



## 標準パターンDP合成法による多数話者の単語音声認識 (2)

|       |  |
|-------|--|
| メタデータ | 言語: jpn<br>出版者:<br>公開日: 2013-11-21<br>キーワード (Ja):<br>キーワード (En):<br>作成者: 高橋, 文彦, 矢田, 茂郎<br>メールアドレス:<br>所属: |
| URL   | <a href="https://doi.org/10.24729/00007879">https://doi.org/10.24729/00007879</a>                          |

# 標準パターン DP 合成法による 多数話者の単語音声認識 (II)

高橋 文彦\* 矢田 茂郎\*\*

Speech Recognition for Multi-speakers by Standard Patterns Synthesized  
by DP Matching Method (II)

Fumihiko TAKAHASHI\* Shigerou YADA\*\*

## ABSTRACT

This paper is a report for a research of speech recognition in case of unspecified multi-speakers. In the former report, we developed a new Synthesis Method for Standard Pattern, averaging many patterns by utilizing DP Matching Method. In this research, peculiar parameter of voice consists of two components, main component is time series of PARCOR Coefficient and subcomponent is time series of voice power. Experimental results of speech recognition are fairly improved by adding voice power component, the score of recognition is about 98% in case of 50 words (Japanese cities) by unspecified male voice.

Key Words: Recognition, Speech, Standard Pattern, Multispeakers

### 1. まえがき

音声認識や音声合成などの音声関連技術は、ここ数年で急速に進歩してきた。特に合成音声は、自動販売機から家電製品、電話の番号案内に至るまで広範囲に利用されるようになった。これらは音声の情報圧縮技術の進歩と、その LSI 化の成功によるものである。

音声認識技術に関しても実用化は進んでいる。簡単な音声認識用ボードが市販され、ホビータン的な用途にまで使われるようになってきている。しかし、これらはすべて使用者が予め自分の音声を登録しておくタイプのもので特定話者音声認識と呼ばれる分野のものである。音声認識が広く利用されるためには、誰の声でもすぐ聞き分けられるような認識装置が要求される。不特定多数話者の音声認識の研究が待たれるわけである。

筆者らは、この方向に沿って研究を進め、さきに標準パターン DP 合成法<sup>1)</sup>を提唱した。これは DP マッチングの手法を利用して、多数のパターンを平均化して標準パターンを合成し、個人特徴のバラツキを抑えた汎用の標準パターンを得るものである。前報では特徴パラメータに PARCOR 係数を用い、50 単語で 92% 程度の認識率 (成年男子) を得た。

しかしながら、PARCOR 係数には本質的に音の強弱や高低に関する情報は含まれない。したがって今回は、強弱に関する情報、すなわち音声パワーの時間変化を特徴パラメータとして追加し認識率の改善をはかった。その結果、認識率を 98% 程度まで向上させることができたので報告する。

### 2. 音声認識システムの概要

このシステムの概要を図 1 に示す。入力音声を予め登録された標準パターンと比較することは、特定話者の場合と同じであるが標準パターンの作成に工夫がほどこされている。標準パターン作成用に入力された音声信号は分析過程で特徴抽出され、各特徴パラメータの時系列パターンとなる。これを同一語彙毎に平均化して標準パタ

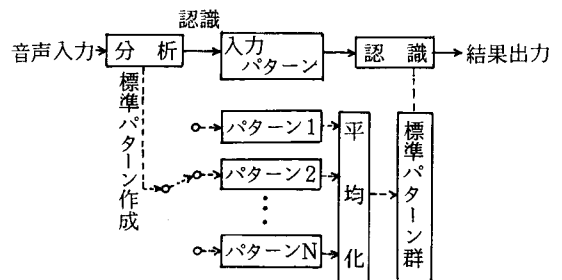


図 1 音声認識システムのブロック図

平成 2 年 4 月 9 日受理

\* 電気工学科(Department of Electrical Engineering)

\*\* 東京大学工学部(Dep. of Engineering, Tokyo Univ.)

ーンを作成し登録する。この平均化の際に、各時系列パターンを単語グループ毎に DP マッチング法で最適伸縮させて重ね合わせ平均化して、標準パターンを作成するわけである。

### 3. 音声の収録と特徴パラメータ

音声収録装置の外観を図 2 に示す。日本電気製パーソナルコンピュータ・PC 9801 を主体とするもので、これに音声入力用 A/D コンバータと音声出力用の D/A コンバータを附加したシステムである。音声は 10kHz のサンプリング間隔で A/D 変換 (12 bit) し、ハードディスク (20MB) に格納するようになっている。



図 2 音声入出力装置

入力音声は 50 個の単語 (日本の都市名) で 11 名の男性が発声したものである。50 単語は連続して発声するが、単語間の区切りと単語内の無音部を機械が混同しないように、区切りの休止期間は十分長くとるよう注意して発声した。

認識に用いる特徴パラメータは PARCOR 係数 (1 ~ 15 次) と音声パワーで、共に 25.6 msec 毎 (データ個数 256) にフレーム化して計算している。したがって単語音声は、25.6 msec 毎の時系列パターンに変換されることになる。

### 4. パターンマッチングと時間軸整合

マッチング法はパターン認識の汎用的手法である。あらかじめ登録された標準パターンと別の未知パターンを重ね合わせ、パターン間のズレ (距離) を計算し、最もズレの小さいものを選び出して、そのパターンに属すると判定する手法がパターンマッチング法である。

音声信号の入力パターンは、PARCOR 係数群と音声パワーの時系列として次式で表わされる。

$$\mathbf{A} = [ a_1, a_2, a_3 \dots a_i \dots a_l ] \quad (4.1)$$

$l$  は入力単語音声の長さ (フレーム数) で、 $a_i$  は第

$i$  フレームの PARCOR 係数群 (1 次 ~  $M$  次) と音声パワーパラメータからなる、 $(M+1)$  次元ベクトルである。

また、予め登録されている  $N$  個の標準パターンの中で、第  $n$  番目のパターンを次の式で表わすことにする。

$$\mathbf{B}^n = [ b_1^n, b_2^n, b_3^n, \dots, b_j^n, \dots, b_l^n ] \quad (4.2)$$

ここで入力パターン  $\mathbf{A}$  と標準パターン  $\mathbf{B}^n$  との間の距離関数としてスカラー量  $D(\mathbf{A}, \mathbf{B}^n)$  を定義すると、 $N$  個のパターンの中で、この値を最小にするパターンに属すると判定するのがパターンマッチング法である。当然の事ながら、判定結果はこの距離関数  $D(\mathbf{A}, \mathbf{B}^n)$  の定義の仕方によって大きく影響を受ける。

さて音声パターン  $\mathbf{A}, \mathbf{B}$  を較するとき、まず発声時間の長短が問題になる。たとえば、オハヨゴザイマスとオハヨゴザイマスは後者の方が長いが同じ言葉である。したがってパターンマッチングの際の時間軸整合は単純で一様な線形伸縮ではなく、上記のことを考えた非線形伸縮で重ね合わせる必要がある。この整合法は時間軸の変動を吸収するのに好都合であるが、あまりに伸縮の幅を広くすると異なる単語でも同一単語と誤認識する恐れがあり、一定の制限を設ける必要がある。

### 5. 距離関数と DP マッチング

ここでは距離関数として、重みつきユークリッド距離を採用した。これに時間軸伸縮の概念をとり入れて関数を定義するわけであるが、図 4 に時間軸整合の模様を示す。 $\mathbf{A}, \mathbf{B}$  両パターンの対応部分は図中の曲線で示される。時間軸の対応は格子点  $(i, j)$  の系列  $F$  で表わすことができる。

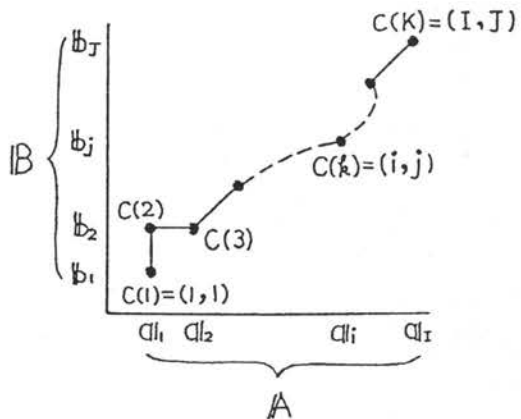


図 3 時間軸整合と DP パス

$$F = [C(1), C(2), \dots, C(k) \dots C(K)] \quad (5.1)$$

ただし,  $C(k) = (i(k), j(k)) \dots k$  番目の格子点

ここで対応する部分の特徴ベクトル  $\mathbf{a}_i$  と  $\mathbf{b}_j$  の距離をベクトルのユークリッド距離として定義すると, その部分の距離  $d(c)$  は

$$d(c) = d(i, j) = | \mathbf{a}_i - \mathbf{b}_j | \quad (5.2)$$

対応経路  $F$  に沿って計算したパターン間距離  $E$  は  $d(c)$  の荷重平均となり

$$E(F) = \frac{\sum_{k=1}^K d(C(k)) W(k)}{\sum_{k=1}^K W(k)} \quad (5.3)$$

$W(k)$  は重みを表わし, その値は対応の仕方によって変わるものである. すなわち経路  $F$  が斜め (伸縮のない部分) では  $W(k) = 2$ , 垂直又は水平 (伸縮のある部分) では  $W(k) = 1$  とすると格子点の数に関係なくパターン間距離は正規化されることになる.

この距離  $E(F)$  を最小にするよう部分伸縮するのが前節の非線形マッチングであり, マッチングの経路  $F$  を DP パス (path) という. 計算の効率化のため動的プログラミング (DP) を使うためである.

この DP マッチング<sup>1)</sup> は伸縮の幅を大きくしすぎると異種パターンにまでマッチングし誤認識の恐れが生じる. したがって一定の制限条件を設ける必要がある. 許容されるマッチングの範囲を整合窓といい, 窓寸法  $r$  を指定するとマッチングの範囲は制限され  $|i - j| \leq r$  となる. また部分伸縮についてもその度合いを制限する方がよい. たとえば,  $1/2 \sim 2$  などのようにである. 伸縮制限指数  $L$  を設定すると許容伸縮限度は  $(1/L \sim L)$  となる.

### 6. 個人別音声パターンの平均化

不特定多数話者の単語音声認識において, もっとも重要な鍵を握るのが汎用標準パターンの作成法である. 筆

者らは多数話者の特徴パターンを平均化する方法で話者の個性を抑え, 汎用的な標準パターンを合成した.

平均化の際にまず問題になるのが時間軸の取扱いである.  $\mathbf{A}, \mathbf{B}$  2つのパターンの整合には, DP マッチングで時間軸を部分的に伸縮して対応部分を探したが, 平均パターンの作成に際しては, この時間軸の伸縮も2つのパターンの中間的なものにしなければならない. 各部の対応関係を明確にしておくため, DP マッチングと同時に DP パス (経路  $F$ ) を記憶しておく必要がある.

この DP パスを使って, 時系列パターン  $\mathbf{A}, \mathbf{B}$  を時間軸伸縮も含めて平均化したパターン  $\mathbf{C}$  を求めるアルゴリズムについて説明する.

まず  $\mathbf{A}, \mathbf{B}$  パターンとその平均パターン  $\mathbf{C}$  を次のように表わすことにする.

$$\mathbf{A} = [ \mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \dots, \mathbf{a}_i, \dots, \mathbf{a}_I ] \quad (6.1)$$

$$\mathbf{B} = [ \mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3, \dots, \mathbf{b}_j, \dots, \mathbf{b}_J ] \quad (6.2)$$

$$\mathbf{C} = [ \mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_3, \dots, \mathbf{c}_h, \dots, \mathbf{c}_H ] \quad (6.3)$$

ただし

$$H = \text{int} [ (I + J) / 2 ] \quad (6.4)$$

ここで  $\text{int}$  は小数部を除去する演算子で, これは離散系に合わせる処置である.  $\mathbf{a}_i, \mathbf{b}_j, \mathbf{c}_h$  はそれぞれ  $i, j, h$  番目の特徴ベクトルである. 式 (6.4) は  $\mathbf{A}, \mathbf{B}$  両パターンを平均したパターン  $\mathbf{C}$  のフレーム数 (時間長) が両者のほぼ中間になることを示す.

次に平均パターン  $\mathbf{C}$  の各要素  $\mathbf{c}$  を求める方法を考える. パターン  $\mathbf{A}, \mathbf{B}$  の各要素  $\mathbf{a}$  と  $\mathbf{b}$  の最適対応 (DP パス) が図 4 のようであったとする. 両者の平均パターンの要素  $\mathbf{c}$  は, 図上の対応する要素  $\mathbf{a}$  と  $\mathbf{b}$  を結んだ線分の中央に相当する時刻に,  $\mathbf{c} = (\mathbf{a} + \mathbf{b}) / 2$  となることが原則である.

ただし, この時系列はフレーム単位の離散系であるため図の□点のような中途半端な位置はとり得ないので, このような場合は半フレームだけ左に寄せて, すぐ左の・点と考えて近似処理 (図 4 の  $\bullet \leftarrow \square$  の部分) をする.

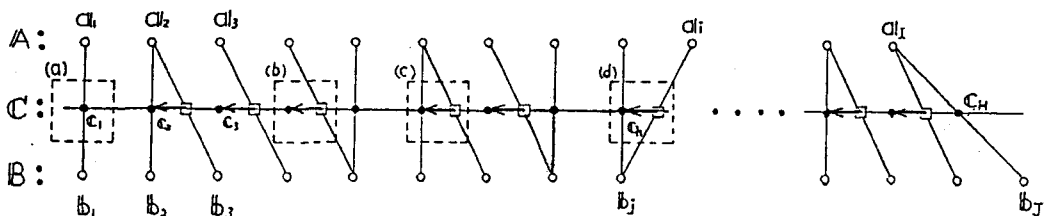


図 4 パターン要素間の対応と時間軸の平均化

以上のことから、要素  $c_h$  の計算には  $a, b$  の対応関係も考慮して、次のように定義すればよい。

(1) 1:1 対応

$a_i \leftrightarrow b_j$  の対応 (図4の [a], [b] 部) では

$$c_h = (a_i + b_j) / 2$$

(2) 1:2 対応

$a_i \leftrightarrow b_j, b_{j+1}$  の対応 (図4の [c] 部) では

$$c_h = (2a_i + b_j + b_{j+1}) / 4$$

$a_i, a_{i+1} \leftrightarrow b_j$  の対応 (図4の [d] 部) では

$$c_h = (a_i + a_{i+1} + 2b_j) / 4$$

以上の手順で、類似した2つの時系列パターンについて時間軸伸縮も含めて平均化したパターンを合成することができる。さらにこの方法を拡張すれば多数パターンから1個の標準パターンが合成できる。まずパターン **A** と **B** を平均し、それと他のパターン **C** と **D** を平均したものをもう一度平均する、という様に平均化を進めれば、 $2^n$  個 ( $n$  整数) の原パターンから1つの標準パターンを合成することができる。

この方法で得られる標準パターンは、話者の特徴を平均化して打ち消し合うため、話者独立性の高い汎用性のあるパターンと考えられる。

## 7. 音声パワーのパラメータ化

さきに述べたように PARCOR 係数には、音声の強弱や高低の情報は含まれていない。これは発声のときの口腔の形状に相当するものである。したがって PARCOR 係数だけで音声認識を行うことは、口の動きだけを見て判断するようなものである。自然言語では常に強弱や高低の変化があり、これらを認識に利用するのが妥当であろう。

図5はこの例を示したものである。話者が異なっても

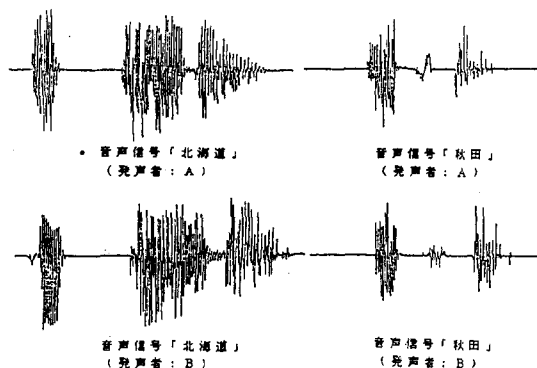


図5 音声パワーの変動

「北海道」「秋田」それぞれの振幅パターンは類似した形状を示す。したがって筆者らは、音声パワーをパラメータ化し、従来の PARCOR 係数に追加することにした。

フレーム単位の音声パワー  $p$  は (振幅)<sup>2</sup> のフレーム長毎の積分値であるから、パワーパラメータの時系列要素  $P_i$  は正規化して次のように定式化される。

$$P_i = \frac{n p_i}{2 \sum_{j=1}^n p_j} \quad (7.1)$$

ここで  $n$  はフレーム数 (単語長)、 $p_i$  は第  $i$  番目のフレームの音声パワーである。

実験の結果、パワーパラメータ  $P_i$  は、通常の単語で 0~6 くらいの範囲で変化することがわかった。

## 8. プログラムの構成

プログラムは前処理段階から認識まで、つぎの4ブロックから構成される。使用言語は FORTRAN である。

(1) 単語切り出しプログラム

一連の音声信号から、それぞれの単語部分を切り出すプログラムである。音声レベルが一定のしきい値より低い場合を無音部とし、基本的にこの部分で単語音声を切り出すようにしている。

ただし、単語内の無音部と混同をさけるため、発声の際に単語間の区切りの休止期間を十分長くとするようにしている。その上このプログラムでは、まずすべての無音部を検出し、それを長さの順に並べ、上位から (単語数 + 1) 番目までを単語間無音部とする方式をとった。切り出した後も誤切断がないか耳で聞いて確認をしている。

(2) 分析プログラム

切り出した音声データを一定長 (256) でフレーム化し、PARCOR 係数と音声パワーパラメータを計算して発声者毎のパターンファイルに登録する。

(3) 平均化プログラム

パターン群を DP マッチングによって平均化し、それぞれの標準パターンとして登録する。

(4) 認識プログラム

入力パターンに対して、登録された標準パターンを順次 DP マッチングで比較し、距離最小のパターンを認識結果として選出する。

## 9. 認識実験

音声の収録段階はオフラインの音声収録装置 (パーソナルコンピュータ) で行い、収録されたデータは、1 MB

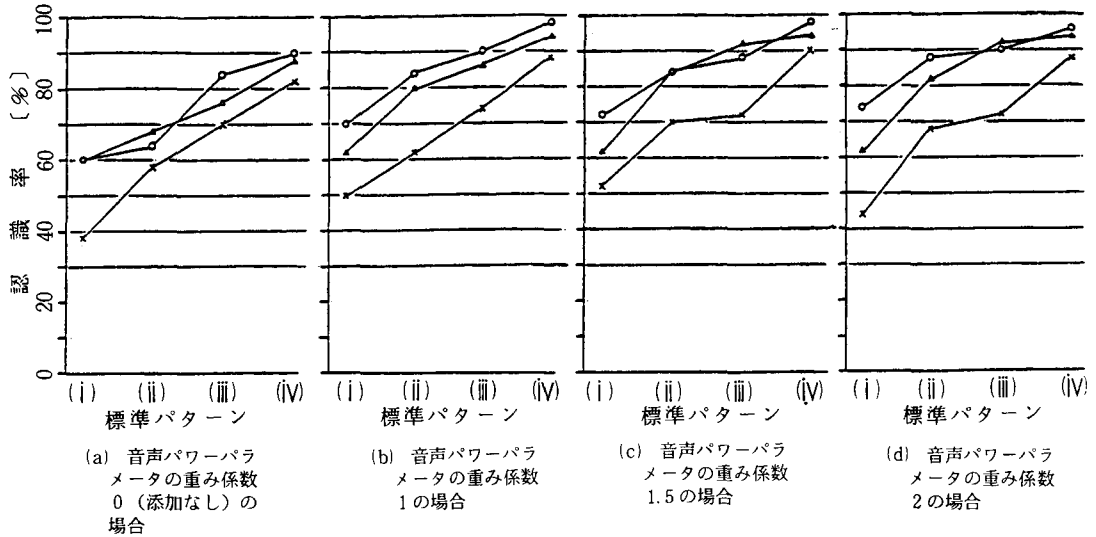


図 6 認識実験

のフロッピーディスクで情報処理センターに持ち込み、中型計算機 FACOM M360R で処理する方式をとった。

音声データは「札幌」「青森」「秋田」など50都市名を11人の男性(20才前後)が発声したものである。このうち8名分は標準パターン作成用とし、他の3名分を独立パターン(A, B, C)とし、それぞれについて認識実験を行った。図6はその結果を示す。図中のマーク○, △, ×はそれぞれ被験者(A, B, C)の音声を示す。

ここで横軸の(i), (ii), (iii), (iv)は標準パターンのレベルを示すもので、(i)は単独パターン、(ii)は2個のパターンを平均化したもの、(iii), (iv)はそれぞれ4個、8個のパターンを平均化した標準パターンを示す。

特徴パラメータは PARCOR 係数(1~15次)に(7.1)の音声パワーパラメータを加え、この重み係数を種々変更して実験したものである。

標準パターンのレベルが高くなるほど(参加人数が多くなるほど)認識率は向上している。また重み係数は1.5附近がもっとも認識率が高くなるようである。

## 10. 考 察

標準パターンの合成に参加する人数が多くなるほど良好な成績が得られる。人数を16か32くらいまで増やすとかなり良質の汎用パターンが得られるであろう。

音声パワーパラメータを加えた効果は明らかで、左端(a)(パワーパラメータなし)と比べると右の3つのグラフは何れも認識率が高くなっており、(c)の重み1.5の場合と比較すると6~8%向上している。

この結果からみると、7・で述べたように PARCOR 係数だけで認識を行う手法は、いわば口の動きだけで言葉を推定するようなもので、本質的に欠陥があると言える。

○, △, ×の個人差によるバラつきもかなり認められるが、標準パターンのレベルが上がるほどバラつきも少なくなっているようである。

パワーパラメータの重みは1.5程度が適当となっているが、このパラメータ自身の変動振幅が大きいので、PARCOR 係数に置きかえると、3~4個分くらいと思われる。パワーパラメータのとり得る値は0~6くらいで平均値は0.8程度、一方 PARCOR 係数のとり得る値は-1~+1で絶対値の平均は0.3~0.4となっているからである。

## 11. むすび

前報<sup>1)</sup> につづいて、筆者らの提唱した標準パターンDP 合成法が、不得定多数話者の単語音声認識に大変有効であることが再確認できた。

また、従来の PARCOR 係数に音声パワーパラメータを加えることによって、認識率が飛躍的に向上することがわかった。標準パターンの合成に参加する人数をもう1~2段階ふやして16名か32名にすれば、50単語で100%近い認識率が得られると思う。

音声データ収録に協力して頂いた電気工学科の学生諸君や便宜をはかって頂いた情報処理センターの職員各位に感謝の意を表す。

参 考 文 献

- 1) 高橋ほか 標準パターン DP 合成法による多数話者の単語音声認識 府立高専研究紀要 (1987)
- 2) 板倉 統計的手法による音声の特徴抽出 音声処理シンポジウム (1971)
- 3) 梅崎, 中川 単音節を標準パターンとした連続音声認識における発声速度の問題 音声研究資料 (1985-6)
- 4) 藤崎博也 連続音声認識のための音節の変動に関する検討 音声研資料 (1984-12)