

ベイズ的アプローチに基づく HMM 音声合成

南角 吉彦[†] 全 炳河[†] 徳田 恵一[†] 北村 正[†] 益子 貴史^{††}

[†] 名古屋工業大学知能情報システム学科

^{††} 東京工業大学大学院総合理工学研究所

あらまし 本論文では、HMM 音声合成に変分ベイズ法を適用し、ベイズ基準による音声のパラメータ生成アルゴリズムを導出する。HMM 音声合成では、これまで、HMM の学習と音声パラメータの生成に ML 基準が用いられてきた。しかし、ML 推定はモデルパラメータを確定的変数として点推定するため、学習データが十分に得られない場合、モデルの推定精度が低下するという問題がある。提案法は、ML 基準に比べて、音声特徴をモデル化する統計モデルの汎化能力を高めることが可能であり、合成音声の品質改善が期待できる。

キーワード 変分ベイズ法、隠れマルコフモデル、音声合成

A Bayesian Approach to HMM-Based Speech Synthesis

Yoshihiko NANKAKU[†], Heiga ZEN[†], Keiichi TOKUDA[†], Tadashi KITAMURA[†],
and Takashi MASUKO^{††}

[†] Department of Computer Science, Nagoya Inst. of Tech. Gokiso-cho, Showa-ku, Nagoya, 466-8555 Japan

^{††} Interdisciplinary Graduate School of Science and Engineering, Tokyo Inst. of Tech. 4259, Nagatsuta,
Midori-ku, Yokohama, 226-8502 Japan

Abstract In this paper, we applied the variational Bayes method to HMM-based speech synthesis and derived a speech parameter generation algorithm based on the Bayesian criterion. In HMM-based speech synthesis, the ML criterion has been used for training HMM and generating speech parameters. However, the ML estimator produces a point estimate of HMM parameters and the accuracy of estimation may be reduced when little training data is available. By using the proposed algorithm, it is expected that the higher generalization ability of the statistical model for speech features is achieved and leads to improving the quality of synthesized speech.

Key words variational Bayes method, hidden Markov models, speech synthesis

1. ま え が き

コーパスに基づく音声合成では、データを多く集めることにより合成音声の品質を高めることができるが、その一方で限られたデータ量で品質の良い音声合成を実現することが重要な課題である。HMM 音声合成は、メルケプストラムや基本周波数などの音声特徴量を用いた分析合成系により音声合成されるため、波形接続型の音声合成に比べて、少ない学習データで滑らかな音声合成が可能である。また、音声特徴の動的特徴量を考慮したパラメータ生成アルゴリズムにより音声単位境界において歪みの少ない音声パラメータを得ることができる。HMM 音声合成における音声品質の上限は、ポストフィルタの使用などの方法もあるが、基本的には音声特徴量を用いた分析合成系に依存している。しかし、学習データが十分に得られない場合、統計モデルの精度が低下する可能性があり、統計モデルの汎化能力を高めることにより、合成音声の品質を改善する余地があると考えられる。また、この点において、統計的学習の様々な

研究成果を応用できることが、HMM 音声合成の大きな利点であると言える。

HMM 音声合成では、これまで、HMM の学習や音声パラメータの生成に ML(Maximum Likelihood) 基準が用いられてきた。しかし、ML 推定は、モデルパラメータを確定的変数として点推定するため、学習データが十分に得られない場合、モデルの推定精度が低下するという問題がある。これに対し、ベイズ学習では、モデルパラメータを確率変数としてとらえ、学習データに対する事後確率分布を推定する。事後分布を用いて、すべての可能なモデルパラメータを考慮することにより、学習データ量が少ない場合においても、高い汎化性能が得られる。しかし、事後確率分布の推定は困難な積分計算を伴うため、何らかの近似を必要とする。従来、マルコフ連鎖モンテカルロ法などのサンプリング手法が用いられてきたが、計算量が膨大となり、大規模な問題への適用は困難であった。この問題に対し、近年、変分ベイズ法 [1], [2] が提案され、ML 基準により学習されていた様々なモデルにベイズ学習が適用可能となり、

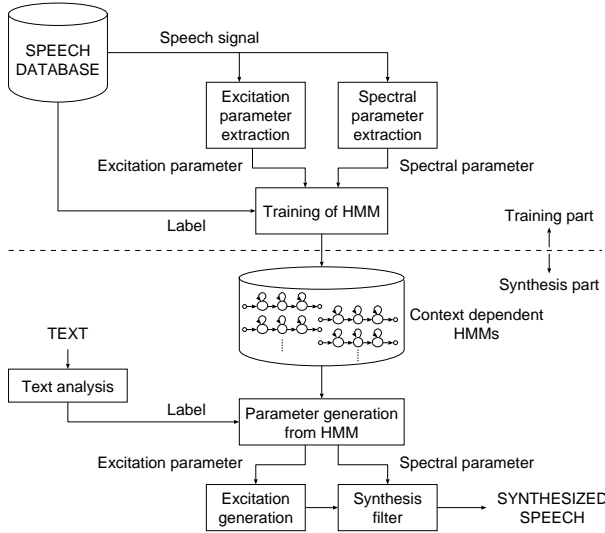


図1 HMMに基づいた音声合成システムの構成

HMMに基づいた音声認識にも適用されている[3]。本研究では、HMM音声合成に変分ベイズ法を適用し、ベイズ基準による音声パラメータの生成アルゴリズムを導出する。また、コンテキストクラスタリングにおいても、ベイズ基準を適用することにより、データ量に応じた最適なモデル構造を決定する。

2. ML 基準による HMM 音声合成

まず、ML 基準による音声合成システムについて説明する。図1にシステムの概要を示す[4]。学習部では、音声データベースからスペクトルパラメータと励振源パラメータを抽出し、音素などの音声単位で HMM を学習する。本論文では、スペクトルパラメータとしてメルケプストラム、励振源パラメータとして基本周波数を用いることとする。合成部では、合成すべきテキストに対応するラベル列に従い、音素 HMM を連結した文章 HMM を作成し、この HMM に対する尤度が最大となる音声パラメータを生成する。基本周波数に従って励振源モデルから励振源を生成し、生成したメルケプストラムから構築された合成フィルタ (MLSA フィルタ) を励振することにより合成音声を得られる。

2.1 隠れマルコフモデルの定義

音声合成に用いる隠れマルコフモデルを定義する。音声特徴量の観測ベクトル列 $O = (o_1, o_2, \dots, o_T)$ と状態遷移を表す隠れ変数列 $Z = (z_1, z_2, \dots, z_T)$, $z_t \in \{1, \dots, N\}$ から成る完全データ $\{O, Z\}$ の対数尤度は次式で定義される。

$$\log P(O, Z|\Lambda) = \sum_{t=1}^{T-1} \sum_{i=1}^N \sum_{j=1}^N z_t^i z_{t+1}^j \log a_{ij} + \sum_{t=1}^T \sum_{i=1}^N z_t^i \log \mathcal{N}(o_t | \mu_i, S_i^{-1}) \quad (1)$$

ただし、 $z_t^i = \delta(z_t, i)$ であり、 $z_t = i$ のとき 1、それ以外は 0 とする。 N は HMM の状態数、HMM のパラメータ $\Lambda = \{a, b\}$ は、状態遷移確率 $a = \{a_{ij}\}_{i,j=1}^N$, $a_{ij} = P(z_{t+1} = j | z_t = i)$ と出力確率 $b = \{b_i(o_t)\}_{i=1}^N$, $b_i(o_t) = P(o_t | z_t = i) = \mathcal{N}(o_t | \mu_i, S_i^{-1})$ により与えられる。ただし、 $\mathcal{N}(\cdot | \mu_i, S_i^{-1})$ は、 μ_i と S_i を、そ

れぞれ、平均ベクトルと共分散行列の逆行列として持つガウス分布である。出力確率分布としては、多次元ガウス分布の重み付き和で表される多次元ガウス混合分布が用いられることが多いが、ここでは、簡単のため単一の多次元ガウス分布を仮定する。

2.2 ML 基準によるパラメータ生成

ML 基準による音声パラメータの生成では、動的特徴量を考慮することにより、状態遷移やモデル接続部分で歪みの少ない滑らかな音声パラメータを生成することができる[5],[6]。以下にそのアルゴリズムを説明する。まず、学習部では与えられた学習データ O の尤度を最大にするモデルパラメータ $\Lambda^{(ML)}$ を推定する。学習データ O は、音声認識などで良く用いられるように、音声パラメータの静的特徴量に動的特徴量を連結したベクトルを特徴ベクトルとして用いる。

$$\Lambda^{(ML)} = \arg \max_{\Lambda} P(O|S, \Lambda) \quad (2)$$

ここで、 S は学習データのラベル列である。合成部では、推定されたモデルから尤度最大となる音声パラメータ $x = [x_1^T, x_2^T, \dots, x_{T_x}^T]^T$ を生成する。

$$x^{(ML)} = \arg \max_x P(Wx|s, \Lambda^{(ML)}) \quad (3)$$

ここで、 s は合成音声のラベル列である。また、 W は学習データと同様に、音声の静的特徴量 x に動的特徴量を加える行列であり、 Wx の時刻 t における要素は、次式で表される。

$$(Wx)_t = [x_t^T, \Delta x_t^T, \Delta^2 x_t^T]^T \quad (4)$$

また、動的特徴ベクトル Δx_t , $\Delta^2 x_t$ は、静的特徴ベクトル x_t から

$$\Delta x_t = \sum_{\tau=-L^{(1)}}^{L^{(1)}} w_1(\tau) x_{t+\tau} \quad (5)$$

$$\Delta^2 x_t = \sum_{\tau=-L^{(2)}}^{L^{(2)}} w_2(\tau) x_{t+\tau} \quad (6)$$

により計算されるものとする。ただし、 $w_1(\tau)$, $w_2(\tau)$ は動的特徴量を計算するための重み係数である。 x_t が D 次元とすれば、 x , Wx は、それぞれ、 TD 次元、 $3TD$ 次元である。 W は、 $3TD \times TD$ の行列であり、1部の要素に係数 1, $w_1(\tau)$, $w_2(\tau)$ をもち、他の多くの要素は 0 とする。

式(3)において、 $P(Wx|s, \Lambda^{(ML)})$ は隠れ変数を含むため、すべての隠れ変数列を考慮する必要があるが、ここでは簡単のため、あらかじめ最適な状態系列が与えられると仮定する^(注1)。

$$x^{(ML)} = \arg \max_x \sum_q P(Wx|q, \Lambda^{(ML)}) P(q|s, \Lambda^{(ML)}) \simeq \arg \max_x P(Wx|q', \Lambda^{(ML)}) \quad (7)$$

(注1): 式(3)を EM アルゴリズムのような繰返しアルゴリズムで解く手法[7]が提案されている。また、隠れ変数である状態系列を確率モデルで最適化する場合、ここで定義した状態遷移確率では不十分であり、状態継続長モデルなどの学習が必要となる。

ここで, $q = (q_1, q_2, \dots, q_{T_x})$ は x の状態系列, q' は与えられた最適な状態系列を表す. 状態系列 q' が与えられたときの x の尤度 $P(\mathbf{W}x|q', \Lambda^{(ML)})$ は, 高次元のガウス分布と見なすことができる.

$$\begin{aligned} \log P(\mathbf{W}x|q', \Lambda) &= -\frac{DT}{2} \log(2\pi) \\ &+ \frac{1}{2} \log |S| - \frac{1}{2} (\mathbf{W}x - \boldsymbol{\mu})^\top S (\mathbf{W}x - \boldsymbol{\mu}) \end{aligned} \quad (8)$$

ただし, $\boldsymbol{\mu}$, S は, 状態系列 q' に従って並べたガウス分布の平均ベクトルと共分散行列を一つにまとめたものである.

$$\boldsymbol{\mu} = [\boldsymbol{\mu}_{z_1}^\top, \boldsymbol{\mu}_{z_2}^\top, \dots, \boldsymbol{\mu}_{z_T}^\top]^\top \quad (9)$$

$$S = \text{diag}[S_{z_1}, S_{z_2}, \dots, S_{z_T}] \quad (10)$$

もし, 動的特徴量を用いない場合, 尤度を最大にする音声パラメータは平均ベクトル列と等しくなり, 滑らかな音声パラメータは得られない. しかし, 動的特徴量を考慮することにより, $P(\mathbf{W}x|q', \Lambda)$ を最大にする x は,

$$\frac{\partial \log P(\mathbf{W}x|q', \Lambda)}{\partial x} = \mathbf{0} \quad (11)$$

とおくことによって得られる線形方程式

$$\mathbf{W}^\top S \mathbf{W} x^{(ML)} = \mathbf{W}^\top S \boldsymbol{\mu} \quad (12)$$

により定められる. $\mathbf{W}^\top S \mathbf{W}$ は, $TD \times TD$ の巨大な行列であるが, コレスキー分解あるいは QR 分解を用いて少ない演算量で解くことができる [8], [9]. また, 時間方向に再帰的な形で計算を行うアルゴリズム [10] も提案されている.

3. ベイズ基準による HMM 音声合成

3.1 変分法による予測分布の近似

ML 基準による音声合成では, モデルパラメータを点推定するのに対し, ベイズ音声合成では, 学習データの事後確率分布を計算し, すべてのモデルパラメータを重み付けした予測分布から音声パラメータが生成される. ベイズ基準による合成音声は次式で得られる.

$$\begin{aligned} x^{(Bayes)} &= \arg \max_{\mathbf{x}} P(\mathbf{W}x|s, \mathbf{O}, S) \\ &= \arg \max_{\mathbf{x}} P(\mathbf{W}x, \mathbf{O}|s, S) \end{aligned} \quad (13)$$

ここで, $P(\mathbf{W}x|s, \mathbf{O}, S)$ は合成音声の予測分布である. また, 動的特徴量を含む合成パラメータ $\mathbf{W}x$ と学習データ \mathbf{O} の周辺尤度は,

$$\begin{aligned} P(\mathbf{W}x, \mathbf{O}|s, S) &= \sum_q \sum_Z \int P(\mathbf{W}x, \mathbf{q}, \mathbf{O}, \mathbf{Z}, \Lambda|s, S) d\Lambda \\ &= \sum_q \sum_Z \int P(\mathbf{W}x, \mathbf{q}|s, \Lambda) P(\mathbf{O}, \mathbf{Z}|S, \Lambda) P(\Lambda) d\Lambda \end{aligned} \quad (14)$$

と書くことができる. ここで, $P(\Lambda)$ は事前確率であり, あらかじめ得られる学習対象に関する知識を与えることができる. また, $P(\mathbf{W}x, \mathbf{q}|s, \Lambda)$ と $P(\mathbf{O}, \mathbf{Z}|S, \Lambda)$ は, それぞれ, 合成データと学習データの完全な尤度を表す. 対数周辺尤度

$\log P(\mathbf{W}x, \mathbf{O}|s, S)$ に対し, 任意の分布 $Q(q, \mathbf{Z}, \Lambda)$ を導入し, Jensen の不等式を用いて, 下限 \mathcal{F} を定義する.

$$\begin{aligned} \log P(\mathbf{W}x, \mathbf{O}|s, S) &= \log \sum_q \sum_Z \int P(\mathbf{W}x, \mathbf{O}, \mathbf{q}, \mathbf{Z}, \Lambda|s, S) d\Lambda \\ &= \sum_q \sum_Z \int \log Q(\mathbf{q}, \mathbf{Z}, \Lambda) \frac{P(\mathbf{W}x, \mathbf{O}, \mathbf{q}, \mathbf{Z}, \Lambda|s, S)}{Q(\mathbf{q}, \mathbf{Z}, \Lambda)} d\Lambda \\ &\geq \sum_q \sum_Z \int Q(\mathbf{q}, \mathbf{Z}, \Lambda) \log \frac{P(\mathbf{W}x, \mathbf{O}, \mathbf{q}, \mathbf{Z}, \Lambda|s, S)}{Q(\mathbf{q}, \mathbf{Z}, \Lambda)} d\Lambda \\ &= \left\langle \log \frac{P(\mathbf{W}x, \mathbf{O}, \mathbf{q}, \mathbf{Z}, \Lambda|s, S)}{Q(\mathbf{q}, \mathbf{Z}, \Lambda)} \right\rangle_{Q(\mathbf{q}, \mathbf{Z}, \Lambda)} \\ &= \mathcal{F} \end{aligned} \quad (15)$$

ここで, 下限 \mathcal{F} は, $Q(\mathbf{q}, \mathbf{Z}, \Lambda)$ を変関数とする汎関数であり, 不等式における $\log P(\mathbf{W}x, \mathbf{O}|s, S)$ と \mathcal{F} の差は, $Q(\mathbf{q}, \mathbf{Z}, \Lambda)$ と $P(\mathbf{q}, \mathbf{Z}, \Lambda|\mathbf{W}x, \mathbf{O}, s, S)$ の KL-divergence で表される.

$$\begin{aligned} \mathcal{F} &= \left\langle \log \frac{P(\mathbf{q}, \mathbf{Z}, \Lambda|\mathbf{W}x, \mathbf{O}, s, S) P(\mathbf{W}x, \mathbf{O}|s, S)}{Q(\mathbf{q}, \mathbf{Z}, \Lambda)} \right\rangle \\ &= \log P(\mathbf{W}x, \mathbf{O}|s, S) \\ &\quad - KL(Q(\mathbf{q}, \mathbf{Z}, \Lambda) || P(\mathbf{q}, \mathbf{Z}, \Lambda|\mathbf{W}x, \mathbf{O}, s, S)) \end{aligned} \quad (16)$$

ここで, $P(\mathbf{W}x, \mathbf{O}|s, S)$ は $\mathbf{q}, \mathbf{Z}, \Lambda$ に関して定数であるため, $Q(\mathbf{q}, \mathbf{Z}, \Lambda)$ に関して \mathcal{F} を最大化することは KL-divergence を最小にすることと等価である. また, $Q(\mathbf{q}, \mathbf{Z}, \Lambda) = P(\mathbf{q}, \mathbf{Z}, \Lambda|\mathbf{W}x, \mathbf{O}, s, S)$ のとき, KL-divergence は 0 となり, $\log P(\mathbf{W}x, \mathbf{O}|s, S) = \mathcal{F}$ となる. よって, 変関数 $Q(\mathbf{q}, \mathbf{Z}, \Lambda)$ に, 積分計算が可能となるような拘束条件を与えた上で, 下限 \mathcal{F} を最大化することにより, 最適な $P(\mathbf{W}x, \mathbf{O}|s, S)$ の近似分布を求めることができる. ここでは, 以下の拘束条件を与える.

$$Q(\mathbf{q}, \mathbf{Z}, \Lambda) = Q(\mathbf{q})Q(\mathbf{Z})Q(\Lambda) \quad (17)$$

このとき, 下限 \mathcal{F} は次式で表される.

$$\begin{aligned} \mathcal{F} &= \langle \log P(\mathbf{W}x, \mathbf{q}|s, \Lambda) \rangle_{Q(\mathbf{q})Q(\Lambda)} \\ &\quad + \langle \log P(\mathbf{O}, \mathbf{Z}|S, \Lambda) \rangle_{Q(\mathbf{Z})Q(\Lambda)} + \log P(\Lambda) \\ &\quad + E[Q(\mathbf{q})] + E[Q(\mathbf{Z})] + E[Q(\Lambda)] \end{aligned} \quad (18)$$

ただし, $E[\cdot]$ はエントロピーを表す. まず, $Q(\mathbf{Z})$ について, \mathcal{F} を最大化する. 変関数 $Q(\mathbf{Z})$ は, 式 (15) の不等式が成り立つため, $\sum_Z Q(\mathbf{Z}) = 1$ を満たす必要があり, ラグランジュ係数 λ_Z を用いて, 次式を最大化する.

$$\begin{aligned} \mathcal{F}_Z &= \mathcal{F} - \lambda_Z \left(\sum_Z Q(\mathbf{Z}) - 1 \right) \\ &= \sum_Z \left(f_Z(\mathbf{Z}, Q(\mathbf{Z})) - \lambda_Z Q(\mathbf{Z}) \right) + Const \end{aligned} \quad (19)$$

ただし,

$$\begin{aligned} f_Z(\mathbf{Z}, Q(\mathbf{Z})) &= Q(\mathbf{Z}) \langle \log P(\mathbf{O}, \mathbf{Z}|S, \Lambda) \rangle_{Q(\Lambda)} - Q(\mathbf{Z}) \log Q(\mathbf{Z}) \end{aligned} \quad (20)$$

ここで、 \mathcal{F}_Z を最大にする関数 $Q(\mathbf{Z})$ を得るため、変分法を適用する．変分法では、次式で表される変分導関数 $\delta\mathcal{F}_Z$ を 0 と置くことにより、関数の最適化問題を解く．

$$\begin{aligned}\delta\mathcal{F}_Z &= \sum_{\mathbf{Z}} \frac{\partial}{\partial Q(\mathbf{Z})} \left\{ f_Z(\mathbf{Z}, Q(\mathbf{Z})) - \lambda_Z Q(\mathbf{Z}) \right\} \delta Q(\mathbf{Z}) \\ &= 0\end{aligned}\quad (21)$$

ただし、 $\delta Q(\mathbf{Z})$ は $Q(\mathbf{Z})$ の変分であり、 $Q(\mathbf{Z})$ の微細なずれを表す任意の関数である．ここで、 \mathcal{F}_Z は、 $Q(\mathbf{Z})$ の微分項を含まない単純な形であるため、すべての \mathbf{Z} において、次式が成り立てば良い．

$$\frac{\partial}{\partial Q(\mathbf{Z})} \left\{ f_Z(\mathbf{Z}, Q(\mathbf{Z})) - \lambda_Z Q(\mathbf{Z}) \right\} = 0 \quad (22)$$

よって、最適な近似分布は、

$$Q(\mathbf{Z}) = C_Z \exp \langle \log P(\mathbf{O}, \mathbf{Z} | S, \Lambda) \rangle_{Q(\Lambda)} \quad (23)$$

となる．ただし、 C_Z は $\sum_{\mathbf{Z}} Q(\mathbf{Z}) = 1$ を満たすための正規化定数である．同様の導出方法により、

$$Q(\mathbf{q}) = C_q \exp \langle \log P(\mathbf{W}\mathbf{x}, \mathbf{q} | s, \Lambda) \rangle_{Q(\Lambda)} \quad (24)$$

$$\begin{aligned}Q(\Lambda) &= C_\Lambda P(\Lambda) \exp \langle \log P(\mathbf{O}, \mathbf{Z} | S, \Lambda) \rangle_{Q(\mathbf{Z})} \\ &\quad \times \exp \langle \log P(\mathbf{W}\mathbf{x}, \mathbf{q} | s, \Lambda) \rangle_{Q(\mathbf{q})}\end{aligned}\quad (25)$$

が得られる．これらの近似分布は相互に関係しているため、個別に更新を繰り返すことにより、結果的に \mathcal{F} を極大に導くことができる．また、各更新において、 \mathcal{F} は必ず増加するため収束性が保証されている．

さらに、 $Q(\Lambda)$ に関して、 $\mathbf{a}_i = \{a_{ij}\}_{j=1}^N$ 、 $\mathbf{b}_i = \{\boldsymbol{\mu}_i, \mathbf{S}_i\}$ とし、事前確率を $P(\Lambda) = \prod_{i=1}^N P(\mathbf{a}_i) \prod_{i=1}^N P(\mathbf{b}_i)$ と設定すると、直ちに、

$$Q(\Lambda) = \prod_{i=1}^N Q(\mathbf{a}_i) \prod_{i=1}^N Q(\mathbf{b}_i) \quad (26)$$

$$\begin{aligned}Q(\mathbf{a}_i) &= C_{\mathbf{a}_i} P(\mathbf{a}_i) \exp \left\{ \sum_{j=1}^N \sum_{t=1}^{T_x-1} \langle z_t^i z_{t+1}^j \rangle \log a_{ij} \right\} \\ &\quad \times \exp \left\{ \sum_{j=1}^N \sum_{t=1}^{T_x-1} \langle q_t^i q_{t+1}^j \rangle \log a_{ij} \right\}\end{aligned}\quad (27)$$

$$\begin{aligned}Q(\mathbf{b}_i) &= C_{\mathbf{b}_i} P(\mathbf{b}_i) \exp \left\{ \sum_{t=1}^T \langle z_t^i \rangle \log \mathcal{N}(\mathbf{o}_t | \boldsymbol{\mu}_i, \mathbf{S}_i^{-1}) \right\} \\ &\quad \times \exp \left\{ \sum_{t=1}^{T_x} \langle q_t^i \rangle \log \mathcal{N}((\mathbf{W}\mathbf{x})_t | \boldsymbol{\mu}_i, \mathbf{S}_i^{-1}) \right\}\end{aligned}\quad (28)$$

となる．ただし、

$$\langle z_t^i \rangle = \sum_{\mathbf{Z}} Q(\mathbf{Z}) z_t^i, \quad \langle z_t^i z_{t+1}^j \rangle = \sum_{\mathbf{Z}} Q(\mathbf{Z}) z_t^i z_{t+1}^j \quad (29)$$

$$\langle q_t^i \rangle = \sum_{\mathbf{q}} Q(\mathbf{q}) q_t^i, \quad \langle q_t^i q_{t+1}^j \rangle = \sum_{\mathbf{q}} Q(\mathbf{q}) q_t^i q_{t+1}^j \quad (30)$$

である．また、 $C_{\mathbf{a}_i}$ と $C_{\mathbf{b}_i}$ は、それぞれ $Q(\mathbf{a}_i)$ と $Q(\mathbf{b}_i)$ を正規化する定数であり、 $C_\Lambda = \prod_{i=1}^N C_{\mathbf{a}_i} \prod_{i=1}^N C_{\mathbf{b}_i}$ である． $Q(\mathbf{Z})$ 、

$Q(\mathbf{q})$ についても、これらの分布を用いて、

$$\begin{aligned}Q(\mathbf{Z}) &= C_Z \prod_{t=1}^{T-1} \exp(\log a_{z_t z_{t+1}})_{Q(\mathbf{a}_{z_t})} \\ &\quad \times \prod_{t=1}^T \exp(\log \mathcal{N}(\mathbf{o}_t | \boldsymbol{\mu}_{z_t}, \mathbf{S}_{z_t}^{-1}))_{Q(\mathbf{b}_{z_t})}\end{aligned}\quad (31)$$

$$\begin{aligned}Q(\mathbf{q}) &= C_q \prod_{t=1}^{T_x-1} \exp(\log a_{q_t q_{t+1}})_{Q(\mathbf{a}_{q_t})} \\ &\quad \times \prod_{t=1}^{T_x} \exp(\log \mathcal{N}((\mathbf{W}\mathbf{x})_t | \boldsymbol{\mu}_{q_t}, \mathbf{S}_{q_t}^{-1}))_{Q(\mathbf{b}_{q_t})}\end{aligned}\quad (32)$$

と書ける．ここで、 $Q(\mathbf{Z})$ 、 $Q(\mathbf{q})$ は HMM の尤度関数と同じ形となっており、式 (29)、(30) の期待値は、Forward-Backward アルゴリズムを用いて計算することができる．

3.2 事前分布

モデルパラメータの事前確率 $P(\Lambda)$ として、共役事前分布を用いる．共役事前分布とは、事前分布と事後分布が同じ分布族となる事前分布のことである．尤度関数が式 (1) の HMM では、遷移確率は Dirichlet 分布、出力確率は Gauss-Wishart 分布となる．つまり、事後分布の近似分布である $Q(\mathbf{a}_i)$ と $Q(\mathbf{b}_i)$ がそれぞれ Dirichlet 分布と Gauss-Wishart 分布となる．以下に事前分布を定義する．

$$P(\mathbf{a}_i) = \mathcal{D}(\mathbf{a}_i | \phi_i) = \frac{\Gamma\left(\sum_{j=1}^N \phi_{ij}\right)}{\prod_{j=1}^N \Gamma(\phi_{ij})} \prod_{j=1}^N a_{ij}^{\phi_{ij}-1} \quad (33)$$

$$P(\mathbf{b}_i) = \mathcal{N}(\boldsymbol{\mu}_i | \boldsymbol{\nu}_i, (\xi_i \mathbf{S}_i)^{-1}) \mathcal{W}(\mathbf{S}_i | \eta_i, \mathbf{B}_i) \quad (34)$$

$$\begin{aligned}\mathcal{N}(\boldsymbol{\mu}_i | \boldsymbol{\nu}_i, (\xi_i \mathbf{S}_i)^{-1}) &= \\ C_{N_i} |\mathbf{S}_i|^{\frac{D}{2}} \exp \left\{ -\frac{1}{2} \text{Tr}(\xi_i \mathbf{S}_i (\boldsymbol{\mu}_i - \boldsymbol{\nu}_i)(\boldsymbol{\mu}_i - \boldsymbol{\nu}_i)^\top) \right\}\end{aligned}\quad (35)$$

$$\begin{aligned}\mathcal{W}(\mathbf{S}_i | \eta_i, \mathbf{B}_i) &= \\ C_{W_i} |\mathbf{S}_i|^{\frac{\eta_i - D - 1}{2}} \exp \left\{ -\frac{1}{2} \text{Tr}(\mathbf{S}_i \mathbf{B}_i) \right\}\end{aligned}\quad (36)$$

$$C_{N_i} = (2\pi)^{-\frac{D}{2}} \xi_i^{\frac{D}{2}} \quad (37)$$

$$C_{W_i} = \frac{|\mathbf{B}_i|^{\frac{\eta_i}{2}}}{2^{\frac{\eta_i D}{2}} \pi^{\frac{D(D-1)}{4}} \prod_{d=1}^D \Gamma\left(\frac{\eta_i + 1 - d}{2}\right)} \quad (38)$$

とする．ここで、 $\Gamma(\cdot)$ はガンマ関数、 D は特徴ベクトルの次元数である．事前分布を表すパラメータをまとめると、 $\{\phi_{ij}, \xi_i, \eta_i, \boldsymbol{\nu}_i, \mathbf{B}_i\}_{i,j=1}^N$ となる．また、共役事前分布を用いているため、事後分布も同様のパラメータセットで表すことができ、 $\{\bar{\phi}_{ij}, \bar{\xi}_i, \bar{\eta}_i, \bar{\boldsymbol{\nu}}_i, \bar{\mathbf{B}}_i\}_{i,j=1}^N$ で表すこととする．

3.3 事後分布パラメータの更新

具体的な更新式を以下に示す．まず、事前分布のパラメータと $Q(\mathbf{Z})$ 、 $Q(\mathbf{q})$ から計算される期待値を以下に定義する．

$$\bar{N}_i = \sum_{t=1}^T \langle z_t^i \rangle + \sum_{t=1}^{T_x} \langle q_t^i \rangle \quad (39)$$

$$\bar{N}_{ij} = \sum_{t=1}^{T-1} \langle z_t^i z_{t+1}^j \rangle + \sum_{t=1}^{T_x-1} \langle q_t^i q_{t+1}^j \rangle \quad (40)$$

$$\bar{\mathbf{o}}_i = \frac{1}{\bar{N}_i} \left(\sum_{t=1}^T \langle z_t^i \rangle \mathbf{o}_t + \sum_{t=1}^{T_x} \langle q_t^i \rangle (\mathbf{W}\mathbf{x})_t \right) \quad (41)$$

$$\bar{C}_i = \frac{1}{\bar{N}_i} \left(\sum_{t=1}^T \langle z_t^i \rangle (\mathbf{o}_t - \bar{\mathbf{o}}_i) (\mathbf{o}_t - \bar{\mathbf{o}}_i)^\top + \sum_{t=1}^{T_x} \langle q_t^i \rangle ((\mathbf{W}\mathbf{x})_t - \bar{\mathbf{o}}_i) ((\mathbf{W}\mathbf{x})_t - \bar{\mathbf{o}}_i)^\top \right) \quad (42)$$

これらの期待値を用いて，モデルパラメータの事後確率分布 $Q(\Lambda)$ は，次式で更新される．

$$\bar{\phi}_{ij} = \phi_{ij} + \bar{N}_{ij}, \quad \bar{\xi}_i = \xi_i + \bar{N}_i, \quad \bar{\eta}_i = \eta_i + \bar{N}_i \quad (43)$$

$$\bar{\nu}_i = \frac{\bar{N}_i \bar{\mathbf{o}}_i + \xi_i \boldsymbol{\nu}_i}{\bar{N}_i + \xi_i} \quad (44)$$

$$\bar{B}_i = \bar{N}_i \bar{C}_i + \mathbf{B}_i + \frac{\bar{N}_i \xi_i}{\bar{N}_i + \xi_i} (\bar{\mathbf{o}}_i - \boldsymbol{\nu}_i) (\bar{\mathbf{o}}_i - \boldsymbol{\nu}_i)^\top \quad (45)$$

また， $Q(\mathbf{Z})$ ， $Q(\mathbf{q})$ の更新に用いられる期待値は，次式で表される．

$$\langle \log a_{ij} \rangle_{Q(a_i)} = \Psi(\bar{\phi}_{ij}) - \Psi \left(\sum_{k=1}^N \bar{\phi}_{kj} \right) \quad (46)$$

$$\begin{aligned} & \langle \log \mathcal{N}(\mathbf{o}_t | \boldsymbol{\mu}_i, \mathbf{S}_i^{-1}) \rangle_{Q(b_i)} \\ &= -\frac{1}{2} \left\{ D \log \pi + \frac{D}{\bar{\xi}_i} - \sum_{d=1}^D \Psi \left(\frac{\bar{\eta}_i + 1 - d}{2} \right) + \log |\bar{B}_i| \right. \\ & \quad \left. + \text{Tr} \left\{ \bar{\eta}_i \bar{B}_i^{-1} (\mathbf{o}_t - \bar{\boldsymbol{\nu}}_i) (\mathbf{o}_t - \bar{\boldsymbol{\nu}}_i)^\top \right\} \right\} \quad (47) \end{aligned}$$

$$\begin{aligned} & \langle \log \mathcal{N}((\mathbf{W}\mathbf{x})_t | \boldsymbol{\mu}_i, \mathbf{S}_i^{-1}) \rangle_{Q(b_i)} \\ &= -\frac{1}{2} \left\{ D \log \pi + \frac{D}{\bar{\xi}_i} - \sum_{d=1}^D \Psi \left(\frac{\bar{\eta}_i + 1 - d}{2} \right) + \log |\bar{B}_i| \right. \\ & \quad \left. + \text{Tr} \left\{ \bar{\eta}_i \bar{B}_i^{-1} ((\mathbf{W}\mathbf{x})_t - \bar{\boldsymbol{\nu}}_i) ((\mathbf{W}\mathbf{x})_t - \bar{\boldsymbol{\nu}}_i)^\top \right\} \right\} \quad (48) \end{aligned}$$

となる．ここで， $\Psi(\cdot)$ は digamma 関数である．これらの期待値を用いて，Forward-Backward アルゴリズムを実行することにより， $\bar{N}_i, \bar{N}_{ij}, \bar{\mathbf{o}}_i, \bar{C}_i$ を更新する．

3.4 ベイズ基準による音声パラメータの生成

変分法による近似事後分布の最適化を行うことにより，下限 \mathcal{F} は周辺尤度 $P(\mathbf{W}\mathbf{x}, \mathbf{O} | s, S)$ を良く近似することができる．よって，最適な音声パラメータは， \mathcal{F} を最大化することで得られる．下限 \mathcal{F} を \mathbf{x} について，

$$\begin{aligned} \frac{\partial \mathcal{F}}{\partial \mathbf{x}} &= \frac{\partial}{\partial \mathbf{x}} \langle \log P(\mathbf{W}\mathbf{x} | \mathbf{q}, \Lambda) P(\mathbf{q} | s, \Lambda) \rangle_{Q(\mathbf{q})Q(\Lambda)} \\ &= \left\langle \frac{\partial}{\partial \mathbf{x}} \log P(\mathbf{W}\mathbf{x} | \mathbf{q}, \Lambda) \right\rangle_{Q(\mathbf{q})Q(\Lambda)} = \mathbf{0} \quad (49) \end{aligned}$$

とおくことにより，

$$\mathbf{W}^\top \langle S \rangle \mathbf{W} \mathbf{x}^{(Bayes)} = \mathbf{W}^\top \langle S \boldsymbol{\mu} \rangle \quad (50)$$

が得られる．ただし，

$$\begin{aligned} \langle S \rangle &= \sum_{\mathbf{q}} \int Q(\mathbf{q}) Q(\Lambda) S d\Lambda \\ &= \text{diag}[\langle S_{q_1} \rangle, \langle S_{q_2} \rangle, \dots, \langle S_{q_T} \rangle] \quad (51) \end{aligned}$$

である．また，行列の各成分は独立に期待値計算ができるため，

$$\langle S_{q_t} \rangle = \sum_{i=1}^N \langle q_t^i \rangle \int Q(\mathbf{b}_i) S_i d\Lambda = \sum_{i=1}^N \langle q_t^i \rangle \bar{\eta}_i \bar{B}_i^{-1} \quad (52)$$

となる．同様に，

$$\begin{aligned} \langle S \boldsymbol{\mu} \rangle &= \sum_{\mathbf{q}} \int Q(\mathbf{q}) Q(\Lambda) S \boldsymbol{\mu} d\Lambda \\ &= [\langle S_{q_1} \boldsymbol{\mu}_{q_1} \rangle^\top, \langle S_{q_2} \boldsymbol{\mu}_{q_2} \rangle^\top, \dots, \langle S_{q_T} \boldsymbol{\mu}_{q_T} \rangle^\top]^\top \quad (53) \end{aligned}$$

$$\begin{aligned} \langle S_{q_t} \boldsymbol{\mu}_{q_t} \rangle &= \sum_{i=1}^N \langle q_t^i \rangle \int Q(\mathbf{b}_i) S_i \boldsymbol{\mu}_i d\boldsymbol{\mu}_i dS_i \\ &= \sum_{i=1}^N \langle q_t^i \rangle \bar{\eta}_i \bar{B}_i^{-1} \bar{\boldsymbol{\nu}}_i \quad (54) \end{aligned}$$

である．式 (50) は，ML 基準によるパラメータ生成 (式 (12)) と同じ形をしており， S ， $S \boldsymbol{\mu}$ を，それぞれ， $\langle S \rangle$ ， $\langle S \boldsymbol{\mu} \rangle$ に置き換えることにより容易に実装可能である．また，事後分布 $Q(\Lambda)$ ， $Q(\mathbf{q})$ ， $Q(\mathbf{Z})$ の更新には， \mathbf{x} を用いるため， \mathbf{x} の更新と事後分布パラメータの更新を繰り返すことにより，より最適な音声パラメータを求めることができる．この繰り返し手順には，様々な組合せが考えられるが，例えば以下のような手順となる．

- (1) 事前分布パラメータと \mathbf{x} の初期値を与える．
- (2) $Q(\mathbf{Z})$ の更新 (式 (46), (47))
- (3) $Q(\mathbf{q})$ の更新 (式 (46), (48))．
- (4) Forward-Backward アルゴリズム (式 (39)–(42))
- (5) $Q(\Lambda)$ の更新 (式 (43)–(45))
- (6) \mathbf{x} の更新 (式 (50))．
- (7) \mathcal{F} が増加しなければ終了．それ以外は (2) へ．

3.5 コンテキストクラスタリング

HMM 音声合成では，HMM による音声認識に比べ，韻律やアクセントなどの様々な情報をコンテキストとして HMM を学習するため，コンテキストクラスタリングが重要な役割を果たす．これまで，コンテキストクラスタリングにおいても，HMM の学習と同様，ML 基準が用いられてきた．また，クラスタリングの分割停止条件として，MDL 基準が良く用いられる．しかし，MDL 基準はその導出に，学習データが十分に得られるという仮定を置いており，データ量が少ない場合，十分な精度が得られない可能性が指摘されている．この問題に対し，音声認識において，ベイズ基準によるコンテキストクラスタリング [11] が提案されており，有効性が確認されている．本研究でも同様に，決定木によるコンテキストクラスタリングにベイズ基準を用いる．

ベイズ基準のコンテキストクラスタリングでは， $Q(\mathbf{Z})$ ， $Q(\mathbf{q})$ ， $Q(a)$ を固定した状態で， \mathcal{F} を最大にする $Q(\mathbf{b})$ の共有構造を決定する． \mathcal{F} に近似事後確率を代入し，クラスタリングに関係のある項だけを考慮すると，

$$\begin{aligned} \mathcal{F} &= -\log C_\Lambda + E[Q(\mathbf{q})] + E[Q(\mathbf{Z})] \\ &= -\sum_{i=1}^N \log C_{b_i} + \text{Const} \quad (55) \end{aligned}$$

となる．ただし，

$$\log C_{b_i} = \log \frac{\bar{C}_{N_i} \bar{C}_{W_i}}{C_{N_i} C_{W_i}} (2\pi)^{\frac{N_i D}{2}} \quad (56)$$

である．簡単のため事前確率はすべてのガウス分布で等しいと仮定し，その正規化定数を C_{N_0} , C_{W_0} とすると，あるクラスタが分割された時の \mathcal{F} の変化量は，

$$\begin{aligned}\delta\mathcal{F} &= -\log C_{b_y} - \log C_{b_n} + \log C_{b_p} \\ &= -\log \bar{C}_{N_y} \bar{C}_{W_y} - \log \bar{C}_{N_n} \bar{C}_{W_n} \\ &\quad + \log \bar{C}_{N_p} \bar{C}_{W_p} + \log C_{N_0} C_{W_0} \\ &= \delta f_y + \delta f_n - \delta f_p - \delta f_0\end{aligned}\quad (57)$$

と書ける．ただし，

$$\delta f_i = -\frac{D}{2} \log \xi_i - \frac{\eta_i}{2} \log |\mathbf{B}_i| + \sum_{d=1}^D \log \Gamma\left(\frac{\eta_i + 1 - d}{2}\right)\quad (58)$$

である．ここで， C_{b_p} は分割前， C_{b_y} , C_{b_n} は分割後の事後確率分布の正規化定数であり，各クラスタの統計量 \bar{N}_i , $\bar{\sigma}_i$, \bar{C}_i から計算される． $\delta\mathcal{F}$ を最大にする質問を選択しながら分割を繰返し， $\delta\mathcal{F} \leq 0$ で分割を停止することにより，ベイズ基準によるクラスタリングを実現できる．

3.6 実装における近似

これまでに説明したアルゴリズムでは， $Q(\mathbf{Z})$ は， $Q(\Lambda)$ を介して x に依存する．よって，合成音声のラベル列 s が決定した後で，すべての学習データに対する Forward-Backward アルゴリズムを計算する必要がある．しかし，実際に，この計算を行うことは計算量の観点から現実的ではない．そこで，実装においては， $Q(\mathbf{Z})$ は x に依存しないと仮定し，あらかじめ学習データから計算した $Q(\mathbf{Z})$ で固定することを考える．

合成音声のラベル列 s が与えられる前は，式 (14) における $P(\mathbf{W}x, q|s, \Lambda)$ を計算することができないため，この項を考慮しないで $Q(\mathbf{Z})$, $Q(\Lambda)$ を最適化する．すなわち， $P(\mathbf{W}x, \mathbf{O}|s, S)$ を $P(\mathbf{O}|S)$ に置き換える．このとき，最適な事後分布は，

$$Q(\mathbf{Z}) = C_{\mathbf{Z}} \exp\langle \log P(\mathbf{O}, \mathbf{Z}|S, \Lambda) \rangle_{Q(\Lambda)}\quad (59)$$

$$Q(\Lambda) = C_{\Lambda} P(\Lambda) \exp\langle \log P(\mathbf{O}, \mathbf{Z}|S, \Lambda) \rangle_{Q(\mathbf{Z})}\quad (60)$$

で得られる．また，この分布を用いた期待値計算は，式 (39)–(42) を次式に置き換えることにより，計算可能である．

$$\bar{N}_i = \sum_{t=1}^T \langle z_t^i \rangle, \quad \bar{N}_{ij} = \sum_{t=1}^{T-1} \langle z_t^i z_{t+1}^j \rangle\quad (61)$$

$$\bar{\sigma}_i = \frac{1}{\bar{N}_i} \sum_{t=1}^T \langle z_t^i \rangle \mathbf{o}_t\quad (62)$$

$$\bar{C}_i = \frac{1}{\bar{N}_i} \sum_{t=1}^T \langle z_t^i \rangle (\mathbf{o}_t - \bar{\sigma}_i) (\mathbf{o}_t - \bar{\sigma}_i)^\top\quad (63)$$

合成システムを実現するための手順は，以下ようになる．

• 学習部

- (1) 事前分布パラメータを与える．
- (2) $Q(\mathbf{Z})$ の更新 (式 (46), (47)) ．
- (3) Forward-Backward アルゴリズム (式 (61)–(63)) ．
- (4) $Q(\Lambda)$ の更新 (式 (43)–(45)) ．
- (5) \mathcal{F} が増加しなければ終了．それ以外は (2) へ ．

• 合成部

- (1) ラベル列 s に基づく x の初期値を与える ．
- (2) $Q(q)$ の更新 (式 (46), (48))
- (3) Forward-Backward アルゴリズム (式 (39)–(42))
- (4) $Q(\Lambda)$ の更新 (式 (43)–(45))
- (5) x の更新 (式 (50)) ．
- (6) \mathcal{F} が増加しなければ終了．それ以外は (2) へ ．

合成部では， $Q(\Lambda)$ の更新を行わないなどの様々な組合せが考えられるが，どのような手順が適当であるかは，今後，評価実験により調査する必要がある．現在，実験として，学習部の実装を終えており，合成部において状態系列 q を外部から与え， $Q(q)$ と $Q(\Lambda)$ を更新しない簡単な場合について，合成音声の生成を確認している．また，2.2 で述べたように，ここで定義した HMM は状態系列の最適化には不十分であり，状態継続長モデルが必要となるが，本論文では，その点は考慮していない．そのため，実装上は， $Q(q)$ を状態継続長モデルを用いて計算するなど，何らかの対処が必要である．HMM の定義として，遷移系列を最適化できるモデルを定義し，学習も含めた最適化を行うのが理想的であるが，計算量の問題があり，今後の課題である ．

4. ま と め

本論文では，ベイズ基準による HMM 音声合成を導出した．今後の課題は，提案手法の評価と状態継続長モデルを用いた状態系列と音声パラメータの同時最適化である ．

文 献

- [1] H. Attias, “Inferring parameters and structure of latent variable models by variational Bayes,” Proc. UAI, 1999.
- [2] 上田修功, “ベイズ学習,” (全4回) 電子情報通信学会誌, Vol. 85, No. 4,6,7,8, 2002.
- [3] S. Watanabe, Y. Minami, A. Nakamura, and N. Ueda, “An application of variational Bayesian approach to speech recognition,” NIPS15, 2002.
- [4] 吉村貴克, 徳田恵一, 益子貴史, 小林隆夫, 北村正, “HMM に基づく音声合成におけるスペクトル・ピッチ・継続長の同時モデル化,” 信学論 (D-II), vol. J83-D-II, no.11, Nov. 2000.
- [5] 益子貴史, 徳田恵一, 小林隆夫, 今井 聖, “動的特徴を用いた HMM に基づく音声合成,” 信学論 (D-II), vol. J79-D-II, no.12, pp.2184–2190, Dec. 1996.
- [6] 徳田恵一, 益子貴史, 小林隆夫, 今井 聖, “動的特徴を用いた HMM からの音声パラメータ生成アルゴリズム,” 日本音響学会誌, vol.53, no.3, pp.192–200, Mar. 1997.
- [7] K. Tokuda, Takayoshi Yoshimura, T. Masuko, T. Kobayashi and T. Kitamura, “Speech parameter generation algorithms for HMM-based speech synthesis,” Proc. ICASSP, pp.1315–1318, June 2000.
- [8] K. Tokuda, T. Kobayashi and S. Imai, “Speech parameter generation from HMM using dynamic features,” Proc. ICASSP, pp.660–663, 1995.
- [9] K. Tokuda, T. Masuko, T. Yamada, T. Kobayashi and Satoshi Imai, “An algorithm for speech parameter generation from continuous mixture HMMs with dynamic features,” Proc. EUROSPEECH, pp.757–760, 1995.
- [10] K. Koishida, K. Tokuda, T. Masuko and T. Kobayashi, “Vector quantization of speech spectral parameters using statistics of dynamic features,” Proc. ICSP, vol.1, pp.247–252, Aug. 1997.
- [11] 渡辺 晋治, 南 康浩, 中村 篤, 上田修功, “ベイズの基準を用いた状態共有型 HMM 構造の選択,” 信学論 (D-II), vol. J86-D-II, no.6, pp.776–786, 2003.