

Trajectory Modeling based on HMMs with the Explicit Relationship between Static and Dynamic Features

Keiichi Tokuda, Heiga Zen, Tadashi Kitamura

Department of Computer Science
Nagoya Institute of Technology, Japan
{tokuda, zen, kitamura}@ics.nitech.ac.jp

Abstract

This paper shows that the HMM whose state output vector includes static and dynamic feature parameters can be reformulated as a trajectory model by imposing the explicit relationship between the static and dynamic features. The derived model, named trajectory HMM, can alleviate the limitations of HMMs: i) constant statistics within an HMM state and ii) independence assumption of state output probabilities. We also derive a Viterbi-type training algorithm for the trajectory HMM. A preliminary speech recognition experiment based on N -best rescoring demonstrates that the training algorithm can improve the recognition performance significantly even though the trajectory HMM has the same parameterization as the standard HMM.

1. Introduction

The HMM has successfully been applied to modeling the sequence of speech spectra in speech recognition systems. Its tractable and efficient implementations are achieved by assumptions: i) constant statistics within an HMM state and ii) independence assumption of state output probabilities. To overcome these limitations in the standard HMM framework, alternative models have been proposed, e.g., [1]–[10]. Although these models can improve the speech recognition performance, they generally require an increase in model parameters and computational complexity. The use of the dynamic features (delta and delta-delta features) [11] also improves the performance of HMM-based speech recognizers. However, it has been thought of as an ad hoc rather than an essential solution.

This paper derives a trajectory model by reformulating the standard HMM whose state output vector includes static and dynamic feature parameters. The standard HMM allows inconsistent static and dynamic features. By imposing the explicit relationship between static features and dynamic features, the standard HMM is naturally translated into a trajectory model, referred to as ‘‘trajectory HMM’’ in this paper. We also derive a Viterbi-type training algorithm for the trajectory model, and evaluate it in a continuous speech recognition task.

The formulation of the trajectory HMM is closely related to a technique for parameter generation from HMM [12]–[14], in which the speech parameter sequence is determined so as to maximize its output probability for the HMM under the constraints between static and dynamic features. While we derived the speech parameter generation algorithm in order to construct HMM-based speech synthesizers [15] which can synthesize speech with various voice characteristics [16]–[18], the generation algorithm was also applied to speech recognition in [19]. The proposed training algorithm can be viewed as a train-

ing algorithm for the method in [19]. The model derived in [20] based on the maximum entropy framework is also closely related to the proposed trajectory model.

The rest of this paper is organized as follows. Section 2 defines the trajectory model. In Section 3 the training algorithm for the model is derived. Results of a continuous speech recognition experiment are shown in Section 4. Concluding remarks and future plans are presented in the final section.

2. Reformulating HMM as trajectory model

The output probability of a speech parameter vector sequence $\mathbf{o} = [\mathbf{o}_1^\top, \mathbf{o}_2^\top, \dots, \mathbf{o}_T^\top]^\top$ for the standard HMM is given by

$$P(\mathbf{o} | \lambda) = \sum_{\text{all } \mathbf{q}} P(\mathbf{o} | \mathbf{q}, \lambda) P(\mathbf{q} | \lambda), \quad (1)$$

where $\mathbf{q} = \{q_1, q_2, \dots, q_T\}$ is a state sequence. In most of speech recognition systems, the speech parameter vector \mathbf{o}_t is assumed to consist of the static feature vector $\mathbf{c}_t = [c_t(1), c_t(2), \dots, c_t(M)]^\top$ (e.g., cepstral coefficients), and dynamic feature vectors $\Delta \mathbf{c}_t, \Delta^2 \mathbf{c}_t$ (e.g., delta and delta-delta cepstral coefficients), that is $\mathbf{o}_t = [\mathbf{c}_t^\top, \Delta \mathbf{c}_t^\top, \Delta^2 \mathbf{c}_t^\top]^\top$. The dynamic features calculated by

$$\Delta \mathbf{c}_t = \sum_{\tau=-L_-^{(1)}}^{L_+^{(1)}} w^{(1)}(\tau) \mathbf{c}_{t+\tau} \quad (2)$$

$$\Delta^2 \mathbf{c}_t = \sum_{\tau=-L_-^{(2)}}^{L_+^{(2)}} w^{(2)}(\tau) \mathbf{c}_{t+\tau}. \quad (3)$$

correspond to the first and second time-derivative of the static feature \mathbf{c}_t , respectively. Conditions (2) and (3) can be arranged in a matrix form:

$$\mathbf{o} = \mathbf{W} \mathbf{c}, \quad (4)$$

where

$$\mathbf{c} = [\mathbf{c}_1^\top, \mathbf{c}_2^\top, \dots, \mathbf{c}_T^\top]^\top \quad (5)$$

$$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T]^\top \quad (6)$$

$$\mathbf{w}_t = [\mathbf{w}_t^{(0)}, \mathbf{w}_t^{(1)}, \mathbf{w}_t^{(2)}] \quad (7)$$

$$\mathbf{w}_t^{(n)} = \left[\mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}, \underset{\text{1st}}{w^{(n)}(-L_-^{(n)}) \mathbf{I}_{M \times M}}, \dots, w^{(n)}(0) \mathbf{I}_{M \times M}, \dots, w^{(n)}(L_+^{(n)}) \mathbf{I}_{M \times M}, \right. \\ \left. \underset{(t-L_-^{(n)})\text{-th}}{\phantom{w^{(n)}(-L_-^{(n)}) \mathbf{I}_{M \times M}}}, \dots, \underset{(t+L_+^{(n)})\text{-th}}{\phantom{w^{(n)}(L_+^{(n)}) \mathbf{I}_{M \times M}}} \right]$$

$$\mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}^{\top}, \quad n = 0, 1, 2 \quad (8)$$

$L_-^{(0)} = L_+^{(0)} = 0$, and $w^{(0)}(0) = 1$. The above model is improper in the sense of statistical modeling: it allows inconsistent static and dynamic feature vector sequences even though they are constrained by (4). The statistical model should be defined as a function of \mathbf{c} because the original observation is \mathbf{c} rather than the augmented variable \mathbf{o} .

When each state output probability distribution is assumed to be single Gaussian, $P(\mathbf{o} | \mathbf{q}, \lambda)$ is given by

$$P(\mathbf{o} | \mathbf{q}, \lambda) = \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t | \boldsymbol{\mu}_{q_t}, \mathbf{U}_{q_t}) = \mathcal{N}(\mathbf{o} | \boldsymbol{\mu}_{\mathbf{q}}, \mathbf{U}_{\mathbf{q}}), \quad (9)$$

where $\boldsymbol{\mu}_{q_t}$ and \mathbf{U}_{q_t} are the $3M \times 1$ mean vector and the $3M \times 3M$ covariance matrix, respectively, associated with q_t -th state, and

$$\boldsymbol{\mu}_{\mathbf{q}} = \left[\boldsymbol{\mu}_{q_1}^{\top}, \boldsymbol{\mu}_{q_2}^{\top}, \dots, \boldsymbol{\mu}_{q_T}^{\top} \right]^{\top} \quad (10)$$

$$\mathbf{U}_{\mathbf{q}} = \text{diag} [\mathbf{U}_{q_1}, \mathbf{U}_{q_2}, \dots, \mathbf{U}_{q_T}]. \quad (11)$$

By substituting (4) for (9), we can rewrite (9) as follows:

$$\begin{aligned} P(\mathbf{W}\mathbf{c} | \mathbf{q}, \lambda) &= \mathcal{N}(\mathbf{W}\mathbf{c} | \boldsymbol{\mu}_{\mathbf{q}}, \mathbf{U}_{\mathbf{q}}) \\ &= K_{\mathbf{q}} \cdot \mathcal{N}(\mathbf{c} | \bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}), \end{aligned} \quad (12)$$

where $\bar{\mathbf{c}}_{\mathbf{q}}$ is given by

$$\mathbf{R}_{\mathbf{q}} \bar{\mathbf{c}}_{\mathbf{q}} = \mathbf{r}_{\mathbf{q}} \quad (13)$$

and

$$\mathbf{R}_{\mathbf{q}} = \mathbf{W}^{\top} \mathbf{U}_{\mathbf{q}}^{-1} \mathbf{W} \quad (14)$$

$$\mathbf{r}_{\mathbf{q}} = \mathbf{W}^{\top} \mathbf{U}_{\mathbf{q}}^{-1} \boldsymbol{\mu}_{\mathbf{q}} \quad (15)$$

$$\mathbf{P}_{\mathbf{q}} = \mathbf{R}_{\mathbf{q}}^{-1} \quad (16)$$

$$\begin{aligned} K_{\mathbf{q}} &= \frac{\sqrt{(2\pi)^{MT} |\mathbf{P}_{\mathbf{q}}|}}{\sqrt{(2\pi)^{3MT} |\mathbf{U}_{\mathbf{q}}|}} \\ &\cdot \exp \left\{ -\frac{1}{2} \left(\boldsymbol{\mu}_{\mathbf{q}}^{\top} \mathbf{U}_{\mathbf{q}}^{-1} \boldsymbol{\mu}_{\mathbf{q}} - \mathbf{r}_{\mathbf{q}}^{\top} \mathbf{P}_{\mathbf{q}} \mathbf{r}_{\mathbf{q}} \right) \right\}. \end{aligned} \quad (17)$$

By omitting the normalization constant $K_{\mathbf{q}}$ in the above expression (12), a new statistical model can be defined:

$$P(\mathbf{c} | \lambda) = \sum_{\text{all } \mathbf{q}} P(\mathbf{c} | \mathbf{q}, \lambda) P(\mathbf{q} | \lambda), \quad (18)$$

where

$$P(\mathbf{c} | \mathbf{q}, \lambda) = \mathcal{N}(\mathbf{c} | \bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}). \quad (19)$$

It is interesting to note that the mean $\bar{\mathbf{c}}_{\mathbf{q}}$ is exactly the same as the speech parameter trajectory obtained by the speech parameter generation technique (Case 1 in [14], see Appendix also), that is,

$$\bar{\mathbf{c}}_{\mathbf{q}} = \arg \max_{\mathbf{c}} P(\mathbf{o} | \mathbf{q}, \lambda) = \arg \max_{\mathbf{c}} P(\mathbf{W}\mathbf{c} | \mathbf{q}, \lambda). \quad (20)$$

This means that by assuming the parameter trajectory $\bar{\mathbf{c}}_{\mathbf{q}}$ as the mean for the spectral parameter vector sequence \mathbf{c} corresponding to an utterance, the standard HMM can naturally be translated into a trajectory model: the state output probability of observing the static part of the output vector changes during a state, and is affected by statistics of neighboring states. Note that the spectral parameter vector sequence \mathbf{c} is modeled by a mixture of Gaussians whose dimensionality is TM , and their covariances $\mathbf{P}_{\mathbf{q}}$ are generally full.

3. Training Algorithm

In this section, we derive a training algorithm for the trajectory model. It should be noted that although the model has the same parameterization as the standard HMM, the output probability is defined by (18) rather than by (1). Accordingly, the model parameters should be trained based on (18).

An auxiliary function of current parameter set λ and new parameter set λ' is defined by

$$\mathcal{Q}(\lambda, \lambda') = \sum_{\text{all } \mathbf{q}} P(\mathbf{q} | \mathbf{c}, \lambda) \log P(\mathbf{c}, \mathbf{q} | \lambda'). \quad (21)$$

It can be shown that by substituting λ' which maximizes $\mathcal{Q}(\lambda, \lambda')$ for λ , the likelihood increases unless λ is a critical point of the likelihood. However, it is not tractable since we have to evaluate all possible state sequences. To avoid this difficulty, we apply the single Viterbi path approximation¹. As a result, the problem is broken down into the following maximization problems:

$$\mathbf{q}_{\max} = \arg \max_{\mathbf{q}} P(\mathbf{c}, \mathbf{q} | \lambda) \quad (22)$$

$$\lambda' = \arg \max_{\lambda} P(\mathbf{c}, \mathbf{q}_{\max} | \lambda) \quad (23)$$

3.1. Model parameter update

First, we solve the maximization problem of (23). The problem is equivalent to maximizing

$$\begin{aligned} \log P(\mathbf{c} | \mathbf{q}, \lambda) &= -\frac{1}{2} [MT \log(2\pi) - \log |\mathbf{R}_{\mathbf{q}}| \\ &\quad + \mathbf{c}^{\top} \mathbf{P}_{\mathbf{q}}^{-1} \mathbf{c} + \mathbf{r}_{\mathbf{q}}^{\top} \mathbf{P}_{\mathbf{q}} \mathbf{r}_{\mathbf{q}} - 2 \mathbf{r}_{\mathbf{q}}^{\top} \mathbf{c}] \end{aligned} \quad (24)$$

with respect to

$$\mathbf{m} = \left[\boldsymbol{\mu}_1^{\top}, \boldsymbol{\mu}_2^{\top}, \dots, \boldsymbol{\mu}_N^{\top} \right]^{\top} \quad (25)$$

$$\boldsymbol{\sigma} = \left[\mathbf{U}_1^{-1}, \mathbf{U}_2^{-1}, \dots, \mathbf{U}_N^{-1} \right]^{\top}, \quad (26)$$

where N is the total number of HMM states. In this paper, we refer to $P(\mathbf{c} | \mathbf{q}, \lambda)$ as ‘‘trajectory likelihood.’’

By setting $\partial \log P(\mathbf{c} | \mathbf{q}, \lambda) / \partial \mathbf{m} = \mathbf{0}$, we obtain a set of linear equations

$$\mathbf{S}_{\mathbf{q}}^{\top} \mathbf{W} \mathbf{P}_{\mathbf{q}} \mathbf{W}^{\top} \mathbf{S}_{\mathbf{q}} \boldsymbol{\Sigma}^{-1} \mathbf{m} = \mathbf{S}_{\mathbf{q}}^{\top} \mathbf{W} \mathbf{c} \quad (27)$$

for determination of \mathbf{m} which maximizes $\log P(\mathbf{c} | \mathbf{q}, \lambda)$, where

$$\boldsymbol{\Sigma}^{-1} = \text{diag}(\boldsymbol{\sigma}) \quad (28)$$

$$\boldsymbol{\mu}_{\mathbf{q}} = \mathbf{S}_{\mathbf{q}} \mathbf{m} \quad (29)$$

$$\mathbf{U}_{\mathbf{q}}^{-1} = \text{diag}(\mathbf{S}_{\mathbf{q}} \boldsymbol{\sigma}), \quad (30)$$

and $\mathbf{S}_{\mathbf{q}}$ is a $3T \times 3MN$ matrix whose elements are 0 or 1 determined according to the state sequence \mathbf{q} . The dimensionality of (27) is $3MN$: although it could be tens of thousands, it is still small enough to solve the set of linear equations using currently available computational resources.

For maximizing $\log P(\mathbf{c} | \mathbf{q}, \lambda)$ with respect to $\boldsymbol{\sigma}$, we apply a steepest descent algorithm using the first derivative

$$\begin{aligned} \frac{\partial \log P(\mathbf{c} | \mathbf{q}, \lambda)}{\partial \boldsymbol{\sigma}} &= \frac{1}{2} \mathbf{S}_{\mathbf{q}}^{\top} \text{diag}^{-1}(\mathbf{W} \mathbf{P}_{\mathbf{q}} \mathbf{W}^{\top} \\ &\quad - \mathbf{W} \mathbf{c} \mathbf{c}^{\top} \mathbf{W}^{\top} + 2 \boldsymbol{\mu}_{\mathbf{q}} \mathbf{c}^{\top} \mathbf{W}^{\top} \\ &\quad + \mathbf{W} \bar{\mathbf{c}}_{\mathbf{q}} \bar{\mathbf{c}}_{\mathbf{q}}^{\top} \mathbf{W}^{\top} - 2 \boldsymbol{\mu}_{\mathbf{q}} \bar{\mathbf{c}}_{\mathbf{q}}^{\top} \mathbf{W}^{\top}) \end{aligned} \quad (31)$$

because (24) is not a quadratic function of $\boldsymbol{\sigma}$.

¹We can also use the N -best approximation.

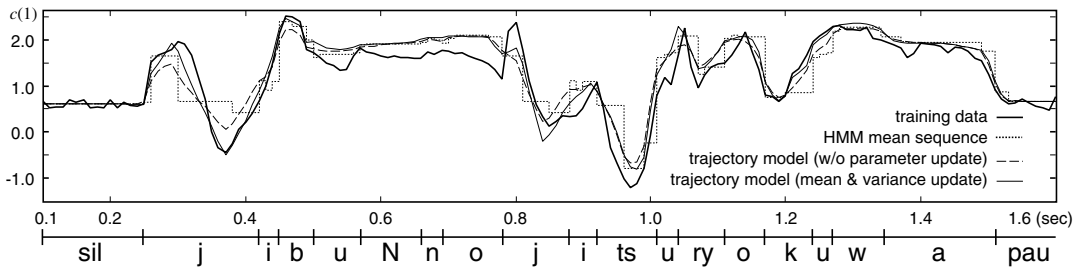


Figure 1: Generated trajectories \bar{c}_q for the trajectory HMM (thin line) and the standard HMM (broken line), the state mean sequence μ_q (dotted line), and one of training data c (thick line).

Table 1: Average log “trajectory likelihood” $P(c | q, \lambda)$.

	trajectory HMM			
	w/o update	m	σ	$m \& \sigma$
training data	6.86	7.91	12.0	14.8
test data	6.93	7.83	12.1	14.8

3.2. Optimal path search

Next, we discuss the optimal path search problem of (22). Based on the approximation

$$\begin{aligned} q_{\max} &= \arg \max_q P(q, c | \lambda) \\ &= \arg \max_q \frac{1}{K_q} P(q, o | \lambda) \end{aligned} \quad (32)$$

$$\simeq \arg \max_q P(q, o | \lambda), \quad (33)$$

we can use the Viterbi algorithm for the standard HMM. To reduce the inaccuracy caused by the approximation, the obtained state boundaries can be adjusted by a recursive algorithm (case 2 in [14]) based on the “trajectory likelihood.” By incorporating time-recursive calculation of the normalization factor K_q in the Viterbi algorithm, we can obtain an optimal path directly by a modified Viterbi algorithm. This would be possible by using the speech parameter generation algorithm of [21] which can operate in a time-recursive manner with some look-ahead.

Noted that extensions of the training algorithm to the multiple observations and multi-mixture distributions are straightforward.

3.3. Decoding

For decoding, we can use i) rescoring scheme, ii) recursive calculation of normalization factor K_q in the Viterbi algorithm. In approach i), N -best candidates are rescored by the “trajectory likelihood.” The state boundaries of each candidate can also be adjusted by the recursive algorithm (case 2 in [14]). On the other hand, in approach ii), K_q is calculated in the Viterbi decoding using the time-recursive algorithm of [21] with some look-ahead.

4. Experiment

A preliminary experiment was carried out. We used phonetically balanced 450 sentences from ATR Japanese speech database for training speaker-dependent monophone HMMs and their trajectory versions. Diagonal covariances were used for both types of models. Speech signals were sampled at 16

Table 2: Recognition error rates (%) for the trajectory HMM and the standard HMM.

	trajectory HMM				standard HMM
	w/o update	m	σ	$m \& \sigma$	
error rate	19.4	18.5	19.0	18.1	19.5
rel. imp.	0.5%	5%	2.6%	7%	ref

kHz and windowed by a 25.6-ms Blackman window with a 10-ms shift, and then mel-cepstral coefficients were obtained by a mel-cepstral analysis technique [22]. The feature vector consists of 19 mel-cepstral coefficients including the zeroth coefficient as well as their delta and delta-delta coefficients. We used 3-state left-to-right HMM structure with no skips. For a given Viterbi path obtained by the standard HMM, parameters of the trajectory HMM were optimized by the algorithm described in Section 3.1 (i.e., the proposed training algorithm was not iterated in this experiment). The trained trajectory HMM was evaluated in a continuous phoneme recognition task based on the N -best rescoring approach: first, a 10-best list was generated for each test utterance by using the conventional Viterbi algorithm with the standard HMM, and then each candidate was rescored by the trajectory HMM.

Table 1 shows average log likelihoods of the trajectory HMMs (“w/o update”: model parameters were not updated, i.e., those of the standard HMM were used, “ m ”: only means were updated, “ σ ”: only variances were updated, “ $m \& \sigma$ ”: both means and variances were updated). It is shown that the training algorithm improves the “trajectory likelihoods” considerably for both training and test data sets.

Figure 1 shows the mean trajectories (the first coefficient in \bar{c}_q) for the trajectory HMM (“w/o parameter update”: model parameter of the standard HMM were used, “mean and variance update”: means and variances were updated). The state mean sequence of the standard HMM (the first coefficient in static part of μ_q) and the trajectory extracted from one of the training data (the first coefficient in c) are also plotted in the figure. It can be seen that the trajectory generated from the trained trajectory HMM is closer to the trajectory of the training data. It should be noted that the trajectory corresponding to a phoneme varies according to its neighboring phonemes (see phoneme /j/ in the figure). This means that the trajectory HMM has the capability to capture the coarticulation effects naturally.

Table 2 shows the phoneme recognition error rates for the trajectory HMM and the standard HMM. The trajectory HMM achieves a relative error reduction of 7% over the standard HMM. Given the fact that without parameter update the perfor-

mance of the trajectory HMM did not improve, model parameter training based on the “trajectory likelihood” is essential in the trajectory HMM approach.

5. Conclusion

This paper defined a new kind of trajectory model (trajectory HMM) for acoustic modeling of speech by reformulating the standard HMM whose observation vector consists of static and dynamic features. We also derived a Viterbi-type training algorithm for the trajectory model. It was shown that the model can work as a trajectory model while it maintains the basic structure of the standard HMM. In a preliminary speech recognition experiment, a relative error rate reduction of 7% over the standard HMM was achieved. One of the advantages of the trajectory HMM is that huge amounts of software resources for the standard HMM could be easily reused.

In the near future, results of a complete speech recognition experiment and application to Text-to-Speech synthesis will be presented. Future work also includes derivation of a fast algorithm for solving (27) and a Baum-Welch-type training algorithm.

6. Acknowledgment

Authors would like to thank Prof. Takao Kobayashi and Dr. Takashi Masuko for discussions.

7. References

- [1] M. J. F. Gales and S. J. Young, “Segmental Hidden Markov Models,” Proc. EUROSPEECH, pp.1579–1582, 1993.
- [2] H. Gish and K. Ng, “Parametric trajectory models for speech recognition,” Proc. ICSLP, pp.I-466–I-469, 1996.
- [3] W. J. Holmes and M. J. Russell, “Speech recognition using a linear dynamic segmental HMM,” Proc. ICASSP, pp.1611–1614, 1995.
- [4] M. Ostendorf, V. Digalakis, and O. A. Kimball, “From HMMs to segment models: a unified view of stochastic modeling for speech recognition,” IEEE Trans. on Speech and Audio Processing, vol.4, no.5, pp.360–378, 1996.
- [5] L. Deng, M. Aksmanovic, X. Sun, J. Wu, “Speech recognition using hidden Markov models with polynomial regression functions as nonstationary states,” IEEE Trans. on Speech and Audio Processing, vol.2, no.4, pp.507–520, Oct. 1994.
- [6] Y. Gong and J. P. Haton, “Stochastic trajectory modeling for speech recognition,” Proc. ICASSP, vol.I, pp57–60, 1994.
- [7] S. Takahashi, “Phoneme HMM’s constrained by frame correlations,” Proc. ICASSP, pp. 219-222, 1993.
- [8] K. K. Paliwal, “Use of temporal correlation between successive frames in hidden Markov model based Speech recognizer,” Proc. ICASSP, pp.215-218, 1993.
- [9] M. Ostendorf and S.Roukos, “A stochastic segment model for phoneme-based continuous speech recognition,” IEEE Trans. On Acoustics, Speech and Signal Processing, vol.37, no.12, pp.1857–1869, 1989.
- [10] C. J. Wellekens, “Explicit correlation in hidden Markov model for Speech Recognition,” Proc. ICASSP, pp.383-386, 1987.
- [11] S. Furui, “Speaker independent isolated word recognition using dynamic features of speech spectrum,” IEEE Trans. Acoust., Speech, Signal Processing, vol.ASSP-34, pp.52–59, 1986.
- [12] K. Tokuda, T. Kobayashi and S. Imai, “Speech parameter generation from HMM using dynamic features,” Proc. ICASSP, pp.660–663, 1995.
- [13] K. Tokuda, T. Masuko, T. Yamada, T. Kobayashi and S. Imai, “An algorithm for speech parameter generation from continuous mixture HMMs with dynamic features,” Proc. EUROSPEECH, pp.757–760, 1995.

- [14] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, T. Kitamura, “Speech parameter generation algorithms for HMM-based speech synthesis,” Proc. ICASSP, vol.3, pp.1315–1318, June 2000.
- [15] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi and T. Kitamura, “Simultaneous modeling of spectrum, pitch and duration in HMM-based speech synthesis,” Proc. EUROSPEECH, pp.2347–2350, 1999.
- [16] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi and T. Kitamura, “Speaker Interpolation in HMM-Based Speech Synthesis System,” Proc. EUROSPEECH, pp.2523–2526, 1997.
- [17] M. Tamura, T. Masuko, K. Tokuda and T. Kobayashi, “Adaptation of pitch and spectrum for HMM-based speech synthesis using MLLR,” Proc. ICASSP, vol.2, pp.805–808, May 2001.
- [18] K. Shichiri, A. Sawabe, K. Tokuda, T. Masuko, T. Kobayashi, T. Kitamura, “Eigenvoices for HMM-based speech synthesis,” Proc. ICSLP, pp.1269–1272, 2002.
- [19] Y. Minami, E. McDermott, A. Nakamura, S. Katagiri, “A recognition method with parametric trajectory synthesized using direct relations between static and dynamic feature vector time series,” Proc. ICASSP, vol.I, pp.957–960, 2002.
- [20] K. S. Van Horn, “A Maximum-entropy solution to the frame-dependency problem in speech recognition,” Tech. Rep., Dept. of Computer Science, North Dakota State University, Nov. 2001.
- [21] K. Koishida, K. Tokuda, T. Masuko and T. Kobayashi, “Vector quantization of speech spectral parameters using statistics of dynamic features,” Proc. ICSP, pp.247–252, 1997.
- [22] T. Fukada, K. Tokuda, T. Kobayashi and S. Imai, “An adaptive algorithm for mel-cepstral analysis of speech,” in Proc. ICASSP, 1992, pp.137–140.
- [23] A. Acero, “Formant analysis and synthesis using hidden Markov models,” Proc. EUROSPEECH, pp.1047–1050, 1999.

A. Appendix

A.1. Maximizing $P(\mathbf{o} | \mathbf{q}, \lambda)$ with respect to \mathbf{o} under $\mathbf{o} = \mathbf{W}\mathbf{c}$

The logarithm of $P(\mathbf{o} | \mathbf{q}, \lambda)$ can be written as

$$\log P(\mathbf{o} | \mathbf{q}, \lambda) = -\frac{1}{2} \mathbf{o}^\top \mathbf{U}_q^{-1} \mathbf{o} + \mathbf{o}^\top \mathbf{U}_q^{-1} \boldsymbol{\mu}_q + \text{Const}, \quad (34)$$

where Const is independent of \mathbf{o} .

It is obvious that $P(\mathbf{o} | \mathbf{q}, \lambda)$ is maximized when $\mathbf{o} = \boldsymbol{\mu}_q$ without the condition (4), that is, the speech parameter vector sequence becomes a sequence of the mean vectors. With the condition (4), maximizing $P(\mathbf{o} | \mathbf{q}, \lambda)$ with respect to \mathbf{o} is equivalent to that with respect to \mathbf{c} . By setting

$$\frac{\partial \log P(\mathbf{W}\mathbf{c} | \mathbf{q}, \lambda)}{\partial \mathbf{c}} = \mathbf{0}, \quad (35)$$

we obtain a set of equations

$$\mathbf{R}_q \mathbf{c} = \mathbf{r}_q. \quad (36)$$

For direct solution of (36), we need $O(T^3 M^3)$ operations² because \mathbf{R}_q is a $TM \times TM$ matrix. By utilizing the special structure of \mathbf{R}_q , (36) can be solved by the Cholesky decomposition or the QR decomposition with $O(TM^3 L^2)$ operations³, where $L = \max_{n \in \{1, 2\}, s \in \{-, +\}} L_s^{(n)}$. Equation (36) can also be solved by an algorithm derived in [12]–[14], which can operate in a time-recursive manner [21].

²When $\mathbf{U}_{q,i}$ is diagonal, it is reduced to $O(T^3 M)$ since each of the M -dimensions can be calculated independently.

³When $\mathbf{U}_{q,i}$ is diagonal, it is reduced to $O(TML^2)$. Furthermore, when $L_-^{(1)} = -1$, $L_+^{(1)} = 0$, and $w^{(2)}(i) \equiv 0$, it is reduced to $O(TM)$ as described in [23].