

6N-7

音声テンプレート統合法

菅原一秀

日本アイ・ビー・エム株式会社 サイエンス・インスティチュート

1.はじめに

テンプレートを用いる音声認識ではその良否が直接認識結果に反映される。テンプレートの最も簡単な作成法は、ある一つの発声をもってテンプレートとする方法である。しかしこの方法は一つの発声のみに頼るので偏りが避けられないし、発声の変動が大きくなったり、テンプレートを適応化させようとした時にはうまくいかない。これらの問題点を解決するためにテンプレートを複数個使う方式(マルチ・テンプレート法)があるが、与えられた学習用発声全部をそのままテンプレートとして持つのは記憶域及び計算量が大きくなってしまい、得策ではない。そこで、与えられたテンプレートから一つあるいは少数のテンプレートを作りだす方法が要求される。本稿ではその方法の一つとして、「重み付きDPマッチング」を利用した新しいテンプレート作成法(AWDP法:Average by Weighted DP)を示す。

2.従来技術の問題点とそれらの解決法

いくつかのテンプレートをDPマッチングにより統合する代表的な方法として次の二つがある。(1)どれか一つのテンプレートを選びこれを軸にして他のテンプレートを最適パスを介しこれに投影し、平均化などの操作により新しいテンプレートを得る。(PDP法:Projection by DP)(2)最適パス上で二つのテンプレートのそれぞれの重みに基づき内分して新しいテンプレートを得る。(新美[2] ADP法:Average by DP)

(1)のPDP法には「どのテンプレートを軸として選ぶか」、「特徴量は平均化されるが、時間長は軸として選ばれたものが使われる」等テンプレートの扱いが不均等であるという問題点がある。また、(2)のADP法にはテンプレートが3個以上のときそれらを組合せる順序により結果が異なってくるという問題がある。

これらを解決するためには、現在二つのテンプレートに対して行われているDPマッチングを任意個数N個のテンプレートに対して行えるよう拡張することが考えられる。しかし、この場合計算量はNに対し指数関数的に増大し少量のデータに対しても実行が困難になる。

ADP法が通常のDPで最適パスを求めた後、二つのテンプレートに対する重みで内分するのに対し、本稿で提案するAWDP法はDPマッチングの段階で重みを考慮して最適パスを求めるので、各テンプレートの重みをより良く反映できることが期待される。

3.テンプレート統合法

以下に説明する1~3の手順による。

1.重み付きDPマッチング

二つのテンプレートを $w:(1-w)$ ($0 \leq w \leq 1$) の重みでマッチングする

2.最適パスの追跡

1.のマッチングで得られた最適パスは正規化時間に対応しており、これをたどって新しいテンプレートを得る

3.複数テンプレートの統合

上記1, 2の方法を組合せて、N個 ($2 \leq N$) のテンプレートから一つのテンプレートを作る

3.1重み付きDPマッチング

統合すべきテンプレートを

$$A : a_1 \dots a_m$$

$$B : b_1 \dots b_n$$

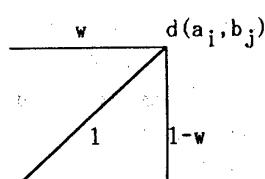
とする。DPマッチングを図1の様な荷重で行う。つまり、下の漸化式により累積距離gを計算する。

$$g(i,1) = d(a_i, b_1)$$

$$g(i,j) = \min \begin{cases} g(i-1,j) + w * d(a_i, b_j) \\ g(i, j-1) + (1-w) * d(a_i, b_j) \\ g(i-1, j-1) + d(a_i, b_j) \end{cases}$$

($1 \leq i \leq m, 1 \leq j \leq n$)

ここでdは局所距離



また、パスの長さを

やはり図1の様に定義する。

図1.荷重とパス長

補足： 前項の説明で $w=p/(p+q)$, $1-w=q/(p+q)$ (p, q は自然数) に選ぶと、上の定義は $p+q$ 次元での、次のような制限付のマッチングを縮退させたものと考えることができる。

1. $1 \sim p$ 次元のテンプレートは A の p 個のコピー
2. $p+1 \sim p+q$ 次元のテンプレートは B の q 個のコピー
3. パスの荷重はその両端の市街地距離
4. マッチングは $\{(x_1, \dots, x_{p+q}) \mid x_1 = \dots = x_p = X, x_{p+1} = \dots = x_{p+q} = Y\}$ で定まる空間で行われる。

このマッチングを 2 次元に縮退させる方法は明らか。パスの重みを正規化することにより図 1. の定義を得る。台形型 DP マッチング（大河内[1]）を用いればより性質の良い最適パス（正規化時間）が得られる。

3.2 最適パスの追跡

3.1 のマッチングで得た最適パスの長さは $w*m + (1-w)*n$ である。これを原点から出発してたどり、パスの長さの累積が整数 k になる点 p_k ($1 \leq k \leq w*m + (1-w)*n$) で新しい特徴量 c_k を求める。例えば、 p_k に一番近い格子点上で A と B の $(1-w):w$ の内分をとることが考えられる。

3.3 複数テンプレートの統合

3.1, 3.2 の方法を使い、二個のテンプレート A, B から、それらを $w:(1-w)$ の割合で反映した新しいテンプレート C ができる。N 個のテンプレートが与えられた時にはこの方法を組合せて、一つのテンプレートを作ることができる。そのとき、中間的にできるテンプレートに対してはそれをするために使用したテンプレート数に比例した重みを掛ける。

5. 実験

本稿で説明した AWDP 法と従来から提案されていた PDP 法及び ADP 法による単語音声認識の比較実験を行った。表 1. に 2 発声（発声 1, 2）、及び 3 発声（発声 1, 2, 3）からこれらの方針により作成したテンプレートによる単語認識実験の結果を示す。（2 発声のときは ADP 法と AWDP 法は同じ）話者 A, B とも、発声 1~5 は同時に収録し、それ以外は 1~2 週間をおいて収録した。認識対象とした発声はテンプレートの作成に使用されなかったもの全部である。

語彙 : 約 6,000 語から選んだ類似 150 単語
話者 : 成人男性 2 名
発声 : 話者 A 7 回
話者 B 6 回

基になるテンプレート数 (N) : 2 及び 3
重み (w) : $1/N$
DP パス : 対称型（図 1.）

表 1. 誤認識率 (%)

	方式	話者	
		A	B
2 発声	单一	1	2.8 1.5
		2	1.87 1.3
	PDP	1→2	1.87 0.5
		2→1	0.8 1.2
	AWDP		1.1 0.67
	ADP		0.8 0.67
3 発声	AWDP		0.77 0.72

5. 考察

統合テンプレートを使った場合は一つのテンプレートを使った場合に比べて、話者 A, B ともに誤認識がかなり減少した。このことは本稿で述べたテンプレートの統合方式の有効性を示している。なお、他の DP マッチング方式（「補間型」や「スタガード・アレイ方式」）においてパスに重みを付けたものについても実験を行い、上の実験と同様の結果を得ている。

PDP 法では、その非対称性により、投影の方向により結果が異なってくることが上の実験においても示されている。ADP 法及び AWDP 法によればこのような事は少なくなり、（ $N=2$ の時はまったくない）更に、上の実験では認識率でみても PDP 法によるものより概ね良い結果を得ている。この実験では ADP 法と AWDP 法はほぼ同等の結果を示していた。不特定話者のためのテンプレート作成時等、統合すべきテンプレートが多くなった場合に両方法間にどのような差が出てくるか今後の検討課題である。また、両方法ともテンプレートの数が 3 以上になると組合せる順序により得られるテンプレートが異なる点は解決されていないので、この点も検討を要する。

[参考文献]

- [1] 大河内正明, "台形型 DP マッチングに関する考察", 音声研究会資料 S81-64 (1981-12)
- [2] Niimi, Y., "A Method for Forming Universal Reference Patterns in an Isolated Word Recognition System," Proc. of the 4-th Int. Joint Conf. on Pattern Recognition, 1978, pp. 1022-1024
- [3] 新美康永, "音声認識", 共立出版, 東京, 1979
- [4] Niimi, Y. and Kobayashi, Y., "Synthesis of Speaker-Adaptive Word Templates By Concatenation of the Monosyllabic Sounds," Proc. of the ICASSP 86, Tokyo, 1986, pp. 2651-2654