

音声認識における時系列パターン照合アルゴリズムの展開

Time-Sequential-Pattern Matching Algorithm on Automatic Speech Recognition

中川 聖一*
Seiichi Nakagawa

* 豊橋技術科学大学情報工学系
Dept. of Information and Computer Sciences, Toyohashi University of Technology.

1988年4月8日 受理

Keywords: speech recognition, DP matching, HMM, time warping, dynamic programming, one pass DP, one stage DP, forward-backward algorithm, Baum-Welch estimation algorithm.

1. 音声認識におけるパターン照合の問題

単語や文は、音韻とか音節とかの言語の基本構成単位から構成されているので、音韻認識や音節認識が確実にできれば単語音声認識や連続音声認識も確実にできる。しかし、実際には、話者や前後のコンテキストによって音声パターンが変動するので(個人差、調音結合の現象)、音韻や音節の認識は容易でない。また、たとえ発声者と発声内容が同一であっても、発声ごとに音声パターンは異なる。特に時間軸上でのパターンの変動がみられる。

時間軸上での変動に対しては、時系列パターンの時間軸上での非線形伸縮を許すパターン照合アルゴリズム、いわゆる DP マッチング法が開発され実用化されている。個人差の問題に対して最も安易な対処方法は、話者を限定することであり、この場合は特定話者方式と呼ばれている。一方、話者を限定しない場合は個人差の正規化や適応化手法が試みられてきた。しかし、満足できる正規化手法は開発されていない。通常は、統計的手法で個人差の問題に対処している。調音結合の問題に対しては、調音結合の影響を受けたままの単語音声パターンを一つの観測パターンとみなして、単語音声どうして直接照合する方法が成功を収めてきた。しかし、この考え方を大語彙の単語音声や連続音声認識に適用することはできない。なぜなら、単語パターンの作成が容易でないことや単語間にわたる調音結合には対処できないからである。

本稿では、音声認識の基本的アルゴリズムとしては確立されている DP マッチング法と HMM (Hidden Markov Model, 隠れマルコフモデル) について解説する。DP マッチング法は、二つの異なる単語音声の時系列パターンを非線形に伸縮しながら照合する手法である。一方、HMM は、個人差や調音結合(単語単位の認識の場合は、単語内の調音結合は問題ないが)によるパターンの変動を確率モデルとしてとらえる方法である。もちろん、時間長の変動も確率モデルで吸収できる。これらの点から原理的には DP マッチン

表1 DP マッチングによる単語音声認識結果

語彙	話者モード	認識率 [%]
10 数字	不特定	99.2
39 数字, アルファベット	特定	79.5
	不特定	79.0
54 コマンド語	不特定	96.5
129 航空座席, 予約単語	特定	88.0
	不特定	91.0
1109 基本単語	特定	79.2

表2 HMM による不特定話者の数字音声認識結果

方法	モデル	認識率 [%]
HMM	離散分布	97.1
	連続分布 (1モデル/数字)	99.1
	連続分布 (2モデル/数字)	99.4
	5 混合対角分布 (2モデル/数字)	99.3
DP	12 テンプレート/数字	99.4

グ法よりも HMM 法のほうが精度よく標準パターンを構成できる。しかし、実際には確率モデルのパラメータ数が多くなるとパラメータの精度よい推定が困難になるから、結局、単語音声を粗っぽいモデル（たとえば、単語音声を時間的に 5~10 状態）で表現せざるを得ないので、マルチテンプレートによる DP マッチング法と比べて認識精度は現在のところ大差はない。

表 1, 2 は AT&T ベル電話研究所で行われてきた一連の音声認識研究の結果をまとめたものである⁽¹⁾⁽²⁾。

2. DP マッチングによる音声認識アルゴリズム⁽³⁾⁽⁴⁾

2.1 孤立単語音声の認識⁽⁴⁾⁽⁶⁾

(a) 単語音声認識の定式化

一つの単語は発声されるたびに継続時間が変わり、その時間構造も線形ではない。発声速度が速くなるに従って、母音の定常部分は短くなるのに対して、子音部やわたりの部分は比較的固有の長さを保っている。このような性質をもった二つの不等長なパターンを比較する場合、時間長の違いによって生ずる本質的でない差を極力排除した距離尺度が望ましい。このような考え方から時間軸に非線形な変換（これを時間変換関数：time warping function と呼ぶ）を施すことにより、パターン長の差の影響をできるかぎり吸収しようとする方式が DP マッチングと呼ばれているものである。

二つの音声パターン A と B が、時間標本化された特徴ベクトルの時系列として表現されている場合に問題を定式化しよう。

$$A = a_1 a_2 \dots a_i \dots a_I$$

$$B = b_1 b_2 \dots b_j \dots b_J$$

とする。ここで、 a_i は A の第 i フレームを示すものとする ($a_i = A(i)$)。 I, J はそれぞれ A および B の時間長でフレーム総数である。

ここで、 A, B 両パターンの時間の対応すなわち時間変換関数は、 $i-j$ 平面上の格子点 $c = (i, j)$ の系列 F で表現することができる。

$$F = c(1) c(2) \dots c(k) \dots c(K)$$

ただし、系列点 $c(k) = (i(k), j(k))$ は i 軸では $i(k)$ 、 j 軸では $j(k)$ の格子点である。二つの特徴ベクトル a_i と b_j との間では、いくつかの距離を定義できるが、その定義の仕方は一応問題にせず、便宜的にこの二者間の距離を格子点を仲介にして $d(c) = d(i, j)$ で表すと、点列 F に沿ってとった距離の荷重平均（時間正規化距離）において点列 F を変化させたときの最小値で、二つのパターン A と B との距離 $D(A, B)$ を定義する。

$$D(A, B) = \min_F \frac{\sum_{k=1}^K d(c(k)) \cdot w(k)}{\sum_{k=1}^K w(k)} \quad (1)$$

式(1)の分母 $\sum_{k=1}^K w(k)$ は点列 F を構成する点の数による距離の値の差を補正するための正規化係数であり、その値は F に依存する。このとき式(1)の最小化問題は複雑になるので、式(1)の分母を F に依存しない定数とすることによりこれを簡単化する。

$$N = \sum_{k=1}^K w(k)$$

とおくと、(1)式は、

$$D(A, B) = \frac{1}{N} \min_F \sum_{k=1}^K d(c(k)) \cdot w(k) \quad (2)$$

となり、最小化する目的関数が加法的になるので、次の(b)で述べるように、その最小化は動的計画法の手法を用いて効率よく解くことができる。

N を定数化する方法は種々ありうるが、次の二つがよく試みられている。

- ① 対称形： $w(k) = (i(k) - i(k-1)) + (j(k) - j(k-1))$,

このとき、 $N = I + J$ (図 1 (a))

- ② 非対称形： $w(k) = i(k) - i(k-1) = 1$,

$$0 \leq j(k) - j(k-1) \leq 2,$$

このとき、 $N = I$ (図 1 (c))

(b) 動的計画法による最小化問題の解法

ここでは、式(2)の最小化問題を動的計画法 (Dynamic Programming) を用いて効率よく解く方法を

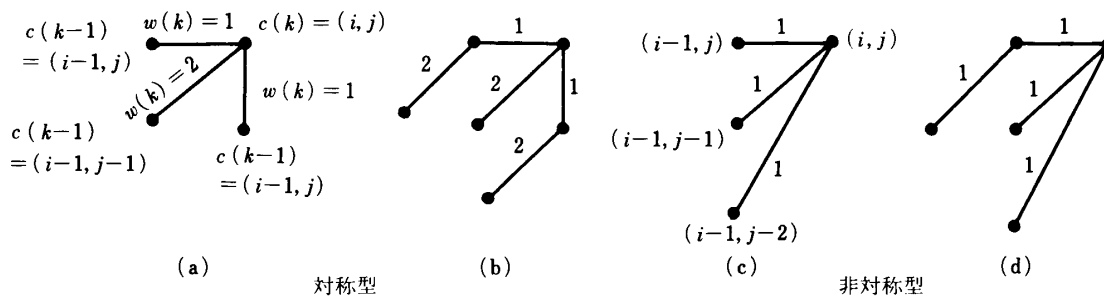


図 1 各種 DP パスと荷重係数の例

説明する。まず部分点列 $(c(1) c(2) \dots c(k) (c(k)=i(k), j(k))$ であるが、便宜上 $c(k) = (i, j)$ と略記する) に対する次の部分積 $g(c(k))$ を考える。

$$g(c(k)) = g(i, j) = \min_{c(l) \dots c(k)} \left(\sum_{l=1}^k d(c(l)) \cdot w(l) \right) \quad (3)$$

上式を変形すると、

$$g(c(k)) = \min_{c(l) \dots c(k)} \left(\sum_{l=1}^{k-1} d(c(l)) \cdot w(l) + d(c(k)) \cdot w(k) \right) \\ = \min_{c(k-1)} \left(g(c(k-1)) + d(c(k)) \cdot w(k) \right) \quad (4)$$

$c(k) = (i, j)$ であるから、対称形の $w(k)$ を用いることにすると、点列に対する単調連続性の条件から $c(k-1)$ は $(i, j-1)$, $(i-1, j)$ および $(i-1, j-1)$ のいずれかである。故に、式(4)は次のように書ける。

$$g(i, j) = \min \left\{ \begin{array}{l} g(i, j-1) + d(i, j) \\ g(i-1, j-1) + 2d(i, j) \\ g(i-1, j) + d(i, j) \end{array} \right\} \quad (5)$$

したがって、式(2)の最小化問題は次の3段階の計算によって解くことができる。

- ① 初期設定： $g(1, 1) = 2d(1, 1) = 2d(a_1, b_1)$ および $i = 1$
- ② $j_1 = \max(1, i-r), j_2 = \min(i+r, J)$ として j を j_1 から j_2 まで一つずつ増やしながら式(5)によって $g(i, j)$ を計算する ($2r$ は整合窓長)。
- ③ i が I に等しくなければ i を1だけ増して②へいく。等しければ、 $D(A, B) = g(I, J)/N$ とする。式(5)を用いる場合は $N = I + J$ である。整合窓は極端な非線形変換を防ぐためのものである。

上記②の計算がフレーム周期 (通常 10~20 ms) 以内で処理できれば、 A が入力パターンの場合、実時間で $D(A, B)$ が計算できることになる。

DP マッチング法を単語音声の認識に適用するには次のようにする。認識の対象となっている単語を n ($n=1 \sim N$)、その標準パターンを $B(n)$ で表す。入力 A と各標準パターン $B(n)$ との距離 $D^n = D(A, B(n))$ を上記の方法で計算し、 $D^n = \min_n (D^n)$ を与える単語 \hat{n} を A に対する認識結果とする。

このほか、非線形伸縮比を局所的にいつも $1/2 \sim 2$ に押さえる方法として傾斜制限つき DP パスを用いる方法がある。例えば、図 1 (d) の非対称形パスの場合は、

$$g(i, j) = \min \left\{ \begin{array}{l} g(i-2, j-1) + d(i-1, j) + d(i, j) \\ g(i-1, j-1) + d(i, j) \\ g(i-1, j-2) + d(i, j) \end{array} \right\} \quad (6)$$

$$D(A, B) = g(I, J)/I$$

となる。

2.2 連続単語音声の認識

(a) 連続単語音声認識の定式化⁽⁷⁾⁽⁸⁾

次に動的計画法に基づいた DP マッチング法の連続単語への拡張について述べる。

認識の対象となっている単語を n ($n=1 \sim N$) で、その標準パターンを $B(n) = b_1^n b_2^n \dots b_{J^n}$ (J^n は標準パターン $B(n)$ の長さ) で示す。また入力パターンを $A = a_1 a_2 \dots a_I$ とする。 x 個の単語からなる単語列 $n = n_1 n_2 \dots n_x$ を発声したときのパターンを、標準パターン $B(n_1), B(n_2), \dots, B(n_x)$ を接続して得られるパターン、 $B = B(n_1) \oplus B(n_2) \oplus \dots \oplus B(n_x)$ (\oplus でパターンの接続を示す) によって近似する。

この B を連続単語の標準パターンと考えて、前節で述べた DP マッチング法によって入力パターン A との間で時間正規化マッチングを実行すると、両者の間の距離が求められると同時に、入力パターン A の単語単位のセグメンテーションも、最適な照合パスをバックトレースすることによって得られる。単語列 n (したがって個数 x も) を変化させたとき、両者の間の距離を最小にするような単語列を認識結果とするわけである。このときの距離の最小値を $D'(I)$ (I は A の長さ) と書くと、

$$D'(I) = \min_{n, x} (D(A, B(n_1) \oplus B(n_2) \oplus \dots \oplus B(n_x))) \quad (7)$$

を満たす単語数 x および単語列 $n = n_1 n_2 \dots n_x$ を求めようとするのである。

式(7)を総当たり法で解くと N^X (X : 入力パターンとマッチングをとる単語列中の最大の単語数) のオーダーの DP マッチングの回数を必要とする。そこで、この最小化問題を2段階に分割して計算量の軽減を図ることを考える。

入力パターン A の部分パターン $a_{l+1} a_{l+2} \dots a_m$ を $A(l, m)$ 、 D および D' に対する正規化前の累積距離を T および T' と書く。非対称形の DP パス ($w(k) = i(k) - i(k-1) = 1$) の性質を用いれば、途中の導出過程は略するが、次の漸化的な表現を得る。

$$T'(I) = \min_{1 < l < I} (T'(l) + \min_n \{ T(A(l, I), B(n)) \}) \quad (8)$$

上式の誘導は A の任意の部分パターン $A(0, i)$ に対して成立するから、上式の I (入力パターン全体の長さ) は任意の $i (< I)$ に置き換えることができ、

$$T'(i) = \min_{1 \leq l < i} (T'(l) + \min_n \{T(A(l, i), B(n))\})$$

$$(i = 1, 2, \dots, I) \quad (9)$$

ここで、 $T'(0) = 0$ とする。上式を満たす l, n をそれぞれ $C(i), N(i)$ と定義する。すなわち、 $N(i)$ は入力音声 A の第 i フレームで終る全区間で DP マッチングがとれた最適単語列の最後尾の単語名、 $C(i)$ は $N(i)$ の入力音声パターンとのマッチング開始フレームの値 ($l+1$) に対し、一つ前の l を指し単語境界フレームに該当する。式(9)を解釈すれば、入力音声パターンが i フレームで終端すると仮定した場合、その最適な単語列 (単語数 x) は入力パターンの第 l フレームまでの最適な単語列 (単語数 $x-1$) と、入力パターンのそれに続く $l+1 \sim i$ フレームと照合がとれる最適単語を結合させたものであり、式(9)の $T'(i)$ は l を可変とした際の最適値を示している。連続単語としての全体の認識結果は、 $N(I), N(C(I)), N(C(C(I))), \dots$ として、最適単語列の最後尾単語名から逆順に求めることができる。

(b) 2段 DP マッチング法⁽⁸⁾

式(9)は式(7)の最小化問題がやはり動的計画法によって解けることを示している。式(9)の $T(A, B)$ の計算は前節で述べた DP マッチング法を使用するので、結局、式(9)は動的計画法を2段階に分けて使用することになる。この意味から、式(9)を直接解く方法は2段 DP 法と呼ばれる。

式(9)において l の範囲は原理的には $1 < l < i$ であるが、前項で述べた整合窓の条件により、 $i - r - \max(J^n) \leq l \leq i + r - \min(J^n)$ (r は窓の大きさ、 J^n は単語 n の標準パターンの長さ) なる範囲をとればよいことがわかる。さらに、1回の T の計算によって多くの部分パターンとの距離計算が並行的に行われる。すなわち、単語 n の標準パターンと入力パターンの第 $l+1$ フレームからの入力の部分パターンとの累積距離を求める際に、並行的に $T(A(l, i), B(n))$ ($l+J^n-r \leq i \leq l+J^n+r$) を一度に求めることができるという事実注意到すると、2段 DP マッチング法は式(9)から予想されるよりかなり少ない計算量で連続単語の認識を可能にしていることがわかる。

(c) $O(n)$ DP マッチング法⁽⁹⁾, One Pass DP 法⁽¹⁰⁾, One Stage DP 法⁽¹¹⁾

式(9)は次のように変形できる。

$$T'(i) = \min_{1 \leq l < i} (\min_n \{ T'(l) + T(A(l, i), B(n)) \}) \quad (10)$$

$T(A(l, i), B(n))$ は、孤立単語の認識の項で述べたように、初期設定を、 $g^n(l+1, 1) = d^n(l+1, 1)$, $g^n(l'+1, 1) = \infty$ ($l' > l$) として以下の漸化式を求めていくことによって得られる (ただし、 $w(k)$ は非対称形)。

$$g^n(h, k) = \min \left\{ \begin{array}{l} g^n(h-1, k) \\ g^n(h-1, k-1) \\ g^n(h-1, k-2) \end{array} \right\} + d^n(h, k) \quad (11)$$

$$T(A(l, i), B(n)) = g^n(i, J^n)$$

しかし、この場合、初期設定を $g^n(l+1, 1) = T'(l) + d^n(l+1, 1)$, $0 < l < i$ として、 $T'(l)$ を g の中へ組み込めば式(10)は、

$$T'(i) = \min_n (g^n(i, J^n)) \quad (12)$$

となり1段の動的計画法を表現することになる。しかも i に関しては2段 DP 法の項で述べたように、一度の T の計算に並行して種々の i について並行的に計算していくことができる。こう考えることによって、2段 DP マッチングより効率のよいアルゴリズムが誘導できる。図2はこのマッチング法のようなすを示している。

この方法の計算量は、入力パターン長 I , 平均標準パターン長 J , 標準パターン数 N に比例することと DP (T の計算) の回数が単語数に比例することから $O(n)$ DP 法⁽⁹⁾ と呼ばれている。この方法は、迫江らによるクロック伝搬法⁽¹²⁾ の特殊な場合とも解釈できる。また、Bridleによっても独立に考案され、Neyによって明確なアルゴリズムが示され、それぞれ One Pass DP 法⁽¹⁰⁾, One Stage DP 法⁽¹¹⁾ とも呼ばれている。

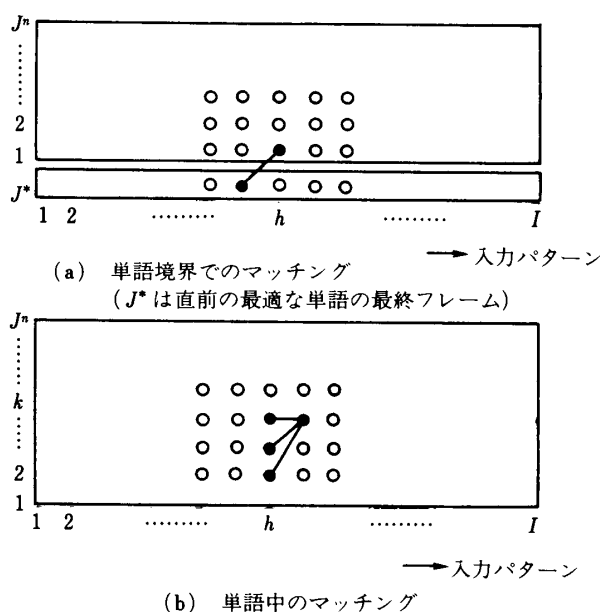


図2 $O(n)$ DP マッチング法, One Pass DP 法におけるマッチング処理

つまり、孤立単語の DP マッチングによる認識の計算量と同じオーダで連続単語の認識が可能となる。

このほかにも種々のアルゴリズムが開発されている。Level Building 法は、第 n レベルの処理として、入力パターンと $(n-1)$ 単語列との最適な照合結果を初期値とし各単語の照合を連結していき、 n 単語列としての最適な照合結果を得るもので、2 段 DP 法よりも計算量は少ない⁽¹³⁾。また、拡張連続 DP 法は、式(9)の第 2 項をワードスポッティングアルゴリズムで能率よく求める方法で、近似解(ただし、厳密解と認識精度は同等)となるが計算量はさらに少ない特徴がある⁽¹⁴⁾。

2.3 オートマトン制御 DP マッチング⁽¹⁵⁾

自然言語文の音声で代表される連続音声は、協力的で明確に発声された場合は、上述の連続単語音声の場合と同様にパターンマッチング法が適用可能である。ただ単語間の接続に関して構文的・意味的制約がある点が異なる。

連続単語音声の場合は、単語間の接続に関しては何も制約がなくあらゆる単語が接続可能であった。連続音声の場合は、接続できる単語と接続できない単語を区別して処理しなければならない。この関係は状態遷移図で表すことができ、一般に有限状態オートマトンで表現できる。

図 3 は簡単な英文を生成する有限状態オートマトンである。たとえば、状態 ART では、boy, girl, little, big の 4 単語が接続可能である。二重丸印の状態は最終状態を表し、この状態に達すると、入力は終了することを示す。

有限状態オートマトン α は次のように定義できる。

- $\alpha = \langle K, \Sigma, \delta, q_0, F \rangle$
- K : 状態 q_i の有限集合, $\{q_i\}$
- Σ : 入力単語 n の有限集合 $\{n\}$
- δ : 状態遷移関数, $K \times \Sigma \rightarrow K, \{\delta(q_i, n) = q_j\}$
- q_0 : 初期状態, $q_0 \in K$
- F : 最終状態の集合, $F \subset K$

単語の接続条件がオートマトン α で表現される場合、単語列 n_1, n_2, \dots がこのオートマトン α に受理さ

れるという条件のもとに、式(9)と類似した以下の式を最小化する最適な単語列 n を求めるとよい。

$$T'(q_k, i) = \min_{l, q_j, n} (T'(q_j, l) + T(A(l, i), B(n))) \tag{13}$$

ただし、 $q_k = \delta(q_j, n)$ 。上式を満たす $\hat{n}, \hat{l}, \hat{q}_j$ によって $N(q_k, i), C(q_k, i)$ が定義され、新たに状態遷移の結果保存のために $Q(q_k, i)$ を定義算出しておく必要がある ($N(q_k, i) = \hat{n}, C(q_k, i) = \hat{l}, Q(q_k, i) = \hat{q}_j$)。

認識結果は、 $q' = \arg \min (T'(q, I))$ として $N(q', I), N(Q(q', I), C(q', I)), \dots$ のようにしてオートマトンで受理される単語列のうち最適単語列が最後尾単語名から逆順に求められる。

この定式化の具体的なアルゴリズムは前述の 2 段 DP 法や $O(n)$ DP 法に容易に組み込むことができる(オートマトン制御 2 段 DP 法⁽¹⁵⁾、オートマトン制御クロック同期伝播形 DP 法⁽¹²⁾ などと呼ばれている)。

3. Hidden Markov モデルによる音声認識アルゴリズム⁽¹⁶⁾⁻⁽²¹⁾

3.1 基本モデルと認識アルゴリズム

(a) 基本モデル

DP マッチングは、時系列パターンの時間的構造の変動に対しては強力であるが、話者の個人差などに起因するスペクトルそのものの変動に対しては弱い。パターンマッチングに基づく音声認識の立場からは、これらの変動成分を統計量によって表現する方法がある。この考え方は、時系列パターンの時間的構造の変動に対しても適用できる。DP マッチングにおいては、単語に対して標準的な時系列パターンを標準パターンとしているが、HMM(Hidden Markov Model)による方法は、各単語を標準的な確率状態遷移機械で表現する⁽¹⁷⁾。このモデルのパラメータは、スペクトルや時間的な変動を最もよく吸収するように学習サンプルによって求められる。

図 4 は、Bakis モデルと呼ばれ、最もよく用いられている HMM の例を示している。もちろん、状態数は適当に決めておく必要がある。状態数を多くすれば、きめ細かい単語パターンの表現ができるが、モデルの

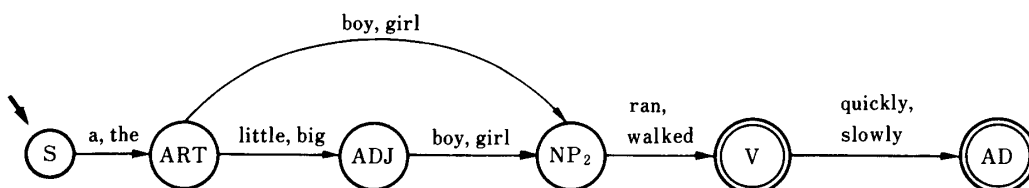


図 3 簡単な英文を生成する有限状態オートマトン

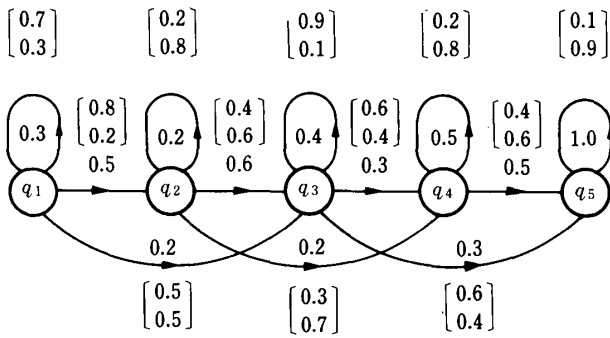


図4 Bakisモデルの例
(□内の数値は、シンボル a, b の出力確率)

パラメータ数が多くなりパラメータ推定の精度が悪くなる。図4の状態遷移のアーク上の数値は状態遷移確率を表している。これは、状態数を S とすると $S \times S$ の行列によって表現できる。状態 q_i から状態 q_j への遷移確率を a_{ij} と書くと図4の例では、

$$A = \{a_{ij}\} = \begin{pmatrix} 0.3 & 0.5 & 0.2 & 0.0 & 0.0 \\ 0.0 & 0.2 & 0.6 & 0.2 & 0.0 \\ 0.0 & 0.0 & 0.4 & 0.3 & 0.3 \\ 0.0 & 0.0 & 0.0 & 0.5 & 0.5 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \end{pmatrix}$$

となる。通常、音声パターンには、時間的な非可逆性の性質があるから、 $i > j$ なら $a_{ij} = 0$ である。

観測パターンが、有限個 (K 個) のシンボルの一つとして表現できる場合は (通常ベクトル量子化などの手法を用いる)、離散分布モデルと呼ばれる。状態 q_i から状態 q_j の遷移でシンボル k のパターンが観測 (出力) される確率を $b_{ij}(k)$ と表す。この $\{b_{ij}(k)\}$ をシンボル出力確率行列と呼ぶ。

簡単な例として、出力シンボルを $\{a, b\}$ の二つに限られた場合の例を図4の状態遷移のアーク上の [] 内に示す。この出力確率が、遷移元の状態によってのみ決定され状態遷移に独立なモデルも考えられる。すなわち、 $b_{ij}(k) = b_{i,k}(k)$ 。この場合は、出力確率行列は $S \times K$ の行列で表現できる。

各状態 q_i の初期確率を π_i で表す。図4の例では、 $\pi_1 = 1$ 、 $\pi_i = 0$ ($i > 1$) である。また、最終状態の集合を F と記す。図4の例では $F = \{q_5\}$ である。

(b) 認識アルゴリズム

今、 $\mathbf{y} = y_1 y_2 \dots y_T$ を出力シンボルの観測系列としよう。たとえば、10ms ごとの音声特徴パラメータ (スペクトルとかケプストラム) の時系列パターンが考えられる。このとき、この HMM の M によって \mathbf{y} が生起する確率 $P(\mathbf{y} | M)$ を求める。これを各単語に対応する HMM に対して求め、最大確率を与えるモデルに対応する単語を認識結果とする。

$\mathbf{q} = q_{i_0} q_{i_1} \dots q_{i_T}$ を状態遷移系列 (ただし、 $q_{i_T} \in F$) とすれば、

$$P(\mathbf{y} | M) = \sum_{i_0, i_1, \dots, i_T} P(\mathbf{y} | \mathbf{q}, M) \cdot P(\mathbf{q} | M) \quad (14)$$

$$\begin{aligned} P(\mathbf{q} | M) &= \prod_t P(q_{i_t} | q_{i_{t-1}}, M) \\ &= \prod_t P(q_{i_t} | q_{i_{t-1}}, M) \end{aligned} \quad (15)$$

(マルコフモデルの性質)

$$\begin{aligned} P(\mathbf{y} | \mathbf{q}, M) &= \prod_t P(y_t | q_{i_{t-1}}, q_{i_t}, M) \\ &= \prod_t P(y_t | q_{i_{t-1}}, q_{i_t}, M) \end{aligned} \quad (16)$$

$$\begin{aligned} P(\mathbf{y} | M) &= \sum_{i_0, i_1, \dots, i_T} \prod_t P(q_{i_t} | q_{i_{t-1}}, M) \cdot P(y_t | q_{i_{t-1}}, q_{i_t}, M) \\ &= \sum_{i_0, i_1, \dots, i_T} \pi_{i_0} a_{i_0 i_1} b_{i_0 i_1}(y_1) \\ &\quad a_{i_1 i_2} b_{i_1 i_2}(y_2) \dots a_{i_{T-1} i_T} b_{i_{T-1} i_T}(y_T) \end{aligned}$$

ここで、 $q_{i_0}^{i_{t-1}} = q_{i_0}, q_{i_1}, \dots, q_{i_{t-1}}$ 、

$$\mathbf{y}_i = y_1 y_2 \dots y_i \quad (17)$$

となる。なお、最後の式中のそれぞれの確率は、モデル M での条件確率である。簡単な例として図4でシンボル系列 aba が出力される確率を求めてみよう。この場合は、考えるべき状態遷移系列は、図4からもわかるように、左上隅から右下隅に至る経路に対応するもので次の6通りである。 $q_1 \rightarrow q_1 \rightarrow q_3 \rightarrow q_5$ 、 $q_1 \rightarrow q_2 \rightarrow q_3 \rightarrow q_5$ 、 $q_1 \rightarrow q_2 \rightarrow q_4 \rightarrow q_5$ 、 $q_1 \rightarrow q_3 \rightarrow q_3 \rightarrow q_5$ 、 $q_1 \rightarrow q_3 \rightarrow q_4 \rightarrow q_5$ 、 $q_1 \rightarrow q_3 \rightarrow q_5 \rightarrow q_5$ 。それぞれの確率は、

$$P_1 = 0.3 \times 0.7 \times 0.2 \times 0.5 \times 0.3 \times 0.6 = 0.00378$$

$$P_2 = 0.5 \times 0.8 \times 0.6 \times 0.6 \times 0.3 \times 0.6 = 0.02592$$

$$P_3 = 0.5 \times 0.8 \times 0.2 \times 0.7 \times 0.5 \times 0.4 = 0.0112$$

$$P_4 = 0.2 \times 0.5 \times 0.4 \times 0.1 \times 0.3 \times 0.6 = 0.00072$$

$$P_5 = 0.2 \times 0.5 \times 0.3 \times 0.4 \times 0.5 \times 0.4 = 0.0024$$

$$P_6 = 0.2 \times 0.5 \times 0.3 \times 0.4 \times 1.0 \times 0.1 = 0.0012$$

となる。故に、 $P(aba | M) = P_1 + P_2 + P_3 + P_4 + P_5 + P_6 = 0.04522$ となる。

一般に $P(\mathbf{y} | M)$ の値は、以下に述べる Forward アルゴリズムで求められる。

Forward 変数 $\alpha(i, t)$ を観測シンボル \mathbf{y}_i を出力して状態 q_i にある確率とすれば、

$$\alpha(i, 0) = \pi_i \text{ for } i = 1, 2, \dots, S$$

$$\alpha(i, t) = \sum_j \alpha(j, t-1) \cdot a_{ji} \cdot b_{ji}(y_t) \quad (18)$$

$$\text{for } t = 1, 2, \dots, T; i = 1, 2, \dots, S$$

$$P(\mathbf{y} | M) = \sum_{i, q_i \in F} \alpha(i, T) \quad (19)$$

一方、 $P(\mathbf{y} | M)$ を厳密に求めないで、モデル M が観測シンボル系列 \mathbf{y} を出力するときの最も可能性の高い状態系列 $\mathbf{q} = q_{i_0} q_{i_1} \dots q_{i_T}$ 上での出力確率を使うことも考えられる。図4の例では、 $q_1 \rightarrow q_2 \rightarrow q_3 \rightarrow q_5$ (P_2) がこれに相当する。この対数尤度 L は、次の動的計画法を適用することによって求められる。

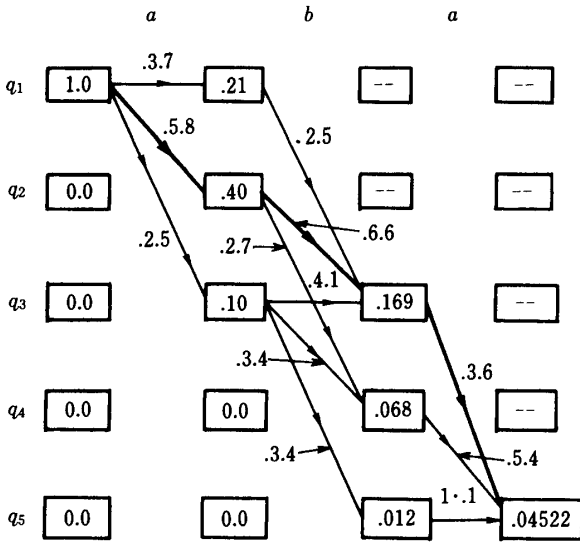


図5 図4に対する $\alpha(i, t)$ の計算例
(--は $\alpha(5, 3)$ を求めるのに無関係なので省略した)

$$f'(i, 0) = \log \pi_i \text{ for } i = 1, 2, \dots, S$$

$$f'(i, t) = \max_j \{ f'(j, t-1) + \log a_{ij} \cdot b_{ij}(y_t) \}$$

for $t = 1, 2, \dots, T; i = 1, 2, \dots, S$

(20)

$$L = \max_{i: q_i \in F} f'(i, T) \tag{21}$$

これは Viterbi アルゴリズムと呼ばれている。この方法は、厳密に確率を求める方法と比較して、計算量が少ないにもかかわらず認識精度は同等であることが実験的に確かめられている⁽¹⁸⁾。

図5に図4に対する $\alpha(i, t)$ の計算例を示す。図中太線は最適状態遷移系列を示している。

3.2 Baum-Welch のパラメータ推定法 (Forward-Backward アルゴリズム)⁽¹⁸⁾⁻⁽²¹⁾

新たに次のパラメータを導入する。Backward 変数 $\beta(i, t)$ を、時刻 t に状態 q_i にいて以後観測シンボル y_{t+1}^T を出力する確率、 $\gamma(i, j, t)$ をモデル M が y_t^T を出力する場合において、状態 q_i からシンボル y_t を出力し、時刻 t で状態 q_j へ遷移する確率とする。

このとき、

$$\beta(i, T) = \begin{cases} 1.0 & q_i \in F \\ 0 & q_i \notin F \end{cases}$$

$$\beta(i, t) = \sum_j a_{ij} \cdot b_{ij}(y_t) \cdot \beta(j, t+1)$$

for $t = T, T-1, \dots, 1; i = 1, 2, \dots, S$

(22)

$$\gamma(i, j, t) = \frac{a(i, t-1) \cdot a_{ij} \cdot b_{ij}(y_t) \cdot \beta(j, t) / P(y|M)}{P(y|M)} \tag{23}$$

図6に図4に対する $\beta(i, t)$ の図4に対する計算例

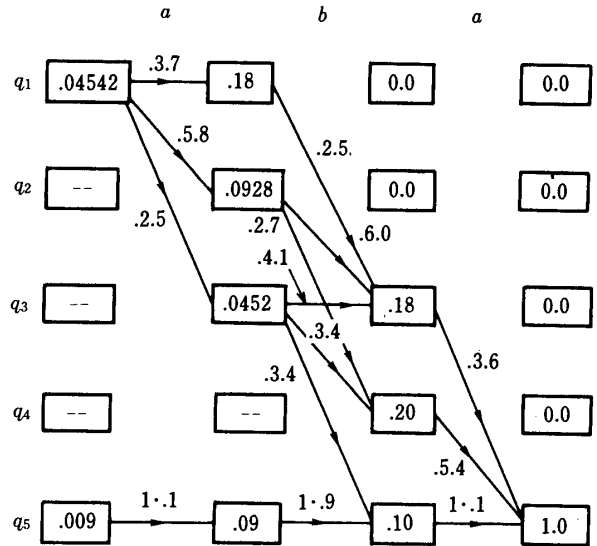


図6 図4に対する $\beta(i, t)$ の計算例
(--は $\beta(1, 1)$ を求めるのに無関係なので省略した)

を示す。これから、パラメータ $\pi_i, a_{ij}, b_{ij}(k)$ を次の再推定を繰り返すことによって求めることができる。

$$\hat{\pi}_i = \frac{\sum_j \gamma(i, j, 1)}{\sum_i \sum_j \gamma(i, j, 1)}$$

$$\hat{a}_{ij} = \frac{\sum_{t=1}^T \alpha(i, t-1) \cdot a_{ij} \cdot b_{ij}(y_t) \cdot \beta(j, t)}{\sum_i \alpha(i, t) \cdot \beta(i, t)}$$

$$= \frac{\sum_j \gamma(i, j, t)}{\sum_i \sum_j \gamma(i, j, t)} \tag{24}$$

$$\hat{b}_{ij}(k) = \frac{\sum_{t: y_t=k} \gamma(i, j, t)}{\sum_j \gamma(i, j, t)} \tag{25}$$

ただし、ここで注意を要することは、訓練サンプル数が n 個ある場合、これらの全サンプルについて上式を用いて1回パラメータを更新し、これを収束するまで更新を繰り返さなければならないことである。1サンプルごとによるパラメータ更新ではパラメータの収束が保証されない。

Baum-Welch の再推定プログラムが正常に動作しているかどうかのデバッグのためには次の項目のチェックが有効である。

- ① $\sum_j a_{ij} = 1, \sum_k b_{ij}(k) = 1$
- ② $\sum_{i \in F} \alpha(i, T) = \sum_i \beta(i, 0) \pi_i$ (26)
- ③ $\hat{P}(y|\hat{\lambda}, M) \geq P(y|\lambda, M)$ (27)

ここで、 $\hat{\lambda}$ は HMM の M のパラメータで、 y による λ からの再推定値である。また、上述の再推定アルゴリズムの $\alpha(i, t), \beta(i, t)$ は、計算途中で非常に小さな値となりアンダーフローの問題が生じるので、スケールリングを行う必要がでてくる⁽¹⁶⁾。

3.3 基本モデルの拡張⁽¹⁶⁾

(a) 観測シンボルが連続分布をとる場合 (連続分布モデル)

今までは、図4の例のように、観測シンボルが有限個のうちの 하나가出力される場合を述べてきた。しかし、特徴パラメータをベクトル量子化などの手法で有限個のシンボルに変換するのは量子化誤差などを伴い好ましくない場合が多い。特に不特定話者の認識の場合は、認識精度の低下を招く。

出力シンボルが、無限個の連続分布をなす場合、 $b_{ij}(k)$ を離散的確率分布から連続的分布に変更する必要がある。一般に $b_{ij}(k)$ を多次元正規分布 $N(\mu_{ij}, \Sigma_{ij})$ で近似される場合が多い。この場合、

$$\hat{\mu}_{ij} = \frac{\sum_t \gamma(i, j, t) \cdot y_t}{\sum_t \gamma(i, j, t)} \quad (28)$$

$$\hat{\Sigma}_{ij} = \frac{\sum_t \gamma(i, j, t) (y_t - \mu_{ij})(y_t - \mu_{ij})^T}{\sum_t \gamma(i, j, t)} \quad (29)$$

として再推定できる。ただし、実際には $b_{ij}(k)$ は一つの多次元正規分布で表現できないことも多いので混合分布 (mixed distribution) で近似することが試みられている⁽²²⁾。

(b) 継続時間分布モデルの導入⁽²⁴⁾⁽²⁵⁾

基本 HMM では状態 i に対して状態 i に τ 時刻 (τ フレーム入力) とどまる確率は、

$$d_i(\tau) = a_{ii}^{-\tau} (1 - a_{ii})$$

で与えられ、 τ と共に指数関数的に減少し、通常の音声現象と一致しなくなる。そこで、状態 i に τ 時刻とどまる確率 $d_i(\tau)$ を離散分布モデルとして推定する方法が提案されている⁽²⁵⁾。このとき、

$$\hat{a}_{ij} = \frac{\sum_{\tau \leq t} \sum_{i \leq j} \alpha(i, t - \tau) a_{ij} d_j(\tau) \prod_{k=1}^{\tau} b_{ij}(y_{t+1-k}) \beta(j, t)}{\sum_i \alpha(i, t) \beta(i, t)}$$

$$\hat{b}_{ij}(h) = \frac{\sum_{\tau \leq t} \sum_{i \leq j} \alpha(i, t - \tau) a_{ij} d_j(\tau) c_h(\tau)}{\sum_i \sum_{\tau \leq t} \alpha(i, t - \tau) a_{ij} d_j(\tau) \cdot \tau}$$

$$\times \frac{\prod_{k=1}^{\tau} b_{ij}(y_{t+1-k}) \beta(j, t)}{\prod_{k=1}^{\tau} b_{ij}(y_{t+1-k}) \beta(j, t)}$$

$$\hat{d}_i(\tau) = \frac{\sum_i \sum_{\tau} \alpha(i, t - \tau) a_{ij} d_j(\tau) \prod_{k=1}^{\tau} b_{ij}(y_{t+1-k}) \beta(j, t)}{\sum_i \alpha(i, t) \beta(i, t)}$$

で与えられる。ただし、 $c_h(\tau)$ は $\{y_{t+1-\tau}, \dots, y_t\}$ の τ 個の出力シンボルのうち h と等しい個数を示す。

3.4 連続単語音声認識への適用⁽¹⁶⁾⁽²⁶⁾

HMM による連続音声認識には、DP マッチングによる連続音声認識アルゴリズムがそのまま適用できる。ただし、注意を要することは、連続単語へ適用するためには、2 段 DP 法以外は厳密な確率計算法 (Baum-Welch のスコアリング) は用いることができず、Viterbi アルゴリズムによるスコアリングを使用する必要があるということである⁽³⁾⁽¹⁶⁾。

図7は、HMMによる連続単語認識のモデルを示している。このモデルでは単語間の接続に際しては何ら制限はない。もし、単語間の接続にオートマトン (正規文法) による制限を課す場合は、図8のようなモデルを考えればよい。ここで、状態 A, B, C は、合流 (confluent) 状態と呼ばれるもので、正規分布の非終端記号に対応している。図の例では、 $A \rightarrow V_1 B, C \rightarrow V_2 B$ の書換え規則がある場合を示している。

このとき、前向きパスアルゴリズムは次のように表現できる。

$$\alpha_V^v(1, t+1) = \max_{\{A \rightarrow vY\}} \{\alpha_V^v(1, t+1), \alpha(A, t)\}$$

$$\alpha_V^v(j, t+1) = \max_{\substack{1 \leq i \leq j \\ (Y_{t+1})}} \alpha_V^v(i, t) a_{ij}^v b_{ij}^v$$

$$\alpha(B, t+1) = \max_{\{X \rightarrow vB\}} \alpha_B^v(j^v, t+1)$$

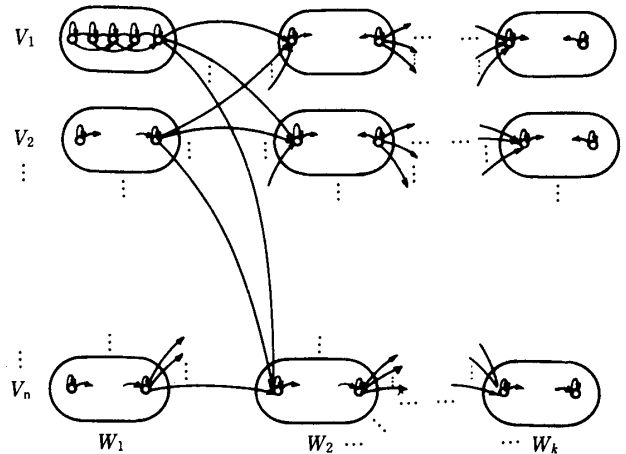


図7 HMMによる連続単語音声認識モデル

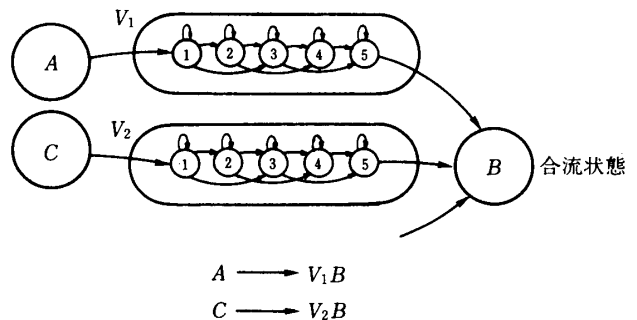


図8 HMMによるオートマトン制御連続単語音声認識モデル

ここで、 J^n は各単語の HMM の最終状態を示している。もちろん、このほかに、単語の接続時点の記憶、前状態の記憶、単語名の記憶が必要であり、これらの記憶方法は DP マッチング法における場合と同一である。

4. 今後の展開

4.1 文脈自由文法および係り受け規則の組み込み⁽¹⁶⁾

2・3節で述べたオートマトン制御 DP マッチング法を文脈自由文法(プッシュダウンオートマトン)制御に拡張することが試みられている。文頭から文末方向への Earley のアルゴリズムによる下降型構文解析法⁽²⁷⁾や CYK アルゴリズムによる上昇型構文解析法⁽²⁸⁾を DP マッチングによる連続音声認識に組み込む方法が提案されている。

また、音声認識特有の解析法として最も確からしいところから左右両方向へ構文解析を進める方法(島駆動方式)も DP マッチングに組み込まれ、文頭から文末に向かって(left-to-right)構文解析する方法との比較が行われている⁽²⁹⁾。

係り受け規則を満たす最適な文節列を選択する音声認識アルゴリズムも提案されている⁽³⁰⁾。しかし、二文節間の係り受け規則は二分木の構文解析木で表現でき、CYK アルゴリズムと本質的に同じである⁽³¹⁾。

また、二文節間の係り受け規則だけではあいまいな構文解析結果を多く生じるので、格構造などの導入が必要であろう⁽³²⁾。

4.2 言語の確率モデルの導入⁽¹⁶⁾

われわれは自然言語音声を理解する場合、話題についての膨大な知識を用いて、発話内容や発話単語を推定し、理解の焦点を絞っている。また、文法的に許される文はどれも同じ確率で生じるものではない。このような現象を確率モデルでとらえる方法が試みられている。一つは音韻・音節・単語連鎖などの三つ組の生起確率を用いる方法、他の一つは確率文法を用いる方法である。これらは、話題の推移確率などを含めて一つの確率モデルとして組み込まれていくであろう。

4.3 ニューラルネットワークの導入⁽¹⁶⁾

音声認識などのパターン認識は、記号処理とか知識処理とかでなく神経回路に組み込まれたメカニズムで行われている可能性が強い。その一つの例が、DP マッチングや HMM 法でありニューラルネットワークである。ニューラルネットワークも音声認識に適用されているが、時系列パターンの照合に関してはこれからの研究に待たなければならない。また、言語処理部においても意味主導型のシステムはパターン間の連想を中心とするものが多く、ここにもニューラルネットワークの適用がなされよう⁽³³⁾。

◇ 参 考 文 献 ◇

- (1) Atal, B. S. and Rabiner, L. R. : Speech Research Direction, *AT&T Tech. Journal*, Vol. 65, Issue5, pp. 75-88 (1986).
- (2) Rabiner, L. R., Juang, B. H., Levinson, S. E. and Sondhi, M. M. : Recognition of isolated digits using Hidden Markov Models with continuous mixture densities, *AT&T Tech. Journal*, Vol. 64, No. 6, pp. 1211-1234 (1985).
- (3) 坂井利之(編): 情報基礎学詳説, コロナ社(1983).
- (4) Vintsyuk, T. K. : Speech discrimination by dynamic programming, *Cybernetics*, Vol. 4, No. 1, pp. 52-57 (1968).
- (5) 迫江, 千葉: 動的計画法を利用した時間軸正規化に基づく連続音声認識, *音響誌*, Vol. 21, No. 9, pp. 438-490 (1971).
- (6) Sakoe, H. and Chiba, S. : Dynamic programming optimization for spoken word recognition, *IEEE Trans.*, Vol. ASSP-26, No. 1, pp. 43-49 (1978).
- (7) Vintsyuk, T. K. : Elementwise recognition of continuous speech composed of words from a specified dictionary, *Cybernetics*, No. 2, pp. 361-372 (1971).
- (8) Sakoe, H. : Two-level DP-matching—A dynamic programming based pattern matching algorithm for connected word recognition, *IEEE Trans.*, Vol. ASSP-27, No. 6, pp. 588-595 (1979).
- (9) 中川聖一: パターンマッチング法による連続単語および連続音節の音声認識アルゴリズム, *信学論誌*, Vol. 66D, No. 6, pp. 637-644 (1983).
- (10) Bridle, J. S., et al. : An algorithm for connected word recognition, *Proc. Int. Conf. ASSP*, pp. 899-902 (1982).
- (11) Ney, H. : The use of a one-stage dynamic programming algorithm for connected word recognition, *IEEE Trans.*, Vol. ASSP-32, No.3, pp. 263-271 (1984).
- (12) 迫江, 巨理: クロック同期伝搬型 DP 法による連続音声認識の検討, *音響学会音声研資*, S81-65 (1981.12).
- (13) Myers, C. S. and Rabiner, L. R. : A level building dynamic time warping algorithm for connected word recognition, *IEEE Trans.*, Vol. ASSP-29, No. 2, pp. 284-297 (1981).
- (14) 中川聖一: 拡張連続 DP 法による連続認識アルゴリズム, *信学論誌*, Vol. 67D, No. 10, pp. 1242-1249 (1984).
- (15) Sakoe, H. : A generalized two-level DP matching algorithm for continuous speech recognition, *信学論誌*, Vol. 65E, No. 11, pp. 649-656 (1982).
- (16) 中川聖一: 確率モデルによる音声認識, *電子情報通信学会*

- (1988).
- (17) Jelinek, J. : Continuous speech recognition by statistical methods, *IEEE Proceedings*, Vol. 64, No. 4, pp. 532-556 (1976).
- (18) Rabiner, L. R., Levinson, S. E. and Sondhi, M. M. : On the application of vector quantization and Hidden Markov Models to speaker-independent isolated word recognition, *BSTJ*, Vol. 62, No. 4, pp. 1075-1105 (1983).
- (19) Levinson, S. E. : Structural methods in automatic speech recognition, *IEEE Proceedings*, Vol. 73, No. 11, pp. 1625-1650 (1985).
- (20) Rabiner, L. R. and Juang, B. H. : An introduction to hidden Markov models, *IEEE ASSP Magazine* (1986. 1).
- (21) 大河内正明 : Hidden Markov Model に基づく音声認識, 音響学会誌, Vol. 42, No. 12, pp. 936-941 (1986).
- (22) Liporace, L. A. : Maximum-likelihood estimation for multivariate observations of Markov process, *IEEE Trans.*, Vol. IT-28, No. 5, pp. 729-734 (1982).
- (23) Juang, B. H. and Rabiner, L. R. : Mixture autoregressive hidden Markov models for speech signal, *IEEE Trans.*, Vol. ASSP-33, pp. 1404-1413 (1985).
- (24) Levinson, S. E. : Continuously variable duration hidden Markov models for automatic speech recognition, *Computer Speech and Language*, 1, pp. 29-45 (1986).
- (25) 橋本, 中川 : HMM 法による連続音声のセグメンテーションの検討, 音響学会講論集, 3-P-2 (1988.3).
- (26) Fallside, F. and Woods, W. A. (ed.) : *Computer Speech Processing*, Prentice Hall (1985).
- (27) 中川聖一 : 文脈自由文法のフレーム同期型構文解析法による連続音声認識, 信学論誌, Vol. 70D, No. 5 (1987).
- (28) Ney, H. : Dynamic programming speech recognition using a context-free grammar, *Proc. ICASSP*, pp. 69-72 (1987).
- (29) 大黒, 中川 : Left-to-right & top down 構文解析法と island-driven & bottom-up 構文解析法による連続音声認識の比較・検討, 信学会言語処理とコミュニケーション技報, MLC 87-12 (1987).
- (30) 尾関和彦 : 文節ラチスから最適係り受け構造を選択する多段決定アルゴリズム, 信学論 (D), J70-D, 12, pp. 2621-2629 (1987).
- (31) 中川聖一 : 連続音声認識のための係り受け解析法と文脈自由文法解析法との統一化, 音響学会講論集, 3-5-10 (1987. 10)
- (32) 中川, 伊藤 : 音節標準パターンと逆時間向き係り受け解析法を用いた日本語文音声の認識, 信学論 (D), J70-D, 12, pp. 2469-2478 (1987).
- (33) 西田豊明 : コネクショニストモデルによる自然言語処理一可能性を探る一, 情報処理学会, 自然言語処理シンポジウム論文集 (1988. 1).
- (34) Sankoff, D. and Kruskal, J. (ed.) : *Time Warps, String Edits and Macromolecules ; The Theory and Practice of Sequence Comparison*, Addison-Wesley (1983).

著者紹介



中川 聖一

昭和 51 年京都大学大学院博士課程修了。同年, 京都大学工学部情報工学科助手。昭和 55 年豊橋技術科学大学情報工学系講師, 昭和 58 年助教授。昭和 60~61 年カーネギーメロン大学客員研究員, 工学博士。音声情報処理, 自然言語処理の研究に従事。昭和 52 年度電子通信学会論文賞受賞。著書「情報基礎学詳説」(共著) コロナ社 (昭和 58 年), 「確率モデルによる音声認識」電子情報通信学会 (昭和 63 年)。