

# $F_0$ パターン生成過程の統計的モデルによる 音声信号からのフレーズ・アクセント指令の推定\*

吉里幸太<sup>1</sup>, 亀岡弘和<sup>1,2</sup>, 齋藤大輔<sup>1</sup>, 嵯峨山茂樹<sup>1</sup>  
(<sup>1</sup> 東大院・情報理工, <sup>2</sup>NTT CS 研)

## 1 はじめに

音声の抑揚を表す物理量である基本周波数 ( $F_0$ ) の時間変化パターンは、構文や意図の伝達に関する様々な種類の非言語情報を含んでおり、発話を用いたコミュニケーションにおいて重要な役割を果たしている。 $F_0$  パターンは、甲状軟骨の運動によって生じる声帯の伸びに応じて決まる量である。ゆえに、音声の  $F_0$  パターンにはある種の物理的制約が課せられる。音声アプリケーションで抑揚を扱う場合、この物理的制約を考慮に入れた  $F_0$  パターンの生成モデルを考えることは非常に有用である。

藤崎の  $F_0$  パターン生成過程モデル (藤崎モデル) とは、甲状軟骨の運動に注目して  $F_0$  パターンの生成過程を説明する、力学的モデルである [1]。藤崎モデルは、発話の言語学的意図と密接に関わる少数のパラメータを与えることで、実測の  $F_0$  パターンに極めて近い  $F_0$  パターンを表現することが可能である。

近年、利用可能な音声データベースが増加してきたことに伴って、統計的手法を用いた音声処理に関して盛んに研究がおこなわれるようになり、多くの成果を挙げている。そうした成果に加えて、例えば、ある音声の集合から抑揚の付き方の傾向を学習でき、音声合成や音声認識に反映できるような統計的枠組が確立されれば極めて有用であろう。

そこで我々は、統計的モデルをベースにした音声認識、音声合成、言語識別、話者認識といった現代の多くの音声アプリケーションに、将来的に韻律モデルを組み込んでいくことを目標にして、 $F_0$  パターンを統計的に扱う枠組みを作る研究を進めている。以前我々は、隠れマルコフモデル (Hidden Markov Model; HMM) を用いて藤崎モデルのパラメータの生成過程を確率的に扱うためのモデルを提案した [2]。本稿ではまず、このモデルを実音声に適用する上で考慮すべき点について議論し、それに基づきモデルの改良を行う。そして、その提案モデルが藤崎モデルのパラメータを統計的に扱ううえで十分な性能を持っていることを確認するために、実音声の  $F_0$  パターンを入力として藤崎モデルのパラメータを推定する問題を解いて定量評価を行う。

## 2 藤崎の $F_0$ パターン生成過程モデル

藤崎の  $F_0$  パターン生成過程モデル (藤崎モデル) とは、甲状軟骨は平行移動と回転運動の二つの独立な運動を行い、それぞれの運動によって伸びた声帯の長さの和が  $\log F_0$  の値に比例する、という仮定をもとに  $F_0$  パターンの生成過程を表現した力学的モデルである [1]。藤崎モデルでは、甲状軟骨の平行移動によって生じる  $F_0$  パターンをフレーズ成分  $y_p(t)$ 、回転運動によって生じる  $F_0$  パターンをアクセント成分  $y_a(t)$  と呼び ( $t$  は時刻)、これらの和に声帯の物理的性質によって決まる定数値であるベースライン成分  $y_b$  を加えたものが  $F_0$  パターンの対数値  $y(t)$  であるとする。 $y_p(t)$  と  $y_a(t)$  は、それぞれフレーズ指令と

呼ばれるパルス波の列  $u_p(t)$  と、アクセント指令と呼ばれる矩形波の列  $u_a(t)$  を入力とした臨界制動の二次線形系により表現される。任意の時刻においてフレーズ指令とアクセント指令が同時に正の値をとることはないという制約のもとで、これらの値の関係は

$$y(t) = y_p(t) + y_a(t) + y_b, \quad (1)$$

$$y_p(t) = G_p(t) * u_p(t), \quad y_a(t) = G_a(t) * u_a(t), \quad (2)$$

$$G_p(t) = \begin{cases} \alpha^2 t e^{-\alpha t} & (t \geq 0) \\ 0 & (t < 0) \end{cases}, \quad (3)$$

$$G_a(t) = \begin{cases} \beta^2 t e^{-\beta t} & (t \geq 0) \\ 0 & (t < 0) \end{cases} \quad (4)$$

と書ける。ここで、 $\alpha$  と  $\beta$  は、それぞれフレーズ制御機構、アクセント制御機構の固有角周波数であり、話者や発話内容によらず、おおよそ  $\alpha = 3 \text{ rad/s}$ ,  $\beta = 20 \text{ rad/s}$  程度であることが経験的に知られている。

なお、フレーズ成分は急激な上昇のあと緩やかに下降していく成分で、アクセント成分は急激な上昇下降をなす成分であり、日本語を含む多くの言語では、前者が句単位での大局的な音調を、後者が語や音節単位での局所的な音調を表現する役割を担っていると考えられている。

## 3 統計的 $F_0$ パターン生成過程モデル

我々は以前、隠れマルコフモデル (Hidden Markov Model; HMM) を用いて藤崎モデルのパラメータの生成過程を確率的に扱うためのモデルを提案した [2]。本章では、このモデルを実音声に適用する上で考慮すべき点を 3 つ挙げ、これらを解決する改良モデルを定式化する。

### 3.1 藤崎モデルの離散時間表現

連続時間システムのフレーズ制御機構とアクセント制御機構に対して後退差分変換を施すことで、藤崎モデルの離散時間表現を得ることができる。 $k$  を離散時刻のインデックスとし、 $y_p[k]$ ,  $u_p[k]$ ,  $y_a[k]$ ,  $u_a[k]$  をそれぞれ  $y_p(t)$ ,  $u_p(t)$ ,  $y_a(t)$ ,  $u_a(t)$  の離散時間表現とすると、これらの値の関係は

$$u_p[k] = a_0 y_p[k] + a_1 y_p[k-1] + a_2 y_p[k-2], \quad (5)$$

$$a_2 = (\psi - 1)^2, \quad a_1 = -2\psi(\psi - 1), \quad a_0 = \psi^2, \quad (6)$$

$$u_a[k] = b_0 y_a[k] + b_1 y_a[k-1] + b_2 y_a[k-2], \quad (7)$$

$$b_2 = (\varphi - 1)^2, \quad b_1 = -2\varphi(\varphi - 1), \quad b_0 = \varphi^2 \quad (8)$$

と書ける。ただし  $\psi = 1 + 1/(\alpha t_0)$ ,  $\varphi = 1 + 1/(\beta t_0)$  であり、 $t_0$  はサンプリング周期である。

### 3.2 藤崎モデルの統計モデル化

フレーズ指令  $u_p[k]$  とアクセント指令  $u_a[k]$  は同時に正の値をとらないという制約を満たしつつ、これ

\* Estimation of phrase and accent commands from speech signals using statistical model of speech  $F_0$  contours. by YOSHIZATO Kota, KAMEOKA Hirokazu, SAITO Daisuke, SAGAYAMA Shigeki (The University of Tokyo)

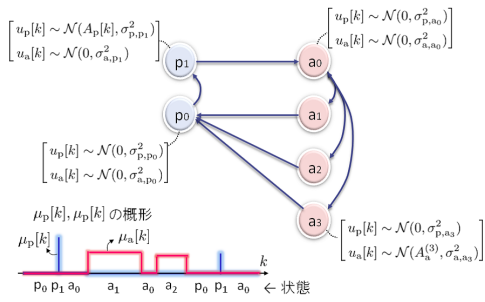


Fig. 1 HMMの構成

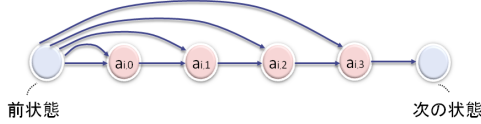


Fig. 2 HMMの状態の分割

らの指令を生成する確率モデルを得るため、藤崎モデルの指令列が、HMMの各状態からの出力  $o[k] = (u_p[k], u_a[k])^T$  として表せると仮定する。提案する統計的  $F_0$  パターン生成過程モデルのHMMの構成を、Fig. 1に示す。出力  $o[k]$  は正規分布に従い、

$$o[k] \sim \mathcal{N}(\nu[k], \Upsilon[k]), \quad (9)$$

$$\nu[k] = \begin{bmatrix} \mu_p[k] \\ \mu_a[k] \end{bmatrix}, \quad \Upsilon[k] = \begin{bmatrix} \sigma_p^2[k] & 0 \\ 0 & \sigma_a^2[k] \end{bmatrix}, \quad (10)$$

と書けるとする。平均ベクトル  $\nu[k]$  と分散共分散行列  $\Upsilon[k]$  は、HMMの状態遷移の結果として決まる値である。

また、フレーズ指令とアクセント指令が生起するタイミングやアクセント指令の持続長には音声の重要な韻律特徴が含まれており、これらを統計的に学習できることが望ましい。そこで、HMMの各状態における停留時間長をより柔軟にパラメトライズできるようにするため、HMMの各々の状態を、全く同じ出力分布を持ついくつかの小状態に分割する。分割の例として、状態  $a_i$  の分割をFig. 2に示す。任意の  $l$  に対して、小状態  $a_{i,l}$  から小状態  $a_{i,l+1}$  への状態遷移確率を1に設定すると、アクセント指令が持続長  $j$  を持つ確率が、前状態から  $a_{i,j}$  への状態遷移確率で表現できる。 $a_i$  と同様の分割を、 $p_0, a_0$  に対しても行う。なお、この状態分割は、[2]のモデルに新たに施された1つ目の改良点である。

提案モデルのHMMを定式化すると次のように書ける。

出力値系列:  $\{o[k]\}_{k=1}^K$   
 状態集合:  $S = \{p_0, p_1, a_0, \dots, a_N\}$   
 状態系列:  $\{s_k\}_{k=1}^K$   
 状態出力分布:  $P(o[k] | s_k = i) = \mathcal{N}(\nu[k], \Upsilon[k])$

$$\nu[k] = \begin{cases} (0, 0)^T & (i = p_0, a_0) \\ (A_p[k], 0)^T & (i = p_1) \\ (0, A_a^{(n)})^T & (i = a_n) \end{cases}$$

$$\Upsilon[k] = \begin{bmatrix} \sigma_p^2 & 0 \\ 0 & \sigma_a^2 \end{bmatrix}$$

状態遷移確率:  $\phi_{i',i} = \log P(s_k = i' | s_{k-1} = i)$

また、ベースライン成分の離散時間表現も、 $y_b[k] \sim \mathcal{N}(\mu_b, \sigma_b^2)$  と正規分布に従うと仮定する。

さらに、実音声には  $F_0$  の値が定義されない無声区間が存在するので、これに対処するために、観測  $F_0$  パターンの離散時刻  $k$  における不確かさの程度  $\sigma_n^2[k]$  という変数を導入する。 $\sigma_n^2[k]$  の値を無声区間では大きく、それ以外の区間では小さく設定することで、提案モデルで有声区間と無声区間とを統一的に扱うために、 $y[k]$  を、真の  $F_0$  の値と  $y_n[k] \sim \mathcal{N}(0, \sigma_n^2[k])$  の和  $y[k] = y_p[k] + y_a[k] + y_b[k] + y_n[k]$  であると考え。この不確かさの導入が、[2]のモデルに新たに施された2つ目の改良点である。

指令や成分の時系列データを、まとめて  $u_p = (u_p[1], \dots, u_p[K])^T$ ,  $u_a = (u_a[1], \dots, u_a[K])^T$ ,  $\mu_p = (\mu_p[1], \dots, \mu_p[K])^T$ ,  $\mu_a = (\mu_a[1], \dots, \mu_a[K])^T$ ,  $y_p = (y_p[1], \dots, y_p[K])^T$ ,  $y_a = (y_a[1], \dots, y_a[K])^T$ ,  $y_b = (y_b[1], \dots, y_b[K])^T$ ,  $y_n = (y_n[1], \dots, y_n[K])^T$ ,  $y = (y[1], \dots, y[K])^T$  と書くと、 $u_p$  と  $u_a$  は、

$$u_p = Ay_p, \quad u_a = By_a \quad (11)$$

と表すことができる。ここで、

$$A = \begin{bmatrix} a_0 & & & & O \\ a_1 & a_0 & & & \\ a_2 & a_1 & a_0 & & \\ & \ddots & \ddots & \ddots & \\ O & & a_2 & a_1 & a_0 \end{bmatrix}, \quad B = \begin{bmatrix} b_0 & & & & O \\ b_1 & b_0 & & & \\ b_2 & b_1 & a_0 & & \\ & \ddots & \ddots & \ddots & \\ O & & b_2 & b_1 & b_0 \end{bmatrix} \quad (12)$$

である。

簡単のため  $\phi_{i',i}$ ,  $\mu_b$ ,  $\sigma_p^2, \sigma_a^2, \sigma_b^2, \sigma_n^2[k]$ ,  $\alpha, \beta$  を全て、観測  $F_0$  パターンに応じて決まる定数とすると、藤崎モデルのパラメータ  $\Theta = \{\{A_p[k], s[k]\}_{k=1}^K, \{A_a^{(n)}\}_{n=1}^N\}$  が与えられたとき  $y = (y[1], \dots, y[K])^T$  を出力する条件付確率  $P(y|\Theta)$  は

$$P(y|\Theta) = \frac{|\Sigma^{-1}|^{1/2}}{(2\pi)^{K/2}} \exp \left\{ -\frac{1}{2} (y - \mu)^T \Sigma^{-1} (y - \mu) \right\},$$

$$\mu = A^{-1} \mu_p + B^{-1} \mu_a + \mu_b \mathbf{1}, \quad (13)$$

$$\Sigma = A^{-1} \Sigma_p (A^T)^{-1} + B^{-1} \Sigma_a (B^T)^{-1} + \Sigma_b + \Sigma_n$$

と書ける。ここで、

$$\Sigma_p = \begin{bmatrix} \sigma_p^2[1] & & & & O \\ & \ddots & & & \\ O & & \sigma_p^2[K] & & \\ & & & \ddots & \\ & & & & \sigma_p^2[K] \end{bmatrix}, \quad \Sigma_a = \begin{bmatrix} \sigma_a^2[1] & & & & O \\ & \ddots & & & \\ O & & \sigma_a^2[K] & & \\ & & & \ddots & \\ & & & & \sigma_a^2[K] \end{bmatrix},$$

$$\Sigma_b = \begin{bmatrix} \sigma_b^2 & & & & O \\ & \ddots & & & \\ O & & \sigma_b^2 & & \\ & & & \ddots & \\ & & & & \sigma_b^2 \end{bmatrix}, \quad \Sigma_n = \begin{bmatrix} \sigma_n^2[1] & & & & O \\ & \ddots & & & \\ O & & \sigma_n^2[K] & & \\ & & & \ddots & \\ & & & & \sigma_n^2[K] \end{bmatrix} \quad (14)$$

である。

#### 4 パラメータ推定アルゴリズム

本章では、3章で提案した統計的  $F_0$  パターン生成過程モデルを用いて、観測  $F_0$  パターン  $y$  が与えられたときに  $P(y|\Theta)$  を最大化する  $\Theta$  を求める問題を解くことを考える。この問題の大域的最適解を求めるのは難しいが、 $x = (y_p^T, y_a^T, y_b^T, y_n^T)^T$  を完全データと見なしてEMアルゴリズムによる不完全データ問題に帰着することで、局所最適解を求めることはできる。このとき、 $x$  が与えられたもとの  $\Theta$  の対数尤度関数は

$$\log P(x|\Theta) \doteq \frac{1}{2} \log |\Lambda^{-1}| - \frac{1}{2} (x - m)^T \Lambda^{-1} (x - m), \quad (15)$$

と書ける．ただし，

$$\mathbf{x} = \begin{bmatrix} \mathbf{y}_p \\ \mathbf{y}_a \\ \mathbf{y}_b \\ \mathbf{y}_n \end{bmatrix}, \quad \mathbf{m} = \begin{bmatrix} \mathbf{A}^{-1} \boldsymbol{\mu}_p \\ \mathbf{B}^{-1} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \mathbf{1} \\ 0 \end{bmatrix}, \quad (16)$$

$$\boldsymbol{\Lambda}^{-1} = \begin{bmatrix} \mathbf{A}^T \boldsymbol{\Sigma}_p^{-1} \mathbf{A} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{B}^T \boldsymbol{\Sigma}_a^{-1} \mathbf{B} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \boldsymbol{\Sigma}_b^{-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \boldsymbol{\Sigma}_n^{-1} \end{bmatrix}$$

である．また，Q 関数  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}')$  は，

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}') \stackrel{c}{=} \frac{1}{2} \left[ \log |\boldsymbol{\Lambda}^{-1}| - \text{tr}(\boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}^T | \boldsymbol{\theta}']) \right. \\ \left. + 2\mathbf{m}^T \boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x} | \boldsymbol{\theta}'] - \mathbf{m}^T \boldsymbol{\Lambda}^{-1} \mathbf{m} \right] + \log P(\boldsymbol{\theta}) \quad (17)$$

と書ける．ここで， $\mathbb{E}[\mathbf{x} | \boldsymbol{\theta}']$  と  $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \boldsymbol{\theta}']$  は， $\mathbf{H} = [\mathbf{I}, \mathbf{I}, \mathbf{I}, \mathbf{I}]$  を用いると， $\mathbf{y} = \mathbf{H}\mathbf{x}$  より，

$$\mathbb{E}[\mathbf{x} | \boldsymbol{\theta}'] = \mathbf{m} + \boldsymbol{\Lambda} \mathbf{H}^T (\mathbf{H} \boldsymbol{\Lambda} \mathbf{H}^T)^{-1} (\mathbf{y} - \mathbf{H}\mathbf{m}), \quad (18)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T | \boldsymbol{\theta}'] = \boldsymbol{\Lambda} - \boldsymbol{\Lambda} \mathbf{H}^T (\mathbf{H} \boldsymbol{\Lambda} \mathbf{H}^T)^{-1} \mathbf{H} \boldsymbol{\Lambda} \\ + \mathbb{E}[\mathbf{x} | \boldsymbol{\theta}'] \mathbb{E}[\mathbf{x} | \boldsymbol{\theta}']^T \quad (19)$$

と書ける．これらが E ステップで更新すべき値である．

$\mathbb{E}[\mathbf{x} | \boldsymbol{\theta}']$  を 4 つの  $K \times 1$  行列に分割して，各小行列を  $\mathbb{E}[\mathbf{x} | \boldsymbol{\theta}'] = (\bar{\mathbf{x}}_p^T, \bar{\mathbf{x}}_a^T, \bar{\mathbf{x}}_b^T, \bar{\mathbf{x}}_n^T)^T$  とすると，M ステップの更新式は次のように書ける．

1) 状態系列: Q 関数の中から  $s := \{s_k\}_{k=1}^K$  に関する項だけを抜き出すと，

$$\mathcal{I}(s) := -\frac{1}{2} \sum_{k=1}^K (\mathbf{o}[k] - \boldsymbol{\nu}[k])^T \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[k] - \boldsymbol{\nu}[k]) \\ + \log P(s_1) + \sum_{k=2}^K \log P(s_k | s_{k-1}) \quad (20)$$

となる．ここで， $\mathbf{o}[k] := ([\mathbf{A}\bar{\mathbf{x}}_p]_k, [\mathbf{B}\bar{\mathbf{x}}_a]_k)^T$  である．ただし  $[\cdot]_k$  はベクトルの  $k$  番目の要素を表す． $\mathcal{I}(s)$  を最大化する状態系列  $\{s_k\}_{k=1}^K$  は，動的計画法を用いて効率的に求めることができる．

2) フレーズ・アクセント指令の大きさ: Q 関数の値を最大化する  $A_p[k]$  と  $A_a^{(n)}$  の値は，

$$A_p[k] = [\mathbf{A}\bar{\mathbf{x}}_p]_k \quad (k \in \mathcal{T}_{p_1}), \quad \mathcal{T}_{p_1} = \{k | s_k = p_1\}, \quad (21)$$

$$A_a^{(n)} = \frac{1}{|\mathcal{T}_{a_n}|} \sum_{k \in \mathcal{T}_{a_n}} [\mathbf{B}\bar{\mathbf{x}}_a]_k, \quad \mathcal{T}_{a_n} = \{k | s_k = a_n\} \quad (22)$$

と書ける．

式 (20) から (22) の中には， $\mathbf{A}\bar{\mathbf{x}}_p$  と  $\mathbf{B}\bar{\mathbf{x}}_a$  の計算が含まれている．この計算は，フレーズ成分  $\bar{\mathbf{x}}_p$  とアクセント成分  $\bar{\mathbf{x}}_a$  から，それを生成するフレーズ指令  $\bar{\mathbf{u}}_p$  とアクセント指令  $\bar{\mathbf{u}}_a$  を求めることに相当する．しかし，例えば日本語のように，指令の大きさは必ず 0 以上でなくてはならないという制約を持つ言語も存在する．こうした言語を扱う際に，単純に  $\mathbf{A}\bar{\mathbf{x}}_p$  や  $\mathbf{B}\bar{\mathbf{x}}_a$  を計算するだけでは，計算結果とし

て得られるベクトルは非負の要素を含まう．これを避けるために， $\mathbf{A}\bar{\mathbf{x}}_p$  と  $\mathbf{B}\bar{\mathbf{x}}_a$  を直接計算する代わりに， $\|\mathbf{A}^{-1}\bar{\mathbf{u}}_p - \bar{\mathbf{x}}_p\|_2^2$  と  $\|\mathbf{B}^{-1}\bar{\mathbf{u}}_a - \bar{\mathbf{x}}_a\|_2^2$  を最小化する  $\bar{\mathbf{u}}_p$  と  $\bar{\mathbf{u}}_a$  の値を非負制約のもとで求める問題を考える．これらは， $G_p[k]$  と  $G_a[k]$  をフレーズ制御機構とアクセント制御機構のインパルス応答の離散時間表現とすると， $\sum_k |G_p[k] * \bar{\mathbf{u}}_p[k] - \bar{\mathbf{x}}_p[k]|^2$  と  $\sum_k |G_a[k] * \bar{\mathbf{u}}_a[k] - \bar{\mathbf{x}}_a[k]|^2$  と書ける．この非負制約つき逆畳み込み問題は，[3] の方法で非常に効率的に解くことができる．なお，こうして指令の大きさに非負制約を導入する方法は，[2] に新たに追加された 3 つ目の改良点である．

## 5 実験

統計的  $F_0$  パターン生成過程モデルが，藤崎モデルのパラメータを扱う上で十分な性能を持っていることを確かめるため，4 章で述べたアルゴリズムを使って，観測  $F_0$  パターンから藤崎モデルのパラメータを推定する問題を解き，その結果を定量評価する．実験には，音素のバランスがとれた 503 文を集めたデータベースである，ATR 日本語発話データベースの B セットを用いる [4]．話者は MHT(男性話者) を選択した．また，[?] の方法を用いて入力音声から  $F_0$  パターンを推定した．加えて，推定パラメータの妥当性を評価する目的といくつかのパラメータを学習する目的で，韻律研究の専門家の手によって発話音声に付与された藤崎モデルのパラメータを用いた．なお，付与されたこのパラメータのベースライン成分は，常に  $\log(60 \text{ Hz})$  である．

この実験において，定数値のパラメータは， $N = 20$ ， $t_0 = 8 \text{ ms}$ ， $\alpha = 3.0$ ， $\beta = 20.0$ ， $\sigma_p^2[k] = 0.2^2$ ， $\sigma_a^2[k] = 0.02^2$ ， $\sigma_b^2 = 0.001^2$ ，無声区間では  $\sigma_n^2[k] = 10^{15}$ ，それ以外の区間では  $\sigma_n^2[k] = 10^{-15}$  に設定し， $\mu_b$  は有声区間の  $\log F_0$  の最小値とした． $\boldsymbol{\theta}$  の初期値は [?] の方法で決定し，EM アルゴリズムの反復回数は 20 回とした．HMM の小状態数と状態遷移確率  $\phi_{i,i}$  は，ATR の No.1 から No.200 までの 200 文に専門家が手動で付与した藤崎モデルのパラメータから学習した．定量評価は，残りの 303 文を用いて行った．

手動で決定された指令列と，提案手法により推定された指令列を，3 つの例について比較した図が Fig. 3 である．より詳しく，推定性能を定量評価するため，次のようにした．まず，Fig. 4 のように， $S$  秒以下のずれかないフレーズ指令同士，アクセント指令同士がマッチングされる可能性があるとして，手動で決定された指令列と推定された指令列との間のマッチング数が最大になるようにマッチングをとった．ここで，アクセント指令の時間のずれとは，アクセント指令の立ち上がり時刻のずれと立ち下がり時刻のずれの平均であるとした． $N_E$  を推定された指令列の指令数， $N_A$  を手動で決定された指令列の指令数， $N_M$  をマッチング数， $N_{E\text{sum}}$ ， $N_{A\text{sum}}$ ， $N_{M\text{sum}}$  をそれぞれ  $N_E$ ， $N_A$ ， $N_M$  を 303 文について足し合わせた値としたとき，挿入エラーを  $E_I = (N_{E\text{sum}} - N_{M\text{sum}})/N_{A\text{sum}}$ ，脱落エラーを  $E_O = (N_{A\text{sum}} - N_{M\text{sum}})/N_{A\text{sum}}$  と定義した．また，正解率を  $A = 1 - E_I - E_O$  と定義し，これを定量評価の基準とした．なお，この実験においては，ベースライン成分の値のずれを考慮して，フレーズ・アクセント指令の大きさの違いは評価基準に含めなかった．

$S = 0.3$  秒としたときの実験結果を Table 1 に記した．これらの表は，フレーズ・アクセント両指令，フレーズ指令のみ，アクセント指令のみに注目したときの，正解率  $A$ ，挿入エラー  $E_I$ ，脱落エラー  $E_O$  を，

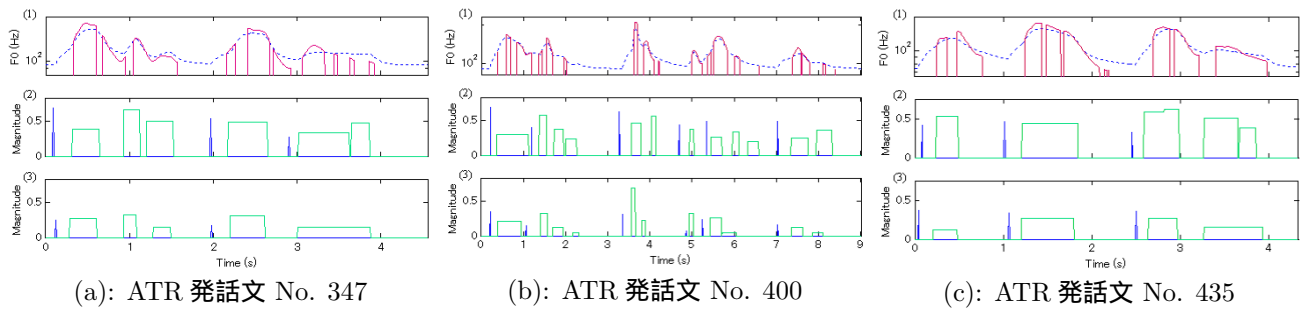


Fig. 3 パラメータ推定結果の例: (1) 実線は有声区間の観測  $F_0$  パターン, 点線は推定された  $y_p + y_a + y_b$  の値. (2) 手で付与されたフレーズ・アクセント指令 (3) 提案手法で推定されたフレーズ・アクセント指令

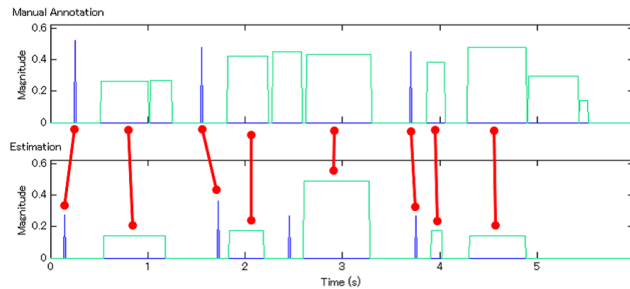


Fig. 4 指令列のマッチングの例

Table 1 定量評価の結果 ( $S=0.3$  秒).

(a) 両指令			
	$A$	$E_I$	$E_O$
初期状態	0.688	0.088	0.224
反復後	0.697	0.127	0.177
(b) フレーズ指令			
	$A$	$E_I$	$E_O$
初期状態	0.647	0.109	0.244
反復後	0.680	0.207	0.112
(c) アクセント指令			
	$A$	$E_I$	$E_O$
初期状態	0.711	0.076	0.213
反復後	0.708	0.083	0.207

それぞれ初期状態と EM アルゴリズムの反復後に計算した結果である. この表より, 初期状態の正解率と反復後の正解率には, 大きな差がないことが分かる. パラメータの初期値を定めるのに使った [?] の方法は,  $F_0$  パターンに関する様々な知見を利用して観測  $F_0$  パターンから藤崎モデルのパラメータを推定するアルゴリズムであり, この実験結果は, 提案した統計的手法によるパラメータ推定がそれと同程度の性能を持っていることを示している. ゆえに我々は, 提案した統計的  $F_0$  パターン生成過程モデルは, 藤崎モデルのパラメータを統計的に扱ううえで十分な性能を持っていると結論づけた.

## 6 おわりに

本研究では, 統計的モデルをベースにした様々な音声アプリケーションに  $F_0$  パターンを組み入れるために統計的  $F_0$  パターン生成過程モデルを提案し, そのモデルが藤崎モデルのパラメータを統計的に扱う

ことができることを実験によって確認した. 今後は, 統計的  $F_0$  パターン生成過程モデルを用いて大量の音声データからフレーズ・アクセント指令の傾向を学習することを通して,  $F_0$  パターンを適切に制御した表情豊かな音声合成や, 抑揚の特徴をふまえた話者認識など, 実際に様々な音声アプリケーションに韻律モデルを組み込んでいきたい.

謝辞 本研究の一部は, 文部科学省科学研究費補助金 (23240021) の助成を受けて行われた. また, 本研究の実験では, 東京大学の広瀬啓吉教授が ATR の音声データに手で付与した藤崎モデルのパラメータを用いた. これを作成した同氏に, 強い感謝の意を表する.

## 参考文献

- [1] H. Fujisaki, *In Vocal Physiology: Voice Production, Mechanisms and Functions*, Raven Press, 1988.
- [2] H. Kameoka *et al.*, “A statistical model of speech  $F_0$  contours,” *ISCA Tutorial and Research Workshop on Statistical And Perceptual Audition (SAPA)*, pp. 43-48, 2010.
- [3] H. Kameoka *et al.*, “Robust speech dereverberation based on non-negativity and sparse nature of speech spectrograms,” *In Proc. 2008 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP2009)*, pp. 45-48, 2009.
- [4] A. Kurematsu *et al.*, “ATR japanese speech database as a tool of speech recognition and synthesis,” *Speech Communication*, vol. 27, pp. 187-207, 1999.
- [5] 亀岡, “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一的音声生成モデル,” 日本音響学会 2010 年秋季研究発表会, pp. 211-214, 2010.
- [6] 成澤他, “音声の基本周波数パターン生成過程モデルのパラメータ自動抽出法の評価,” 情報処理学会音声言語情報処理研究会, pp. 1-6, 2003.