

全極型声道モデルと F_0 パターン生成過程モデルを内部にもつ 統計的音声スペクトルモデル

亀岡弘和[†]

[†] 日本電信電話株式会社 NTT コミュニケーション科学基礎研究所

〒 243-0198 神奈川県厚木市森の里若宮 3-1

E-mail: †kameoka@cs.brl.ntt.co.jp

あらまし 音声情報処理研究の歴史において、音声生成過程に関する重要な数理モデルとして、声帯振動の制御機構を模した藤崎モデルと、縦続連結した等長音響管で声道断面積関数をモデル化したことと等価な全極型声道モデルが有名である。本発表では、ある特定条件のウェーブレット変換のスペクトル領域において、これら両生成過程モデルを同時に内包する確率モデルが立てられることを示し、そのパラメータの推論法について述べる。提案法は、全極型声道モデルと藤崎モデルのパラメータをスペクトルから同時に直接推定することが可能な初めての音声分析手法であり、様々な興味深い応用展開が期待される。

キーワード 全極型声道モデル, F_0 パターン生成過程モデル (藤崎モデル), ウェーブレット変換, EM アルゴリズム

Statistical speech spectrum model incorporating all-pole vocal tract model and F_0 contour generating process model

Hirokazu KAMEOKA[†]

[†] NTT Communication Science Laboratories

3-1 Morinosato Wakamiya, Atsugi, Kanagawa, 243-0198 Japan

E-mail: †kameoka@cs.brl.ntt.co.jp

Abstract In this paper, we propose to introduce a statistical speech spectrum model, simultaneously incorporating the all-pole vocal tract model and the F_0 contour generating process model.

Key words All-pole vocal tract model, F_0 contour generating process model, wavelet transform, EM algorithm

1. まえがき

音声は、肺から供給された空気が声帯振動によって周期的な疎密波となり、さらにその音波が声道で共振されることによって「声」となり、口から発せられる。人は音声対話において、声帯の固有振動数を巧みに制御しながら音声に適切な抑揚を与え、同時に、声道の形状を巧みに変化させながら異なる種類の母音を生成することで、聴き手に意図した情報を伝達している。音声情報処理研究の歴史において、音声生成過程に関する重要な数理モデルとして、声帯振動の制御機構を模した藤崎モデルと、縦続連結した等長音響管で声道断面積関数をモデル化したことと等価な全極型声道モデルが有名である。本稿では、全極型声道モデルと藤崎モデルを同時に内包する統計的音声スペクトルモデルを提案する。

声道スペクトルと基本周波数 (F_0) は、分析、合成、認識、強調、符号化などの幅広い音声アプリケーションにおいて根幹となる音声の重要な特徴量であり、音声信号からこれらを自動的に推定する技術は長年にわたって研究され、発展を遂げてきた。声道スペクトルを推定する問題と F_0 を推定する問題は、一見

すると別の問題のようで実は密接に関係している。例えば、以下、音声信号が声道のインパルス応答と声帯による周期的なパルス波との畳み込みで表現できる場合を考えると、声道スペクトルは調波成分のピークをつないだ包絡線に対応するので、 F_0 がもし既知であればどの成分が調波成分であるかが特定でき、声道スペクトルを精度良く推定することができる。一方で、 F_0 推定では、実際の F_0 の整数分の一倍と推定し誤ってしまうことがしばしば問題となるが、この半ピッチ誤りの問題は、声道スペクトルが滑らかであるという先験的な知識を使えば効果的に回避することができる (半ピッチ誤りの状況では滑らかでないスペクトル包絡を仮定してしまっていることに相当するため)。また、音声の F_0 の時間軌跡 (F_0 パターンと呼ぶ。) も通常は滑らかであり、これを制約とすれば同様の推定誤りを回避することが多い。以上のように、声道スペクトルと F_0 パターンは各々の推定において互いに有用な情報であるため、別個に求めるより、同時最適化問題として推定できるようにした方が多分に合理的である。以上の動機に基づき、スペクトル包絡構造と F_0 パターンを同時モデル化したことが提案モデルにおける第一のポイントである。

提案モデルの第二の重要なポイントは、音声生成過程の確率モデル化を通して、与えられた任意の信号がどれだけ音声のものらしいかを定量的に評価することが可能になる点である。音声生成過程から生成しえない信号には低いスコアを、音声生成過程から生成しうる信号には高いスコアを与えるような規準がもしあれば、音声分離、雑音除去、残響除去などの音声を対象としたあらゆる逆問題において音声を抽出する手掛かりとして大いに活用できることが期待されるからである。

以上より、スペクトル包絡構造と F_0 パターンをいかにして同時モデル化するかという点と、いかに音声の生成過程を正確に統計モデルの形で表現できるかという点に主眼を置き、以下で音声スペクトルのモデル化を行う。

2. 音声パワースペクトルモデル

本章では音声信号モデルとして、 n 次調波成分の瞬時振幅が $w_n(u)$ 、瞬時位相が $n\theta(u) + \rho_n$ の擬似周期信号の解析信号表現

$$s(u) = \sum_{n=1}^N w_n(u) e^{j(n\theta(u) + \rho_n)}, \quad w_n(u) \geq 0 \quad (1)$$

を考え (u は時刻、 N は調波成分の最大次数)、この信号のウェーブレット変換スペクトルを導く。周波数が 1 の、アドミッシブル条件を満たす適当なマザーウェーブレットを $\psi(u)$ とし、ウェーブレット基底関数を

$$\psi_{x,t}(u) := \frac{1}{\sqrt{2\pi}e^{-x}} \psi\left(\frac{u-t}{e^{-x}}\right) \quad (2)$$

と定義する。 e^{-x} と t はスケールとシフトパラメータに該当し、 $\psi_{x,t}(u)$ は時刻 t 周辺に局在する周期 e^{-x} の成分を測るためのウェーブレットである。 x は対数周波数に対応する。これを用いて、 s の連続ウェーブレット変換を

$$S_w(x, t) = \langle s, \psi_{x,t} \rangle = \int_{-\infty}^{\infty} s(u) \psi_{x,t}^*(u) du \quad (3)$$

と定義する。ただし、 $\langle \cdot, \cdot \rangle$ は内積を表す。

ここで、時刻 t におけるウェーブレット変換スペクトルについて着目する。音声信号 s において、 θ と w_n の時間変化はいずれも十分緩やかであると仮定し、これらを時刻 t 周辺で

$$w_n(u) \simeq w_n(t), \quad \theta(u) \simeq \theta(t) + \dot{\theta}(t)(u-t) \quad (4)$$

のように 0 次近似および 1 次近似する。 $\dot{\theta}$ は θ の導関数 (時間微分) であり、 s の瞬時基本周波数を表す。式 (4) を式 (1) に代入すると、時刻 t 周辺の s を近似した信号

$$s_t(u) = \sum_{n=1}^N w_n(t) e^{j(n\dot{\theta}(t)(u-t) + \xi_n(t))} \quad (5)$$

を得る。ただし、 $\xi_n(t) := n\theta(t) + \rho_n$ である。 s と s_t は大域的にはまったく異なる信号であるが、時刻 t 周辺においてのみ似通っている点がポイントである。通常ウェーブレット基底関数 $\psi_{x,t}$ は時刻 t 周辺でのみ優勢であるため、時刻 t におけるウェーブレット変換スペクトルは、式 (5) を式 (3) に代入して得られる s_t のウェーブレット変換

$$S_w(x, t) \simeq \langle s_t, \psi_{x,t} \rangle \quad (6)$$

$$= \sum_{n=1}^N w_n(t) e^{j\xi_n(t)} \Psi^*(ne^{-x+\Omega(t)}) \quad (7)$$

により近似できる。ただし、 Ψ はマザーウェーブレット ψ のフーリエ変換である。また、 $\Omega(t) := \log \dot{\theta}(t)$ であり、 $\Omega(t)$ は対数瞬時基本周波数を表す。以上より、ウェーブレット変換のスペクトルはマザーウェーブレット $\psi(t)$ を具体的にどう定めるかによって決まる。ここで、周波数特性 $\Psi(\omega)$ が、次のような $\omega = 1$ で最大値をとる単峰的な実関数

$$\Psi(\omega) := \begin{cases} e^{-\frac{1}{4\sigma^2}(\log \omega)^2} & (\omega > 0) \\ 0 & (\omega \leq 0) \end{cases} \quad (8)$$

となるようなマザーウェーブレットを選ぶと、 S_w は

$$S_w(x, t) = \sum_{n=1}^N w_n(t) e^{j\xi_n(t)} e^{-\frac{1}{4\sigma^2}(x-\Omega(t)-\log n)^2} \quad (9)$$

と具体的に書ける。このマザーウェーブレットによって構成される定 Q フィルタバンクを定 Q 対数正規分布型フィルタバンクと呼ぶことにする。ここで、 $|S_w(x, t)|^2$ は

$$|S_w(x, t)|^2 = \lambda(x, t) + o(x, t) \quad (10)$$

$$\lambda(x, t) := \sum_{n=1}^N w_n(t)^2 e^{-\frac{1}{2\sigma^2}(x-\Omega(t)-\log n)^2} \quad (11)$$

と書ける。 $o(x, t)$ は調波間干渉を表す交差項であり、調波成分同士の重なりが少ないほど小さくなる。

周期 u_0 でサンプリングされた実際の観測音声信号に対して上記と同条件で算出したウェーブレット変換のパワースペクトルを $Y[k, l] := Y(x_k, t_l)$ (ただし k は対数周波数インデックス、 l は離散時刻インデックス) とすると、 $w_n(t)$ と $\Omega(t)$ が真の値と等しい場合には以上より

$$Y[k, l] \simeq \lambda[k, l] \quad (12)$$

となることが期待される。ただし、 $\lambda[k, l] := \lambda(x_k, t_l)$ である。 $w_n(t)$ と $\Omega(t)$ が真値と等しくともなお生じる $Y[k, l]$ と $\lambda[k, l]$ との間の誤差にはさまざまな要因があり、例えば、式 (1) の形で定義した擬似周期性の仮定からの逸脱、式 (4) の近似、交差項 o の存在、などが挙げられる。提案法の枠組では、これらの複合的な誤差要因を詳細にモデル化することはせず、まとめて一挙に確率的な現象と見なし、

$$Y[k, l] \sim \text{Pois}(\lambda[k, l]) \quad (13)$$

と具体的に仮定する。ただし、 Pois はポアソン分布を表す。この仮定の下での λ の最尤推定問題は、スペクトル間の近さを測る尺度の一つとして近年音響信号処理分野で多用される I ダイバージェンスと呼ぶ歪み尺度を規準とした Y と λ の最適フィッティング問題と等価となる。

以上のパワースペクトルモデル $\lambda[k, l]$ において推定すべき未知変数は、瞬時振幅 $w_n[l] := w_n(t_l)$ および対数基本周波数 $\Omega[l] := \Omega(t_l)$ である。提案法のポイントは、これら両未知変数の事前分布を、それぞれ音声生成過程のモデルとして古典的で有名な全極型声道モデル [1] と F_0 パターン生成過程モデル [2] に基づいて設計し、スペクトルの生成過程を階層モデルにより表現する点にある。以下、各事前分布の設計方法を順に述べる。

3. 全極モデルに基づく瞬時振幅の確率分布

本章では、全極型声道モデルに基づき $w_n[l]$ の事前分布を導く。時刻 t_l における式 (5) の音声信号モデルの、観測信号と同じサンプリング周期 u_0 の下での離散時間表現を $s_l[i] := s_{t_l}(iu_0)$

とし、これを全極型システムからの出力

$$s_l[i] = \sum_{m=0}^M a_l[m] s_l[i-m] + \epsilon_l[i] \quad (14)$$

と考えよう。 $a_l[0], \dots, a_l[M]$ は全極モデルのパラメータであり、線形予測モデルにおいて予測係数と呼ぶものに該当する。ここで、 $s_l[i]$ の基本周波数は $e^{\Omega[l]}$ であるため、声帯音源に対応する入力 $\epsilon_l[i]$ も基本周波数が $e^{\Omega[l]}$ の周期信号

$$\epsilon_l[i] = \sum_{n=1}^N v_{l,n} e^{j n e^{\Omega[l]} i u_0} \quad (15)$$

と仮定しておく必要がある（線形フィルタには入力の周波数成分の周波数値を変える効力はないため）。また、上記声帯音源モデルにおける複素振幅 $v_{l,n}$ について

$$v_{l,n} \sim \mathcal{N}_C(0, 1), \quad n = 1, \dots, N \quad (16)$$

を仮定する。ただし、 \mathcal{N}_C は複素正規分布を表す。よって、この仮定は、駆動音源に関してすべての調波成分が等しいパワースペクトル密度をもち、位相を一様にとりうることを仮定していることに相当する。式 (14) の両辺に離散時間フーリエ変換 (DTFT) を行い、 $A_l(z) := a_l[0] - a_l[1]z^{-1} \dots - a_l[M]z^{-M}$ と置いて整理すると、

$$S_l(\omega) = \frac{\sqrt{2\pi}}{A_l(e^{j\omega})} \sum_{n=1}^N v_{l,n} \delta(\omega - n e^{\Omega[l]} u_0) \quad (17)$$

となる。ただし、 S_l は $s_l[1], \dots, s_l[T]$ の DTFT、 ω は規格化角周波数である。上式に対し逆 DTFT を行くと $s_l[i]$ の陽な表現

$$\begin{aligned} s_l[i] &= \frac{1}{\sqrt{2\pi}} \int_{-\pi}^{\pi} \frac{\sqrt{2\pi}}{A_l(e^{j\omega})} \sum_{n=1}^N v_{l,n} \delta(\omega - n e^{\Omega[l]} u_0) e^{j\omega i} d\omega \\ &= \sum_{n=1}^N \frac{v_{l,n}}{A_l(e^{j n e^{\Omega[l]} u_0})} e^{j n e^{\Omega[l]} i u_0} \end{aligned} \quad (18)$$

を得る。式 (5) を離散時間表現したもの

$$s_l[i] = \sum_{n=1}^N w_n[l] e^{j(n e^{\Omega[l]}(i u_0 - t_l) + c_n(t_l))} \quad (19)$$

と式 (18) を比較すると、

$$w_n[l] = \left| \frac{v_{l,n}}{A_l(e^{j n e^{\Omega[l]} u_0})} \right|, \quad n = 1, \dots, N \quad (20)$$

なる関係式を導ける。 $e^{\Omega[l]} u_0$ は基本周波数 $e^{\Omega[l]}$ を規格化角周波数の単位で表わしたもので、 $1/|A_l(e^{j e^{\Omega[l]} u_0})|, \dots, 1/|A_l(e^{j N e^{\Omega[l]} u_0})|$ は全極スペクトルを基本周波数の間隔でサンプリングした値に他ならない。以上より、全極モデルパラメータおよび $\Omega[l]$ が与えられた下での $w_n[l]$ が従う分布を具体的に導くことができる。式 (16) より

$$\frac{v_{l,n}}{A_l(e^{j n e^{\Omega[l]} u_0})} \sim \mathcal{N}_C(0, 1/|A_l(e^{j n e^{\Omega[l]} u_0})|^2) \quad (21)$$

が言え、実部と虚部が独立でかつ等しい分散の正規分布に従う複素数の絶対値はレイリー分布に従うため、式 (20), (21) より

$$w_n[l] \sim \text{Rayleigh}(1/|A_l(e^{j n e^{\Omega[l]} u_0})|) \quad (22)$$

が得られる。この事前分布のパラメータ (ハイパーパラメータ) は、 $\pi_w = \{a_l[m] \mid 1 \leq l \leq L, 0 \leq m \leq M\}$ である。この事前分布は、 $\Omega[l]$ に依存する点に注意が必要である。

4. 藤崎モデルに基づく F_0 パターンの確率分布

本章では、藤崎モデルに基づいて $\Omega := (\Omega[1], \dots, \Omega[L])^T$ の事前分布を導く。藤崎モデルは喉頭の声帯振動制御機構の物理モデルであるが、これを離散時間表現の確率過程として記述した確率モデルが最近筆者らによって提案された [4]。本章の内容はこの新モデルがベースになっている。

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、声帯の伸びと対数 F_0 の変化が比例関係にあるという仮定に基づき F_0 パターンがモデル化される。甲状軟骨の平行移動運動に関係する F_0 パターンをフレーズ成分、回転運動に関係する F_0 パターンをアクセント成分と呼び、それぞれ $\Omega_p(t)$, $\Omega_a(t)$ とする。 $\Omega_p(t)$ の生成過程 (フレーズ制御機構) はフレーズ指令と呼ぶパルス波を入力とした臨界制動の二次線形系、 $\Omega_a(t)$ の生成過程 (アクセント制御機構) はアクセント指令と呼ぶ矩形波を入力とした臨界制動の二次線形系により表現される。以上の二つの成分と、声帯の物理的性質によって決まるベースライン成分 Ω_b の和 $\Omega_p(t) + \Omega_a(t) + \Omega_b$ で F_0 パターン $\Omega(t)$ を表したものが藤崎モデルである。フレーズ制御機構およびアクセント制御機構は

$$G_p(t) = \alpha^2 t e^{-\alpha t} \quad (t \geq 0) \quad (23)$$

$$G_a(t) = \beta^2 t e^{-\beta t} \quad (t \geq 0) \quad (24)$$

と与えられるインパルス応答 $G_p(t)$ と $G_a(t)$ によって特徴づけられる。 α と β はそれぞれの制御機構の固有角周波数である。フレーズ指令とアクセント指令をそれぞれ $u_p(t)$ と $u_a(t)$ のように時間の関数として表せば、 $\Omega(t)$ は $G_p(t)$ と $G_a(t)$ を用いて

$$\Omega(t) = \Omega_b + \Omega_p(t) + \Omega_a(t) \quad (25)$$

$$\Omega_p(t) = (G_p * u_p)(t) \quad (26)$$

$$\Omega_a(t) = (G_a * u_a)(t) \quad (27)$$

と表される。ただし、 $*$ は畳み込みを表す。

藤崎モデルの離散時間表現を得るため、連続時間システム of フレーズ制御機構およびアクセント制御機構に対して後退差分変換により離散化を行う。まず、フレーズ制御機構の伝達関数 (ラプラス変換) は $\mathcal{G}_p(s) = \mathcal{L}[G_p(t)] = \alpha^2 / (s + \alpha)^2$ である。後退差分変換は、時間微分演算子 s を z 領域における後退差分演算子 $s \simeq (1 - z^{-1})/t_0$ に置き換える変換であり (ただし、 $t_0 = t_l - t_{l-1}$)、この変換により、 $\mathcal{G}_p^{-1}(s)$ の逆システムの伝達関数を z 領域で $\mathcal{G}_p^{-1}(z) = g_{p2} z^{-2} + g_{p1} z^{-1} + g_{p0}$ と書ける。ただし、 $g_{p2} = (\phi - 1)^2$, $g_{p1} = -2\phi(\phi - 1)$, $g_{p0} = \phi^2$, $\phi = 1 + 1/(\alpha t_0)$ である。ここで、時刻 t_l におけるフレーズ指令関数およびフレーズ成分をそれぞれ $u_p[l]$, $\Omega_p[l]$ と表すと、 $\Omega_p[l]$ は単一のパラメータ ϕ (または α) によって特性が決まる拘束つき全極モデルからの出力

$$u_p[l] = g_{p0} \Omega_p[l] + g_{p1} \Omega_p[l-1] + g_{p2} \Omega_p[l-2] \quad (28)$$

と見なすことができる。同様に、アクセント指令関数 $u_a[l]$ とアクセント成分 $\Omega_a[l]$ の関係も

$$u_a[l] = g_{a0} \Omega_a[l] + g_{a1} \Omega_a[l-1] + g_{a2} \Omega_a[l-2] \quad (29)$$

と書くことができる。ただし、 $g_{a2} = (\varphi - 1)^2$, $g_{a1} = -2\varphi(\varphi - 1)$, $g_{a0} = \varphi^2$, $\varphi = 1 + 1/(\beta t_0)$ である。ベースライン成分 $\Omega_b(t)$ の離散時間表現を $\Omega_b[l]$ とすると、藤崎モデルの離散時間表現はこれらの三成分の和 $\Omega[l] = \Omega_p[l] + \Omega_a[l] + \Omega_b[l]$ で与えられる。

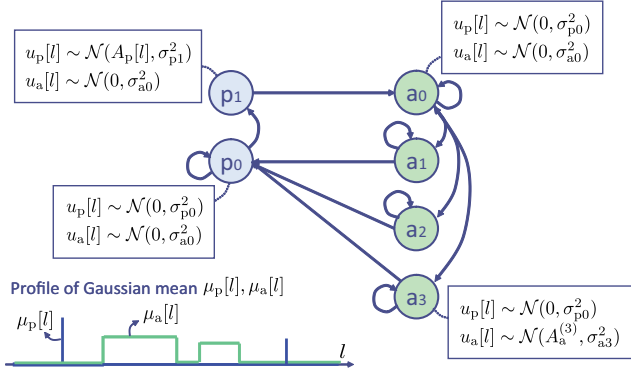


図1 Command function modeling with HMM.

次に、 $u_p[l]$ と $u_a[l]$ を確率モデル化する。そこで、藤崎モデルにおいてフレーズ指令とアクセント指令に関して以下のような規則が定められていることに注意する。

(A1) フレーズ指令はパルス波、アクセント指令は矩形波である。(A2) 複数の指令が同時に発生することはない。

藤崎モデルのパラメータ推定における難しさの一つは、これらの制約の下で最適推定をいかに行えるかという点にある。提案法では、これを解決するため隠れマルコフモデル (HMM) を用いて $u_p[l]$ と $u_a[l]$ をモデル化する。 $u_p[l]$ と $u_a[l]$ を組にしたベクトル $\mathbf{o}[l] := (u_p[l], u_a[l])^T$ を考え、 $\mathbf{o}[l]$ を

$$\mathbf{o}[l] = \boldsymbol{\nu}[l] + \boldsymbol{\epsilon}[l] \quad (30)$$

のようにフレーズ指令からなる時系列 $\mu_p[l]$ とアクセント指令からなる時系列 $\mu_a[l]$ を組にしたベクトル $\boldsymbol{\nu}[l] = (\mu_p[l], \mu_a[l])^T$ に雑音 $\boldsymbol{\epsilon}[l]$ が加わったものと見なす。 $\boldsymbol{\epsilon}[l]$ を共分散行列が $\mathbf{\Upsilon}[l] = \text{diag}(v_p^2[l], v_a^2[l])$ のガウス性雑音とすると、以上より

$$\mathbf{o}[l] \sim \mathcal{N}(\boldsymbol{\nu}[l], \mathbf{\Upsilon}[l]) \quad (31)$$

が言える。ここで、 $\boldsymbol{\nu}[l]$ および $\mathbf{\Upsilon}[l]$ が状態遷移に伴って変化する図1のようなHMMを考える。HMMによる $\mathbf{o}[l]$ のモデル化により、状態遷移の経路制限を通して両指令の組 $\boldsymbol{\nu}[l]$ に対して(A1), (A2)を満たすような制約を巧く与えることが可能となる。提案するHMMの構成は以下のとおりである。

出力系列: $\{\mathbf{o}[l]\}_{l=1}^L$
 状態集合: $\mathcal{S} := \{p_0, p_1, a_0, \dots, a_N\}$
 状態系列: $\{s_l\}_{l=1}^L$
 出力分布: $P(\mathbf{o}[l] | s_l = j) = \mathcal{N}(\mathbf{c}_j[l], \mathbf{D}_j)$
 $\mathbf{c}_{p_0}[l] = \mathbf{c}_{a_0}[l] = (0, 0)^T$
 $\mathbf{c}_{p_1}[l] = (A_p[l], 0)^T$, $\mathbf{c}_{a_n}[l] = (0, A_a^{(n)})^T$
 $\mathbf{D}_{p_0} = \mathbf{D}_{a_0} = \text{diag}(\sigma_{p_0}^2, \sigma_{a_0}^2)$
 $\mathbf{D}_{p_1} = \text{diag}(\sigma_{p_1}^2, \sigma_{a_0}^2)$, $\mathbf{D}_{a_n} = \text{diag}(\sigma_{p_0}^2, \sigma_{a_n}^2)$
 遷移確率: $\phi_{j', j} := \log P(s_l = j | s_{l-1} = j')$

簡単のため状態遷移確率および分散 $\sigma_{p_0}^2, \sigma_{p_1}^2, \sigma_{a_0}^2, \dots, \sigma_{a_N}^2$ を定数とすると、以上より指令入力モデルにおいて推定すべき未知パラメータは、状態遷移系列 s_l 、フレーズ指令の大きさ $A_p[l]$ 、アクセント指令の大きさ $\{A_a^{(n)}\}_{n=1}^N$ であり、これらをまとめて π_w と記す。また、平均値系列 $\{\mu_p[l]\}_{l=1}^L$ 、 $\{\mu_a[l]\}_{l=1}^L$ および分散値系列 $\{v_p[l]\}_{l=1}^L$ 、 $\{v_a[l]\}_{l=1}^L$ は、状態遷移系列 $\{s_l\}_{l=1}^L$ がひとたび決まれば $(\mu_p[l], \mu_a[l])^T \leftarrow \mathbf{c}_{s_l}[l]$, $\text{diag}(v_p[l], v_a[l]) \leftarrow \mathbf{D}_{s_l}$ ($l = 1, \dots, L$) により一挙に決まる。

前述の指令入力モデルに基づき $\boldsymbol{\Omega} = (\boldsymbol{\Omega}[1], \dots, \boldsymbol{\Omega}[L])^T$ の確

率密度関数を導く。式 (31) より、 $\mathbf{u}_p := (u_p[1], \dots, u_p[L])^T$, $\mathbf{u}_a := (u_a[1], \dots, u_a[L])^T$, $\boldsymbol{\mu}_p := (\mu_p[1], \dots, \mu_p[L])^T$, $\boldsymbol{\mu}_a := (\mu_a[1], \dots, \mu_a[L])^T$ とすると、

$$\mathbf{u}_p \sim \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p), \quad \boldsymbol{\Sigma}_p = \text{diag}(v_p^2[1], \dots, v_p^2[L]) \quad (32)$$

$$\mathbf{u}_a \sim \mathcal{N}(\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a), \quad \boldsymbol{\Sigma}_a = \text{diag}(v_a^2[1], \dots, v_a^2[L]) \quad (33)$$

が言える。式 (28), (29) の関係式より、フレーズ成分 $\boldsymbol{\Omega}_p := (\boldsymbol{\Omega}_p[1], \dots, \boldsymbol{\Omega}_p[L])^T$ と \mathbf{u}_p の関係、および、アクセント成分 $\boldsymbol{\Omega}_a := (\boldsymbol{\Omega}_a[1], \dots, \boldsymbol{\Omega}_a[L])^T$ と \mathbf{u}_a の関係は、

$$\mathbf{G}_p := \begin{bmatrix} g_{p_0} & & & & & & O \\ g_{p_1} & g_{p_0} & & & & & \\ g_{p_2} & g_{p_1} & g_{p_0} & & & & \\ & \ddots & \ddots & \ddots & & & \\ O & & & g_{p_2} & g_{p_1} & g_{p_0} & \end{bmatrix}, \quad \mathbf{G}_a := \begin{bmatrix} & & & & & & O \\ g_{a_0} & & & & & & \\ g_{a_1} & g_{a_0} & & & & & \\ g_{a_2} & g_{a_1} & g_{a_0} & & & & \\ & \ddots & \ddots & \ddots & & & \\ O & & & g_{a_2} & g_{a_1} & g_{a_0} & \end{bmatrix}$$

と置くと、それぞれ $\mathbf{u}_p = \mathbf{G}_p \boldsymbol{\Omega}_p$, $\mathbf{u}_a = \mathbf{G}_a \boldsymbol{\Omega}_a$ のように表現できることから、式 (32), (33) より

$$\boldsymbol{\Omega}_p \sim \mathcal{N}(\mathbf{G}_p^{-1} \boldsymbol{\mu}_p, \mathbf{G}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{G}_p^{-1})^T) \quad (34)$$

$$\boldsymbol{\Omega}_a \sim \mathcal{N}(\mathbf{G}_a^{-1} \boldsymbol{\mu}_a, \mathbf{G}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{G}_a^{-1})^T) \quad (35)$$

が導かれる。ベース成分 $\boldsymbol{\Omega}_b[l]$ についても、同様に $\boldsymbol{\Omega}_b \sim \mathcal{N}(\boldsymbol{\mu}_b \mathbf{1}, \boldsymbol{\Sigma}_b)$ とし、 $\boldsymbol{\Omega}_p$, $\boldsymbol{\Omega}_a$ とは独立と仮定する。ただし、 $\boldsymbol{\Sigma}_b$ は定数とする。仮定より、 $\boldsymbol{\Omega}_p$, $\boldsymbol{\Omega}_a$, $\boldsymbol{\Omega}_b$ は独立なので、藤崎モデルパラメータ $\pi_w := \{\pi_w, \alpha, \beta, \mu_b\}$ が与えられたもとの F_0 パターン $\boldsymbol{\Omega} = \boldsymbol{\Omega}_p + \boldsymbol{\Omega}_a + \boldsymbol{\Omega}_b$ の確率密度関数は、

$$\boldsymbol{\Omega} \sim \mathcal{N}(\mathbf{G}_p^{-1} \boldsymbol{\mu}_p + \mathbf{G}_a^{-1} \boldsymbol{\mu}_a + \mu_b \mathbf{1}, \mathbf{G}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{G}_p^{-1})^T + \mathbf{G}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{G}_a^{-1})^T + \boldsymbol{\Sigma}_b) \quad (36)$$

で与えられる。以上より、以下を得る。

$$P(\boldsymbol{\Omega} | \pi_w) = \frac{|\boldsymbol{\Sigma}^{-1}|^{1/2}}{(2\pi)^{T/2}} \exp \left\{ -\frac{1}{2} (\boldsymbol{\Omega} - \boldsymbol{\mu}_\Omega)^T \boldsymbol{\Sigma}_\Omega^{-1} (\boldsymbol{\Omega} - \boldsymbol{\mu}_\Omega) \right\}$$

$$\boldsymbol{\mu}_\Omega = \mathbf{G}_p^{-1} \boldsymbol{\mu}_p + \mathbf{G}_a^{-1} \boldsymbol{\mu}_a + \mu_b \mathbf{1} \quad (37)$$

$$\boldsymbol{\Sigma}_\Omega = \mathbf{G}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{G}_p^T)^{-1} + \mathbf{G}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{G}_a^T)^{-1} + \boldsymbol{\Sigma}_b$$

5. パラメータ推定アルゴリズム

以上のモデルをグラフィカルモデルで表わすと図2のようになる。 $\boldsymbol{\Omega}$, $w := \{w_n[l] | 1 \leq n \leq N, 1 \leq l \leq L\}$, π_w , π_w が未知のパラメータであり、観測スペクトル $Y := \{Y[k, l] | 1 \leq k \leq K, 1 \leq l \leq L\}$ を最も良く説明するような $\boldsymbol{\Omega}$, w , π_Ω , π_w の値を推定することが目的である。そこで、モデルパラメータの事後確率

$$p(\boldsymbol{\Omega}, w, \pi_\Omega, \pi_w | Y) \propto p(Y | \boldsymbol{\Omega}, w) p(\boldsymbol{\Omega}, w, \pi_\Omega, \pi_w) \\ = p(Y | \boldsymbol{\Omega}, w) p(\boldsymbol{\Omega} | \pi_\Omega) p(w | \boldsymbol{\Omega}, \pi_w) p(\pi_\Omega) \quad (38)$$

を $\boldsymbol{\Omega}$, w , π_Ω , π_w に関して最大化する最大事後確率 (MAP) 推定問題を考える。2~4. より各確率は

$$p(Y | \boldsymbol{\Omega}, w) = \prod_{k, l} \frac{\lambda[k, l]^{Y[k, l]} e^{-\lambda[k, l]}}{Y[k, l]!} \quad (39)$$

$$p(\boldsymbol{\Omega} | \pi_\Omega) = \frac{|\boldsymbol{\Sigma}_\Omega^{-1}|^{1/2}}{(2\pi)^{L/2}} e^{-\frac{1}{2} (\boldsymbol{\Omega} - \boldsymbol{\mu}_\Omega)^T \boldsymbol{\Sigma}_\Omega^{-1} (\boldsymbol{\Omega} - \boldsymbol{\mu}_\Omega)} \quad (40)$$

$$p(w | \boldsymbol{\Omega}, \pi_w) = \prod_{n, l} \frac{w_n[l]}{P_n[l]} e^{-\frac{w_n[l]}{P_n[l]}} \quad (41)$$

で与えられる。ただし、 $\lambda[k, l]$ は式 (11) のとおりであり、 $P_n[l] := 1/|A_l(e^{jne^{\Omega l} u_0})|^2$ である。以上の MAP 推定問題の大域最適解は解析的に解くことはできないが、以下に示すように Expectation-Maximization(EM) アルゴリズムにより局所解を得ることができる。

EM アルゴリズムの導出では、潜在データと呼ぶ、観測不可な架空のデータをどのように定義するかがポイントとなる。ここでは、観測スペクトル $Y[k, l]$ のうち n 番目の高調波に対応しているであろう成分 $C_n[k, l]$ を潜在データと見なす。従って、

$$Y[k, l] = \sum_n C_n[k, l] \quad (42)$$

を満たす。今、我々は $Y[k, l] \sim \text{Pois}(\lambda[k, l])$ を仮定しているため、これと矛盾しないように $C_n[k, l]$ が従う分布を仮定する必要がある。スペクトルモデル $\lambda[k, l]$ が n 個の調波成分

$$\lambda_n[k, l] = w_n[l]^2 e^{-\frac{1}{2\sigma^2}(x_k - \Omega[l] - \log n)^2} \quad (43)$$

の和で表されていることを思い出せば、

$$C_n[k, l] \sim \text{Pois}(\lambda_n[k, l]), \quad (n = 1, \dots, N) \quad (44)$$

を仮定すれば、ポアソン分布の再生性より、矛盾することなく $Y[k, l] \sim \text{Pois}(\lambda[k, l])$ が成り立つことが分かる。以上より、観測データ $Y := \{Y[k, l]\}_{1 \leq k \leq K, 1 \leq l \leq L}$ と潜在データ $C := \{C_n[k, l]\}_{1 \leq n \leq N, 1 \leq k \leq K, 1 \leq l \leq L}$ を組にした完全データ $\{Y, C\}$ の確率密度関数が

$$p(Y, C|\Omega, w) = p(Y|C)p(C|\Omega, w) \quad (45)$$

$$p(Y|C) = \prod_{k,l} \delta\left(Y[k, l] - \sum_n C_n[k, l]\right) \quad (46)$$

$$p(C|\Omega, w) = \prod_n \prod_{k,l} \frac{\lambda_n[k, l]^{C_n[k, l]} e^{-\lambda_n[k, l]}}{C_n[k, l]!} \quad (47)$$

で与えられる。完全データの対数確率の、観測データとモデルパラメータが与えられた下での条件つき期待値を Q 関数と呼び、Q 関数をもとに EM アルゴリズムの各ステップの更新則が導かれる。以下で本問題における Q 関数を導く。 $\log p(Y, C|\Omega, w)$ の中で Ω, w に関係する項に対し、 $C_n[k, l]$ の期待値をとったものが Q 関数となるので、Q 関数は

$$Q(\Omega, w) = \sum_n \sum_{k,l} (\langle C_n[k, l] \rangle_{p(C_n[k, l]|Y[k, l], \Omega', w')}) \log \lambda_n[k, l] - \lambda_n[k, l] \quad (48)$$

で与えられる。ただし $\langle f(x) \rangle_{p(x|y)}$ は条件つき期待値 $\mathbb{E}_x[f(x)|y]$ を表す。ここで、ポアソン分布の性質より、

$$C_n[k, l]|Y[k, l] \sim \text{Binom}\left(Y[k, l], \frac{\lambda'_n[k, l]}{\lambda'[k, l]}\right) \quad (49)$$

が言えるため、期待値は具体的に

$$\langle C_n[k, l] \rangle_{p(C_n[k, l]|Y[k, l], \Omega', w')} = Y[k, l] \frac{\lambda'_n[k, l]}{\lambda'[k, l]} \quad (50)$$

と求まる。ただし、Binom は二項分布を表し、 $\lambda'_n[k, l], \lambda'[k, l]$ は Ω', w' を代入した $\lambda_n[k, l], \lambda[k, l]$ を表す。以上より Q 関数の具体形を得た。あとは、Q 関数に対数事前分布を加えた

$$I(\Omega, w, \pi_\Omega, \pi_w) := \sum_{n,k,l} \left(Y[k, l] \frac{\lambda'_n[k, l]}{\lambda'[k, l]} \log \lambda_n[k, l] - \lambda_n[k, l] \right) + \log p(\Omega|\pi_\Omega) + \log p(w|\Omega, \pi_w) \quad (51)$$

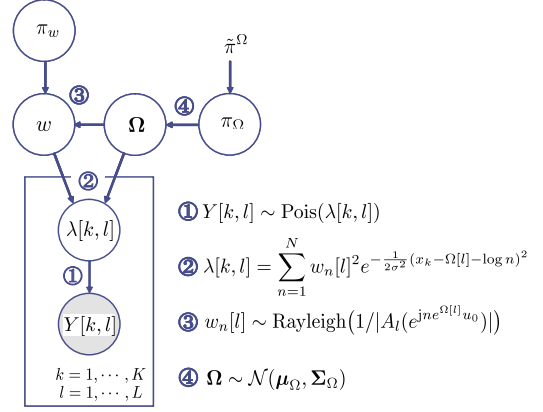


図 2 Graphical model representation of the proposed model.

を最大化する $\Omega, w, \pi_\Omega, \pi_w$ の更新則を導出できれば、以下のとおり EM アルゴリズムが実現できる。

(E ステップ) 更新された Ω, w を Ω', w' に代入する。

(M ステップ) $I(\Omega, w, \pi_\Omega, \pi_w)$ を最大化 (または増加) するように $\Omega, w, \pi_\Omega, \pi_w$ を更新し、E ステップに戻る。

M ステップの更新則を導出するにあたり、以下を仮定する。

(B1) 対数周波数のサンプリング間隔 $x_0 := x_{k+1} - x_k$ が十分小さく、以下の近似が成り立つ。

$$\int_{-\infty}^{\infty} e^{-x^2/2\sigma^2} dx \simeq x_0 \sum_k e^{-x_k^2/2\sigma^2} \quad (52)$$

よって、 $\sum_k \lambda_n[k, l] \simeq \sqrt{2\pi}\sigma w_n[l]^2/x_0$ である。

(B2) 事前分布の信頼度を調節できるようにする目的で、式 (41) 右辺を 2γ 乗したものを $p(w|\Omega, \Theta_w)$ とする。

(B3) $p(w|\Omega, \pi_w)$ において Ω を 1 ステップ前の更新値 Ω' に置き換える。すなわち、 $p(w|\Omega, \pi_w) \simeq p(w|\Omega', \pi_w)$ と近似する。

以上の準備により、M ステップの更新式が導ける。

1) F_0 パターン Ω : $I(\Omega, w, \pi_\Omega, \pi_w)$ で Ω に関係する項は

$$I(\Omega, w, \pi_\Omega, \pi_w) \stackrel{c}{=} -\frac{1}{2}(\Omega - \hat{\Omega})^T \Sigma^{-1}(\Omega - \hat{\Omega}) - \frac{1}{2}(\Omega - \mu_\Omega)^T \Sigma_\Omega^{-1}(\Omega - \mu_\Omega) \quad (53)$$

であるので、 I を最大化する Ω は

$$\Omega = (\Sigma + \Sigma_\Omega)^{-1}(\Sigma \hat{\Omega} + \Sigma_\Omega \mu_\Omega) \quad (54)$$

で与えられる。ただし、 $\hat{\Omega} := (\hat{\Omega}[1], \dots, \hat{\Omega}[L])^T$ は

$$\hat{\Omega}_l = \sum_k \frac{Y[k, l]}{\lambda'[k, l]} \sum_n \lambda'_n[k, l](x_k - \log n) / \sum_k Y[k, l] \quad (55)$$

を要素とするベクトル、 Σ は対角要素 Σ_l が

$$\Sigma_l = \sigma^2 / \sum_k Y[k, l] \quad (56)$$

の対角行列である。

2) n 次調波成分の瞬間パワー $w_n[l]$: $\partial I / \partial w_n[l] = 0$ を解くと、 I を最大化する $w_n[l]$ は

$$w_n[l]^2 = \frac{\sum_k Y[k, l] \frac{\lambda'_n[k, l]}{\lambda'[k, l]} + \gamma}{\sqrt{2\pi}\sigma/x_0 + \gamma/P_n[l]} \quad (57)$$

で与えられる。この更新式は $\gamma \rightarrow \infty, \gamma \rightarrow 0$ の極限をとると

$$\lim_{\gamma \rightarrow \infty} \frac{\sum_k Y[k, l] \frac{\lambda'_n[k, l]}{\lambda[k, l]} + \gamma}{\sqrt{2\pi\sigma/x_0 + \gamma/P_n[l]}} = P_n[l]$$

$$\lim_{\gamma \rightarrow 0} \frac{\sum_k Y[k, l] \frac{\lambda'_n[k, l]}{\lambda[k, l]} + \gamma}{\sqrt{2\pi\sigma/x_0 + \gamma/P_n[l]}} = \frac{x_0}{\sqrt{2\pi\sigma}} \sum_k Y[k, l] \frac{\lambda'_n[k, l]}{\lambda[k, l]}$$

となる。上段の式は直前に更新された全極スペクトルを基本周波数の間隔でサンプリングした値、下段の式はスペクトルモデルが観測スペクトルに最も良くフィットする $w_n[l]^2$ の値をそれぞれ表しており、式 (57) は両者の間を補間する値に相当する。 γ はその補間係数と見ることができる。

3) 藤崎モデルパラメータ π_Ω : I において π_Ω に関する項は

$$I(\Omega, w, \pi_\Omega, \pi_w) \stackrel{c}{=} -\frac{1}{2} \log |\Sigma_\Omega| - \frac{1}{2} (\Omega - \mu_\Omega)^T \Sigma_\Omega^{-1} (\Omega - \mu_\Omega) \quad (58)$$

であり、[4] の目的関数と同形となるので、 π_Ω の更新方法には当該文献と同様の反復アルゴリズムが適用できる。詳細は [4] を参照されたい。

4) 全極モデルパラメータ π_w : I において π_w に関する項は、

$$I(\Omega, w, \pi_\Omega, \pi_w) \stackrel{c}{=} \gamma \sum_l \sum_n \left(\log \frac{w_n[l]^2}{P_n[l]^2} - \frac{w_n[l]^2}{P_n[l]} \right) \quad (59)$$

であり、[3] で設定されている目的関数と類似した形となるので、 π_w の更新方法には当該文献と同様のアイデアが適用できる。詳細は省略するが、 $\mathbf{a}_l := (a_l[0], a_l[1], \dots, a_l[M])^T$ とすると、

$$\mathbf{a}_l \leftarrow \mathbf{R}_l^{-1} \mathbf{h}_l \quad (60)$$

$$\mathbf{h}_l \leftarrow \hat{\mathbf{R}}_l \mathbf{a}_l \quad (61)$$

を繰り返すことで R を増加させる \mathbf{a}_l を得ることができる。ただし、 \mathbf{R}_l および $\hat{\mathbf{R}}_l$ はそれぞれ

$$R_l(m' - m) = \frac{1}{N} \sum_{n=1}^N w_n[l]^2 \cos ne^{\Omega[l]} u_0(m' - m) \quad (62)$$

$$\hat{R}_l(m' - m) = \frac{2}{N} \sum_{n=1}^N P_n[l] \cos ne^{\Omega[l]} u_0(m' - m) \quad (63)$$

を m' 行 m 列成分とする Toeplitz 行列である。

6. 動作実験

ATR 音声データベースの女性話者音声データ (サンプリング周波数 16kHz のデジタル信号) を用いて提案法の基本動作の確認を行った。図 3 にその動作実験例を示す。図 3 上段に上記の音声信号に対して $x_0 = 12\text{cent}$, $\sigma = 0.02$, $t_0 = t_l - t_{l-1} = 8\text{ms}$ の条件下で算出したウェーブレット変換スペクトログラム、中段に 5. で述べた最適化アルゴリズムにより得た $\lambda[k, l]$ の推定値を示す。また、下段に F_0 パラメータ Ω (赤線) および推定した藤崎モデルパラメータにより表現される F_0 パターン (青線) を示す。 F_0 パラメータの推定結果を見ると、基本周波数推定においてしばしば問題となる倍ピッチ・半ピッチ誤りがほとんどないことが分かる。これはスペクトル包絡の滑らかさと、 F_0 パターンの滑らかさの制約が効果的に働いていることを示して

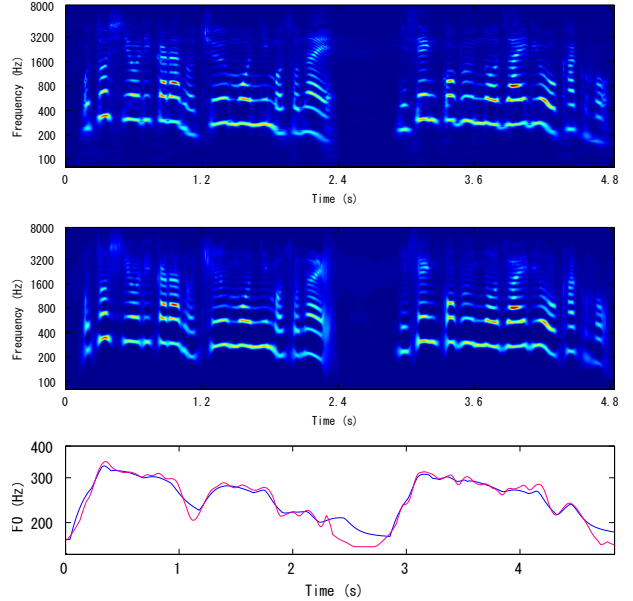


図 3 観測スペクトログラム $Y[k, l]$ (上段) と $\lambda[k, l]$ の推定値 (中段) と Ω の推定値 (下段・赤線) および μ_Ω の推定値 (下段・青線)

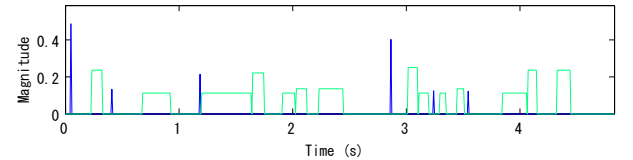


図 4 フレーズ指令 (青線) とアクセント指令 (緑線) の推定結果

いる。また、図 4 に示すように藤崎モデルパラメータがスペクトルから直接得られる点も提案法の大きな特徴の一つである。

現段階では提案法の具体的な性能を定量的に評価したわけではないが、動作実験の過程で既にいくつかの課題が明らかになった。具体的には、 Ω がまだ明らかに正しい F_0 の解に至っていない段階にもかかわらず収束してしまうケースがしばしば見られた。これは、 Ω と μ_Ω は互いに近くなるようにいずれも更新されるため、 Ω の初期値に強い拘束を受けてしまうのが原因であると考えられる。この問題は w と π_w についても同様であった。また、反復計算の初期段階において藤崎モデルのパラメータ推定が不安定な動作をすることが多かった。これは、藤崎モデルのパラメータ更新ステップにおいては、有音区間の $\Omega[l]$ の推定値と本来は無視するべきは無音区間における $\Omega[l]$ の推定値も同等の値を置いて推定を行っていることに原因があると考えられる。以上のことから、 Ω と w を周辺化して Y から直接 π_Ω と π_w を推定できる方法をいずれ検討していく必要があるだろうという感触を得た。筆者は既に、変分ベイズ法に基づき Ω と w の事後分布を近似計算するアルゴリズムを導出しており、これをもとに、今後上記の課題解決に向けて検討を行っていく予定である。

7. まとめ

本稿では、全極型声道スペクトルモデルと藤崎の F_0 パターン生成過程モデルを事前分布の形として内包する音声ウェーブレットスペクトルの統計モデルを提案し、パラメータの効率的な最適化方法を EM アルゴリズムに基づいて導出した。この枠組により、全極型声道モデルと藤崎モデルのパラメータを直接スペクトルから同時推定することが可能となった。

文 献

- [1] F. Itakura, S. Saito, "Analysis synthesis telephony based upon the maximum likelihood method," In *Proc. 6th Int'l Cong. Acoust. (ICA '68)*, C-5-5, C17-20, 1968.
- [2] H. Fujisaki, "A note on the physiological and physical basis for the phrase and accent components in the voice fundamental frequency contour," In *Vocal Physiology: Voice Production, Mechanisms and Functions*, (O. Fujimura, ed.) Raven Press, pp. 347-355, 1988.
- [3] A. El-Jaroudi, J. Makhoul, "Discrete all-pole modeling," *IEEE trans. Signal Process.*, **39**(2), pp. 411-423, 1991.
- [4] 亀岡, ルルー, 大石, "音声 F_0 パターン生成過程の確率モデル," *音講論 (秋)*, 1-1-3, pp.207-210, Sep. 2010.