

# A NONSTATIONARY AUTOREGRESSIVE HMM WITH GAIN ADAPTATION FOR SPEECH RECOGNITION

Ki Yong Lee and \*Jooahn Lee

School of Electronics Engr., Soongsil University, 1-1, Sangdo 5Dong, Dongjak-Ku, Seoul 156-743, Korea

Tel: +82-2-820-0908 Fax: +82-2-814-3627 E-mail: [kylee@saint.soongsil.ac.kr](mailto:kylee@saint.soongsil.ac.kr)

\*Dept of Information and Telecommunication, Dong-Ah Broadcasting College, Ansung, Korea

## ABSTRACT

In this paper, a time domain approach for speech recognition is developed. The nonstationary autoregressive (AR) hidden markov model (HMM) with gain contour is proposed for modeling the statistical characteristics of the speech signal. The parameter of nonstationary AR model was modeled by the polynomial function with linear combination of  $M$  known basis functions. In this proposed model, speech signal is blocked by samples into fixed-length frames and modeled by nonstationary AR model controlled by markov switching sequences at each frame. Given the HMM parameter set of the speech, the gain-adapted recognition algorithm is developed for speech recognition.

## 1. Introduction

The autoregressive hidden markov model (ARHMM) [1,2] is useful methods to represent the statistical characteristics of the clean speech in speech recognition and speech enhancement. In the conventional ARHMM, individual states are assumed to be stationary stochastic sequences. This stationary-state assumption appears to be reasonable when a state is intended to represent piece-wise stationary segment of speech. Since speech sounds, such as fricative, glides, liquids, and transition regions between phones, reveal the most notable nonstationary nature [3,4], we can not expect to obtain the better performance by the conventional methods based on the above assumption. Another basic issue arising from ARHMM for speech recognition is matching problem of the energy contour of the signal to the energy contour of the model for that signal. The energy matching is usually attained by appropriate normalization but is not applicable when only noisy speech signals are available [5]. In this letter, for energy matching we propose a training algorithm of nonstationary ARHMM with gain adaptation and gain adapted speech recognition algorithm.

To overcome these problems, the nonstationary ARHMM with gain adaptation for gain normalized clean signals are design using maximum likelihood (ML) estimates of the gain contours of the clean training sequences. The parameter of nonstationary AR model was modelled by the polynomial function with linear combination of  $M$  known basis functions. Then, the speech signal is blocked by samples into fixed-length frames and modeled by nonstationary AR model with frame varying polynomial function controlled by markov switching sequences at each frame. Our model is formally very similar to the trend HMM [3,4], but it is designed to handle the speech signal at the frame level, where it is represented by the signal, rather than dealing with feature vectors directly. Also, for  $M=0$ , the

proposed model become to conventional ARHMM [5]. The proposed method for speech recognition is combined with ML estimates of the gain contours of the clean test signals, obtained from the given clean signals, in performing recognition using the maximum *a posteriori* decision rule.

We have evaluated our nonstationary ARHMM method with  $M=1$  on a base of ten isolated Korean digits with three versions of each digit pronounced by seven male speakers. In speech recognition, the proposed method was compared with a conventional ARHMM method with gain-normalized approach and with gain adaptation.

## 2. Nonstationary AR-HMM with gain adatation for clean speech

Let  $y = y_n, n=1, \dots, T$  be the sequence of clean signal vectors, where  $y_n = \{y(t), (n-1)N+1 \leq t \leq nN\}$  and  $s_n \in \{1, \dots, L\}$ , be a sequence of states corresponding to  $y$ . Let  $g = \{g_n, n=1, \dots, T\}$ , be a sequence of gain factors, or a gain contour, for the signal  $y$ .

Then, at  $n$ -th frame speech signal conditioned on state  $i$  is expressed as a linear combination of its past values plus an excitation source with gain contour, as

$$y(t) = \sum_{k=1}^p \sum_{m=0}^M B_k^i(m) y(t-k) + g_n \cdot e_i(t), \quad n-1N+1 \leq t \leq nN. \quad (1)$$

where  $B_k^i(m)$  is the state-dependent time-varying coefficients.

$y(t-1) = y(t-1) \quad y(t-2) \quad \vdots \quad y(t-p)^T$ ,  $e_i \cdot$  is the excitation source with state-dependent variance  $\sigma_i^2$ ,  $N$  and  $n$  is the frame length and number, respectively, and  $g_n > 0$  for all  $n$  is a gain term to take into account the mismatch between training data and testing data for the clean speech models.

We now turn our attention to the problem of estimating the time-varying coefficients in our model. In order to gain insight into the behavior of the coefficients and to make the estimation problem tractable. We choose to model them as a linear combination of  $M$  known basis functions:

$$B_k^i(n) = \sum_{m=0}^M B_{k,m}^i f_m(n) \quad (2)$$

where  $f_m(n)$  represents the  $m$ th basis function and  $B_{k,m}^i$  the weight associated with the basis function.

Here, we choose  $M=1$  and our basis functions to be such that

$$\begin{aligned} f_0 \quad n = 1, \quad & 1 \leq n \leq T, \\ f_1 \quad n = n, \quad & 1 \leq n \leq T. \end{aligned} \quad (3)$$

Therefore, (1) may be rewritten by vector form as

$$y(t) = \mathbf{B}^i y(t-1) + g_n \cdot e_i(t), \quad n-1 \leq t \leq nN \quad (4)$$

where  $\mathbf{B}^i = B_{1,1}^i \ B_{1,2}^i \ B_{2,1}^i \ B_{2,2}^i \ \dots \ B_{P,1}^i \ B_{P,2}^i$  and

$$\begin{aligned} y(t-1) = & \left[ y((n-1)N+t-1), ny((n-1)N+t-1)L, \right. \\ & \left. y((n-1)N+t-p), ny((n-1)N+t-p) \right]^T. \end{aligned}$$

In essence, the model with time-varying coefficients has been transformed into one with time-invariant weights. The problem is now reduced to one of estimating  $2P$  time-invariant parameters that completely characterize the behavior of the coefficients. It should be noted that the choice of basis functions is by no means limited to polynomials.

The likelihood of the observation sequence  $y$  under the model  $\lambda$  and gain contour  $g$  is calculated as

$$\begin{aligned} p_\lambda(y|g) &= \sum_s p_\lambda(s, y|g) \\ &= \sum_s p_\lambda(s|g) p_\lambda(y|s, g), \end{aligned} \quad (5)$$

where  $p_\lambda(s|g)$  denotes the probability of sequence of states  $s$ , and  $p_\lambda(y|s, g)$  is the probability density function (pdf) of the sequence of out vectors  $y$  given  $g$  and  $s$ . For first-order HMM's,  $p_\lambda(s|g)$  is given by

$$p_\lambda(s|g) = \prod_{n=1}^T a_{s_n \rightarrow s_n}, \quad (6)$$

where  $a_{s_n \rightarrow s_n}$  denotes the transition probability from state at time  $t-1$  to state at time  $t$ .

The  $p_\lambda(y|s, g)$  is

$$p_\lambda(y|s, g) = \prod_{n=1}^T p_\lambda(y_n|s_n, g_n) \quad (7)$$

where the pdf  $p_\lambda(y_n|s_n, g_n)$  of the vector  $y_n$  given that this vector was generated from state  $s_n$  and gain contour  $g_n$ .

Then, from (4)  $p_\lambda(y_n|s_n, g_n)$  is given by

$$p_\lambda(y_n|s_n, g_n) = \prod_{t=(n-1)N}^{nN} \frac{\exp\left\{-\frac{(y(t) - \mathbf{B}^{s_n} y(t-1))^2}{2g_n^2 \sigma_{s_n}^2}\right\}}{\sqrt{2\pi g_n^2 \sigma_{s_n}^2}}. \quad (8.a)$$

or in matrix form

$$p_\lambda(y_n|s_n, g_n) = \frac{\exp\left(-\frac{1}{2} y_n^T C_{s_n}^{-1} y_n\right)}{(2\pi)^{N/2} \det^{1/2}(C_{s_n})}. \quad (8.b)$$

where  $C_{s_n} = g_n^2 \sigma_{s_n}^2 (A_{s_n}^T A_{s_n})^{-1}$  and  $A_{s_n}$  is a  $N \times N$  lower triangular Toeplitz matrix in which the first elements  $2P+1$  of the first column constitute the coefficients of the AR process. The parameter set  $\lambda = \{a_{ij}, \mathbf{B}^j, \sigma_j^2, i, j = 1, \dots, L\}$  of the nonstationary ARHMM and gain contour  $g$  for the clean speech is estimated from training sequences of clean speech signals. Note that  $\lambda$  denotes the parameter set of the ARHMM for the gain-normalized signal.

### 3. Gain adapted training algorithm

Gain-adapted training of the nonstationary ARHMM for the word results from ML estimation of the parameter set  $\lambda$  from a training sequence  $y$  from word using an ML estimate of the gain contour  $g$ . Then,  $\lambda$  can be estimated from

$$\max_{\lambda} \max_g p_\lambda(y|g) \quad (9)$$

However, the gradient equations of  $p_\lambda(y|g)$  with respect to  $\{\lambda, g\}$  are nonlinear and therefore have no simple solution. Hence, the estimation of  $\{\lambda, g\}$  is performed here iteratively using the expectation-maximization (EM) approach [6,7]. Then, each iteration constitutes one EM iteration for estimating a value of  $\lambda$  given  $g$  and one EM iteration for estimating  $g$  using the resulting  $\lambda$ .

The training algorithm for nonstationary ARHMM with gain contour can be summarized as follows.

#### 3.1 Estimation of $\lambda$ ;

Assuming that the gain contour  $g$  for all  $n$  is known, the estimation formulas for  $\lambda$  can be derived from maximizing Baum's auxiliary function [3,4]. Each iteration of the Baum algorithm starts with an old set of parameters, say  $\lambda_{l-1}$ , and estimates a new set of parameters, say  $\lambda_l$ , by maximizing the following auxiliary function:

$$Q_l(\lambda, g) = \sum_s p_{\lambda_{l-1}}(s|y, g) \log p_{\lambda_l}(y, s|g). \quad (10)$$

The objective function can be simplified to

$$\begin{aligned} Q_l(\lambda, g) &= \sum_{i=1}^L \sum_{n=1}^T p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j|y, g), \\ &\quad \times [\log a_{ij} + \log p_{\lambda_l}(y_n|s_n, g_n)] \end{aligned} \quad (11)$$

where

$$\begin{aligned} \log p_{\lambda_l}(y_n|s_n, g_n) &= -\frac{N}{2} \log(2\pi g_n^2 \sigma_j^2) \\ &\quad - \sum_{t=(n-1)N+1}^{nN} \frac{(y(t) - \mathbf{B}^j y(t-1))^2}{g_n^2 \sigma_j^2} \end{aligned}$$

and  $p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j|y, g)$  is the *posteriori probability* of

the transition from state  $i$  to state  $j$  given the signal  $y$  and the gain contour  $g_n$  at the  $(l-1)$ th iteration.

The auxiliary function  $Q_1(\lambda, g)$  is optimized by differentiation with respect to each of  $a_{ij}, \mathbf{B}^j, \sigma_j^2, (i, j = 1, \dots, L)$ , respectively. As in standard HMM, the reestimation formulas of the Markov chain and nonstationary AR parameter are established by setting

$$\begin{aligned}\frac{\partial Q_1}{\partial a_{ij}} &= 0, \\ \frac{\partial Q_1}{\partial \mathbf{B}^j} &= 0, \\ \frac{\partial Q_1}{\partial \sigma_j^2} &= 0,\end{aligned}\quad (12)$$

for  $i, j = 1, \dots, L$ , subject to the constraint  $\sum_{j=1}^L a_{ij} = 1$ . We obtain the reestimation formula

$$a_{ij} = \frac{\sum_{n=1}^T p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j | y, g)}{\sum_{j=1}^L \sum_{n=1}^T p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j | y, g)} \quad (13)$$

where the probability  $p_{\lambda_l}(s_{n-1} = i, s_n = j | y, g)$  can be calculated efficiently by the forward-backward algorithm. Also, the reestimation formulas for nonstationary AR parameter are obtained by solving

$$\begin{aligned}B^j &= \left[ \sum_{i=1}^L \sum_{n=1}^T p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j | y, g) \sum_{t=(n-1)N+1}^{nN} \hat{y}(t-1) \right. \\ &\quad \left. \hat{y}^T(t-1) \right]^{-1} \left[ \sum_{i=1}^L \sum_{n=1}^T p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j | y, g) \right. \\ &\quad \left. \times \sum_{t=(n-1)N+1}^{nN} \hat{y}(t) \hat{y}(t-1) \right] \quad (14)\end{aligned}$$

and

$$\begin{aligned}\sigma_j^2 &= \frac{1}{\sum_{i=1}^L \sum_{n=1}^T p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j | y, g)} \\ &\quad \times \left[ \sum_{i=1}^L \sum_{n=1}^T p_{\lambda_{l-1}}(s_{n-1} = i, s_n = j | y, g) \right. \\ &\quad \left. \times \sum_{t=(n-1)N+1}^{nN} (\hat{y}(t) - B^j \hat{y}(t-1))^2 \right] \quad (15)\end{aligned}$$

where  $\hat{y}(t) = \frac{y(t)}{g_n}$  and  $\hat{y}(t-1) = \left[ \frac{y(t-1)}{g_n} \frac{y(t-1)}{g_n} \frac{y(t-2)}{g_n} \dots \frac{y(t-p)}{g_n} \right]^T$  are normalized observation sequence by gain contour.

If  $M$  is set to zero, then (14)-(15) reverts to the reestimation formula for the Gaussian mean vectors in the standard

ARHMM. This is expected since, as the time-varying component in the trend function is removed, the resulting degenerated nonstationary AR HMM is no different from the standard ARHMM.

### 3.2 Estimation of gain contour $g$ :

Next, for estimating the gain contour  $g$ , assume that  $\lambda$  is known. Let  $g_{l-1}$  and  $g_l$  be a current and a new estimate of the gain contour of the signal  $y$ , respectively. Then, the objective function for  $g_l$  is given by

$$\begin{aligned}Q_2(\lambda, g) &= \sum_s p_\lambda(s | y, g_{l-1}) \log p_\lambda(s | y, g_{l-1}) \\ &= \sum_{i,j=1}^L \sum_{n=1}^T p_\lambda(s_{n-1} = i, s_n = j | y, g_{l-1}) \\ &\quad \times [\log a_{ij} + \log p_\lambda(s | y, g_{l-1})].\end{aligned}\quad (16)$$

Similarly to what we have seen before, maximization of the auxiliary function over  $g$  results in an estimate of the gain contour  $g$  by setting the gradient of  $Q_2(\lambda, g)$  with respect to  $g$

$$\frac{\partial Q_2}{\partial g^2} = 0.$$

We arrive at the following gain reestimation formula

$$g_{l,n}^2 = \sum_{j=1}^L p_\lambda(s_n = j | y, g_{l-1}) \sum_{t=(n-1)N+1}^{nN} \frac{(y(t) - B^j y(t-1))^2}{\sigma_j^2}. \quad (17)$$

The iteration process start with initial gain contour  $g_{l=0,n} = 1$  for all  $n$ , and repeated until either a fixed point  $\lambda_{l+1} = \lambda_l, g_{l+1,n} = g_{l,n}$  is reached or the difference in likelihood in two consecutive iterations is sufficiently small.

## 4. Gain-adapted Speech recognition

In speech recognition, given the word speech sequence  $y$  the decision rule for spoken word  $y$  is

$$\max_{1 \leq l \leq Z} \max_g p_\lambda(y | W_l, g). \quad (18)$$

where  $Z$  is number of the total word for speech recognition.

The algorithm for local maximization of  $p_\lambda(y | g)$  over  $g$  can be summarized as follows;

**Step-0:** Initialization: For given the parameter

$$\lambda = \{a_{ij}, \mathbf{B}^j, \sigma_j^2, i, j = 1, \dots, L\}, g_0 = 1, \text{ and } \varepsilon,$$

evaluate  $p_\lambda(y | g_0)$  and  $l=0$ .

**Step-1:** Gain estimation: Calculate the posterior probabilities  $p_\lambda(s_n | y, g_l)$  for  $s_n = 1, \dots, L$  and  $n = 1, \dots, T$ , and estimate  $g_{l+1}$  using (17).

**Step-2:** If  $p_\lambda(y | g_{l+1}) - p_\lambda(y | g_l) \leq \varepsilon$ ,

assign  $\max_g p_\lambda(y|g) = p_\lambda(y|g_{l+1})$  and stop

Otherwise, set  $l \rightarrow l+1$  and goto **Step-1**.

## 5. Experimental Results

We have evaluated our new method on a base of ten isolated Korean digits with three versions of each digit pronounced by seven male speakers. Only 50 speech data of five speakers have participated in training and other 160 speech data have been used for test. This speech data were sampled at 10kHz and modeled by state  $L=5$  and AR order of 15. Training and recognition were performed on nonoverlapping vectors of the speech word whose dimension was  $N=256$ . The proposed method with gain-adapted recognition was compared with a conventional gain-normalized approach [2]. For the data (A) of speakers who participated in training and those (B) of the speakers who didn't participate in training, the recognition rates of the conventional training algorithm has scored 96 % and 90%, respectively while those of the proposed algorithm has scored 98.8% and 93.2 %, respectively. Table 1 shows a results of comparison of the recognition accuracies. Table 2 shows a results of comparison on conventional ARHMM with gain adaptation [5], and proposed method with gain adaptation (in this case, conventional method is equal to proposed method with  $M=0$ ) and  $M=1$ . From this result, the proposed method had a good performance.

ACKNOWLEDGEMENT: This work was supported in part by KOSEF (971-0917-105-1)

## References

- [1] A.B. Poritz, "Hidden markov models: A guided tour," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, Apr. 1988, pp. 7-13.
- [2] B. Juang and L.R. Rabiner, "Mixture autoregressive hidden Markov models for speech signals," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 33, pp. 1404-1413, Dec. 1986.
- [3] L. Deng, "A generalized hidden Markov model with state-conditioned trend functions of time for speech signal," *Signal Processing*, 27, pp. 65-72, 1992.
- [4] L.Deng, M.Aksmanovic, X. Sun, and C.F. JeffWu, "Speech recognition using HMM with polynomial regression functions as nonstationary states," *IEEE Trans. Speech and Audio Processing*, vol.2, no.4, pp.507-520, Oct. 1994.
- [5] Y. Ephraim, "Gain adapted hidden Markov models for recognition of clean and noisy speech," *IEEE Trans. Signal Processing*, vol. 40, no.6, pp.1303-1316, June 1992.
- [6] A.P. Dempster, N.M. Laird, and D.B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Royal Stat. Soc. B*, vol. 39, no.1, pp. 1-38, 1977
- [7] L.E. Baum, T. Petrie, G. Soules, and N. Weiss, "A maximization technique in the statistical analysis of probabilistic functions of Markov chains," *Ann. Math. Stat.*, vol. 41, pp.164-171, 1970.

Table 1. Recognition results

| Data | Gain-normalized | Gain-adapted |
|------|-----------------|--------------|
| A    | 96%             | 98.8%        |
| B    | 90%             | 93.2%        |

Table 2. Recognition results

| Data | ARHMM with $M=0$ | ARHMM with $M=1$ |
|------|------------------|------------------|
| A    | 98%              | 98.8%            |
| B    | 92.5%            | 93.2%            |