

Subspace and Hypothesis based Effective Segmentation of Co-articulated Basic-units for Concatenative Speech Synthesis

R. Muralishankar, R. Srikanth and A. G. Ramakrishnan
Department of Electrical Engineering
Indian Institute of Science
Bangalore, INDIA
080-293-2935
sripad, srikanth, ramkiag@ee.iisc.ernet.in

Abstract— In this paper, we present two new methods for Vowel-Consonant segmentation of a co-articulated basic-units employed in our Thirukkural Tamil Text-to-Speechsynthesis system [1]. The basic-units considered in [1] are CV, VC, VCV, VCCV and VCCC, where C stands for a consonant and V for any vowel. In the first method, we use subspace-based approach for *vowel-consonant* segmentation. It uses *orientedprincipal component analysis* (OPCA) where the test feature vectors are projected on to the V and C subspaces. The crossover of the norm-contours obtained by projecting test basic-unit onto the V and C subspaces give the segmentation points which in turn helps in identifying the V and C durations of a test basic-unit. In the second method, we use *probabilistic principal component analysis* (PPCA) [2] to get probability models for V and C. We then use Neymen-Pearson (NP) test to segment the basic-unit into V and C. Finally, we show that the hypothesis testing turns out to be an energy detector for V-C segmentation which is similar to the first method.

1. INTRODUCTION

For the purpose of recognition or synthesis, speech often needs to be segmented into phonetic units. Manual segmentation is tedious, time consuming and error prone. Due to variability both in human visual and acoustic perceptual capability, it is almost impossible to reproduce the manual segmentation results. Hence manual segmentation is inherently inconsistent. Automatic segmentation is not faultless, but it is inherently consistent and results are reproducible. Ideally, one likes to have an automatic segmentation which can handle basic-units uttered by different speakers. There are two broad categories of speech segmentation [3] namely, *implicit* and *explicit*. Implicit methods split up the utterance without explicit information, such as the phonetic transcription, and are based on the definition of a segment as a spectrally stable part of a signal.

Motivation for Subspace based Segmentation

In our synthesis scheme [1], concatenation is always performed across identical vowels. Changes in duration, pitch and amplitude are obtained by processing the vowel parts of

the basic-units only. Thus, the segmentation of basic-units into *vowel* and *consonant* parts is needed to keep the consonant portion of the waveform intact. Plosives, affricates and fricatives have a common property of low energy when compared with any of the vowels. Figure 1 shows the performance of energy based segmentation for plosive and co-articulated basic-units. As shown in Figs. 1(a) and (b), accurate segmentation can be obtained for non co-articulated units, and not for co-articulated basic-units. The true consonant part /y/ in the signal /eyo/ is shown in Fig. 1(b) with the boundaries dotted.

We propose two methods for co-articulated basic-unit segmentation. In the first method, we use subspace approach using *orientedprincipal component analysis* (OPCA) for basic-unit segmentation. In the second method, we use *probabilisticprincipal component analysis* (PPCA) [2] to obtain probability models for vowel (V) and consonants (C). We employ Neymen-Pearson (NP) test using the probability models for basic-unit segmentation.

2. SUBSPACE BASED SEGMENTATION

When we consider an individual vowel or consonant, there exist techniques like LPC to model their statistical properties. While segmenting the vowel part of a basic-unit, we can consider the vowel information (VI) in the feature vectors as the signal and the consonant information (CI) as noise. Similarly when the segmentation of the consonant part is required, we can view CI as signal and VI as noise. We present a linear feature transformation that aims at finding a subspace, of the feature space, in which the Signal-to-Noise ratio (SNR) is maximum. Such a decomposition can be arrived at by representing VI and CI by training vectors obtained using manual segmentation. The directions in the feature space where the SNR is maximum can be obtained by the *generalized eigenvalue decomposition* of the covariance matrices of the above vectors. Consider a linear transformation matrix \mathbf{A} that maps the original feature vectors x on to \hat{x} .

$$\hat{x} = \mathbf{W}^T x \quad (1)$$

where x is an n -dimensional vector, \hat{x} is an m -dimensional vector, $m \leq n$, and \mathbf{W} is an $n \times m$ matrix with m linearly independent columns. Let d_v and d_c represent the training vectors containing VI and CI, respectively, in the original feature space. The covariance matrices for these training feature

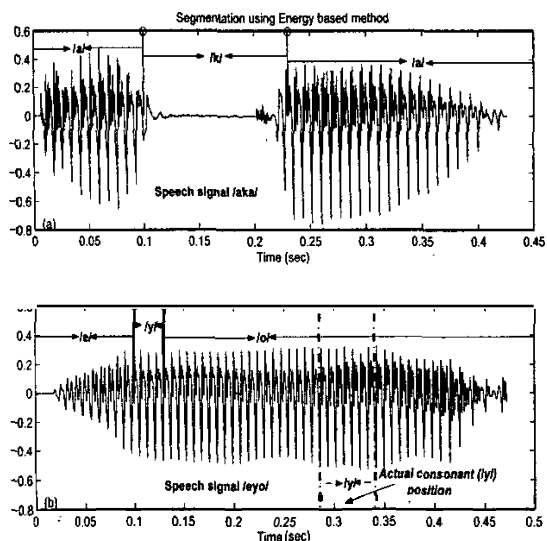


Figure 1. Basic-unit segmentation using energy based method. (a) Speech signal /aka/. (b) Co-articulated signal feyol (continuous vertical line: segmentation using energy based method).

vectors can be written as

$$\begin{aligned} C_v &= E\{(d_v - \bar{d}_v)(d_v - \bar{d}_v)^T\} \\ C_c &= E\{(d_c - \bar{d}_c)(d_c - \bar{d}_c)^T\} \end{aligned} \quad (2)$$

where \bar{d}_v and \bar{d}_c represent the means of d_v and d_c respectively. We collect an ensemble of feature vectors of length N corresponding to different vowels to estimate the $N \times N$ vowel covariance matrix C_v . Similarly we estimate the consonant covariance matrix C_c . We wish to find a W that maximizes the ratio of the variance of VI to that of CI after the transformation. If the density functions of d_v and d_c are assumed to be normally distributed, then their covariance matrices after transformation are given by

$$\begin{aligned} \widehat{C}_v &= W^T C_v W \\ \widehat{C}_c &= W^T C_c W \end{aligned} \quad (3)$$

A simple measure of the variance or the 'scatter' is the determinant of the covariance matrix [4]. Thus, the criterion function to be maximized is given by

$$J(W) = \frac{|\widehat{C}_v|}{|\widehat{C}_c|} = \frac{|W^T C_v W|}{|W^T C_c W|} \quad (4)$$

The columns of the optimum W are obtained as generalized eigenvectors for vowels (GEVV), corresponding to the largest eigenvalues in

$$C_v w_i^{(v)} = \lambda_i C_c w_i^{(v)} \quad (5)$$

Similarly, we obtain generalized eigenvectors for consonants (GEVC) as

$$C_c w_i^{(c)} = \lambda_i C_v w_i^{(c)} \quad (6)$$

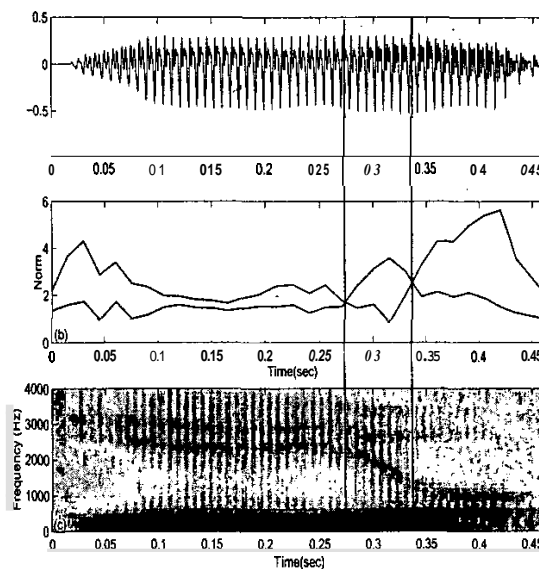


Figure 2. (a) Speech signal /eyo/. (b) Its segmentation into vowel (/e/ and /o/) and consonant (/y/) regions, using both the vowel and consonant norm-contours. (c) Spectrogram of /eyo/.

Evaluation of Norm-Contours

The GEVV and GEVCs are obtained by solving equations 5 and 6. The test signal is divided into overlapping frames and the feature vector x_k corresponding to the k^{th} frame is obtained using LP-Cepstral coefficients (LPCC) or Mel-Cepstral coefficients (MCC). We evaluate the norm-contours as follows.

$$N_v(k) = \sum_{i=1}^M (w_i^{(v)})^T x_k \quad \& \quad N_c(k) = \sum_{i=1}^M (w_i^{(c)})^T x_k \quad (7)$$

N_v and N_c give the norm-contours from V and C subspaces. Norm of the projection of the feature vectors derived from the test basic-unit on GEVV and GEVC give the norm-contours. One of them represents the vowel information and the other, the consonant information. The resulting norm contours obtained for a test signal cross each other at the beginning and at the end of consonant region of a given test basic-unit. The segmentation points are the ones where $N_v(k) \approx N_c(k)$. We found that optimum results were obtained when $M = 3$.

3. STATISTICAL TESTING FOR SEGMENTATION

Here, we view the basic-unit segmentation as classifying frames of a given test basic-unit into one of two classes: vowels and consonants. This indeed is a hypothesis testing (two-class) problem and it requires probability models for vowel & consonant and a robust threshold as well. Tipping and Bishop [2] proposed PPCA, to emphasize the advantages of associating a probability model with principal component analysis (PCA), rather than considering the algorithmic perspective of

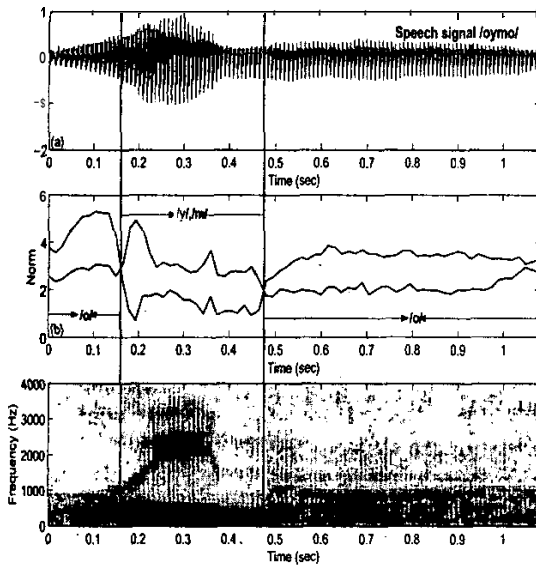


Figure 3. (a) Speech signal /oymo/. (b) Its segmentation into vowel (/o/) and consonant (/y/ and /m/) regions, using both the vowel and consonant norm-contours. (c) Spectrogram of /oymo/.

determining principal components. Using PPCA approach [2], we generate probability model for V and C.

Motivation for Hypothesis based Segmentation

In the section 2, we saw the performance of *oriented principal component analysis* based basic-unit segmentation. Although the performance of OPCA is good, the crossovers (threshold) are not robust. So, segmentation becomes difficult when there are multiple crossovers. Figure 4(a) shows a basic-unit /aumai/, and Fig. 4(b) shows the norm-contours obtained after projecting the basic-unit in Fig. 4(a) on V and C subspaces. One can see from Fig. 4(b) that there are multiple crossovers, making it difficult to choose an appropriate threshold for segmentation. Here, classical decision theory [5] can be used to find the robust threshold. Along with the crossover information, decision regarding vowel or consonant class is used to obtain correct segmentation points by eliminating the false crossovers. The class decision is shown in Fig. 4(d).

Probability Models using PPCA

Unlike in PCA, PPCA gives a probability density for the data through a latent variable model. A latent variable model relates a d -dimensional observed data vector \mathbf{t} to a q -dimensional ($q < d$) latent vector \mathbf{x} by defining a noise model and a prior on the distribution on the latent variables in the form,

$$\mathbf{t} = \mathbf{U}\mathbf{x} + \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad (8)$$

where $\boldsymbol{\epsilon}$ is an \mathbf{x} -independent noise process and \mathbf{U} is the $(d \times q)$ generative matrix. It is common that the prior distribution of the latent variables is a sample Gaussian distribution $\mathbf{x} \sim$

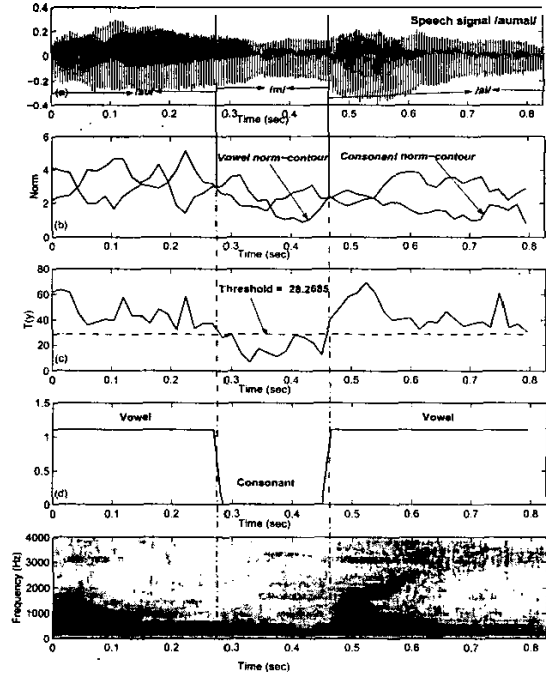


Figure 4. Segmentation of basic-unit using OPCA and PPCA. (a) Test basic-unit /aumai/. (b) Vowel and consonant-norm contours after projecting the test basic-unit into V & C subspaces. (c) Test statistics $T(y)$, along with the threshold for a given $P_{FA} \approx 10^{-4}$. (d) V & C decision after the class testing. (e) Spectrogram of /aumai/.

$N(\mathbf{0}, \mathbf{I})$ over the latent space. The noise model may also be Gaussian with $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \boldsymbol{\psi})$. The model is restricted to having a non-zero mean by the parameter $\boldsymbol{\mu}$. These aspects mean that the observed data vectors are normally distributed by $\mathbf{t} \sim N(\boldsymbol{\mu}, \mathbf{C})$, where the model covariance is given by

$$\mathbf{C} = \boldsymbol{\psi} + \mathbf{U}\mathbf{U}^T \quad (9)$$

If $\boldsymbol{\psi} = \sigma^2 \mathbf{I}$, the latent variable model is called PPCA [2]. In this terminology, conventional PCA is recovered when $\sigma^2 \rightarrow 0$. One of the ways of estimating the parameters of these latent variable models is by the Expectation-Maximization algorithm. Tipping and Bishop [6] have recently formulated PCA within a maximum-likelihood framework based on a specific form of Gaussian latent variable model. They also showed that with $\mathbf{C} = \sigma^2 \mathbf{I} + \mathbf{U}\mathbf{U}^T$, the only stationary points occur when

$$\mathbf{U} = \mathbf{W}_q (\boldsymbol{\Lambda}_q - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{R}. \quad (10)$$

Here \mathbf{W}_q has q eigenvectors of \mathbf{S} (sample covariance matrix, given by $\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{t}_n - \boldsymbol{\mu})(\mathbf{t}_n - \boldsymbol{\mu})^T$) as column vectors, $\boldsymbol{\Lambda}_q$ is a diagonal matrix with the corresponding eigenvalues, and \mathbf{R} is an arbitrary $q \times q$ orthogonal rotation matrix. Tipping and Bishop prove that when \mathbf{W}_q contains the *principal* eigenvectors of \mathbf{S} , global maximum of the likelihood occurs. So with \mathbf{U} in Equation 10, the latent variable model defines

a mapping from the latent space into the *principal subspace* of the observed data. The maximum likelihood estimator for noise variance σ^2 is given by the average variance lost for each discarded dimension, and can be formulated as

$$\sigma_{ML}^2 = \frac{1}{d-q} \sum_{j=q+1}^d \lambda_j. \quad (11)$$

To sum up, a Probabilistic PCA is obtained by finding the q principal eigenvectors and eigenvalues of the sample covariance matrix \mathbf{S} , while Gaussian distribution with sample mean $\boldsymbol{\mu}$ and covariance matrix $\sigma^2 \mathbf{I} + \mathbf{U}\mathbf{U}^T$, gives the density model.

Vowel-Consonant Segmentation as Hypothesis Testing

To bring the segmentation problem into a Hypothesis testing framework, we need to have a probability model for *vowels* and *consonants*. This can be accomplished by calculating \mathbf{U} , using equation 10 and plugging \mathbf{U} into the covariance model given by $\mathbf{C} = \sigma^2 \mathbf{I} + \mathbf{U}\mathbf{U}^T$. Let us call the probability model for *vowels* and *consonants* as $C_{,,}$ and $C_{,,}$, respectively. We have,

$$\begin{aligned} C_{pv} &= \sigma_v^2 \mathbf{I} + U_v U_v^T \\ C_{pc} &= \sigma_c^2 \mathbf{I} + U_c U_c^T \end{aligned} \quad (12)$$

We use GEVV ($W_q^{(v)}$), **GEVC** ($W_q^{(c)}$) and the corresponding eigenvalue matrices ($\Lambda_q^{(v)}$ & $\Lambda_q^{(c)}$) to obtain the parameters U_v, U_c, σ_v and σ_c using equations 10 and 11. These parameters are in turn used to generate probability model C_{pv} and C_{pc} using equation 12. Assuming the \mathbf{V} and \mathbf{C} features to be zero mean Gaussian random processes with covariance matrices C_{pv} and $C_{,,}$, the segmentation problem is to distinguish between the hypotheses

$$\mathbf{x} \sim \begin{cases} N(0, C_{pc}) & \text{under } H_0 \\ N(0, C_{pv}) & \text{under } H_1 \end{cases}$$

Here, \mathbf{x} is a zero mean *test basic-unit*. A NP detector decides H_1 if the likelihood ratio exceeds a threshold or if

$$Lh(\mathbf{x}) = \frac{p(\mathbf{x}; H_1)}{p(\mathbf{x}; H_0)} > \gamma$$

where,

$$p(\mathbf{x}; H_1) = \frac{1}{\sqrt{(2\pi)^d |C_{pv}|}} \exp\left[-\frac{1}{2} \mathbf{x}^T C_{pv}^{-1} \mathbf{x}\right],$$

$$p(\mathbf{x}; H_0) = \frac{1}{\sqrt{(2\pi)^d |C_{pc}|}} \exp\left[-\frac{1}{2} \mathbf{x}^T C_{pc}^{-1} \mathbf{x}\right],$$

$$Lh(\mathbf{x}) = \frac{\frac{1}{\sqrt{(2\pi)^d |C_{pv}|}} \exp\left[-\frac{1}{2} \mathbf{x}^T C_{pv}^{-1} \mathbf{x}\right]}{\frac{1}{\sqrt{(2\pi)^d |C_{pc}|}} \exp\left[-\frac{1}{2} \mathbf{x}^T C_{pc}^{-1} \mathbf{x}\right]}$$

the log-likelihood ratio (LLR) becomes

$$lh(\mathbf{x}) = \left\{ (\mathbf{x}^T C_{pc}^{-1} \mathbf{x}) - (\mathbf{x}^T C_{pv}^{-1} \mathbf{x}) \right\} > \ln \left[\gamma \frac{\sqrt{|C_{pv}|}}{\sqrt{|C_{pc}|}} \right].$$

Hence, we decide H_1 if

$$T(\mathbf{x}) = \mathbf{x}^T (C_{pc}^{-1} - C_{pv}^{-1}) \mathbf{x} > \gamma_{th} \ln \left[\gamma \frac{\sqrt{|C_{pv}|}}{\sqrt{|C_{pc}|}} \right] \quad (13)$$

Evaluation of the Threshold

By simplifying the test statistics $T(\mathbf{x})$ (eq. 13) into a standard distribution form, then the threshold (γ_{th}) can be evaluated [5]. Let $\mathbf{A} = C_{pc}^{-1}$ and $\mathbf{B} = C_{pv}^{-1}$, equation 13 becomes

$$T(\mathbf{x}) = \mathbf{x}^T (\mathbf{A} - \mathbf{B}) \mathbf{x} > \gamma_{th}. \quad (14)$$

Define $\mathbf{G} = \mathbf{A}^{-1} \mathbf{B}$. Let the eigenvalue decomposition of \mathbf{G} , be

$$\mathbf{G}\mathbf{V} = \mathbf{A}\mathbf{V},$$

such that

$$\begin{aligned} V^T \mathbf{A} \mathbf{V} &= \Lambda_1 \\ V^T \mathbf{B} \mathbf{V} &= \Lambda_2 \end{aligned} \quad (15)$$

where, Λ, Λ_1 and Λ_2 are diagonal matrices. Writing equation 15 in terms of \mathbf{A} and \mathbf{B} , we have

$$\begin{aligned} \mathbf{A} &= \mathbf{V} \Lambda_1 \mathbf{V}^T \\ \mathbf{B} &= \mathbf{V} \Lambda_2 \mathbf{V}^T \end{aligned} \quad (16)$$

Substituting eq. 16 into 14, we have

$$T(\mathbf{x}) = \mathbf{x}^T \mathbf{V} [\Lambda_1 - \Lambda_2] \mathbf{V}^T \mathbf{x}. \quad (17)$$

Let $\mathbf{y} = \mathbf{V}^T \mathbf{x}$. Then equation 17 becomes

$$T(\mathbf{y}) = \mathbf{y}^T \tilde{\Lambda} \mathbf{y} \quad (18)$$

where, $\tilde{\Lambda} = \Lambda_1 - \Lambda_2$. We know that \mathbf{x} is zero mean. Therefore \mathbf{y} is also zero mean, i.e., $E\{\mathbf{y}\} = 0$ and the covariance of \mathbf{y} is

$$E\{\mathbf{y}\mathbf{y}^T\} = \begin{cases} V^T E\{\mathbf{x}\mathbf{x}^T\} V \\ V^T C_{xx} V \end{cases}$$

If \mathbf{x} is a *vowel*, then $C_{xx} = C_{,,}$ and $V^T C_{pv} V = \Lambda_1^{-1}$. Similarly, If \mathbf{x} is a *consonant*, then $C_{xx} = C_{,,}$ and $V^T C_{pc} V = \Lambda_2^{-1}$. So, we can write equation 18 as

$$T(\mathbf{y}) = \sum_{n=1}^N \frac{y^2(n)}{\eta^2(n)} \quad (19)$$

where $\left\{ \frac{1}{\eta^2(1)}, \frac{1}{\eta^2(2)}, \dots, \frac{1}{\eta^2(n)} \right\}$ are diagonal elements of $\tilde{\Lambda}$. Because \mathbf{y} is i.i.d., $T(\mathbf{y})$ is χ^2 (Chi-square) distributed and the detector turns out to be an *energy detector*. The required threshold, γ_{th} , for segmentation is computed by using an optimization algorithm for the fixed *probability of false alarm*, P_{FA} [5].

$$P_{FA} = \int_{\gamma_{th}}^{\infty} \int_{-\infty}^{\infty} \prod_{n=1}^N \frac{1}{\sqrt{1-2j\alpha_n\omega}} \exp(-j\omega t) \frac{d\omega}{2\pi} dt \quad (20)$$

$$P_D = \int_{\gamma_{th}}^{\infty} \int_{-\infty}^{\infty} \prod_{n=1}^N \exp(-j\omega t) \frac{d\omega}{2\pi} dt. \quad (21)$$

where,

$$\alpha = \frac{1}{\eta^2(n)}.$$

Now, the decision (segmentation) rule is

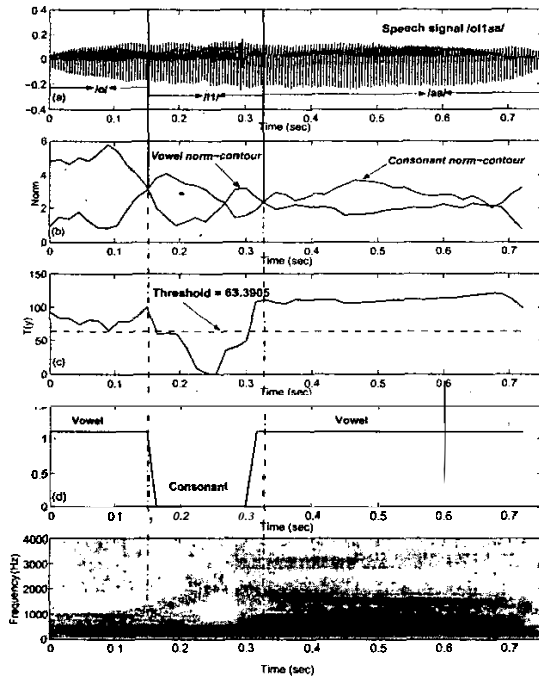


Figure 5. Segmentation of basic-unit using OPCA and PPCA. (a) Test basic-unit /o11aa/. (b) Vowel and consonant norm contours after the projection of test basic-unit into V & C subspaces. (c) Test statistics $T(y)$, along with the threshold for a given $P_{FA} = 10^{-4}$. (d) V & C decision after the class testing. (e) Spectrogram of /o11aa/.

Decide **Vowel** (H_1) if

$$T(y) = \sum_{n=1}^N \frac{y^2(n)}{\eta^2(n)} > \gamma_{ll}$$

otherwise decide **Consonant** (H_0).

4. RESULTS AND DISCUSSION

Speech segmentation experiments were conducted on a Kannada speech database spoken by a female volunteer. GEVV's and GEVC's were obtained from a Tamil speech database spoken by a male volunteer using the method discussed in section 2. Feature vectors were obtained for each frame of a test basic-unit. Duration of each frame of speech was 30 ms, with an overlap of 20 ms between successive frames. Each frame of speech was Hamming windowed and processed to yield a 13-dimensional feature vector. The feature vectors used were MCC and LPCC. We have seen in Fig. 1(b) that energy based segmentation fails to identify the co-articulated consonant region. On the other hand, in Fig. 2, the same consonant region has been correctly identified using subspace-based method. Figure 2(c) also displays the spectrogram of the basic-unit /eyo/. Here, spectrogram shows clear second formant transition from the frontal vowel /e/ to the back vowel /o/ and the transition region corresponds to

the consonant /y/. Basic-unit considered in Fig. 3 is VCCV and both C's are glide and a nasal (/yl and /ml/).

Probability models for V and C were obtained using equations 10, 11 and 12. Here, the observed feature vector \mathbf{t} , and its dimension is 13, the latent variable \mathbf{x} has dimension of 3, i.e., $d = 13$ and $q = 3$. Dimension of the latent variable \mathbf{x} has been chosen based on the dimensions of GEVV and GEVC. The multiple crossovers present some uncertainty in classifying the basic-units into vowel and consonant parts. By employing classical decision theory, we can remove this uncertainty by choosing a threshold on a statistical basis. In plot 4(d), the basic-unit is seen to be clearly segmented. In Fig. 4(b), there are multiple crossovers and using the class decision approach shown in Fig. 4(d), we can eliminate false crossovers. We used fixed $P_{FA} = 10^{-4}$ to calculate the threshold, γ_{ll} . Figure 4(e) shows the spectrogram of /aumai/. Figure 5(a) shows the co-articulated basic-unit /o11aa/. In Fig. 5(b) we can see the multiple crossovers of the norm-contours. Figure 5(c) shows the test statistics along with the threshold and V & C decision is shown in Fig. 5(d). Figure 5(e) shows the spectrogram of /o11aa/.

5. CONCLUSION

We have presented subspace and hypothesis based methods for co-articulated basic-unit segmentation. The first method uses crossovers of norm-contours for segmentation. Norm-contours give a measure of energy in the projected subspaces. There may be an ambiguity when there are multiple crossovers. In the second method this ambiguity is resolved by finding a statistical threshold. The test statistic using NP criterion turns out to be an energy detector as in the first method but without any ambiguity in segmentation. PPCA method can be further extended to classify within the vowel or consonant segments in a basic-unit using mixture-PPCA [6].

REFERENCES

- [1] G. L. Jayavardhana Rama, A. G. Ramakrishnan, R. Muralishankar and P. Prathibha, "A complete text-to-speech synthesis system in Tamil," *IEEE workshop on Speech Synthesis*, 2002.
- [2] M. E. Tipping and C. M. Bishop, "Probabilistic principal component analysers," *Journal of the Royal Statistical Society*, vol. 61, no. 3, pp. 611-622, 1999.
- [3] Jan P. van Hemert, "Automatic segmentation of speech," *IEEE Trans. Signal Proc.*, vol. 39, no. 4, pp. 1008-1012, 1991.
- [4] R. Duda and P. Hart *Pattern Classification and Scene Analysis*, New York: Wiley, 1973.
- [5] Steven M. Kay, *Fundamentals of Statistical Signal Processing - Detection Theory*, Prentice-Hall, Inc. A Simon & Schuster, New Jersey 07458, 1998.
- [6] M.E. Tipping and C. M. Bishop, "Mixture of probabilistic principal component analysers," *Neural Computation*, vol. 11, no. 2, pp. 443-482, 1999.