

FACTOR ANALYZED VOICE MODELS FOR HMM-BASED SPEECH SYNTHESIS

Kyosuke Kazumi, Yoshihiko Nankaku, Keiichi Tokuda

Nagoya Institute of Technology, Nagoya, Japan

ABSTRACT

This paper describes factor analyzed voice models for realizing various voice characteristics in the HMM-based speech synthesis. The eigenvoice method can synthesize speech with arbitrary voice characteristics by interpolating representative HMM sets. However, the objective of PCA is to accurately reconstruct each speaker-dependent HMM set, and this is not equivalent to estimating models which represent training data accurately. To overcome this problem, we propose a general speech model which generates speech utterances with various voice characteristics directly. In the proposed method, the HMM states, factors representing voice characteristics and contextual decision trees are simultaneously optimized within a unified framework.

Index Terms— HMM-based speech synthesis, eigenvoice, factor analysis, expectation maximization algorithm, deterministic annealing EM algorithm

1. INTRODUCTION

To let machines speak naturally like a human, an HMM-based speech synthesis system has been proposed [1, 2]. The system models spectrum, pitch and state duration simultaneously in a unified framework of HMMs and synthesizes speech using parameters generated from HMM sets. One of the advantages of the system is that voice characteristics of synthesized speech can be easily changed by transforming HMM parameters.

The eigenvoice method for HMM based synthesis [3] can synthesize speech with arbitrary voice characteristics by interpolating representative HMM sets. In this method, a set of speaker-dependent HMMs is represented by a single large dimensional vector called a supervector and representative supervectors called eigenvoice vectors are constructed by applying Principal Component Analysis (PCA). Finally, speech with a desired voice characteristic is generated by a target HMM set which is reconstructed by interpolating the eigenvoice vectors using given weights representing factors of voice characteristics. However, the objective of PCA in the eigenvoice method is to accurately reconstruct each speaker-dependent HMM set, and this is not equivalent to estimating models which represent training data accurately. Furthermore, speaker-dependent HMMs for all speakers are required with the same decision trees (parameter tying structures). However, it is sometimes difficult to collect enough speech data for all speakers.

To overcome these problems, we propose a general speech model which generates speech utterances with various voice characteristics directly. In this paper, we name the proposed model factor analyzed voice model (FA-voice model), because the structure of the proposed method is based on the factor analysis (FA) model for representing various speaker characteristics. In the proposed method, the likelihood is directly calculated from speech utterances of training data, therefore speaker-dependent HMMs are not required. In the training of the proposed model, the HMM states,

factors representing voice characteristics and contextual decision trees are simultaneously optimized within a unified maximum likelihood (ML) framework based on a single statistical model. This simultaneous optimization is expected to achieve a better quality of synthesized speech for a desired voice characteristics.

The parameters of the proposed model can be estimated via the expectation maximization (EM) algorithm for approximating the Maximum Likelihood (ML) estimate. However, the exact expectation step (E-step) is computationally intractable due to the combination of hidden variables. To derive a feasible algorithm, we applied the variational EM algorithm [4] to the proposed model. The variational method approximates the posterior distribution over the hidden variables by a tractable distribution. However, the EM algorithm has the problem that the solution converges to a local optimum and the convergence point depends on the initial model parameters. This problem causes that estimated model parameters are not appropriate for training data and the quality of synthesized speech may be degraded. To overcome this problem, we apply the deterministic annealing EM (DAEM) algorithm [5] to the training of the FA voice model.

The rest of the paper is organized as follows. Section 2 describes the eigenvoice method based on PCA. In Section 3, we propose the voice generation model based on factor analysis and estimation method for model parameters using the EM algorithm, and Section 4 describes the DAEM algorithm for estimating the proposed model parameters. Experimental results are presented in Section 5, and concluding remarks and future work are presented in the Section 6.

2. EIGENVOICE METHOD BASED ON PCA

In the eigenvoice method for HMM-based speech synthesis, speaker dependent HMM sets represent voice characteristics of each speaker. Therefore, mean vectors of these HMM sets are generated from eigenvoice vectors and these weights. For each of S speaker-dependent HMM sets, we extract parameters representing all HMMs. Then, we concatenate all parameters for each speaker and create a vector of a large dimensionality which is called “supervector.” By applying PCA, S eigenvoice vectors are calculated from S supervectors. Eigenvoice vectors are called “eigenvoices.” Using the first L eigenvoices e_l , L arbitrary weight coefficients λ_l for eigenvoices $e_l, l = 1, 2, \dots, L$, and a vector $\bar{\mu}$ is a mean vector of S supervectors, a new supervector $\hat{\mu}$ is calculated as follows:

$$\hat{\mu} = \bar{\mu} + \sum_{l=1}^L \lambda_l e_l. \quad (1)$$

Then a new speaker dependent HMM set is reconstructed from the generated supervector $\hat{\mu}$. It is noted that eigenvoices for speech synthesis should be different from those in speech recognition because they should capture not only spectrum but also $F0$ parameters. Various voice qualities can be synthesized by setting arbitrary weight coefficients λ_l .

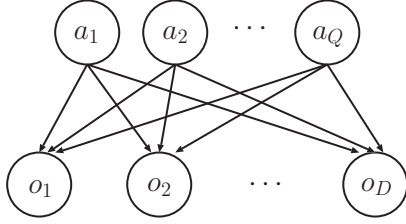


Fig. 1. Factor analysis model

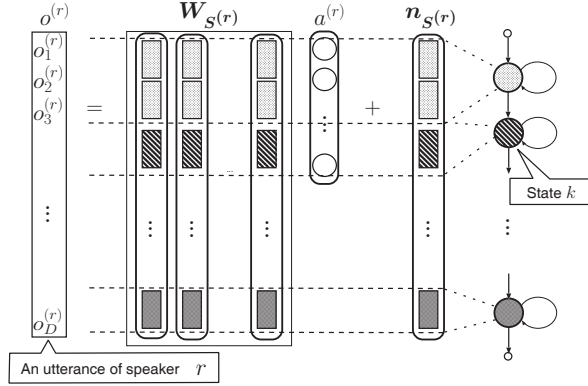


Fig. 2. Structure of the FA voice model

3. FACTOR ANALYZED VOICE MODELS

3.1. Factor analysis

The factor analysis (FA) is one of a statistical method for modeling the covariance structure of high dimensional data using a small number of latent variables. In the FA model, observation variable \mathbf{o} is generated by following equation:

$$\mathbf{o} = \mathbf{W}\mathbf{a} + \mathbf{n}, \quad (2)$$

where \mathbf{a} is a Q dimensional latent variable, \mathbf{n} is a noise vector, and $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_Q]$, $\mathbf{w}_i = [w_{i1}, w_{i2}, \dots, w_{iD}]^\top$ is a $D \times Q$ matrix known as a factor loading matrix. The structure of the FA model is shown in Figure 1.

Applying the FA model as speech generation model, acoustic feature vectors of speaker dependent utterance $\mathbf{o} = [o_1, o_2, \dots, o_D]^\top$ are generated from a common latent variable \mathbf{a} and noise vector \mathbf{n} which are inherent in each variable given by

$$\mathbf{a} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{n} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (3)$$

where \mathbf{I} is identity matrix. In this condition, observation sequence \mathbf{o} is generated from

$$P(\mathbf{o} | \boldsymbol{\Lambda}) = \int P(\mathbf{o} | \mathbf{a}, \boldsymbol{\Lambda}) P(\mathbf{a} | \boldsymbol{\Lambda}) d\mathbf{a}, \quad (4)$$

$$P(\mathbf{o} | \mathbf{a}, \boldsymbol{\Lambda}) = \mathcal{N}(\mathbf{o} | \mathbf{W}\mathbf{a} + \boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (5)$$

where $\boldsymbol{\Lambda}$ is model parameter.

3.2. FA voice models

In the factor analyzed voice models, the loading matrix \mathbf{W} corresponds to eigenvoice vectors and we assume that the loading matrix and noise vector are generated from HMM sets. Figure 2 shows the proposed model structure. In the proposed models, eigenvoice

vectors can convert their structure according to various input data. Hence, this model receives states changing similarly to HMM. Furthermore, HMM parameters are generated from liner combination of eigenvoice vectors. As a result, this model is feature generating model with liner transformation.

In the FA voice model, an utterance observation sequence of speaker r is generated from following expression:

$$\mathbf{o}^{(r)} = \mathbf{W}_{\mathbf{S}^{(r)}} \mathbf{a}^{(r)} + \mathbf{n}_{\mathbf{S}^{(r)}}, \quad (6)$$

where $\mathbf{o}^{(r)}$ is an observation sequence of speaker r , $\mathbf{a}^{(r)}$ is latent variable, and $\mathbf{S}^{(r)}$ is HMM path. The loading matrix $\mathbf{W}_{\mathbf{S}^{(r)}}$ and a noise vector $\mathbf{n}_{\mathbf{S}^{(r)}}$ depend on state transition. The latent variable is prepared for each speaker and represents speaker qualities.

In this model, the likelihood function for utterances of all speakers is written as

$$P(\mathbf{o} | \boldsymbol{\Lambda}) = \prod_r \sum_{\mathbf{S}^{(r)}} \int P(\mathbf{o}^{(r)} | \mathbf{a}^{(r)}, \mathbf{S}^{(r)}, \boldsymbol{\Lambda}) \times P(\mathbf{a}^{(r)} | \boldsymbol{\Lambda}) P(\mathbf{S}^{(r)} | \boldsymbol{\Lambda}) d\mathbf{a}^{(r)}, \quad (7)$$

$$P(\mathbf{o}^{(r)} | \mathbf{a}^{(r)}, \mathbf{S}^{(r)}, \boldsymbol{\Lambda}) = \mathcal{N}(\mathbf{o}^{(r)} | \mathbf{W}_{\mathbf{S}^{(r)}} \mathbf{a}^{(r)} + \boldsymbol{\mu}_{\mathbf{S}^{(r)}}, \boldsymbol{\Sigma}_{\mathbf{S}^{(r)}}), \quad (8)$$

where $\boldsymbol{\Lambda}$ is model parameter and $\boldsymbol{\mu}_{\mathbf{S}^{(r)}}$, $\boldsymbol{\Sigma}_{\mathbf{S}^{(r)}}$ are mean and variance of noise vector $\mathbf{n}_{\mathbf{S}^{(r)}}$.

3.3. EM algorithm for FA voice models

In the proposed, model parameters are estimated by the expectation maximization (EM) algorithm to maximize the log-likelihood for given training data. \mathcal{Q} -function of the EM algorithm is given by

$$\mathcal{Q}(\boldsymbol{\Lambda}, \boldsymbol{\Lambda}') = \sum_r \sum_{\mathbf{S}^{(r)}} \int P(\mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \mathbf{o}^{(r)}, \boldsymbol{\Lambda}') \times \ln P(\mathbf{o}^{(r)}, \mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \boldsymbol{\Lambda}') d\mathbf{a}^{(r)}, \quad (9)$$

where $\boldsymbol{\Lambda}'$ is estimated model parameter. Maximizing \mathcal{Q} -function, the log-likelihood is guaranteed to increase or remain unchanged. However, calculating \mathcal{Q} -function is very hard, because it needs the summation for all paths $\mathbf{S}^{(r)}$ and the integral over the latent variable $\mathbf{a}^{(r)}$. Therefore, we introduce the variational methods to the EM algorithm.

In the variational EM algorithm, we approximate the posterior distribution in Eq. (9) by using an arbitrary distribution $Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)})$ and define the lower bound of the log-likelihood:

$$\begin{aligned} & \ln P(\mathbf{o} | \boldsymbol{\Lambda}) \\ &= \ln \sum_r \sum_{\mathbf{S}^{(r)}} \int P(\mathbf{o}^{(r)}, \mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \boldsymbol{\Lambda}) d\mathbf{a}^{(r)} \\ &= \ln \sum_r \sum_{\mathbf{S}^{(r)}} \int Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) \frac{P(\mathbf{o}^{(r)}, \mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \boldsymbol{\Lambda})}{Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)})} d\mathbf{a}^{(r)} \\ &\geq \sum_r \sum_{\mathbf{S}^{(r)}} \int Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) \ln \frac{P(\mathbf{o}^{(r)}, \mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \boldsymbol{\Lambda})}{Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)})} d\mathbf{a}^{(r)} \\ &= \sum_r \sum_{\mathbf{S}^{(r)}} \int Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) \ln P(\mathbf{o}^{(r)}, \mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \boldsymbol{\Lambda}) d\mathbf{a}^{(r)} \\ &\quad - \sum_r \sum_{\mathbf{S}^{(r)}} \int Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) \ln Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) d\mathbf{a}^{(r)} \\ &= \mathcal{F}(Q, \boldsymbol{\Lambda}), \end{aligned} \quad (10)$$

where Jensen's inequality is applied. The difference between $\ln P(\mathbf{o} | \mathbf{\Lambda})$ and \mathcal{F} is given by Kullback-Leibler divergence between $Q(\mathbf{a}, \mathbf{S})$ and the true posterior distribution $P(\mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \mathbf{o}^{(r)}, \mathbf{\Lambda})$. The E-step computes the posterior probabilities $Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)})$ which maximizes \mathcal{F} .

In this paper, we assume the following constraints:

$$Q(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) = Q(\mathbf{a}^{(r)})Q(\mathbf{S}^{(r)}), \quad (11)$$

$$\sum_{\mathbf{S}^{(r)}} Q(\mathbf{S}^{(r)}) = 1, \quad \int Q(\mathbf{a}^{(r)}) d\mathbf{a} = 1. \quad (12)$$

Based on these constraints, we use Lagrange multiplier method and obtain following $Q(\mathbf{a}^{(r)})$ and $Q(\mathbf{S}^{(r)})$:

$$Q(\mathbf{a}^{(r)}) \propto P(\mathbf{a}^{(r)} | \mathbf{\Lambda}) \exp \left\langle \log P(\mathbf{o}^{(r)} | \mathbf{a}^{(r)}, \mathbf{S}^{(r)}, \mathbf{\Lambda}) \right\rangle_{Q(\mathbf{S}^{(r)})}, \quad (13)$$

$$Q(\mathbf{S}^{(r)}) \propto P(\mathbf{S}^{(r)} | \mathbf{\Lambda}) \exp \left\langle \log P(\mathbf{o}^{(r)} | \mathbf{a}^{(r)}, \mathbf{S}^{(r)}, \mathbf{\Lambda}) \right\rangle_{Q(\mathbf{a}^{(r)})}, \quad (14)$$

where $\langle \cdot \rangle_Q$ represents an expectation with respect to a distribution Q . The structure of Eq. (14) includes an expectation of a latent variable. Furthermore, Eq. (13) can be transformed into a Gaussian distribution. The mean $\boldsymbol{\mu}_a^{(r)}$ and variance $\boldsymbol{\Sigma}_a^{(r)}$ of Eq. (13) are given by

$$\boldsymbol{\mu}_a^{(r)} = \boldsymbol{\Sigma}_a^{(r)} \left(\sum_k \mathbf{W}_k^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{U}_k^{(r)} - N_k^{(r)} \boldsymbol{\mu}_k) \right), \quad (15)$$

$$\boldsymbol{\Sigma}_a^{(r)} = \left(\mathbf{I} + \sum_k N_k^{(r)} \mathbf{W}_k^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{W}_k \right)^{-1}. \quad (16)$$

We define following equations:

$$N_k^{(r)} = \sum_t \langle S_t^{(r)}, k \rangle, \quad \mathbf{U}_k^{(r)} = \sum_t \langle S_t^{(r)}, k \rangle \mathbf{o}_t^{(r)}, \quad (17)$$

where $\langle S_t^{(r)}, k \rangle$ is a probability of staying in state k at time t given observation \mathbf{o} .

In the M-step, the model parameter $\mathbf{\Lambda}$ is calculated for maximizing the lower bound \mathcal{F} . For simplifying, we define the following equations.

$$\bar{\mathbf{a}}^{(r)} = [\mathbf{1} \quad \mathbf{a}^{(r)\top}]^\top, \quad (18)$$

$$\bar{\mathbf{W}}_{\mathbf{S}^{(r)}} = [\boldsymbol{\mu}_{\mathbf{S}^{(r)}} \quad \mathbf{W}_{\mathbf{S}^{(r)}}], \quad (19)$$

The lower bound $\mathcal{F}(Q, \mathbf{\Lambda})$ is partial differentiated about each model parameter and we obtain following equations for re-estimation:

$$\bar{\mathbf{W}}_k = \left(\sum_r \mathbf{U}_k^{(r)} \langle \bar{\mathbf{a}}^{(r)\top} \rangle \right) \left(\sum_r N_k^{(r)} \langle \bar{\mathbf{a}}^{(r)} \bar{\mathbf{a}}^{(r)\top} \rangle \right)^{-1}, \quad (20)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_r N_k^{(r)}} \sum_r \left(\mathbf{V}_k^{(r)} - \mathbf{U}_k^{(r)} \langle \bar{\mathbf{a}}^{(r)} \rangle^\top \bar{\mathbf{W}}_k^\top \right), \quad (21)$$

where

$$\mathbf{V}_k^{(r)} = \sum_t \langle S_t^{(r)}, k \rangle \mathbf{o}_t^{(r)} \mathbf{o}_t^{(r)\top}, \quad (22)$$

$$\langle \bar{\mathbf{a}}^{(r)} \rangle = [\mathbf{1} \quad \boldsymbol{\mu}_a^{(r)\top}]^\top, \quad (23)$$

$$\langle \bar{\mathbf{a}}^{(r)} \bar{\mathbf{a}}^{(r)\top} \rangle = \begin{bmatrix} 1 & \boldsymbol{\mu}_a^{(r)\top} \\ \boldsymbol{\mu}_a^{(r)} & \boldsymbol{\Sigma}_a^{(r)} + \boldsymbol{\mu}_a^{(r)} \boldsymbol{\mu}_a^{(r)\top} \end{bmatrix}. \quad (24)$$

The variational EM algorithm iteratively maximizes \mathcal{F} with respect to the Q and $\mathbf{\Lambda}$ holding the other parameters fixed:

$$(\text{E step}) : Q^{(k+1)} = \arg \max_Q \mathcal{F}(Q, \mathbf{\Lambda}^{(k)}), \quad (25)$$

$$(\text{M step}) : \mathbf{\Lambda}^{(k+1)} = \arg \max_{\mathbf{\Lambda}} \mathcal{F}(Q^{(k+1)}, \mathbf{\Lambda}), \quad (26)$$

In this procedure, the lower bound \mathcal{F} is guaranteed to increase instead of the value of the Q -function.

3.4. Context clustering

The proposed models can make arbitrary shared structure for each model. Therefore, we apply some questions to each node for increasing the log-likelihood and make clusters like HMM. In this paper, we make each model has same structure. In this clustering method, it is easy to re-estimate model parameters simultaneously because each model have a same structure. Calculating the log-likelihood directly costs a lot of computational cost. Hence, we use the lower bound of the log-likelihood \mathcal{F} in Eq. (10).

4. ANNEALING BASED ESTIMATION ALGORITHM FOR FA VOICE MODELS

In this paper, we propose the DAEM algorithm [5] as annealing based estimation algorithm. Applying the DAEM algorithm to the FA voice model, the free energy function \mathcal{L}_β can be defined as

$$\mathcal{L}_\beta = -\frac{1}{\beta} \ln \sum_r \sum_{\mathbf{S}^{(r)}} \int P(\mathbf{o}^{(r)}, \mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \mathbf{\Lambda})^\beta d\mathbf{a}^{(r)}, \quad (27)$$

where $\frac{1}{\beta}$ is called "temperature" and β is temperature parameter. Since the negative free energy function corresponds to the log-likelihood function, the lower bound of $-\mathcal{L}_\beta$ is defined as

$$\begin{aligned} & -\mathcal{L}_\beta \\ & \geq \frac{1}{\beta} \sum_r \sum_{\mathbf{S}^{(r)}} \int Q_\beta(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) \ln P(\mathbf{o}^{(r)}, \mathbf{a}^{(r)}, \mathbf{S}^{(r)} | \mathbf{\Lambda})^\beta d\mathbf{a}^{(r)} \\ & \quad - \frac{1}{\beta} \sum_r \sum_{\mathbf{S}^{(r)}} \int Q_\beta(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) \ln Q_\beta(\mathbf{a}^{(r)}, \mathbf{S}^{(r)}) d\mathbf{a}^{(r)} \\ & = \mathcal{F}_\beta(Q, \mathbf{\Lambda}), \end{aligned} \quad (28)$$

where $Q_\beta(\mathbf{a}^{(r)}, \mathbf{S}^{(r)})$ is the approximated posterior distribution similarly to the variational EM algorithm. $Q_\beta(\mathbf{a}^{(r)})$ can be derived as the Gaussian distribution which has following mean $\boldsymbol{\mu}_{a_\beta}^{(r)}$ and variance $\boldsymbol{\Sigma}_{a_\beta}^{(r)}$:

$$\boldsymbol{\mu}_{a_\beta}^{(r)} = \beta \boldsymbol{\Sigma}_{a_\beta}^{(r)} \left(\sum_k \mathbf{W}_k^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{U}_k^{(r)} - N_k^{(r)} \boldsymbol{\mu}_k) \right), \quad (29)$$

$$\boldsymbol{\Sigma}_{a_\beta}^{(r)} = \frac{1}{\beta} \left(\mathbf{I} + \sum_k N_k^{(r)} \mathbf{W}_k^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{W}_k \right)^{-1}. \quad (30)$$

These equations have the structure which adds the temperature parameter β to the approximated posterior distribution $Q(\mathbf{a}^{(r)})$ in the variational EM algorithm. Furthermore, \mathcal{F}_β is partial differentiated by model parameter and we can obtain a following equation:

$$\frac{\partial \mathcal{F}_\beta}{\partial \mathbf{\Lambda}} = \frac{\partial \mathcal{F}}{\partial \mathbf{\Lambda}} = \mathbf{0}. \quad (31)$$

This equation equivalent to the updating equation of the variational EM algorithm. In the DAEM algorithm, the temperature parameter β is gradually increased while iterating the EM-steps at each temperature. When $\frac{1}{\beta}$ is set to an initial temperature $\beta \simeq 0$, the EM-steps may achieve a single global maximum of \mathcal{F}_β . Finally at the $\beta = 1$, \mathcal{F}_β is identical with the variational EM algorithm.

5. EXPERIMENTS

5.1. Experimental conditions

For training speaker independent HMM sets, we used ATR Japanese speech database B-set uttered by six male and four female speakers. Nine speakers excluding one male speaker were used for the experiment. Each speaker uttered 503 sentences: different 50 sentences were used for each speaker in the training, and the remaining 53 sentences were used for evaluation. The speech data was down-sampled from 20kHz to 16kHz, windowed at a 5-ms frame rate using a 25-ms Blackman window. Feature vectors consisted of spectral and F_0 feature vectors. Each spectral feature vector consisted of 24 mel-cepstral coefficients and their delta and delta-delta coefficients. The F_0 parameter vectors consisted of log F_0 , its delta and delta-delta. A left-to-right, 5-state, MSD-HMM with no skip structure was used. The number of eigenvoice vectors and the number of factors in the proposed method are both two. For the DAEM algorithm, the temperature parameter β was updated by

$$\beta(i) = \left(\frac{i}{20}\right)^2, \quad (i = 0, 1, \dots, 20), \quad (32)$$

where i denotes the iteration number. At each temperature, 20 EM-steps were conducted.

To evaluate the performance of the FA voice models, the following four training methods were compared:

- “PCA”: the PCA based eigenvoice model with the model structure obtained by speaker independent HMMs.
- “PCA_STC”: the PCA based eigenvoice model with the model structure obtained by STC [6].
- “FA_EM”: the FA voice model initialized by speaker independent HMMs, and the EM algorithm was used.
- “FA_DAEM”: the FA voice model initialized by speaker independent HMMs, and the DAEM algorithm was used.

5.2. Experimental results

A subjective listening test was conducted to evaluate quality of synthesized speech. The test compared the naturalness of synthesized speech by the mean opinion score (MOS) test method. The subjects were 10 Japanese graduate students. Speech samples were randomly chosen from the evaluation sentences. Voice characteristics of speech samples were average voice and that of nine speakers in the training data, and two sentences for each characteristic, in total 20 sentences were prepared for each subject. In the MOS test, after listening to each test sample, the subjects were asked to assign it a 5-point naturalness score (5:excellent, 4:good, 3:fair, 2:poor, 1:bad).

Figure 3 plots the experimental results. It can be seen from the figure that the proposed methods “FA_EM” and “FA_DAEM” achieved better subjective scores than the conventional methods “PCA” and “PCA_STC.” Since the FA voice models were estimated from training data directly, appropriate model parameters were obtained for representing speech utterances accurately. It is also the reason of the improvement that the context clustering was performed based on the same criterion as the parameter estimation. It can also be seen that “FA_DAEM” obtained a better subjective score than “FA_EM.” Because the DAEM algorithm can improve the local maxima problem. These results clearly show the effectiveness of the proposed method in speech synthesis and the DAEM algorithm is more effective than the EM algorithm for estimating model parameters.

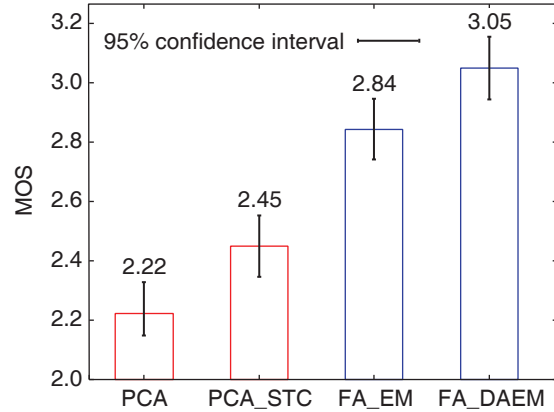


Fig. 3. Mean opinion scores of synthesized speech with 95% confidence intervals by the conventional and proposed methods

6. CONCLUSION

This paper proposed the FA voice model for HMM-based speech synthesis. It can synthesize speech with various voice characteristics. This method can estimate more appropriate model parameters than eigenvoice method based on PCA. We also derived the EM and DAEM algorithm for the proposed method. In the experiments, the proposed method achieved a higher performance than the conventional eigenvoice method. Furthermore, the DAEM algorithm improved the performance of the proposed method. Experiments on larger datasets and evaluation of synthesized speech with various voice characteristics will be future work.

7. ACKNOWLEDGEMENT

The research leading to these results was partly funded from the European Community’s Seventh Framework Programme (FP7/2007-2013) under grant agreement 213845 (the EMIME project), and the Strategic Information and Communications R&D Promotion Programme (SCOPE), Ministry of Internal Affairs and Communication, Japan.

8. REFERENCES

- [1] T. Masuko, K. Tokuda, T. Kobayashi, and S. Imai, “Speech synthesis from HMMs using dynamic features,” Proc. of ICASSP, pp.389–392, 1996.
- [2] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura, “Simultaneous modeling of spectrum, pitch, and duration in HMM-based speech synthesis,” Proc. of EUROSPEECH, pp.2347–2350, 1999.
- [3] K. Shichiri, A. Sawabe, T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura, “Eigenvoices for HMM-based speech synthesis,” Proc. of ICSLP, vol.1, pp.1269–1272, 2002.
- [4] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul, “An introduction to variational methods for Graphical models,” Machine Learning, vol.37, pp.183–233, 1997.
- [5] N. Ueda and R. Nakano, “Deterministic annealing EM algorithm,” Neural Networks, vol.11, no.2, pp.271–282, 1998
- [6] J. Yamagishi, M. Tamura, T. Masuko, K. Tokuda, and T. Kobayashi, “A context clustering technique for average voice model in HMM-based speech synthesis,” Proc. of ICSLP, pp. 133–136, 2002.