

ノート指令と表現指令によって駆動される歌声 F_0 生成過程の統計モデル*

大石康智 (NTT), 亀岡弘和 (NTT/東大), 持橋大地 (統数研), 柏野邦夫 (NTT)

1 はじめに

歌声の音高 (F_0) 軌跡には、発声器官の物理的制約に起因する動的変動成分 (オーバーシュートや微細変動成分) や歌唱者の意図的表現による動的変動成分 (ビブラートやポルタメント) が含まれ、知覚的には前者は人間らしさ・自然性に関係し、後者は巧拙感に関係することがわかってきている [1]。ただし、これらの動的変動成分を F_0 軌跡から特徴抽出し、歌唱者ごとにその特性を精緻に学習することはまだ十分に検討されていない。

本研究の目的は、このような動的変動成分を F_0 軌跡から楽譜情報と分離して特徴抽出し、歌唱者ごとにどのようなパターンをもちうるのか、各パターンが文脈 (楽譜の音符列) にどう依存するかを計算機に学習させることである。ここでは、うろ覚えの状態であつた歌声ではなく、楽曲のメロディまたはその楽譜を既知として、歌唱者なりに表情付けて歌つた歌声を対象とする。本研究は、歌唱者の歌い方や個性、癖を学習することを目指しており、歌唱力評価や歌唱者識別、現在盛んに研究される歌声合成や歌声変換 [2, 3, 4] への応用が期待できる。

本研究では、物理モデルに基づいて F_0 軌跡の生成過程を記述し、そこから動的変動成分の特徴抽出に取り組む。これまで、話声の F_0 パターンを表現する藤崎モデル [5] を参考に、2次系を利用した歌声 F_0 制御モデルが提案された [1, 6]。2次系の伝達関数は

$$G(s) = \frac{\Omega^2}{s^2 + 2\zeta\Omega s + \Omega^2} \quad (1)$$

であり、この減衰率 ζ によって、様々な振動現象を表現できる。文献 [1] では、楽譜の音符列を表す階段状信号に式 (1) のインパルス応答を部分的に畳み込んで得られる F_0 軌跡を利用して、表情豊かな歌声合成音を実現した。制御パラメータ (減衰率 ζ と固有周波数 Ω) は手作業あるいは規則に基づいて決定された。一方、我々は観測される F_0 軌跡から制御パラメータを推定する逆問題の解法を検討してきた [7]。ただし、ビブラートや微細変動成分がすべて白色雑音としてモデル化されたため、意図的表現による動的変動成分を微細変動成分と分離して特徴付けられなかった。

本稿では、歌唱の意図的表現を特徴抽出するために、ノート指令信号と表現指令信号によって駆動される歌声 F_0 軌跡の生成過程を提案する (Fig. 1)。ここで、ノート指令信号は楽譜に記載される音符の並びを

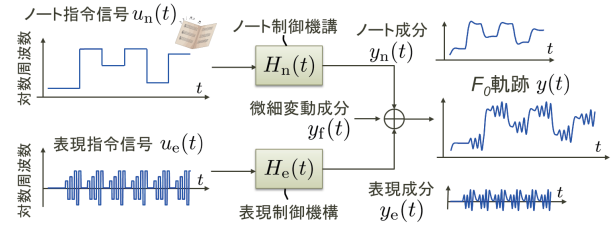


Fig. 1 提案する歌声の F_0 生成過程の概略

表現する。一方、表現指令信号は歌唱者の音楽的な表現意図を矩形状の細かい指令として表現する。ノート成分と表現成分は各指令信号によって駆動される制御機構の出力であり、これらの制御機構は2次系で表現される。ノート制御機構はオーバーシュートなどのノートの立ち上がり方を制御する。表現制御機構は矩形状の細かい指令信号を制御してビブラートやポルタメントを生成する。微細変動成分は10 Hz以上の不規則な振動成分を想定する [1]。最終的に、対数スケールの F_0 軌跡 $y(t)$ (ここで、 t は時間) は、これら3つの成分の重ね合わせと想定する。この F_0 生成過程を想定した理由は2点ある。

理由1 藤崎モデルでは、甲状軟骨の二つの独立な運動 (平行移動と回転) に伴う声帯の長さの変化の合計が F_0 の時間的変化をもたらすと解釈され、これらの運動に関する成分をそれぞれ、フレーズ成分とアクセント成分とした [5]。歌声にも同様の機構が存在すると仮定し、大域的な変化を表すノート成分とアクセントのような歌唱者が意図的に制御できる表現成分から歌声 F_0 軌跡が構成されると考えた。

理由2 従来、ビブラートは正弦波で表現され、そのパラメータは Vibrato rate と Vibrato extent とされた [1]。これらのパラメータが指数的に変化するモデルも提案された [3]。しかし実際は、これらのパラメータは時々刻々と変化するものの、指数的な変化であるとは言えない。この複雑な振動現象をモデル化して歌い方を認識するために、意図を表す矩形状の細かい指令信号が2次系によって制御されるビブラートの生成過程を想定した。

以降の節では、上記の歌声 F_0 軌跡の生成過程を離散時間表現して確率モデル化し、統計的手法に基づいてモデルパラメータの推定アルゴリズムを導出する。そして、提案モデルの有効性を実測の F_0 軌跡を用いて評価する。文献 [8] では、同じ枠組みで藤崎モデルを確率過程に基づいて統計モデル化する試みも行っているのでこちらも参照されたい。

* A Statistical Model of Singing Voice F_0 Contours Driven by Musical Note and Expression Command Functions by OHISHI, Yasunori (NTT), KAMEOKA, Hirokazu (NTT / The University of Tokyo), MOCHIHASHI, Daichi (The Institute of Statistical Mathematics), KASHINO Kunio (NTT)

2 歌声 F_0 生成過程の確率モデル

歌声 F_0 軌跡の生成過程を離散時間表現し、確率過程に基づいて統計モデル化する。連続時間領域で表現される 2 次系制御機構の伝達関数の離散時間表現を得るために、後差分変換を利用する [7]。 t_0 をサンプリング周期とすると、この変換によりノート制御機構の逆システムの伝達関数 $\mathcal{H}_n^{-1}(s)$ は z 領域で、

$$\mathcal{H}_n^{-1}(z) = a_2 z^{-2} + a_1 z^{-1} + a_0 \quad (2)$$

と書くことができる。ただし、 $a_2 = \varphi^2$, $a_1 = -2\varphi(\psi + \varphi)$, $a_0 = 1 + 2\varphi\psi + \varphi^2$, および、 $\varphi = 1/(\Omega t_0)$, $\psi = \zeta$ と表現される。ここで、 k を離散時刻インデックスとし、ノート指令信号およびノート成分の離散時間表現をそれぞれ $u_n[k]$, $y_n[k]$ とすると、 $y_n[k]$ は、ノート制御パラメータ φ , ψ によって特性が決まる拘束つき全極モデルからの出力

$$u_n[k] = a_0 y_n[k] + a_1 y_n[k-1] + a_2 y_n[k-2] \quad (3)$$

と見なすことができる。同様に、表現指令信号 $u_e[k]$ と表現成分 $y_e[k]$ の関係も $u_e[k] = b_0 y_e[k] + b_1 y_e[k-1] + b_2 y_e[k-2]$ と書ける。ただし、表現制御機構は臨界制動系 ($\zeta = 1$ の場合) を想定し、 $b_2 = \xi^2$, $b_1 = -2\xi(1 + \xi)$, $b_0 = 1 + 2\xi + \xi^2$, 表現制御パラメータを $\xi = 1/(\Omega t_0)$ とする。微細変動成分 $y_f(t)$ の離散時間表現を $y_f[k]$ として、歌声 F_0 軌跡の離散時間表現は、 $y[k] = y_n[k] + y_e[k] + y_f[k]$ と想定する。

次に、 F_0 生成過程を確率モデル化する。ノート指令信号と表現指令信号はそれぞれ、楽譜に記載されるメロディの音符の並びと歌唱者の音楽的な表現意図を表す (Fig. 1)。これらの指令信号を表現するために、HMM を利用して、 $u_n[k]$ と $u_e[k]$ を確率モデル化する (Fig. 2)。 $o[k] := (u_n[k], u_e[k])^T$ を

$$o[k] \sim \mathcal{N}(\nu[k], \Upsilon), \nu[k] := \begin{bmatrix} \mu_n[k] \\ \mu_e[k] \end{bmatrix}, \Upsilon := \begin{bmatrix} \sigma_n^2 & 0 \\ 0 & \sigma_e^2 \end{bmatrix}$$

のように正規分布する確率変数と見なし、平均 $\nu[k]$ が状態遷移に伴って変化するモデルを考える。

具体的には、この HMM は $I \times J$ 個の状態集合 $S := \{S_{i,j}\}_{i=1,j=1}^{I,J}$ からなる。状態 $S_{i,j}$ では、 $\mu_n[k]$ は $A_n^{(i)} + d_i$ の値をとる。ここで、 $A_n^{(i)}$ は楽譜に記載されるメロディの i 番目の音符の音高を表し (この値は楽譜から与えられるものとする)、 d_i はその絶対音高からのずれ (音高シフトパラメータと呼ぶ)、 I は音符の総数を表す。Fig. 2 に示す状態遷移により、 $\mu_n[k]$ は I 個の音符区間からなる階段状信号を表現する。一方、 $\mu_e[k]$ は $B_e^{(i,j)}$ の値をとり、これは歌唱者の表現意図を表すための、矩形の指令の大きさを表す。ここで、 i 番目の音符では、 $S_{i,1}$ を通らずして、状態 $S_{i,j}$ から別の状態 $S_{i,j'}$ ($j \neq j'$, $2 \leq j \leq J$, $2 \leq j' \leq J$) へ直接に遷移できない制約を設けることによって、 $\mu_e[k]$ は Fig. 2 に示すよ

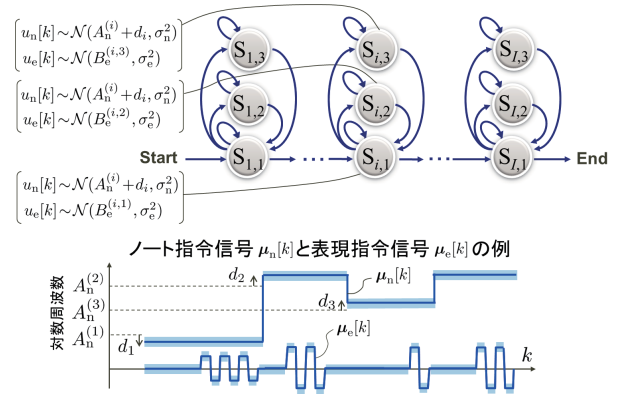


Fig. 2 HMM による指令信号の統計モデル化

うな矩形信号を表現する。 J は表現指令を構成するための状態数を表す (Fig. 2 は $J = 3$ である)。指令信号を生成する HMM (指令信号モデル) の構成をまとめると、出力系列: $\{o[k]\}_{k=1}^K$, 状態集合: $\{S_{i,j}\}_{i=1,j=1}^{I,J}$, 状態系列: $\{s_k\}_{k=1}^K$, 出力確率分布: $P(o[k]|s_k = n) = \mathcal{N}(c_n[k], \Upsilon)$, $c_n[k] = (A_n^{(i)} + d_i, B_e^{(i,j)})^T$, ($n = S_{i,j}$), 状態遷移確率 $\phi_{n',n} := \log P(s_k = n|s_{k-1} = n')$ となる。簡単のため、状態遷移確率 $\phi_{n',n}$ は定数とすると、指令信号モデルにおいて推定すべきパラメータは、HMM の経路を表す状態系列 $\{s_k\}_{k=1}^K$ 、音高シフトパラメータ $\{d_i\}_{i=1}^I$ 、表現指令の大きさパラメータ $\{B_e^{(i,j)}\}_{i=1,j=1}^{I,J}$ 、指令信号の分散パラメータ σ_n^2 , σ_e^2 であり、これらをまとめて θ_n と記述する。また、平均値系列 $\{\mu_n[k]\}_{k=1}^K$ および $\{\mu_e[k]\}_{k=1}^K$ は、状態系列 $\{s_k\}_{k=1}^K$ が与えられたもとで、 $(\mu_n[k], \mu_e[k])^T \leftarrow c_{s_k}[k]$ で与えられる。

$y = (y[1], \dots, y[K])^T$ の確率密度関数を導く。 $u_n := (u_n[1], \dots, u_n[K])^T$, $u_e := (u_e[1], \dots, u_e[K])^T$, $\mu_n := (\mu_n[1], \dots, \mu_n[K])^T$, $\mu_e := (\mu_e[1], \dots, \mu_e[K])^T$ とすると

$$u_n | \theta_u \sim \mathcal{N}(\mu_n, \Sigma_n), \Sigma_n = \sigma_n^2 I_K \quad (4)$$

$$u_e | \theta_u \sim \mathcal{N}(\mu_e, \Sigma_e), \Sigma_e = \sigma_e^2 I_K \quad (5)$$

と書ける。ここで、 I_K は $K \times K$ の単位行列を表す。ノート成分 $y_n := (y_n[1], \dots, y_n[K])^T$ とノート指令信号 u_n の関係、および表現成分 $y_e := (y_e[1], \dots, y_e[K])^T$ と表現指令信号 u_e の関係は、

$$A := \begin{bmatrix} A^{(1)} \\ \vdots \\ A^{(i)} \\ \vdots \\ A^{(I)} \end{bmatrix} = \begin{bmatrix} a_0^{(1)} & & & & O \\ & a_1^{(1)} & & & \\ & & \ddots & & \\ & & & a_2^{(i)} & a_1^{(i)} & a_0^{(i)} \\ & & & & \ddots & \ddots \\ O & & & & & a_2^{(I)} & a_1^{(I)} & a_0^{(I)} \end{bmatrix}$$

$$B := \begin{bmatrix} b_0 & & & & O \\ b_1 & b_0 & & & \\ b_2 & b_1 & b_0 & & \\ & \ddots & \ddots & \ddots & \\ O & & b_2 & b_1 & b_0 \end{bmatrix}$$

と置くと、それぞれ、 $\mathbf{u}_n = \mathbf{A}\mathbf{y}_n, \mathbf{u}_e = \mathbf{B}\mathbf{y}_e$ と表現できる。ここでは、先行研究で得られた知見 [1] に基づいて、ノート制御パラメータがノートごとに異なるものと想定し、 $a_0^{(i)}, a_1^{(i)}, a_2^{(i)}$ は、 $a_2^{(i)} = (\varphi^{(i)})^2, a_1^{(i)} = -2\varphi^{(i)}(\psi^{(i)} + \varphi^{(i)}), a_0^{(i)} = 1 + 2\varphi^{(i)}\psi^{(i)} + (\varphi^{(i)})^2$ と表現され、 $\{\varphi^{(i)}, \psi^{(i)}\}_{i=1}^I$ を制御パラメータとする。そこで、ノート成分と表現成分の確率密度関数は、

$$\mathbf{y}_n \sim \mathcal{N}(\mathbf{A}^{-1}\boldsymbol{\mu}_n, \mathbf{A}^{-1}\boldsymbol{\Sigma}_n(\mathbf{A}^{-1})^T) \quad (6)$$

$$\mathbf{y}_e \sim \mathcal{N}(\mathbf{B}^{-1}\boldsymbol{\mu}_e, \mathbf{B}^{-1}\boldsymbol{\Sigma}_e(\mathbf{B}^{-1})^T) \quad (7)$$

が導出される。微細変動成分 $\mathbf{y}_f := (y_f[1], \dots, y_f[K])^T$ はガウス性白色雑音を想定する。

$$\mathbf{y}_f | \sigma_f^2 \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_f), \quad \boldsymbol{\Sigma}_f = \sigma_f^2 \mathbf{I}_K \quad (8)$$

仮定より、 $\mathbf{y}_n, \mathbf{y}_e, \mathbf{y}_f$ は独立なので、パラメータ $\Theta := \{\theta_u, \{\varphi^{(i)}, \psi^{(i)}\}_{i=1}^I, \xi, \sigma_f^2\}$ が与えられた下での歌声 F_0 軌跡 $\mathbf{y} = \mathbf{y}_n + \mathbf{y}_e + \mathbf{y}_f$ の確率密度関数（尤度関数）は式 (6), (7) と式 (8) より、

$$\mathbf{y} | \Theta \sim \mathcal{N}(\mathbf{A}^{-1}\boldsymbol{\mu}_n + \mathbf{B}^{-1}\boldsymbol{\mu}_e, \mathbf{A}^{-1}\boldsymbol{\Sigma}_n(\mathbf{A}^{-1})^T + \mathbf{B}^{-1}\boldsymbol{\Sigma}_e(\mathbf{B}^{-1})^T + \boldsymbol{\Sigma}_f) \quad (9)$$

となる。 Θ の事前確率は、各パラメータは互いに独立で、パラメータ $\sigma_n^2, \sigma_e^2, \sigma_f^2$ は一様に分布すると仮定し、 $P(\Theta) \propto P(\varphi^{(1)}) \dots P(\varphi^{(I)}) P(\psi^{(1)}) \dots P(\psi^{(I)}) P(\xi) P(d_1) \dots P(d_I) P(B_e^{(1,1)}) \dots P(B_e^{(I,J)}) P(s_1) \prod_{k=2}^K P(s_k | s_{k-1})$ と記述される。

3 パラメータ推定アルゴリズム

Θ の事後確率 $P(\Theta | \mathbf{y})$ を最大化する問題を解析的に解くことはできないが、 $\mathbf{x} := (\mathbf{y}_n^T, \mathbf{y}_e^T, \mathbf{y}_f^T)^T$ を完全データと見なすことで、EM アルゴリズムによる不完全データ問題に帰着できる。不完全データと完全データの関係は $\mathbf{y} = \mathbf{H}\mathbf{x}$ となり、ここで $\mathbf{H} := [\mathbf{I}_K \ \mathbf{I}_K \ \mathbf{I}_K]$ とする。このとき、Q 関数は、

$$Q(\Theta, \Theta') \stackrel{c}{=} \frac{1}{2} [\log |\boldsymbol{\Lambda}^{-1}| - \text{tr}(\boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta'])]$$

$$+ 2\mathbf{m}^T \boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] - \mathbf{m}^T \boldsymbol{\Lambda}^{-1} \mathbf{m}] + \log P(\Theta)$$

$$\mathbf{m} := \begin{bmatrix} \mathbf{A}^{-1}\boldsymbol{\mu}_n \\ \mathbf{B}^{-1}\boldsymbol{\mu}_e \\ \mathbf{0} \end{bmatrix}, \boldsymbol{\Lambda}^{-1} := \begin{bmatrix} \mathbf{A}^T \boldsymbol{\Sigma}_n^{-1} \mathbf{A} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{B}^T \boldsymbol{\Sigma}_e^{-1} \mathbf{B} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \boldsymbol{\Sigma}_f^{-1} \end{bmatrix}$$

となる。ここで、 $\text{tr}(\cdot)$ は行列のトレースを表す。E ステップでは、直前のステップで更新されたパラメータを Θ' に代入し、条件付きガウス分布の性質から $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ と $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$ を計算する [7]。 $\mathbf{y}_n, \mathbf{y}_e, \mathbf{y}_f$ に対応するように、 $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ および $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$ を以下のように区分表現する。

$$\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] = \begin{bmatrix} \bar{\mathbf{x}}_n \\ \bar{\mathbf{x}}_e \\ \bar{\mathbf{x}}_f \end{bmatrix}, \mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta'] = \begin{bmatrix} \mathbf{R}_n & * & * \\ * & \mathbf{R}_e & * \\ * & * & \mathbf{R}_f \end{bmatrix}$$

各パラメータの M ステップ更新式を以下に示す。

1) 状態系列: Q 関数で $s := \{s_k\}_{k=1}^K$ に関する項は

$$\mathcal{I}_1(s) := -\frac{1}{2} \sum_{k=1}^K (\mathbf{o}[k] - \mathbf{c}_{s_k}[k])^T \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[k] - \mathbf{c}_{s_k}[k]) + \log P(s_1) + \sum_{k=2}^K \log P(s_k | s_{k-1}) \quad (10)$$

である。ただし、 $\mathbf{o}[k] := ([\mathbf{A}\bar{\mathbf{x}}_n]_k, [\mathbf{B}\bar{\mathbf{x}}_e]_k)^T$ であり、 $[\cdot]_k$ はベクトルの k 番目の要素を表す。これを最大化する $\{s_k\}_{k=1}^K$ は Viterbi 探索により導出できる。

2) ノート制御パラメータ: $\varphi^{(i)}$ と $\psi^{(i)}$ に関する事前分布を $\varphi^{(i)} \sim \mathcal{N}(\mu_\varphi, \sigma_\varphi^2), \psi^{(i)} \sim \mathcal{N}(\mu_\psi, \sigma_\psi^2)$ とする。Q 関数の中で $\varphi^{(i)}$ と $\psi^{(i)}$ に関する項は、

$$\begin{aligned} \mathcal{I}_2(\varphi^{(i)}, \psi^{(i)}) &= |\mathcal{T}_{S_{i,\cdot}}| \log(1 + 2\varphi^{(i)}\psi^{(i)} + (\varphi^{(i)})^2) \\ &\quad - \frac{1}{2\sigma_n^2} \text{tr}((\mathbf{A}^{(i)})^T \mathbf{A}^{(i)} \mathbf{R}_n) + \frac{1}{\sigma_n^2} ([\boldsymbol{\mu}_n]_{\mathcal{T}_{S_{i,\cdot}}})^T \mathbf{A}^{(i)} \bar{\mathbf{x}}_n \\ &\quad - \frac{1}{2\sigma_\varphi^2} (\varphi^{(i)} - \mu_\varphi)^2 - \frac{1}{2\sigma_\psi^2} (\psi^{(i)} - \mu_\psi)^2 \end{aligned}$$

$$\mathcal{T}_{S_{i,\cdot}} = \{k | s_k \in \{S_{i,1}, \dots, S_{i,J}\}\} \quad (11)$$

となる。ここで、 $|\mathcal{T}|$ は集合 \mathcal{T} の要素数を表す。また、 $[\boldsymbol{\mu}]_{\mathcal{T}}$ は、集合 \mathcal{T} の要素を添え字として、その添え字に相当する $\boldsymbol{\mu}$ の要素を取り出した部分ベクトルを表す。今、定数行列 $\mathbf{U}_0, \mathbf{U}_1, \mathbf{U}_2$ を用いて $\mathbf{A}^{(i)}$ は、

$$\mathbf{A}^{(i)} = (\varphi^{(i)})^2 [\mathbf{U}_2]_{\mathcal{T}_{S_{i,\cdot}}} + \varphi^{(i)}\psi^{(i)} [\mathbf{U}_1]_{\mathcal{T}_{S_{i,\cdot}}} + [\mathbf{U}_0]_{\mathcal{T}_{S_{i,\cdot}}}$$

と表現される。ここで、 $[\mathbf{U}]_{\mathcal{T}}$ は集合 \mathcal{T} の要素を添え字として、行列 \mathbf{U} からその添え字に相当する行ベクトルを取り出して構成される部分行列を意味する。ニュートン・ラフソン法を利用して、 $\mathcal{I}_2(\varphi^{(i)}, \psi^{(i)})$ を最大化する $\varphi^{(i)}$ と $\psi^{(i)}$ が数値的に導出される。

3) 表現制御パラメータ: ノート制御パラメータと同様の手順で更新式を導出できる。ただし、推定すべきパラメータ数が多いため、今回は $\xi = 3$ に固定する。

4) その他のパラメータ: d_i と $B_e^{(i,j)}$ の事前分布をそれぞれ $d_i \sim \mathcal{N}(0, \sigma_d^2)$ と $B_e^{(i,j)} \sim \mathcal{N}(\mu_{B^{(i,j)}}, \sigma_B^2)$ とする。残されたパラメータの更新式は

$$d_i = \frac{1}{|\mathcal{T}_{S_{i,\cdot}}| + \sigma_n^2 / \sigma_d^2} \sum_{k \in \mathcal{T}_{S_{i,\cdot}}} ([\mathbf{A}\bar{\mathbf{x}}_n]_k - A_n^{(i)}) \quad (12)$$

$$B_e^{(i,j)} = \frac{1}{|\mathcal{T}_{S_{i,j}}| / \sigma_e^2 + 1 / \sigma_B^2} \left(\sum_{k \in \mathcal{T}_{S_{i,j}}} \frac{[\mathbf{B}\bar{\mathbf{x}}_e]_k}{\sigma_e^2} + \frac{\mu_{B^{(i,j)}}}{\sigma_B^2} \right)$$

$$\mathcal{T}_{S_{i,j}} = \{k | s_k = S_{i,j}\} \quad (13)$$

$$\sigma_n^2 = \left(\text{tr}(\mathbf{A}^T \mathbf{A} \mathbf{R}_n) - 2\boldsymbol{\mu}_n^T \mathbf{A} \bar{\mathbf{x}}_n + \boldsymbol{\mu}_n^T \boldsymbol{\mu}_n \right) / K \quad (14)$$

$$\sigma_e^2 = \left(\text{tr}(\mathbf{B}^T \mathbf{B} \mathbf{R}_e) - 2\boldsymbol{\mu}_e^T \mathbf{B} \bar{\mathbf{x}}_e + \boldsymbol{\mu}_e^T \boldsymbol{\mu}_e \right) / K \quad (15)$$

$$\sigma_f^2 = \text{tr}(\mathbf{R}_f) / K \quad (16)$$

と導出される。無声音のため F_0 が推定されない区間の取り扱いや各パラメータの実際の更新方法については、文献 [9] を参照されたい。

4 評価実験

実測の F_0 軌跡を用いて、提案手法の性能を評価する。「RWC 研究用音楽データベース：ポピュラー音楽」(RWC-MDB-P-2001)[10]における、歌手名：緒方智美、曲番号：No.07、曲名：PROLOGUE のメロディの F_0 を手作業でラベル付した結果 [11] を実測 F_0 軌跡として利用した。本来ならば音響信号から F_0 を推定すべきであるが、今回は提案手法の性能の上限を調べるためにこのようなデータを利用した。 F_0 は 10ms ごとにラベル付されているので、アップサンプリングによって、5ms ごとの実測 F_0 軌跡を得た。同楽曲の楽譜情報はその MIDI データを利用した。

パラメータ推定アルゴリズムの設定値を以下に示す。 I と $\{A_n^{(i)}\}_{i=1}^I$ は MIDI データから与えられる。表現指令信号を構成するための状態数 J は 5 とした。HMM の状態遷移確率は、 $\phi_{S_{i,1},S_{i,1}} = \log(0.9999 \times (J-1)/J)$, $\phi_{S_{i,1},S_{i,j}} = \log(0.9999/J)$, $\phi_{S_{i,1},S_{i+1,1}} = \log(0.0001)$, $\phi_{S_{i,j},S_{i,j}} = \log(0.9999)$, $\phi_{S_{i,j},S_{i,1}} = \log(0.0001)$, ($1 \leq i \leq I$, $2 \leq j \leq J$) とした。パラメータの事前分布における固定パラメータは、 $\mu_\varphi = 6$, $\sigma_\varphi^2 = 0.1$, $\mu_\psi = 0.6$, $\sigma_\psi^2 = 0.02$, $\sigma_d^2 = 2500$, $\sigma_B^2 = 100$, $\mu_{B^{(i,1)}} = 0$, $\mu_{B^{(i,2)}} = 30$, $\mu_{B^{(i,3)}} = -30$, $\mu_{B^{(i,4)}} = 60$, $\mu_{B^{(i,5)}} = -60$, ($1 \leq i \leq I$) とした。これらは文献 [1] と予備実験に基づいて決定した。

実測 F_0 軌跡からノート指令信号と表現指令信号、制御機構の出力であるノート成分と表現成分の推定結果を Fig. 3 に示す。これらを足し合わせたものを合成 F_0 軌跡とし、定性的には実測 F_0 軌跡に近い軌跡が得られた。また、微細変動成分パラメータである σ_f^2 の値は 187.7 であるため、振幅が 13 cent 程度の変動成分が推定され、先行研究の知見 [1] と整合する。無声音の区間を欠損データとして扱うため、合成 F_0 軌跡ではその部分が周囲の F_0 から補間される。

Fig. 4 は具体的に、ポルタメントとビブラートの生成される様子を示す。楽譜 (MIDI) では一つの音を伸ばし続けるよう記述されているものの、歌手の表現意図によって、半音だけ滑らかに変化させるポルタメントが観測される。このとき、ノート指令信号は同一音符区間として推定され、表現指令信号は音高を段階的に変化させる指令が推定される。一方、実測 F_0 軌跡には必ずしも正弦波とは言えないビブラートが観測される。提案法では、部分的に矩形化する表現指令が推定されてビブラートを生成する。ポルタメントとビブラートの生成のどちらにおいても合成 F_0 軌跡は実測 F_0 軌跡に近い軌跡が得られた。

5 おわりに

歌声 F_0 軌跡に含まれる動的変動成分を楽譜情報と分離して抽出することを目的として、2 つの指令信号によって駆動される F_0 軌跡の生成過程を記述し、モ

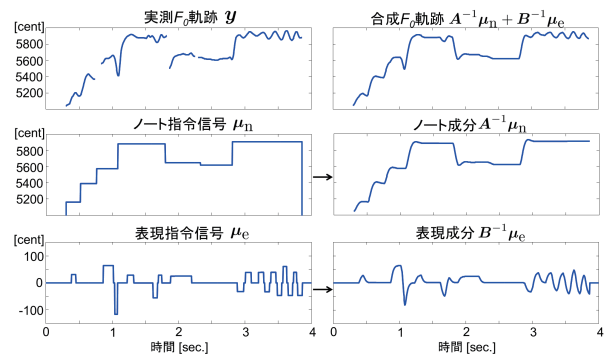


Fig. 3 実測 F_0 軌跡に対するパラメータ推定結果

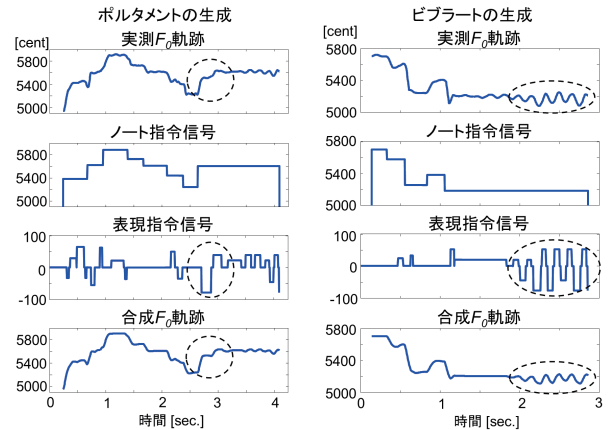


Fig. 4 ポルタメントとビブラートの生成

デルパラメータの推定アルゴリズムを提案した。評価実験では、実測の F_0 軌跡から推定される歌唱者の表現意図を考察した。提案モデルの有効性を主観的に評価する聴取実験も行なっているので、文献 [9] も参照されたい。今後の課題は、抽出された動的変動成分のパターンを歌唱者ごとに学習することである。

参考文献

- [1] Saitou *et al.*, *Proc. EUROSPEECH'09*, pp. 832–835, 2009.
- [2] Mase *et al.*, *Proc. INTERSPEECH'10*, pp. 845–848, 2010.
- [3] 右田ほか, 情処論文誌, Vol.52, No.5, pp. 1910–1922, 2011.
- [4] 中野ほか, 情処論文誌, Vol.52, No.12, pp. 3853–3867, 2011.
- [5] Fujisaki, In *Vocal Physiology: Voice Production, Mechanisms and Functions*, (O. Fujimura, ed.) Raven Press, pp. 347–355, 1988.
- [6] 柏野ほか, 音講論集, 2-9-1, pp. 625–626 1998.
- [7] 大石ほか, 音講論集, 1-8-19 pp. 279–282, 2011.
- [8] Kameoka *et al.*, *Proc. SAPA'10*, pp. 43–48, 2010.
- [9] 大石ほか, 情処研報, 2012–MUS94–12, 2012.
- [10] 後藤ほか, 情処論文誌, Vol.45, No.3, pp. 728–738, 2004.
- [11] Goto, *Proc. ISMIR'06*, 2006.