

SPEAKER VERIFICATION EMPLOYING TESPAP CODING AND DTW

Petre G. POP¹, Eugen LUPU¹

Abstract – TESPAP (Time Encoding Signal Processing and Recognition) is a processing and recognition method in the time domain, proposed by [1]. The key problem in TESPAP is to define the alphabet used for the coding process, alphabet usually generated by a quantization process. To avoid this complicated process, this paper presents an approach in which we used DTW (Dynamic Time Warping) to align sequences of epochs in order to generate a verification decision.

Keywords: TESPAP, DTW, speaker verification.

I. INTRODUCTION

The speaker verification consists of automatically authenticating the identity claimed by a speaker, given only some samples of his voice. There are three categories of approaches in speaker verification. In the first one, the verification system is trained on a particular utterance and the same utterance is later spoken by the speaker who claims that identity, making up **text-dependent speaker verification**. Within the second approach, verification decisions are based on utterances selected by the speaker and not previously known by the verification system, making up the **text-independent speaker verification**. In **text-prompted** approach, the verification system generate the text that each speaker has to utter in both training and testing stage.

A typical approach to the text dependent speaker verification is DTW (Dynamic Time Warping) in which the unknown speaker's utterances are time aligned to the reference stored for the speaker whose identity is claimed, and the decision to accept or reject is based on a measure of similarity between two time series of parameters corresponding to the utterances.

TESPAP is a method based on the approximations to the locations of real and complex zeros, derived from an analysis of a band-limited signal. The key features of the TESPAP coding in the speech processing field are the following:

- the capability to separate and classify many signals that are indistinguishable in the frequency domain;
- an ability to code the time varying speech waveforms into optimum configurations for processing with Neural Networks.

II. TESPAP

A. BASICS

TESPAP method is based on the approximations to the locations of the real and complex zeros, derived from the analysis of a band limited signal. The real zeros correspond to the zero crossings of the signal while the complex zeros are associated with local maxima. Numerical descriptors of the signal waveform may be obtained via the classical Shanon numbers resulted from the analysis [2].

The Shanon model involves detecting the ordinates of a waveform at a series of points equally spaced at $1/2W$. A variety of mainly linear transforms (e.g. Fourier, LPC, Wavelet or Walsh) has been developed for describing and classifying key features of the sampled data set. This coding strategies involve the following requirements:

- the use of amplitude descriptors;
- the use of regular sampling;
- an approximation domain dependent upon the numbers of bits per sample.

TESPAP coding is based on the zero-crossings of the signal analysed.

B. TESPAP CODING

The key in the interpretation of the TESPAP coding possibilities consists in the complex zeros concept. The band-limited signals generated by natural information sources include complex zeros that are not physically detectable. The real zeros of a function (represented the zero crossing) and some complex zeros can be detected by visual inspection (Fig. 1), but the detection of all zeros (real and complex) is a complex task.

Locating all complex zeros involves the numerical factorization of a $2TW^n$ -order polynomial [3]. A signal waveform of bandwidth W and duration T , contains $2TW$ zeros, a number which usually exceeds several thousand [2]. The numerical factorization of a

¹ Technical University of Cluj-Napoca, Faculty of Electronics & Telecommunications Comm. Dept.
email : Petre.Pop@com.utcluj.ro

$2TW^h$ -order polynomial is computationally infeasible for real time.

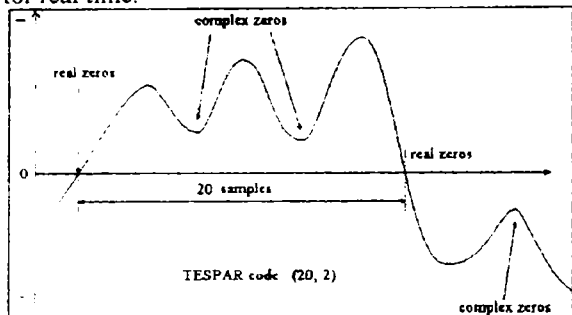


Fig. 1. TESPAP waveform analysis

The key to overcome this drawback is to introduce an approximation in the complex zeros location. Instead of detecting all zeros of the function the following procedure may be used:

- the waveform is segmented between successive real zeros;
- this duration information is combined with simple approximations of the wave shape between these two locations.

Only some of the complex zeros that can be identified directly from the waveform by these approximations detect. In this transformation of signals, from time-domain in the zero-domain, the real zeros, in the time-domain, are identical to the locations of the real zeros in the zero-domain, and the complex zeros occur in conjugate pairs associated with features (minima, maxima, points of inflexion etc.) in the wave shape that appear between the real zeros. In this way an important subset of complex zeros may be identified by examining the features of the wave shape between its successive real zeros.

In the simplest implementation of the TESPAP method [1], two descriptors are associated with every segment or epoch of the waveform, in order to generate the TESPAP symbol alphabet :

- the duration between successive real zeros (in number of samples);
- the shape between two successive real zeros.

In this simple TESPAP model implementation, not all complex zeros can be identified from the wave shape, so the approximation is limited to those zeros that can be so identified.

The band-limitation of the signal imposes significant restrictions upon the maximum and minimum duration of any epoch, and also upon the maximum number of significant waveform extreme points that each epoch may contain. The longest epoch may have a duration approximately equal to half the period of the lowest frequency component allowed by band-limiting; the shortest epoch may have a duration approximately equal to half the period of the highest frequency component allowed within the band of signal. Also, short epochs have no or few features, whilst long epochs may contain few or many features. For the simplest implementation, each epoch may be

classified in terms of its duration (D) - number of samples and the number of minima (S), that it contains.

The TESPAP coding process is presented in Fig.2, using an alphabet (symbol table) to map the duration/shape (D/S) attributes of each epoch to a single descriptor or symbol.

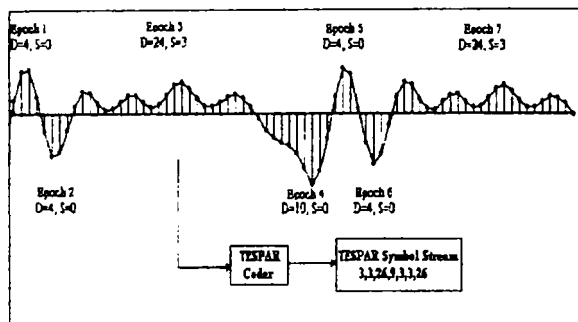


Fig. 2. TESPAP coding process

The TESPAP symbols string may be converted into a variety of fixed-dimension matrices [3]. For example, the S-matrix is a single dimension $1 \times N$ (N - number of symbols of the alphabet) vector, which contains the histogram of symbols that appear in the data stream. Another option is the A-matrix, which is a two dimensional $N \times N$ matrix that contains the number of times each pair of symbols appears with a possible lag of n symbols. The matrices obtained in the training phase are compared to that obtained in the testing phase allowing tasks like verification or recognition. The TESPAP alphabet may be generated in vector quantization process or using neural networks [7].

III. DYNAMIC TIME WARPING

DTW is a recursive recognition algorithm, which is usually used to evaluate a distance between a previously stored reference set, and a test set of speech parameters. The main advantage of this method consists of temporal alignment of the two compared sets of data. In fact, DTW seeks a way between reference and test data so that the cumulated distance to be minimal (Fig. 3).

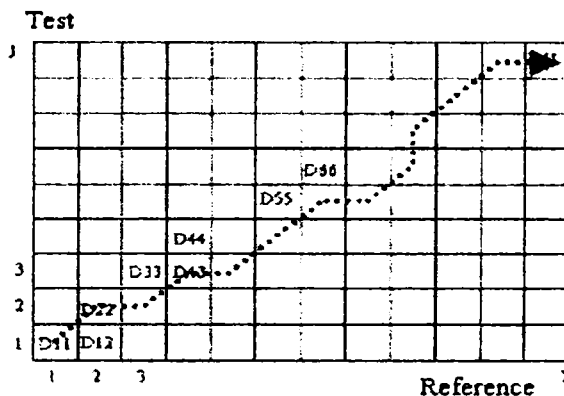


Fig. 3. The optimal path in DTW algorithm

The cumulated distance between the two utterances A and B is given by [6]:

$$D(F) = D(A, B) = \frac{g(i, j)}{(i + j)}, \quad (1)$$

where:

$$g(i, j) = \min \begin{pmatrix} g(i, j-1) + d(i, j) \\ g(i-1, j-1) + 2d(i, j) \\ g(i-1, j) + d(i, j) \end{pmatrix}, \quad (2)$$

and:

$$g(i, 1) = 2 \cdot d(i, 1) \quad (3)$$

In previous equations, $d(i, j)$ is the Euclidean distance between two parameters vectors of the reference and test utterances.

IV. THE THRESHOLD COMPUTING

A major problem in speaker verification consists in formulating a criterion for accepting/rejecting the speaker. To decide whether to reject a speaker or not, for a particular utterance, a threshold is associated to each speaker. An unknown speaker is rejected if its distortion exceeds the threshold.

One way to compute the threshold for a given speaker is to estimate the parameters for two Gaussians distributions: the *in-class* distribution of the distortion obtained by encoding utterances from that speaker in his codebook and an *out-of-class* distribution of the distortion obtained by encoding utterances spoken by other speakers. Equalising the overlapping areas of the two distributions, thus equalising the expected numbers of false acceptances and false rejections, chooses the threshold. The threshold computation involves the following steps:

- compute the mean distortion μ_i^{in} resulted from encoding the training set of the speaker "i" in his codebook and the corresponding standard deviation σ_i^{in} ;
- compute μ_i^{out} , the mean distortion obtained by encoding utterances not spoken by the speaker "i", using the "i" speaker's codebook and the corresponding standard deviation σ_i^{out} .

To equalize the numbers of false rejection and false acceptances, the threshold T_i is chosen to be at an equal number of standard deviations away of each mean (Fig. 4).

$$T_i = \frac{\mu_i^{in} \sigma_i^{out} + \mu_i^{out} \sigma_i^{in}}{\sigma_i^{in} + \sigma_i^{out}} \quad (4)$$

This method for threshold computation assumes Gaussians distributions. As distortion metric, the Euclidean distance was employed.

Another way to compute the decision threshold for each speaker enrolled in the experiment was proposed

in [5] and implies the clustering of the reference and test utterances in the parameters space (Fig. 5).

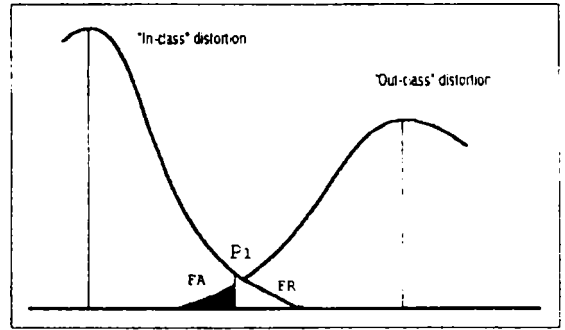


Fig. 4. Threshold computing based on mean distortions and standard deviations

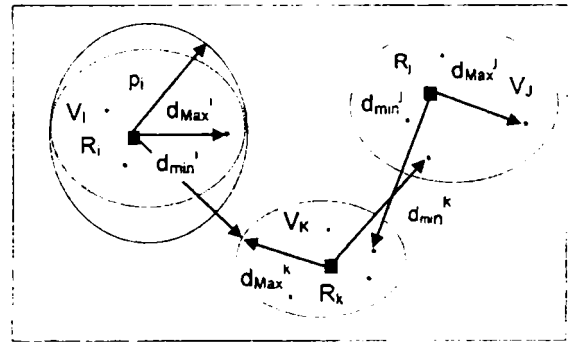


Fig. 5. Threshold estimation based on averaged distances

First, the maximum distance between the optimal reference and the other references for the "i" speaker is computed:

$$d_{MAX}^i = d(R_i, r_{im}), \quad m=1, \dots, M, \quad (5)$$

where d is Euclidean distance and M is the number of references for each speaker.

Second, the minimum distance between the optimal reference of speaker "i" and the other references of the speakers different from "i" is also computed:

$$d_{min}^i = d(R_i, r_{km}), \quad m=1, \dots, M; \quad k=1, \dots, K, \quad (6)$$

where K is the number of the enrolled speakers.

Finally, the threshold p_i for the speaker "i", is computed as the average of the two distances previously defined [4]:

$$T_i = \frac{d_{MAX}^i + d_{min}^i}{2} \quad (7)$$

V. COMBINING TESPAP AND DTW

The key problem in TESPAP is to define the alphabet used for the coding process, alphabet usually generated by a quantization process. This process is influenced by the maximum number for shape and duration that strongly depend on sampling rate. To

avoid this complicated process, in our approach we used DTW (Dynamic Time Warping) to align sequences of epochs in order to generate a verification decision. We used TESPAP to generate speech features (sequences of epochs, each with specific shape and duration) for each speaker utterance, in training and testing stage. DTW is used for successive alignments to compute a speaker model in training stage and in testing stage in order to evaluate a distance between the speaker's models and the test utterance model.

Our speaker verification system works as follows:

- all training utterances are processed with TESPAP method and the resulting epochs (D/S) are saved; before TESPAP coding an endpoint detection process is applied to each utterance;
- DTW is used to build a model for each speaker by successive alignments between epochs files corresponding to each training utterance; each model consists of a sequence of epochs;
- speaker's models and training epochs files are used to compute a threshold for each speaker which will be used in testing stage;
- each test utterance is processed with TESPAP method;
- the resulting epochs string and speaker models are time-aligned with DTW and a distance is computed for each speaker;
- this distance is compared with each speaker threshold and a verification decision is emerged;

We used a weighted Euclidean distance in DTW alignment to discriminate between the contribution of shape and duration.

VI. EXPERIMENTS

Two particular Romanian utterances, "Lămâia ia anemia" (U1) and "Aoleu lăna are molii" (U2) were used for speaker verification experiments. The experiments involved 25 speakers (15 males and 10 females) and 5 utterances for each test phrase were collected from each speaker, 3 utterances were used for training and 2 utterances for testing.

We used Error Recognition Rate (ERR) as verification criteria:

$$ERR = \sqrt{FAR \cdot FRR} \quad (8)$$

where FAR is the False Acceptance Rate and FRR is the False Rejection Rate.

The experiments were carried out for each utterance (U1, U2), using both types of decision threshold (Th₁-eq.4, Th₂-eq.7).

The speaker verification results are presented in the following table.

Table 1

	Th ₁	Th ₂
ERR[%]	11.4	9.7

Certain conclusions can be outlined from these experiments:

- the Th₂ decision threshold leads to better results than Th₁;
- the verification performances are worse than other methods (ERR <= 2%);
- data reduction is very high, because a set of speech samples corresponding to two successive zero-crossings is replaced by two TESPAP values, S(shape) and D(duration);
- all calculations are made in the time domain.

V. CONCLUSIONS

In this paper, we presented an approach to text dependent speaker verification using a combined TESPAP-DTW method in which TESPAP coding is used to generate speech features and DTW is used to generate a model for each speaker from training utterances as well as in the test stage to compute a distance between the test utterance and speaker's models.

The verification experimental results show medium performances. The decision threshold based on averaged distances generates better results than threshold based on mean distortions and standard deviations.

This approach seems to be promising because all calculations are made in the time domain and data reduction is about 15-20 times.

REFERENCES

- [1] R. A. King, T. C. Phipps. "Shannon, TESPAP And Approximation Strategies", ICSPAT 98, Vol. 2, pp. 1204-1212, Toronto, Canada, September 1998.
- [2] A. A. G. Requicha, "The zeros of entire functions, theory and engineering applications", Proceedings of the IEEE, vol. 68 no 3, pp. 308-328, March 1980.
- [3] J. Holbeche, R. D. Hughes and R. A. King, "Time Encoded Speech (TES) descriptors as a symbol feature set for voice recognition systems", IEE Int. Conf. Speech Input/Output, Techniques and Applications, pp. 310-315, March 1996.
- [4] Rabiner, L.R., Juang, B.H. Fundamentals of speech recognition, Prentice-Hall International, Inc., 1993.
- [5] Lupu E., Pop G.P., Todorean G., "Speaker Verification Using Vector Quantization", Proceedings of the First Workshop on Text, Speech, Dialog (TDS 98) Brno, Czech Republic 23-26 Sept. 1998, p. 275-380.
- [6] Furui, S. Digital Speech Processing, synthesis and recognition, Marcel Dekker Publications, 1989, 2001.
- [7] Lupu E., Moca V., Pop G.P., "Environment for Speaker Recognition Using Speech Coding", Proceedings of the International Conference COMMUNICATIONS 2004, June 03-05, 2004, Bucharest, pp. 199-204.