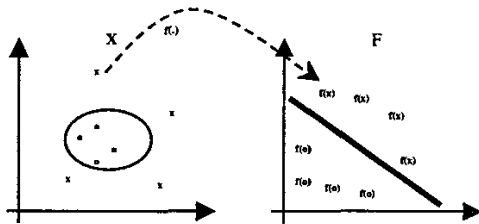


Figure 4 Polynomial kernel

Using this kernel architecture SVM can be seen as a NN based system with 3 layers: first input layer with the dimension equal with the number of features of the pattern, than an hidden layer in the future (kernel) space and finally the output layer which will give the binary decision. This architecture is described in the next figure.

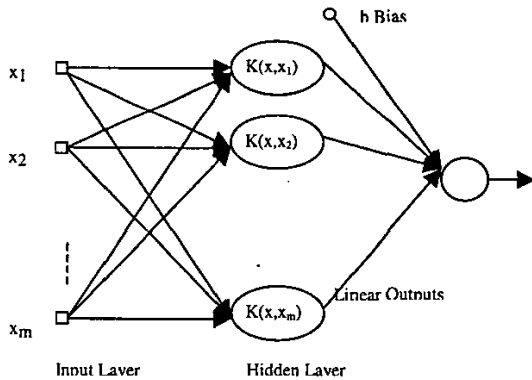


Figure 5 SVM network

4. Multiclass SVM

Like we said in the beginning, SVM is a binary decision method but it can be extended to multiclass task using different algorithms. The most common algorithms use combinations of binary SVMs.: 'one against one' method, 'one against all' method and DAGSVM (Directly Acyclic Graph SVM). The oldest method, 'one against all', consists in building several binary SVMs (equal with the number of classes). In the training phase we will train each SVM with one of the classes against the rest of the classes and in the testing phase we test the test data with all SVMs and the decision is taken based on the distance between data test and the hyperplanes from all SVMs.

'One against one' method consists in building more binary SVMs where we train each class with another class until we trained each class with all the other classes. In the testing phase we test the current data with all SVMs and if for the classes (i,j) binary SVM the decision is that is in the class i for example the index of i is increased with one and the index of j is decreased with one. At the end the decision is taken based on the biggest index. This is why this method is also called 'max wins' method. The DAGSVM method is similar to the 'max wins' method, but we construct only k(k-1) binary SVMs if we have k classes, then form the tree given in Figure 6. In the testing phase, we start at the top of the tree and if the decision is that is in the i class then we go to the left path if not we go to the right path and continuing until the end of the tree where we will have the final decision

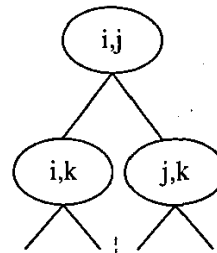


Figure 6 DAGSVM tree

This method is faster than all the other and our test proved that the results obtained with DAGSVM is very closed to the one obtained with 'max wins' method

5. Experimental results

We used for feature extraction from the speech signal two methods: the perceptual linear prediction (PLP)[5] analysis and the mel-cepstral analysis[3]. For each window, we extract 5 PLP[7] coefficients or 13 mel-cepstral coefficients. For the image sequence we use a face-tracking algorithm and we extract geometric features of the speaker face. For each frame we extract 3 geometric features (the mouth width and the height of the upper lip and downer lip)[6]. For synchronization between image and speech, the video sequence was recorded at 30 fps and we made the length of the analysis window for the speech to be 33ms. So for each frame we will have 3 features from image and 5 PLP or 13 MFC coefficients for speech. For fusion we used simple concatenation between the two feature vectors. Then we formed a 'supervector' putting together the features calculated for each window and we construct bimodal models using those 'supervectors'. For classification we used the DAGSVM algorithm.

We tested our system using database created by the Advanced Multimedia Laboratory from the Carnegie Mellon University. The database contains 10 words (digits from one to ten) spoken by 10 peoples each with 10 pronunciations.

We performed two types of tests: first with enrolled speakers, namely speakers involved both in training SVMs and testing SVMs. We used five pronunciations for training and five for testing. For the second type of test with unenrolled speakers we used the leave- one- out method. For each word we trained the SVM net with 9 speakers and tested with the 10th repeating the procedure for each speaker.

The results are presented in the Table 1 and Table 2

Table 1 Recognition rates for unenrolled speakers

Coefficients [Numberof]	PLP(5)	MFCC (13)	PLP+image (8)	MFCC+image (16)
SNR=30dB	84.75%	91.71%	87.73%	92.84%
SNR=25dB	77.71%	90.49%	82.42%	92.49%
SNR=19dB	76.8%	87.89%	80.08%	91.13%

Table 2 Recognition rates for enrolled speakers

Coefficients [Numberof]	PLP (5)	MFCC (13)	PLP+image (8)	MFCC+image (16)
SNR=30dB	91.71%	97.74%	92.13%	97.42%
SNR=25dB	90.85%	96.53%	91.98%	96.85%
SNR=19dB	86.71%	94.13%	91.85%	96.14%

In figure 7 are represented the variations of recognition rates when artificial noise is added over the speech signal.

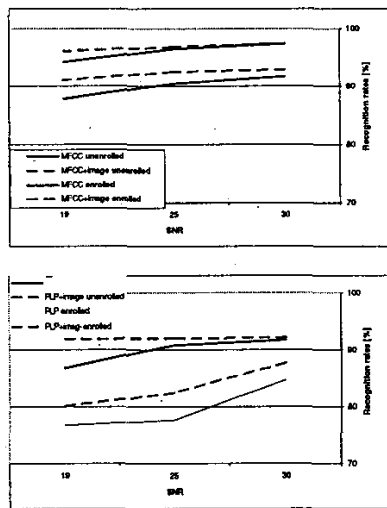


Figure 7 Experimental results

The performance obtained using bimodal recognition compared with classic unimodal recognition based only on the speech signal is sensible higher, especially under difficult conditions, namely when the speech signal is corrupted with noise. It can be observed that when using coefficients both from image and speech the variations of recognition rates are considerably smaller than when using only speech parameters.

5. Conclusions

In this paper a new approach for building robust speech recognizer systems was presented. The robustness was accomplished by using additional features obtained from the speaker image along with the features obtained from the speech signal. We extract features from the speech signal using the PLP and the mel-cepstral technique and from the image of the speaker we extract geometric features. For classification we used the SVM algorithm which we extended to multiclass decision using the DAGSVM algorithm. The experimental results confirmed the stability of the recognition rates when we added artificial noise over the speech signal. Another observation from the experimental results is that when using the MFC coefficients (best 97.42%) the rate of recognition is higher than when using PLP coefficients (best 91.71%). The difference between recognition rates for the enrolled speakers (best 97.42%) and for unenrolled speakers (best 92.84) is not so high which indicate that SVM has a good generalization property.

References:

- [1] C.C. Chibelushi, F. Deravi and J.S.D Mason A *Review of Speech Based Bimodal Recognition* IEEE Trans on Multimedia Vol 4 Nr 1 March 2002 pp 23-38
- [2] G. Costache, I. Gavata *Multimedia Classifier* Proceedings ESA-EUSC 2004 Conference, Madrid, Spain March 17-18, 2004
- [3] I. Gavata, C.O. Dumitru, G. Costache and D. Militaru *Continuous Speech Recognition Based on Statistical Methods* Proceedings of Sped2003 Conference Bucharest, Romania April 10-11 2003 pp115-127
- [4] S. Haykin, *Neural Networks: A Comprehensive Foundation*, Prentice Hall, 1999
- [5] H. Hermansky, *Perceptual Linear Predictive (PLP) Analysis of Speech*, J. Acoust. Soc. America, pp. 1738-1752, vol. 87, no. 4, Apr. 1990.
- [6] F. J. Huang and T. Chen, *Real-Time Lip-Synch Face Animation Driven by Human Voice*, IEEE Workshop on Multimedia Signal Processing, Los Angeles, California, Dec 1998
- [7] J. Koehler, N. Morgan, H. Hermansky Hirsch H. and Tong G.: *Integrating RASTA-PLP into Speech Recognition*, in Proceedings ICASSP'94, pp. 421-424, vol. 1, Adelaide, Australia, April 1994.
- [8] V. Vapnik, *An overview of statistical learning theory* IEEE Transactions on Neural Networks Vol 10, No 5, 1999, pp. 988-1000