

Probabilistic Feature Mapping Based on Trajectory HMMs

Heiga Zen, Yoshihiko Nankaku, Keiichi Tokuda

Department of Computer Science and Engineering, Nagoya Institute of Technology
Gokiso-cho, Showa-ku, Nagoya, 466–8555 Japan

{zen,nankaku,tokuda}@sp.nitech.ac.jp

Abstract

This paper proposes a feature mapping algorithm based on the trajectory GMM or trajectory HMM. Although the GMM or HMM-based feature mapping algorithm works effectively, its conversion quality sometimes degrades due to the inappropriate dynamic characteristics caused by the frame-by-frame conversion. While the use of dynamic features can alleviate this problem, it also introduces an inconsistency between training and mapping. The proposed algorithm can solve this inconsistency while keeping the benefits of the use of dynamic features, and offers an entire sequence-level transformation rather than the frame-by-frame conversion. Experimental results in voice conversion show that the proposed algorithm outperforms the conventional one both in objective and subjective tests.

Index Terms: trajectory HMM, voice conversion

1. Introduction

During these last years, a probabilistic feature mapping algorithm based on Gaussian mixture models (GMMs) or hidden Markov models (HMMs) has been used in various speech applications [1–3]. In this algorithm, GMM or HMM-based mapping functions are estimated using pairs of source and target feature vectors. The resulting mapping functions make it possible to continuously transform of any sample of the source into that of target. Although this algorithm works effectively, its conversion performance is still insufficient. One of major factors which deteriorate the conversion quality could be inappropriate dynamic characteristics caused by the frame-by-frame conversion. To alleviate this problem, Toda et al. introduced the dynamic feature constraints [4]. The use of dynamic features makes it possible to convert a source feature vector sequence to a target one while satisfying the statistics of both static and dynamic characteristics. However, it also introduces an inconsistency between training and mapping: In probabilistic feature mapping, first joint probability density functions (PDFs) of source and target feature vector sequences are modeled by a probabilistic model [5]. Second, conditional PDFs of a target feature vector sequence for a given source one are obtained from the joint PDFs.¹ Finally, the converted target parameter vector sequence is determined so as to minimize its mean-square error (MMSE) [5] or maximize its likelihood (ML) [4] based on the conditional PDFs. In the conventional mapping algorithm using dynamic features, the relationship between static and dynamic features² is ignored on both the joint and conditional PDF estimation stages but utilized on the mapping stage [4, 6]. Ignoring these interdependencies allows inconsistency between the static and dynamic

¹Modeling conditional PDFs directly has also been proposed [6].

²Generally, dynamic features are calculated as regression coefficients from static features of their neighboring frames.

feature vector sequences when the HMM or GMM is used as a generative model in the obvious way. Zen et al. reported that this inconsistency degraded the quality of HMM-based speech synthesis [7]. It suggests that this inconsistency also degrades the performance of GMM or HMM-based feature mapping algorithm with dynamic features.

Recently, a trajectory model, derived from the HMM by imposing the explicit relationships between static and dynamic features, was proposed [7]. This model, named trajectory HMM, can overcome the frame-wise conditional independence assumption of state output probabilities and piecewise constant statistics of the HMM without any additional parameters.

In this paper, we propose a probabilistic feature mapping algorithm based on the trajectory GMM or trajectory HMM. The proposed algorithm can solve the inconsistency introduced by the use of dynamic features while keeping its benefits, and offers an entire sequence-level transformation rather than the frame-by-frame conversion. We evaluate the proposed algorithm in voice conversion, which is one of the most successful applications of probabilistic feature mapping.

This paper is organized as follows: Section 2 describes the joint probability modeling by the trajectory GMM or trajectory HMM. Section 3 derives mapping algorithms based on the MMSE and ML criteria. Section 4 shows objective and subjective experimental results in voice conversion. Concluding remarks and future plans are given in the final section.

2. Probability Density Function

Source and target static feature vector sequences, \mathbf{x} and \mathbf{y} , are written as

$$\mathbf{x} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_T^\top]^\top, \quad \mathbf{y} = [\mathbf{y}_1^\top, \dots, \mathbf{y}_T^\top]^\top, \quad (1)$$

where \mathbf{x}_t and \mathbf{y}_t are the M_x and M_y -dimensional source and target static feature vectors at the t -th frame, respectively, and T is the total number of frames. In this paper, we model the joint probability of \mathbf{x} and \mathbf{y} by a trajectory GMM or a set of trajectory HMMs [7] as follows:

$$P(\mathbf{z} | \lambda) = \sum_{\mathbf{q}} P(\mathbf{q} | \lambda) P(\mathbf{z} | \mathbf{q}, \lambda), \quad (2)$$

$$P(\mathbf{z} | \mathbf{q}, \lambda) = \mathcal{N}(\mathbf{z}; \bar{\mathbf{z}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}), \quad (3)$$

$$P(\mathbf{q} | \lambda) = \begin{cases} \prod_{t=1}^T w_{q_t}, & \text{trajectory GMM} \\ \pi_{q_1} \prod_{t=2}^T a_{q_{t-1}q_t}, & \text{trajectory HMM} \end{cases} \quad (4)$$

$$z = [z_1^\top, \dots, z_T^\top]^\top, \quad z_t = [x_t^\top, y_t^\top]^\top, \quad (5)$$

where λ denotes the set of model parameters, z is the MT -dimensional joint static feature vector sequence, z_t is the the M -dimensional joint static feature vector at the t -th frame, $M = M_x + M_y$ is the dimensionality of joint static feature vectors, $q = \{q_1, \dots, q_T\}$ is a Gaussian component sequence, $q_t \in \{1, \dots, N\}$ is the Gaussian component at the t -th frame, N is the total number of Gaussian components in the model set, w_i is the mixture prior probability of the i -th Gaussian component, π_i is the initial probability of the i -th Gaussian component, and a_{ij} is the transition probability from the i -th Gaussian component to the j -th one. In Eq. (3), \bar{z}_q and P_q are the $MT \times 1$ mean vector and the $MT \times MT$ covariance matrix for q , respectively. They are given by

$$R_q \bar{z}_q = r_q, \quad (6)$$

$$R_q = W^\top \Omega_q W = P_q^{-1}, \quad (7)$$

$$r_q = W^\top \Omega_q \mu_q, \quad (8)$$

$$\mu_q = [\mu_{q_1}^\top, \dots, \mu_{q_T}^\top]^\top, \quad (9)$$

$$\mu_i = [\mu_i^{(x)\top}, \mu_i^{(y)\top}]^\top, \quad (10)$$

$$\Omega_q = \text{diag} [\Omega_{q_1}, \dots, \Omega_{q_T}], \quad (11)$$

$$\Omega_i = \begin{bmatrix} \Omega_i^{(xx)} & \Omega_i^{(xy)} \\ \Omega_i^{(yx)} & \Omega_i^{(yy)} \end{bmatrix}, \quad (12)$$

where μ_i and Ω_i are the $3M \times 1$ mean vector and the $3M \times 3M$ inverse covariance (precision) matrix associated with the i -th Gaussian component, respectively, and W is a $3MT \times MT$ window matrix which appends dynamic features (delta and delta-delta features) to z . Equations (2)–(4) show that the trajectory GMM and trajectory HMM can be interpreted as a MT -dimensional GMM whose mixture weights are given by products of the mixture prior probabilities and the transition probabilities over time, respectively. It should be noted that the intra and inter-frame covariance matrix of Eq. (3) is generally full. Therefore, the trajectory GMM and trajectory HMM can model both intra-frame and inter-frame dependencies between source and target static feature vector sequences without increasing the number of model parameters compared to the GMM and HMM with the same model topology, respectively.

3. Mapping

3.1. Conditional probability

We can rewrite Eq. (3) as follows:

$$P(z | q, \lambda) = \mathcal{N} \left(\begin{bmatrix} x \\ y \end{bmatrix}; \begin{bmatrix} \bar{x}_q \\ \bar{y}_q \end{bmatrix}, \begin{bmatrix} P_q^{(xx)} & P_q^{(xy)} \\ P_q^{(yx)} & P_q^{(yy)} \end{bmatrix} \right), \quad (13)$$

where

$$P_q^{(xx)} = C_q^{(xx)-1}, \quad P_q^{(yy)} = C_q^{(yy)-1}, \quad (14)$$

$$P_q^{(xy)} = -R_q^{(xx)-1} R_q^{(xy)} C_q^{(yy)-1} = P_q^{(yx)\top}, \quad (15)$$

$$C_q^{(xx)} = R_q^{(xx)} - R_q^{(xy)} R_q^{(yy)-1} R_q^{(yx)}, \quad (16)$$

$$C_q^{(yy)} = R_q^{(yy)} - R_q^{(yx)} R_q^{(xx)-1} R_q^{(xy)}, \quad (17)$$

$$R_q^{(xx)} = W^{(x)\top} \Omega_q^{(xx)} W^{(x)}, \quad (18)$$

$$R_q^{(yy)} = W^{(y)\top} \Omega_q^{(yy)} W^{(y)}, \quad (19)$$

$$R_q^{(xy)} = W^{(x)\top} \Omega_q^{(xy)} W^{(y)} = R_q^{(yx)\top}, \quad (20)$$

$$\Omega_q^{(xx)} = \text{diag} [\Omega_{q_1}^{(xx)}, \dots, \Omega_{q_T}^{(xx)}], \quad (21)$$

$$\Omega_q^{(yy)} = \text{diag} [\Omega_{q_1}^{(yy)}, \dots, \Omega_{q_T}^{(yy)}], \quad (22)$$

$$\Omega_q^{(xy)} = \text{diag} [\Omega_{q_1}^{(xy)}, \dots, \Omega_{q_T}^{(xy)}] = \Omega_q^{(yx)\top}, \quad (23)$$

and $W^{(x)}$ and $W^{(y)}$ are $3M_x T \times M_x T$ and $3M_y T \times M_y T$ window matrices, respectively. The mean vectors of Eq. (13), \bar{x}_q and \bar{y}_q , are given as follows:

$$\bar{x}_q = P_q^{(xx)} (r_q^{(x)} - R_q^{(xy)} R_q^{(yy)-1} r_q^{(y)}), \quad (24)$$

$$\bar{y}_q = P_q^{(yy)} (r_q^{(y)} - R_q^{(yx)} R_q^{(xx)-1} r_q^{(x)}), \quad (25)$$

where

$$r_q^{(x)} = W^{(x)\top} (\Omega_q^{(xx)} \mu_q^{(x)} + \Omega_q^{(xy)} \mu_q^{(y)}), \quad (26)$$

$$r_q^{(y)} = W^{(y)\top} (\Omega_q^{(yy)} \mu_q^{(y)} + \Omega_q^{(yx)} \mu_q^{(x)}), \quad (27)$$

$$\mu_q^{(x)} = [\mu_{q_1}^{(x)\top}, \dots, \mu_{q_T}^{(x)\top}]^\top, \quad (28)$$

$$\mu_q^{(y)} = [\mu_{q_1}^{(y)\top}, \dots, \mu_{q_T}^{(y)\top}]^\top. \quad (29)$$

As a result, the conditional PDF of y for given x and λ can be expressed as

$$P(y | x, \lambda) = \sum_{\forall q} P(q | x, \lambda) P(y | x, q, \lambda), \quad (30)$$

$$P(y | x, q, \lambda) = \mathcal{N}(y; \tilde{y}_q, \tilde{P}_q^{(yy)}), \quad (31)$$

$$\tilde{y}_q = \bar{y}_q + P_q^{(yx)} C_q^{(xx)} (x - \bar{x}_q), \quad (32)$$

$$\tilde{P}_q^{(yy)} = P_q^{(yy)} - P_q^{(yx)} C_q^{(xx)} P_q^{(xy)}. \quad (33)$$

3.2. MMSE-based mapping

The MMSE-based mapping [5] determines the estimated target static feature vector sequence \hat{y} as follows:

$$\hat{y} = E[y | x] = \int P(y | x, \lambda) y d y \quad (34)$$

$$= \int \sum_{\forall q} \gamma_q \cdot P(y | x, q, \lambda) y d y \quad (35)$$

$$= \sum_{\forall q} \gamma_q \cdot \tilde{y}_q \quad (36)$$

$$= \sum_{\forall q} \gamma_q \cdot (A_q x + b_q), \quad (37)$$

where $E[\cdot]$ means expectation, γ_q is a posterior probability of q for the given x , and A_q and b_q are given as follows:

$$A_q = P_q^{(yx)} C_q^{(xx)}, \quad (38)$$

$$b_q = \bar{y}_q - P_q^{(yx)} C_q^{(xx)} \bar{x}_q. \quad (39)$$

The estimated static feature vector sequence \hat{y} is defined as the weighted sum of the conditional mean vectors \tilde{y}_q , where the posterior probabilities of q are used as weights.

3.3. ML-based mapping

The ML-based mapping [4] determines the target parameter \hat{y} so as to maximize $P(y | x, \lambda)$ using the EM algorithm. Specifically, we maximize an auxiliary function of a current target static feature vector sequence y and an updated one \hat{y} defined by

$$\mathcal{Q}(y, \hat{y}) = \sum_{\forall q} P(q | x, y, \lambda) \log P(\hat{y}, q | x, \lambda). \quad (40)$$

The mapped target static feature vector sequence maximizing the auxiliary function is given by

$$\hat{y} = \left(\sum_{\forall q} \gamma'_q \cdot \tilde{P}_q^{(yy)-1} \right)^{-1} \sum_{\forall q} \gamma'_q \cdot \tilde{P}_q^{(yy)-1} \tilde{y}_q \quad (41)$$

$$= A'x + b', \quad (42)$$

where γ'_q is a posterior probability of q for given x and y , and A' and b' are as follows:

$$A' = \left(\sum_{\forall q} \gamma'_q \cdot \tilde{P}_q^{(yy)-1} \right)^{-1} \sum_{\forall q} \gamma'_q \cdot \tilde{P}_q^{(yy)-1} P_q^{(yx)} C_q^{(xx)} \quad (43)$$

$$b' = \left(\sum_{\forall q} \gamma'_q \cdot \tilde{P}_q^{(yy)-1} \right)^{-1} \cdot \sum_{\forall q} \gamma'_q \cdot \tilde{P}_q^{(yy)-1} \left(\tilde{y}_q - P_q^{(yx)} C_q^{(xx)} \tilde{x}_q \right). \quad (44)$$

Note that under the Viterbi approximation with the same q , \hat{y} by the ML-based mapping is identical to that by the MMSE-based one ($\hat{y} = \tilde{y}_q$).

3.4. Discussions

The proposed algorithm incorporates the relationship between static and dynamic features explicitly in whole process: joint PDF training, conditional PDF estimation, and feature mapping. Therefore, we can solve the inconsistency between training and mapping. Zen et al. have reported that solving this inconsistency improved the quality of HMM-based speech synthesis [7]. It suggests that the proposed technique may also improve the quality of feature mapping.

The transformation matrices, A_q and A' , are generally full. Hence, the proposed algorithm offers an entire sequence-level transformation rather than the frame-by-frame conversion.

The proposed algorithm is similar to the articulatory-acoustic modeling algorithm proposed by Le et al. [8]. This algorithm models source and target features using two independent data streams:³ interactions between source and target feature vector sequences are represented implicitly by the posterior probability of q . On the other hand, the proposed algorithm models these interactions explicitly via the state output PDFs.

³In this case, the estimated target static vector sequence is determined as $\hat{y} = P_q^{(yy)} r_q^{(y)}$ because $R_q^{(yx)}$ becomes a zero matrix.

4. Experiment

4.1. Experimental conditions

To evaluate the performance of the proposed algorithm, voice conversion experiments on the CMU ARCTIC database were performed. The first 593 sentences uttered by female speakers CLB and SLT were used for training (CLB: source, SLT: target). The last 40 sentences, which were not included in the training data, were used for evaluation. Speech signals were sampled at a rate of 16 kHz and windowed by a 25 ms Blackman window with a 5 ms shift, and then 24-th order mel-cepstral coefficients were obtained by a mel-cepstral analysis technique [9]. For each utterance, $\log F_0$ values were also extracted by ESPS get_f0. Static feature vectors of source and target speech consisted of 25 mel-cepstral coefficients including the zeroth coefficient. To obtain the joint vector sequences, we performed dynamic time warping (DTW) between the source and target vector sequences using an Euclidean norm including both static and dynamic features. In this paper, we calculated dynamic features as $\Delta x_t = 0.5(x_{t+1} - x_{t-1})$, $\Delta^2 x_t = x_{t+1} - 2x_t + x_{t-1}$. Unlike the common DTW method, we needed look-ahead [7] dynamic programming (DP) because dynamic features at each frame depended on the static features of both preceding and succeeding one frames. Using DTW, each target vector sequence was aligned to the corresponding source one. Therefore, the number of frames of the joint vector sequence was equal to that of the source one. Each joint vector was augmented with its dynamic features (150 dimensions in total).

The three-state left-to-right structure with no-skip was used for HMMs and trajectory HMMs to model 41 phonemes including silence (123 states in total). Each state had a single Gaussian distribution for its state output PDF. To compare the performance of each algorithm with the close number of model parameters, we trained a GMM and a trajectory GMM with 128 Gaussian components. The diagonal structure was used for $\{\Omega_i^{(xx)}, \Omega_i^{(xy)}, \Omega_i^{(yx)}, \Omega_i^{(yy)}\}$. After training the GMM and HMMs in the standard way based on the ML criterion, the trajectory GMM and trajectory HMMs were re-estimated by the Viterbi training [7] using the GMM and HMMs as their initial models, respectively. Here, we optimized Cholesky factors of $\{\Omega_i\}$ rather than $\{\Omega_i\}$ themselves because we can easily keep these matrices positive definite. The limited memory BFGS method was used for optimization.

For mapping, we used the Viterbi approximation and did not iterate the EM algorithm.⁴ In this experiment, $\log F_0$ values were converted by a simple linear conversion used in [4]. To avoid the effect of recognition errors, correct transcriptions were given for HMM and trajectory HMM-based mapping. From the converted mel-cepstral coefficients and $\log F_0$ values, a speech waveform was synthesized through the MLSA filter [9, 10] with pulse or noise excitation.

4.2. Experimental results

The mel-cepstral distortion [4] between the target and converted mel-cepstral coefficients in the evaluation set was used as the objective evaluation measure. In [4], the GMM-based mapping algorithm with dynamic features [4] has been shown to obtain superior performance than the GMM-based mapping algorithm without dynamic features [5] on the MOCHA database.

⁴It was found in preliminary experiments that more EM iterations and avoiding Viterbi approximation improved the likelihoods, as measured by the mapping models, but did not improve the quality of mapped speech.

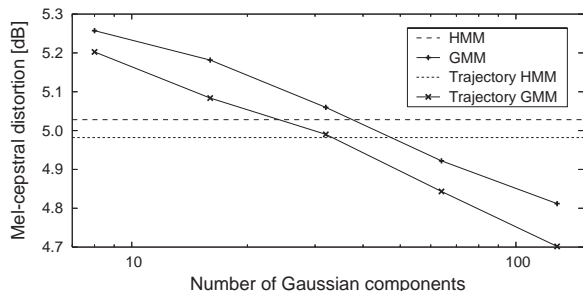


Figure 1: Average mel-cepstral distortions between the target and converted mel-cepstral coefficients in the evaluation set. Without conversion, average mel-cepstral distortion was 6.204.

Therefore, we focus on the performance comparison between the proposed trajectory GMM or trajectory HMM-based mapping algorithm and the GMM or HMM-based mapping with dynamic features in [4]. Figure 1 plots the mel-cepstral distortions. It can be seen from the figure that the proposed algorithm constantly outperformed the conventional one. This is because the inconsistency between training and mapping was solved in the proposed algorithm. Compared with the HMM and trajectory HMM-based mapping, the GMM and trajectory GMM-based ones achieved better objective scores, respectively. We expect that the use of single Gaussian distributions for state output PDFs of the HMMs and trajectory HMMs caused this phenomenon.

Two subjective listening tests were conducted to evaluate the speech quality and speaker similarity. The first test compared the naturalness of converted speech by the mean opinion score (MOS) test method. The second one evaluated the speaker similarity between the analysis-synthesized target speech and converted speech by the differential MOS (DMOS) test method. The subjects were 12 Japanese graduate students, all of whom completed both tests. Fifteen sentences were randomly chosen from the evaluation sentences for each test of each subject. Samples were presented in a random order for each test sentence. Before starting the test, the subjects listened to speech samples of one sentence to become familiar with the task. This sentence was randomly chosen for each subject and excluded from the actual test. In the MOS test, after listening to each test sample, the subjects were asked to assign it a five-point naturalness score (5: natural – 1: poor). In the DMOS test, after listening to the analysis-synthesized target speech and each test sample, the subjects were asked to assign it a five-point similarity score (5: very similar – 1: very far). Both experiments were carried out in a soundproof room using headphones.

Figure 2 plots the experimental results. It can be seen from the figures that the proposed algorithm achieved better subjective scores than the conventional ones in both the MOS and DMOS tests.

5. Conclusions

This paper proposed the probabilistic feature mapping algorithm based on the trajectory GMM or trajectory HMM. This algorithm can solve the inconsistency of the conventional GMM or HMM-based mapping algorithm using dynamic features, and offers the entire sequence-level transformation rather than the frame-by-frame conversion. Experimental results in voice conversion showed that the proposed algorithm outperformed the

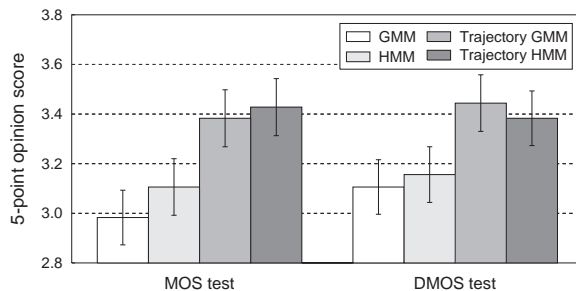


Figure 2: Mean opinion scores and differential mean opinion scores of converted speech by the conventional and proposed algorithms. Error bars show 95% confidence intervals.

conventional GMM or HMM-based one both in the objective and subjective tests.

Future work includes improving the quality of converted speech using global variance [4], investigating the proposed algorithms on other applications, and integrating the DTW procedure into the model definition like DPGMM [11].

6. Acknowledgments

The authors would like to thank Yosuke Uto and Dr. Tomoki Toda for their helpful comments. This work was partly supported by the MEXT e-Society project, the Hori information science promotion foundation, and the Grant-in-Aid for Scientific Research (No. 1880009) of JSPS.

7. References

- [1] Y. Stylianou, O. Cappe, and E. Moulines, "Continuous probabilistic transform for voice conversion," *IEEE Trans. SAP*, vol. 6, no. 2, pp. 131–142, 1998.
- [2] S. Hiroya and M. Honda, "Estimation of articulatory movements from speech acoustics using an HMM-based speech production model," *IEEE Trans. SAP*, vol. 12, no. 2, pp. 175–185, 2004.
- [3] X.-D. Cui, M. Afify, and Y. Gao, "MMSE-based stereo feature stochastic mapping for noise robust speech recognition," in *Proc. ICASSP*, 2008, pp. 4077–4080.
- [4] T. Toda, A. Black, and K. Tokuda, "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory," *IEEE Trans. ASLP*, vol. 15, no. 8, pp. 2222–2235, 2007.
- [5] A. Kain and M. Macon, "Spectral voice conversion for text-to-speech synthesis," in *Proc. ICASSP*, 1998, pp. 285–288.
- [6] K. Richmond, "Trajectory mixture density network with multiple mixtures for acoustic-articulatory inversion," in *Proc. NOLISP*, 2007, pp. 67–70.
- [7] H. Zen, K. Tokuda, and T. Kitamura, "Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic features," *Computer Speech & Language*, vol. 21, no. 1, pp. 153–173, 2007.
- [8] L. Zhang and S. Renals, "Acoustic-articulatory modelling with the trajectory HMM," *IEEE Signal Processing Letters*, vol. 15, pp. 245–248, 2008.
- [9] 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.
- [10] S. Imai, "Cepstral analysis synthesis on the mel frequency scale," in *Proc. ICASSP*, 1983, pp. 93–96.
- [11] Y. Nankaku, K. Nakamura, T. Toda, and K. Tokuda, "Spectral conversion based on statistical models including time-frequency matching," in *Proc. ISCA SSW6*, 2007, pp. 333–338.