

ON THE APPLICATION OF AR MODEL IN SEGMENTING ISOLATED-WORD SPEECH SIGNALS

Xi Xiao, D. Nandagopal & D.A.H. Johnson
Department of Electrical & Electronic Engineering
University of Melbourne
Parkville, 3052
Victoria
Australia

ABSTRACT

In this paper, we present a segmentation method which uses AR modelling of the spectrum of the fullwave rectified speech signal. The FFT of the model coefficients yields a smoothed time domain signal with well defined minima, which locate the segment boundaries. The method is robust in the presence of noise. It is a useful first step in speech processing to segment speech signals into sub-frames which can be treated as time-invariant (stationary) processes.

THE SEGMENTATION ALGORITHM

Let $s(n)$ be the sequence of speech samples representing a word and let $P(n)$ be a modified version of $s(n)$ where

$$P(n) = \text{abs}[s(n)]$$

$P(n)$ can be regarded as the real power spectrum of some 'pseudo' time domain signal. The inverse Fourier transform of $P(n)$ yields a complex 'pseudo' time domain signal, $p(n)$. $p(n)$ is then AR modelled using LPC of order say 26, and the spectrum of the model, $\hat{P}(n)$, found. This 'pseudo' frequency domain corresponds to the real time domain and can therefore be related directly to the original speech signal $s(n)$. In particular the real part of $\hat{P}(n)$ is a smoothed version of the signal and has well defined minima which can be used as the segment boundaries.

FORMULATION

Let $P(n)$ be the finite length fullwave rectified signal ($0 < n < N - 1$) and $x(k)$ be the 'pseudo' time domain signal such that

$$x(z) = \sum_{k=0}^{w-1} x(k) z^{-k} \quad (1)$$

and

$$\begin{aligned} P(n) &= x(z) \Big|_z = e^{j \frac{2\pi n}{N}} = x \left(e^{j \frac{2\pi n}{N}} \right) \\ &= \sum_{k=0}^{N-1} x(k) e^{-jnk2\pi/N} \\ & \quad (0 \leq n, k \leq N - 1) \end{aligned} \quad (2)$$

It is obvious that $P(n)$ is the Discrete Fourier Transform of $x(k)$ by equally evaluating $x(z)$ on unit circle at Z -plane.

Suppose that $x(k)$ can be modelled by Linear Predictive Coding of order P , then the estimate for $x(k)$ is given as

$$\begin{aligned} \hat{x}(k) &= - \sum_{l=1}^P a(l) x(k-l) \\ & \quad (p < k \leq N-1) \end{aligned} \quad (3)$$

where $a(1) \dots a(p)$ are LPC constants to be determined.

The estimate error is thus

$$\begin{aligned} E(k) &= x(k) - \hat{x}(k) = x(k) + \sum_{l=1}^P a(l) x(k-l) \\ & \quad (p < k \leq N-1) \end{aligned} \quad (4)$$

The constants $a(1), a(2) \dots a(p)$ are determined by minimizing the total squared estimate error

$$E_p = \min \left\{ \sum_k (x(k) - \hat{x}(k))^2 \right\} \quad (5)$$

$a(1) \dots a(p)$

Applying Z-transform to equation (4), we obtain an estimate for $X(Z)$,

$$\hat{X}(z) = \frac{E_p}{1 + \sum_{\ell=1}^p a(\ell) z^{-\ell}} \quad (6)$$

Now, let $Z = e^{jnk2\pi/N}, n = 0, 1, \dots, N-1$,

We have

$$\hat{X}(e^{jn \cdot 2\pi/N}) = \frac{E_p}{1 + \sum_{\ell=1}^p a(\ell) e^{-j\ell \cdot n \cdot 2\pi/N}} \quad (7)$$

compare equation (2) and (7), an estimate for $P(n)$ is given as

$$\hat{P}(n) = \hat{X}(e^{jn \cdot 2\pi/N}) = \frac{E_p}{1 + \sum_{\ell=1}^p a(\ell) e^{-j\ell \cdot n \cdot 2\pi/N}} \quad (8)$$

From equation (6) and (8), we can now describe the problem as below:

Given a finite length of signal $P(n)$ with N sample points, find a polynomial of Z with an order P much lower than the sample points such that the reciprocals of the values of the polynomial on unit circle on Z -plane are the best estimate of the whole span of $P(n)$ under the criteria that the spectrum of $P(n)$ can be modelled by LPC constants $a(1), a(2) \dots a(p)$ as defined in equation (3).

A relatively high LPC order is used in the above procedure to ensure that the estimate error is small. As will be shown in next section, this algorithm is robust in detecting the onsets of both speech signal and the voiced parts, but oversegmentation is brought up because of high LPC order.

The next step is to measure the characteristics distance between every two adjacent segments. If the distance is smaller than an empirical threshold, these two segments are merged into one.

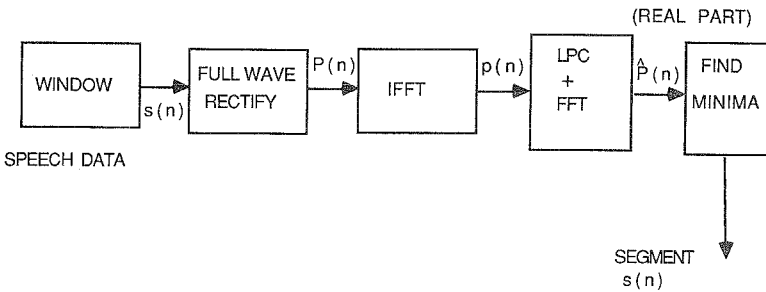
Now let's consider two adjacent segments R and T, and denote the signal within these two segments as X_R of lengths N_R and X_T of lengths N_T and, respectively. The statistic

$$\chi^2 = (\hat{\alpha}_R - \hat{\alpha}_T) \left\{ \sigma_R^{-2} (N_R' D_R)^{-1} + \sigma_T^{-2} (N_T' D_T)^{-1} \right\}^{-1} (\hat{\alpha}_R - \hat{\alpha}_T)$$

has an (asymptotic) χ^2 distribution under the null hypothesis that the two sequences emanate from series with the same LPC vector, α .

(Refer to Peter De Souza and Peter J. Thomson for symbol definition and the proof) [1].

When χ^2 smaller than a threshold, the null hypothesis is accepted, or, the two segments are merged.



RESULTS

Figure 1 shows a speech data segment (Windowed-Hamming), with white noise added resulting in a signal to noise ratio of 12 db. Figure 2 is the 'pseudo' time domain signal obtained by calculating the IFFT of the fullwave rectified signal in Figure 1. An LPC analysis of order 26 on this 'pseudo' time domain signal gives its 'pseudo' LPC spectrum estimate. In particular, the real part of this spectrum, depicted in Figure 3, is a smoothed version of the signal $s(n)$. The vertical lines superimposed on Figure 3 indicate the minima between peaks. These peaks are LPC estimates of the dominant peaks in Figure 1 and thus, the minima correspond approximately to the zero crossings and are used as segment boundaries. The final segmentation is shown in Figure 4. Table 1 lists the segment boundaries extracted from the minima positions.

As shown in Figure 4, the algorithm detects the onset of both the plosives and the voiced portions quite well at sample 213 and 523 respectively (the segment from sample 0 to 213 corresponds to intentionally introduced white noise). Similar segments from sample 785 to 1631 result from the periodical structures of the voiced portions.

Figure 5 is a clear speech signal of word/TWO/ sampled at 10 KHZ. The number of sample points is 4096 including some silence on the front and end of the speech. An LPC of order 26 was used to estimate the contour shown in Figure 6. 16 valleys are found in the contour and are used as segment boundaries represented by the vertical lines in Figure 5. Table 2 lists the LPC distance between every two adjacent segments. The distances measured at boundary 1, 3 and 15 are significantly larger than others and thus, these boundaries are finally chosen.

CONCLUSION

The AR model used in this segmentation method provides stable estimates of the segment boundaries in the presence of noise. The boundaries occur at, or very near the zero crossings of the signal. The number of segments is primarily determined by the order chosen for the LPC process and the oversegmentation can be reduced by LPC distance measures.

REFERENCES

- [1] Peter De Souza and Peter J. Thomson, "LPC Distance Measures and Statistical Tests with Particular Reference to the Likelihood Ratio", IEEE Trans. Acoust., Speech, Signal Processing, Vol. ASSP - 30, pp 304-315, April 1982.

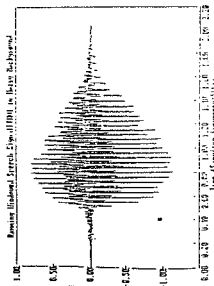


FIG. 1

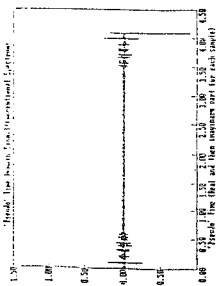


FIG. 2

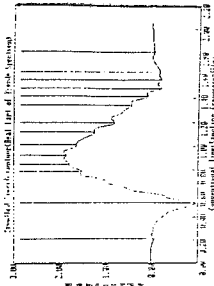


FIG. 3

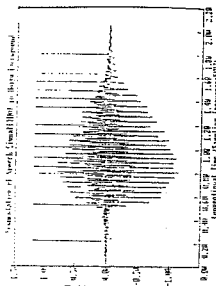


FIG. 4

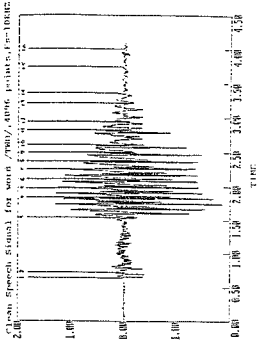


FIG. 5

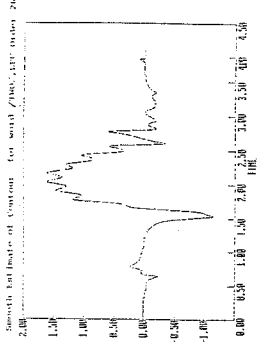


FIG. 6

Table 1. Segment Boundary List

Boundary	Distance	Boundary	Distance
1	10.26	1	10.26
2	18.97	1	18.97
3	28.16	1	28.16
4	34.53	1	34.53
5	45.88	1	45.88
6	54.10	1	54.10
7	62.69	1	62.69
8	73.43	1	73.43

Table 2. Distance Between Adjacent Elements

Boundary	Distance	Boundary	Distance
1	10.26	1	10.26
2	18.97	1	18.97
3	28.16	1	28.16
4	34.53	1	34.53
5	45.88	1	45.88
6	54.10	1	54.10
7	62.69	1	62.69
8	73.43	1	73.43