

ガウス過程回帰の混合エキスパートモデルを用いた 歌声 F_0 軌跡の予測と歌唱表現変換

大石 康智^{1,a)} 持橋 大地² 亀岡 弘和¹ 柏野 邦夫¹

概要： 歌声の声の高さ（基本周波数， F_0 ）の時間変化に表れる歌唱表現を特徴抽出し，任意の楽譜にその歌唱表現を転写して F_0 軌跡を予測できる生成過程モデルを提案する．先行研究より，歌声の F_0 軌跡には歌唱表現に起因する動的変動成分が観測され，歌唱者はこれらを巧みに制御して，歌声に表情をつけていることが明らかとなった．この歌唱動作を計算機上で実現するために，ガウス過程回帰の混合エキスパートモデルを用いて，楽譜における様々なコンテキスト（音符の音高や音長，音符内位置，前後の音符情報など）と F_0 軌跡の関係を学習する．個々のガウス過程回帰によって F_0 軌跡の多様な動特性が特徴づけられ，これらの混合エキスパートモデルによって歌唱者が時々刻々と動特性パターンを使い分ける動作が表現される．評価実験では，単一のガウス過程回帰或多項式回帰を用いるより，ガウス過程の混合エキスパートモデルを用いて F_0 軌跡の動特性を特徴づけた方が，未知の楽譜に対する F_0 軌跡の予測性能が高いことを示す．また，学習データの量やこのモデルを歌唱表現変換に応用することについて議論する．

1. はじめに

本研究の目的は，歌声の基本周波数の時間変化（ F_0 軌跡）に表れる歌唱表現を特徴抽出し，任意の楽譜に対して，その歌唱表現を反映した F_0 軌跡を予測可能な生成過程モデルを構築することである．歌唱表現の定義はまだ十分になされていないが，歌唱表現が歌唱フォルマントや F_0 動的変動成分のような信号特徴と関連付けられることが明らかにされた [1]．特に，ビブラートやポルタメントのような F_0 動的変動成分は知覚的な観点からも，歌唱表現と密接に関係する [2]．ビブラートは，持続音において， F_0 が準周期的に変調される成分であり，速度，振幅，継続長，波形，規則性などによって特徴付けられる．一方，ポルタメントは，音符の切り替わりにおいて， F_0 が緩やかに変化する成分である [3]．歌唱者は，楽譜のコンテキスト（音符の高さや長さなど）が与えられた下で，これらの動的変動成分を巧みに制御して，歌声に表情を付ける．このような生成過程を計算機で実現できれば，歌唱者の個性や表現力を学習することにつながり，歌声の認識や合成 [4-6] への応用を期待できる．例えば，ある歌唱者の大量の歌声から歌唱表現を学習することによって，どんな新しい楽譜に対しても，その表現を付与して合成することが可能となるだ

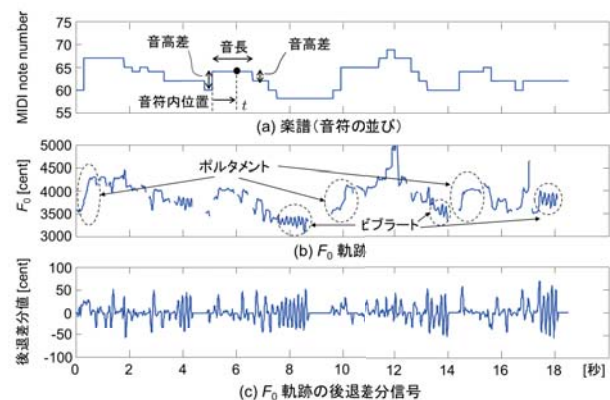


図 1 メロディの楽譜とそれを歌った歌声の F_0 軌跡: MIDI 信号からノートナンバと音長を算出してメロディの楽譜とした．

Fig. 1 Musical note sequence of melody and F_0 contour of singing voice which sang that melody

ろう．

F_0 動的変動成分を特徴抽出する試みとして，文献 [7, 8] は歌唱力の自動評価や歌声合成のためのビブラート特徴抽出手法（速度，振幅，継続長など）を検討した．文献 [3] はビブラートとポルタメントを正弦波の加法モデルで表現し，推定されたモデルパラメータを歌唱者認識に利用した．文献 [2, 9] では，2 次系モデルによって F_0 動的変動成分が表現され，歌声合成のために利用された．このように，個々の F_0 動的変動成分を特徴抽出して，歌声の認識や合成，歌唱力評価に利用する試みはなされてきたものの，動的変動成分を楽譜のコンテキストと結びつけて，体系的にその動特性を学習する方法は十分に検討されていなかった．隠

¹ NTT コミュニケーション科学基礎研究所
NTT Communication Science Laboratories, Atsugi, Kanagawa 243-0198, Japan
² 情報・システム研究機構 統計数理研究所
The Institute of Statistical Mathematics, Tachikawa, Tokyo 190-8562
a) ohishi.yasunori@lab.ntt.co.jp

れマルコフモデル (HMM) に基づく歌声合成のように、コンテキストラベルが付与された区間毎に、例えば、5 状態の HMM を用いて F_0 の動特性を学習することもできるが、状態内で出力確率分布が一定であるという HMM の制約のために、短時間に細かく変動する歌声の F_0 を、固定された状態数と局所的な動的特徴で表現することは難しいとされる [10, 11]. また、未知のコンテキストに対して頑健なモデルを構築するために、木構造に基づくクラスタリングを行うが、実際に生成される F_0 軌跡は過剰に平滑化され、本来の動特性を再現できないという問題もある。

本稿では、楽譜のコンテキストから歌声の F_0 軌跡が生成される入出力関係を、ガウス過程回帰の混合エキスパートモデル [12] を利用して直接的にモデル化する。混合エキスパートモデルは混合モデルの一つの拡張であり、混合係数を入力変数の関数 (ゲート関数と呼ぶ) で表現する。すなわち、入力空間 (楽譜のコンテキストからなる空間) の異なる領域ごとに、個々の混合要素の密度関数 (エキスパートと呼ばれ、ここではガウス過程回帰に相当する) が F_0 軌跡の予測を行う。ゲート関数はいずれの混合要素がどの領域において優勢であるかを判定する。これにより、歌唱者がコンテキストに依存して、歌唱表現に起因する F_0 動的特徴の変動成分を多様に使い分ける動作を表現できると考える。さらに、動的特徴成分を精緻に特徴づけるために、個々のガウス過程回帰におけるカーネル関数の設計を工夫する。評価実験では、新規の楽譜に対する F_0 軌跡の予測の観点から提案法を評価する。また、モデルの学習データ量を議論し、歌唱表現変換への応用を検討する。

2. 歌声 F_0 軌跡の生成過程モデル

メロディの楽譜とその歌声が同期して与えられた下で、楽譜に含まれる様々なコンテキストから F_0 軌跡への回帰問題を考える。入力変数は楽譜に基づいて、例えば下記のように構成される。

$$\mathbf{x}_t = [\text{音符内位置 (発音開始時刻からの時間)}, \text{音長}, \text{先行音符との音高差}, \text{後続音符との音高差}]^T \quad (1)$$

もちろん、歌詞の音素情報やクレッシェンドのような演奏情報の有無などを入力変数に加えることも可能である。一方で、出力変数 $\{y_t\}_{t=1}^T$ には、 F_0 軌跡から音符情報を取り除いた後退差分信号を利用する。差分信号を扱うことで、図 1 に示すように、動的特徴成分を顕著に観測することができる。 F_0 は例えば、YIN [13] を利用して 10 ms ごとに推定され、下記のように Hz で表される周波数単位を cent に変換する。

$$y_{\text{cent}} = 1200 \log_2 \frac{y_{\text{Hz}}}{440 \times 2^{\frac{3}{12} - 5}} \quad (2)$$

学習データ $D = \{\mathbf{x}_t, y_t\}_{t=1}^T$ が与えられ、出力変数を並べ

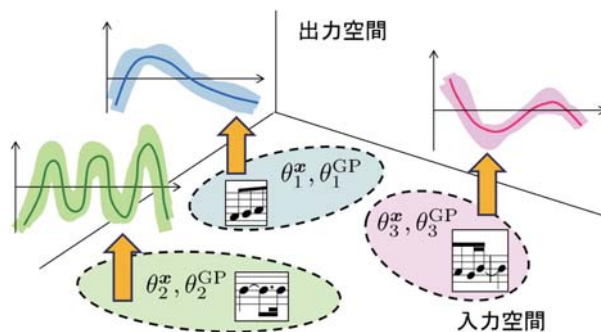


図 2 MoGPEs における入力空間の領域割り当て: 各楕円領域は全共分散型ガウス分布を表し、この領域では同じ GPR によって出力が予測される。

Fig. 2 An example allotment of input space in the MoGPEs. Each elliptical region represents a two-dimensional full covariance Gaussian distribution. The shaded regions are owned by the same GPE.

たベクトル $\mathbf{y} = [y_1, \dots, y_T]^T$ をガウス過程回帰 (Gaussian process regression, GPR) によって回帰する場合、その確率密度関数は

$$\mathbf{y} \sim \mathcal{GP}(\mathbf{y}; \mathbf{0}, \mathbf{K} + \eta^2 \mathbf{I}) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{K} + \eta^2 \mathbf{I}) \quad (3)$$

のような多次元ガウス分布で表現される。ここで、 \mathbf{K} は $K_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ を要素に持つグラム行列、 $k(\mathbf{x}_i, \mathbf{x}_j)$ は 2 変数間の類似性を定義するカーネル関数である。また、 η^2 と \mathbf{I} はそれぞれ、観測雑音の分散値と単位行列を表す。GPR は新しい入力変数 \mathbf{x}_* が与えられたときの出力変数 y_* の予測分布を推論できる。予測分布は、

$$p(y_* | D, \mathbf{x}_*) = \mathcal{N}(y_*; \mu_*, \sigma_*^2), \quad (4)$$

$$\mu_* = \mathbf{k}_*^T (\mathbf{K} + \eta^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\sigma_*^2 = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (\mathbf{K} + \eta^2 \mathbf{I})^{-1} \mathbf{k}_*$$

となり、 \mathbf{k}_* は $k(\mathbf{x}_t, \mathbf{x}_*)$ ($t = 1, \dots, T$) からなるベクトルを表す。学習データが“事例”となって、新たな入力変数に対する予測に貢献することがわかる。

カーネル関数の必要条件はグラム行列が半正定値対称行列となることである。出力信号が定常であることを仮定し、Squared exponential (SE) 共分散関数が一般的に利用されるが、図 1 に示すように、差分信号は必ずしも定常ではない。歌唱者は音符の並びに依存して、様々な動的特徴成分を使い分けながら歌唱するためである。このような非定常な動特性をモデル化するために、ガウス過程回帰の混合エキスパートモデル (Mixture of Gaussian process experts, MoGPEs) を利用する。

MoGPEs では、入力空間がゲート関数によって確率的に分割される (図 2)。それぞれの領域では、異なる GPR によって、動的特徴成分が特徴付けられて、出力が予測される。最終的に、出力は個々のゲート関数と GPR との混合モデルとして表現される。MoGPEs はこれまで、2 つの構

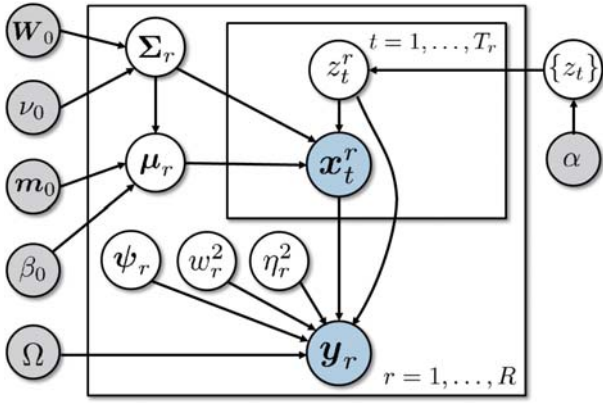


図 3 提案法のグラフィカル表現

Fig. 3 Graphical representation of our model

成方法 [12, 14] が提案されているが、文献 [12] は MoGPEs の完全な生成モデルを定義し、入力変数の欠損や出力変数から入力変数への逆予測が可能であるため、この構成方法を利用する。生成過程の流れを以下に示す。

- (1) ディリクレ多項分布モデルを用いて、 T 個の入力変数を R 個の領域のいずれかに割り当てる。これは指示変数集合 $\{z_t\}_{t=1}^T$ によって表現される。
- (2) 指示変数の集合 $\{z_t : z_t = r\}$ が与えられた下で、領域 r の入力変数の密度分布パラメータ $\theta_r^\alpha \equiv \{\mu_r, \Sigma_r\}$ が生成される。ここで、密度分布は全共分散型ガウス分布を仮定する。
- (3) θ_r^α が与えられた下で、領域 r に属する入力変数集合 $X_r \equiv \{\mathbf{x}_t : z_t = r\}$ が生成される。
- (4) カーネル関数のパラメータ $\theta_r^{\text{GP}} \equiv \{\psi_r, w_r^2, \eta_r^2\}$ が領域毎に生成される。
- (5) 領域毎に X_r と θ_r^{GP} を使ってグラム行列が計算され、出力変数を並べたベクトル $\mathbf{y}_r \equiv \{y_t : z_t = r\}$ が生成される。

図 3 はこのグラフィカル表現であり、同時分布は、

$$p(\{\mathbf{x}_t, y_t\}_{t=1}^T, \{z_t\}_{t=1}^T, \{\theta_r^{\text{GP}}\}_{r=1}^R, \{\theta_r^\alpha\}_{r=1}^R | \Omega) \quad (5)$$

$$= \prod_{r=1}^R [p(\theta_r^\alpha | \Omega) p(X_r | \theta_r^\alpha) p(\mathbf{y}_r | X_r, \theta_r^{\text{GP}}, \Omega)] \times p(\{z_t\}_{t=1}^T | \Omega)$$

と書ける。 R と Ω はそれぞれ、エキスパートの総数と超パラメータ集合を表す。本稿では、指示変数集合 $\{z_t\}_{t=1}^T$ を積分消去せず直接表現して、時間的な依存関係を考慮する。式 (5) における個々の分布を以下に示す。

$$p(\{z_t\}_{t=1}^T | \Omega) = \frac{\Gamma(\alpha)}{\Gamma(\alpha + T)} \prod_{r=1}^R \frac{\Gamma(T_r + \alpha/R)}{\Gamma(\alpha/R)}$$

$$p(\theta_r^\alpha | \Omega) = \mathcal{N}(\mu_r; \mathbf{m}_0, \Sigma_r / \beta_0) \mathcal{W}(\Sigma_r^{-1}; \mathbf{W}_0, \nu_0)$$

$$p(X_r | \theta_r^\alpha) = \mathcal{N}(X_r; \mu_r, \Sigma_r)$$

$$p(\mathbf{y}_r | X_r, \theta_r^{\text{GP}}, \Omega) = \mathcal{GP}(\mathbf{y}_r; \mathbf{0}, \mathbf{K}_r + \eta_r^2 \mathbf{I}_r) \quad (6)$$

ここで、 $\alpha, T_r, \mathbf{I}_r, \eta_r^2, \mathcal{W}$ はそれぞれ、ディリクレ多項分

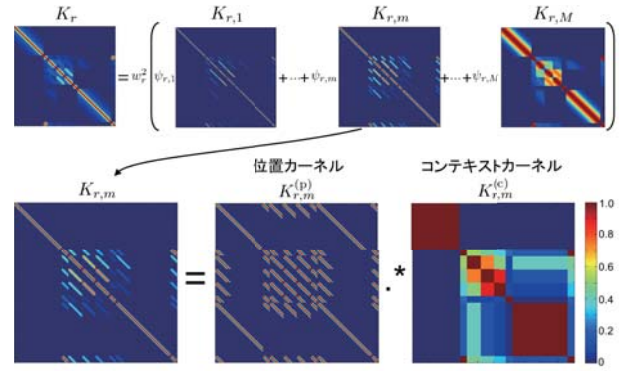


図 4 マルチカーネル学習に基づくグラム行列表現

Fig. 4 Gram matrix based on multiple kernel learning

布の超パラメータ、 X_r の要素数、 $T_r \times T_r$ の単位行列、観測雑音の分散値、ウィシャート分布を示す。

本稿では、マルチカーネル学習に基づき、図 4 のような複数のカーネルの線形結合

$$k_r(\mathbf{x}_i, \mathbf{x}_j) = w_r^2 \sum_{m=1}^M \psi_{r,m} k_{r,m}(\mathbf{x}_i, \mathbf{x}_j) \quad (7)$$

でグラム行列 \mathbf{K}_r を表現し、強度 w_r^2 と各カーネルの優勢度 $\psi_{r,m}$ を推定すべき未知パラメータとみなす [15, 16]。ここで、 $\mathbf{x}_i, \mathbf{x}_j \in X_r$ 、 $\sum_{m=1}^M \psi_{r,m} = 1$ とする。 M は線形結合するカーネル関数の総数である。単位の異なる様々なコンテキスト（音符内位置や音高差、音長など）を扱うため、個々のカーネル関数 $k_{r,m}(\mathbf{x}_i, \mathbf{x}_j)$ を 2 つのカーネル関数の積で表現する [17]。

$$k_{r,m}(\mathbf{x}_i, \mathbf{x}_j) = k_{r,m}^{(p)}(x_i^{(p)}, x_j^{(p)}) k_{r,m}^{(c)}(\mathbf{x}_i^{(c)}, \mathbf{x}_j^{(c)}) \quad (8)$$

ここで、 $k_{r,m}^{(p)}(x_i^{(p)}, x_j^{(p)})$ と $k_{r,m}^{(c)}(\mathbf{x}_i^{(c)}, \mathbf{x}_j^{(c)})$ はそれぞれ、音符内位置の類似を表す位置カーネルと音符の類似度を表すコンテキストカーネルと呼ぶ。入力変数ベクトル \mathbf{x}_i を $x_i^{(p)}, \mathbf{x}_i^{(c)}$ のように位置に関する変数と音符のコンテキストに関する変数に分け、カーネル関数を計算する。位置カーネルは動的変動成分の時間的連続性や周期性を捉えるために、SE 共分散関数と周期共分散関数を利用する。

$$k_{r,m}^{(p)}(x_i^{(p)}, x_j^{(p)}) = \exp\left(-\frac{(x_i^{(p)} - x_j^{(p)})^T (x_i^{(p)} - x_j^{(p)})}{2l_m^{(p)2}}\right)$$

$$k_{r,m}^{(p)}(x_i^{(p)}, x_j^{(p)}) = \exp\left(-2 \sin^2\left(\frac{l_m^{(p)}}{2\pi} (x_i^{(p)} - x_j^{(p)})\right)\right)$$

コンテキストカーネルには SE 共分散関数を利用する。

$$k_{r,m}^{(c)}(\mathbf{x}_i^{(c)}, \mathbf{x}_j^{(c)}) = \exp\left(-\frac{(\mathbf{x}_i^{(c)} - \mathbf{x}_j^{(c)})^T \Lambda (\mathbf{x}_i^{(c)} - \mathbf{x}_j^{(c)})}{2}\right)$$

$$\Lambda^{-1} = \text{diag}(l_{m,1}^{(c)2}, l_{m,2}^{(c)2}, \dots, l_{m,D_c}^{(c)2})$$

まとめると、パラメータと超パラメータをそれぞれ、 $\Theta = \{z_1, \dots, z_T, \theta_1^\alpha, \dots, \theta_R^\alpha, \theta_1^{\text{GP}}, \dots, \theta_R^{\text{GP}}\}$ 、 $\Omega = \{\alpha, \mathbf{m}_0, \mathbf{W}_0, \beta_0, \nu_0, l_1^{(p)}, \dots, l_M^{(p)}, l_{1,1}^{(c)}, \dots, l_{M,D_c}^{(c)}\}$ と定義す

る. D_c は $\mathbf{x}_i^{(c)}$ の次元数を表す.

新たな入力変数 \mathbf{x}_* に対する出力変数 y_* の予測分布を考える. 式 (4) を参考にすると, MoGPEs では,

$$\begin{aligned} p(y_* | \mathcal{D}, \mathbf{x}_*, \Theta, \Omega) & \quad (9) \\ &= \sum_{r=1}^R p(y_* | \mathbf{y}_r, X_r, \mathbf{x}_*, z_* = r, \theta_r^{\text{GP}}) p(z_* = r | \mathbf{x}_*, \theta_r^{\text{x}}) \\ &= \mathcal{N}(y_*; \mu_*, \sigma_*^2) \end{aligned}$$

となり, 右辺の第 1 式が GPR とゲート関数の混合モデルとなることがわかる. ここで,

$$\begin{aligned} \mu_* &= \sum_{r=1}^R c_r \mathbf{k}_{r,*}^T (\mathbf{K}_r + \eta_r^2 \mathbf{I}_r)^{-1} \mathbf{y}_r \\ \sigma_*^2 &= \sum_{r=1}^R c_r^2 (k_r(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_{r,*}^T (\mathbf{K}_r + \eta_r^2 \mathbf{I}_r)^{-1} \mathbf{k}_{r,*}) \\ c_r &= p(z_* = r | \mathbf{x}_*, \theta_r^{\text{x}}) \end{aligned}$$

であり, ゲート関数 c_r はベイズの定理を利用して容易に計算できる. この予測分布に基づいて, F_0 軌跡 $\{f_t\}_{t=1}^{T_m}$ を以下のように音符毎に再合成する.

$$f_t = o_t + c - \bar{o}, \quad o_t = \sum_{n=1}^t \mu_{*,n}, \quad \bar{o} = \sum_{t=1}^{T_m} o_t / T_m \quad (10)$$

ここで, c と T_m はそれぞれ, 楽譜に記載される音符の音高と音長を表す.

GPR はこれまで音声合成 [17, 18] や F_0 生成 [19], 声質変換 [20] などに適用されており, 合成音の観点から性能の向上が確認されている. 提案法はこれらの研究と関連するが, 動的変動成分の動特性を精緻に特徴づけるために, MoGPEs と多様なカーネル関数を利用することが特徴的である. また, MoGPEs を利用すれば, 通常の GPR よりも計算コストを低減できることも特筆すべき点である [12].

3. パラメータの推論

MCMC-EM アルゴリズムを用いて, すべてのパラメータを推論する [21]. 具体的には, $\{z_t\}_{t=1}^T$ と $\{\theta_r^{\text{x}}\}_{r=1}^R$ の推論のためにギブスサンプリングを利用する [12]. また, 式 (6) の多次元ガウス分布を独立なガウス分布の和とみなして, $\{\theta_r^{\text{GP}}\}_{r=1}^R$ の推論のために EM アルゴリズムを利用する [16, 22, 23].

4. 評価実験

新たな入力 (楽譜) に対する出力 (F_0 軌跡の差分信号) の予測性能の観点から提案法を評価する. 「RWC 研究用音楽データベース: ポピュラー音楽」(RWC-MDB-P-2001) [24] における 10 楽曲 (No. 38, 39, 42, 44, 45, 46, 64, 72, 74, 76) の歌声の F_0 軌跡を手作業でラベル付けした結果 [25] を利用する. また, その MIDI 信号を楽譜として利用する.

表 1 RMSE に基づく予測性能の比較

Table 1 Average RMSEs for MoPRs and our model (MoGPEs)

R	10	20	30	40	50
MoPRs	71.8	59.7	73.1	79.6	56.5
MoGPEs	25.0	24.0	23.9	23.1	22.3

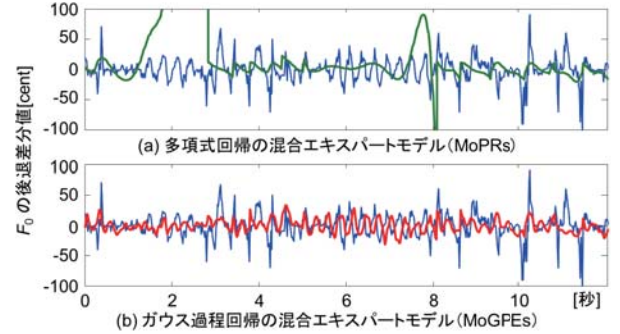


図 5 予測結果の比較

Fig. 5 Comparison of the prediction results: blue, green, and red lines correspond to the actual outputs, the result of the MoPRs, and the result of the MoGPEs, respectively.

これらの楽曲はすべて同一歌手によって歌われたものである. 本来ならばこれらの楽曲の音響信号から F_0 を推定すべきであるが, 今回は提案法の性能の上限を調べるためにこのようなデータを用いた. \mathbf{x}_t と F_0 の後退差分値 y_t は 10 ms ごとに計算され, 無声音のため F_0 が観測されない区間は除去した. 楽曲 No. 38, 39, 42, 44, 45 はモデルパラメータの学習に, 楽曲 No. 46, 64, 72, 74, 76 は予測性能の評価に利用した. 学習データと評価データの量はそれぞれ, 649.3 秒, 625.6 秒である.

指示変数の初期値は, 学習データの入力ベクトル集合 $\{\mathbf{x}_t\}$ を k-means クラスタリングすることによって得られた割当て結果を利用した. 各領域に割当てられた入力ベクトルを利用して, θ_r^{x} を初期化した. $M = 216$ とし, θ_r^{GP} の初期値はそれぞれ, $w_r^2 = 100$, $\psi_{r,1} = 1/M, \dots, \psi_{r,M} = 1/M$, $\eta_r^2 = 10$ ($r = 1, \dots, R$) とした. 超パラメータの設定を脚注に示す*1. パラメータの推論回数は 50 回とした.

評価尺度として, 実際の出力値系列と式 (9) の予測平均値系列間の二乗平均平方根誤差 (Root mean square error, RMSE) を利用した.

$$\text{RMSE} = \sqrt{\sum_{t=1}^{T_t} (y_t - \mu_{*,t})^2 / T_t}, \quad (11)$$

ここで, T_t は評価データの長さを表す. 比較として, 多項式回帰の混合エキスパートモデル (Mixture of polynomial

*1 入力ベクトルの次元数 D は 4 であり, $\alpha = 1$, $\beta_0 = 0.1$, $\nu_0 = D + 1$ に固定する. \mathbf{m}_0 と \mathbf{W}_0 はそれぞれ, 学習データのすべての入力ベクトルの平均ベクトルおよび共分散行列の逆行列を ν_0 で割算した行列に設定する. 位置カーネルの超パラメータとして, SE カーネルは $\{l_m^{(p)}, m = 1, \dots, 108\} \{0.05, 0.11, 0.23, 0.5\}$, 周期カーネルは $\{l_m^{(p)}, m = 109, \dots, 216\} \{0.13, 0.15, 0.17, 0.2\}$ から選択した. コンテキストカーネルの超パラメータは $\{l_{m,1}^{(c)} | 1, 2.2, 5\}$, $\{l_{m,2}^{(c)} | 1, 2.2, 5\}$, $\{l_{m,3}^{(c)} | 0.1, 0.55, 3\}$ の組合せで設定した.

表 2 学習データの量に対する RMSE

Table 2 Average RMSEs for the number of songs in training data

学習データの曲数	1	3	5
MoGPEs	22.3	20.7	20.5

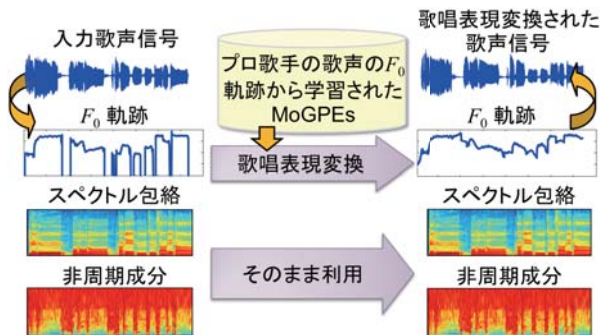


図 6 F_0 軌跡に関する歌唱表現変換

Fig. 6 Singing voice conversion for expression in F_0 contour regression experts, MoPRs) を定義した. 式 (6) の GPR が n 次の多項式回帰に置き換わる.

表 1 はエキスパートの数 R を変化させたときの RMSE の平均値を示す. この結果は楽曲 No. 35 だけを利用してパラメータを学習し, 評価した結果である. MoPRs の多項式の次数は $n = 8$ とした. これより高次の多項式を利用した場合, 多項式の重み係数の推定が安定しなかったためである. RMSE の観点から, 提案法は多項式回帰からなる単純な手法と比較して, 予測性能が上回ることを確認した. また R を増やすにつれて予測性能は向上した. 図 5 は 2 つの手法の予測結果を比較する. 動的変動成分の動特性を表現するために, 周期カーネルを利用する提案法の有効性を確認できた. 表 2 は学習データの量に対する予測性能を検証する. $R = 50$ に固定した. 実験結果より, 学習データを増やすに連れて RMSE が小さくなることを確認した. 予測性能を高めるために, エキスパートの数や学習データ量, 入力コンテキストの種類を今後さらに検討する予定である.

5. 歌唱表現変換への応用

MoGPEs によって学習された歌唱表現を, 別の歌唱者が歌った歌声の F_0 軌跡に転写して, 歌唱表現を変換させることを考える. 図 6 に示すように, 歌声信号における F_0 軌跡のみを加工して再合成する. まず, TANDEM-STRAIGHT [26] を利用して, 歌声信号から F_0 軌跡とスペクトル包絡, 非周期成分を抽出する. 図 7(a) が, 抽出された F_0 軌跡である. 次に, 文献 [9] を利用して, この F_0 軌跡から音符指令列を推定する. 図 7(b) が推定された音符指令列である. 歌唱するメロディの音符の並びが階段状の系列として推定される. この音符指令列から, 式 (1) のコンテキスト情報を 10ms ごとに算出し, 式 (9) における \mathbf{x}_* とする. 式 (9) と (10) を用いて, 音符毎に F_0 軌跡を予測したもの

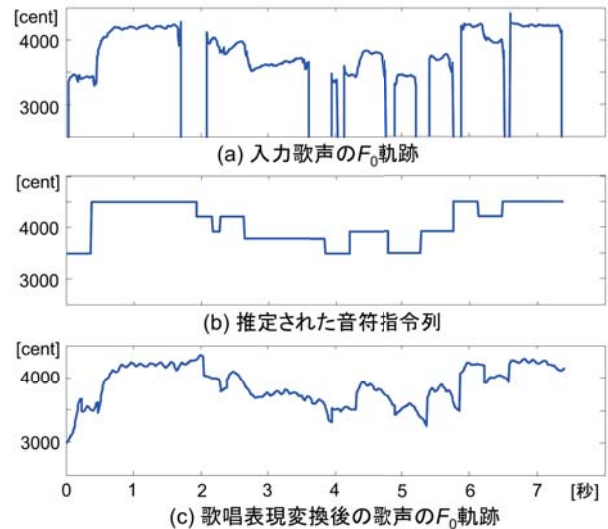


図 7 歌唱表現の変換前後における F_0 軌跡

Fig. 7 F_0 contours before and after singing voice conversion

を変換後の F_0 軌跡とする (図 7(c)). (a) と (c) の軌跡を比較すると, 歌唱表現に起因する動的変動成分が互いに異なることが観測される. 最終的に, TANDEM-STRAIGHT を利用して, 変換された F_0 軌跡と元のスペクトル包絡および非周期成分から歌声信号を再合成すると, 歌唱表現変換が完了する. 学習データの歌唱表現がどの程度転写されたか, 聴取実験を通して検証することが今後の課題である.

6. おわりに

歌声の F_0 軌跡に表れる様々な動的変動成分を特徴抽出して歌唱者の歌唱表現を学習し, 任意の楽譜に対して, その歌唱表現を反映した F_0 軌跡を予測する生成過程モデルを提案した. 機械学習の分野で注目を集めている GPR の混合エキスパートモデルである MoGPEs を利用して, 楽譜のコンテキストと F_0 軌跡の入出力関係を直接的に学習した. 実験結果より, 単純な多項式回帰を用いるよりも GPR の混合エキスパートモデルを利用することの有効性を確認した. 予測性能を高めるために, 学習データを増やすこと, エキスパートの数や超パラメータを調整すること, エキスパートの割当てにディリクレ過程を導入することが挙げられる. 最後に, 提案モデルの応用として, 歌唱表現変換の手順を提案した. 聴取実験を通して, 学習データの歌唱表現がどの程度転写されたか調査することが必要である. 今後の課題は, この枠組を歌唱者認識や歌唱表現認識, 歌声合成に応用することである.

参考文献

- [1] Sundberg, J.: *The Science of the Singing Voice*, Northern Illinois University Press (1987).
- [2] Saitou, T. et al.: Speech-To-Singing synthesis: Converting speaking voices to singing voices by controlling acoustic features unique to singing voices, in *Proc. WASPAA 2007*, pp. 215–218 (2007).

- [3] Regnier, L.: Localization, Characterization and Recognition of Singing Voices, PhD Thesis, IRCAM / UPMC in Paris, France (2013).
- [4] Oura, K. et al.: Pitch adaptive training for HMM-based singing voice synthesis, in *Proc. ICASSP 2012*, pp. 5377–5380 (2012).
- [5] Nose, T. et al.: A Style control technique for singing voice synthesis based on multiple-regression HSMM, in *Proc. INTERSPEECH 2013*, pp. 378–382 (2013).
- [6] Doi, H. et al.: Evaluation of a singing voice conversion method based on many-to-many eigenvoice conversion, in *Proc. INTERSPEECH 2013*, pp. 1067–1071 (2013).
- [7] Nakano, T. et al.: An automatic singing skill evaluation method for unknown melodies using pitch interval accuracy and vibrato features, in *Proc. Interspeech 2006*, pp. 1706–1709 (2006).
- [8] Migita, N. et al.: A study of vibrato features to control singing voices, in *Proc. ICA 2010*, pp. 23–27 (2010).
- [9] Ohishi, Y. et al.: A stochastic model of singing voice F0 contours for characterizing expressive dynamic components, in *Proc. INTERSPEECH 2012* (2012).
- [10] Yoshimura, T. et al.: Simultaneous modeling of spectrum, pitch and duration in HMM-based speech synthesis, in *Proc. EUROSPEECH 1999* (1999).
- [11] Zen, H. et al.: Statistical parametric speech synthesis, *Speech Communication*, Vol. 51, No. 11, pp. 1039–1064 (2009).
- [12] Meeds, E. et al.: An alternative infinite mixture of Gaussian process experts, in *Proc. NIPS 2006* (2006).
- [13] de Cheveigné, A. et al.: YIN, a fundamental frequency estimator for speech and music, *Journal of the Acoustical Society of America*, Vol. 111, No. 4, pp. 1917–1930 (2002).
- [14] Rasmussen, C. E. et al.: Infinite mixtures of Gaussian process experts, in *Proc. NIPS 2002* (2002).
- [15] Bach, F. et al.: Multiple kernel learning, conic duality, and the SMO algorithm, in *Proc. ICML 2004*, pp. 6–13.
- [16] Yoshii, K. et al.: Infinite kernel linear prediction for joint estimation of spectral envelope and fundamental frequency, in *Proc. ICASSP 2013*, pp. 463–467 (2013).
- [17] Koriyama, T. et al.: Statistical nonparametric speech synthesis based on Gaussian process regression, in *Proc. INTERSPEECH 2013*, pp. 912–916 (2013).
- [18] Henter, G. E. et al.: Gaussian process dynamical models for nonparametric speech representation and synthesis, in *Proc. ICASSP 2012*, pp. 4505–4508 (2012).
- [19] Fernandez, R. et al.: F0 contour prediction with a deep belief network-Gaussian process hybrid model, in *Proc. ICASSP 2013*, pp. 6885–6889 (2013).
- [20] Pilkington, N. C. V. et al.: Gaussian process experts for voice conversion, in *Proc. INTERSPEECH 2011*, pp. 2761–2764 (2011).
- [21] Andrieu, C. et al.: An introduction to MCMC for machine learning, *Machine Learning*, Vol. 50, No. 1-2, pp. 5–43 (2003).
- [22] Feder, M. et al.: Parameter estimation of superimposed signals using the EM algorithm, *IEEE Transactions on Acoustics, Speech, and Signal Processing*, Vol. 36, No. 4, pp. 477–489 (1988).
- [23] 亀岡弘和ほか：マルチカーネル線形予測モデルによる音声分析，音講論集，2-Q-24，pp. 499–502 (2010).
- [24] Goto, M. et al.: RWC music database: Popular, classical, and jazz music databases, in *Proc. ISMIR 2002* (2002).
- [25] Goto, M.: AIST annotation for the RWC music database, in *Proc. ISMIR 2006* (2006).
- [26] Kawahara, H. et al.: TANDEM-STRAIGHT: A temporally stable power spectral representation for periodic signals and applications to interference-free spectrum, F0, and aperiodicity estimation, in *Proc. ICASSP 2008*, pp. 3933–3936 (2008).