

1B-4

手書き文字認識における
識別辞書の構成法に関する一検討

池田 正幸 田中 英彦 元岡 達
(東京大学 工学部)

1. はじめに

手書き文字認識の一般的手法として、<step1> 原文字パターンからの特徴抽出、<step2> 得られた特徴と認識用辞書との照合による類似度の算出及び候補の決定といった手順が定着しつつある。従来の研究の多くは step1に工夫をこらし step2に関しては入力パターンの特徴ベクトルと各カテゴリーの特徴平均ベクトルとのユークリッド距離のみで類似性を表現してきた。これに対して本研究ではstep2に着目し、距離算出に関して、自然な拡張を行なうことにより、識別率を改善することを試みたので、以下に報告する。

2. 類似度(距離)の算出

従来の特徴抽出は、対象となるすべての文字カテゴリーの種々の変形を吸収することに主眼をおいてきた。それに対し本方式は、各カテゴリーごとの変形の様相を分散ベクトルや超平面の形で反映させた辞書を作成し、それを考慮して類似度を算出しようとするものである。従来の方式を含めて類似度の算出方式について述べる。

2-1. 単純類似度 S_s - 従来の方式 -

各カテゴリーごとに特徴の平均ベクトル M_i (i はカテゴリーを示す) を求めて辞書とし、入力サンプル X と M_i とのユークリッド距離 D_s をもってカテゴリーと X との間の距離とする。すなわちカテゴリー i ($i = 1 \dots I$)、サンプル j ($j = 1 \dots J$) の特徴ベクトルを $F_i^j = (f_{1i}^j, f_{2i}^j, \dots, f_{ki}^j)^t$ (K は次元数) とすると J をサンプル数として

$$m_{ki} = \frac{1}{J} \sum_{j=1}^J f_{ki}^j \quad \dots\dots\dots ①$$

$$D_s^2(X, i) = \sum_{k=1}^K (x_k - m_{ki})^2 \quad \dots\dots\dots ②$$

2-2. 重み付け類似度 S_v

各カテゴリー毎に特徴ベクトルの分散ベクトル $V_i = (v_{1i}, v_{2i}, \dots, v_{ki})^t$ を求め、 M_i とともに辞書に登録する。重み付け距離は③式を用いて定義する。

$$v_{ki} = \frac{1}{J} \sum_{j=1}^J (f_{ki}^j - m_{ki})^2$$

$$D_v^2(X, i) = \sum_{k=1}^K \frac{(x_k - m_{ki})^2}{v_{ki}} \quad \dots\dots\dots ③$$

2-3. 投影類似度 S_p

2-1 での距離算出が、カテゴリーを平均ベクトルで代表させて 2ベクトル間の距離を基本においていたのに対して、この方式はカテゴリーを超平面として近似表現し、未知サンプルからその超平面への距離をもってカテゴリーへの距離と定義するものである (fig.1)。そのためには各カテゴリー毎に特徴点を散布して主成分を取り出し、それを平均ベクトルとともに辞書として登録すればよい。すなわち

$$C_i = \frac{1}{J-1} \sum_{j=1}^J (F_i^j - M_i)(F_i^j - M_i)^t \quad \dots\dots\dots ④$$

によって、分散共分散行列を求め

$$C_i \phi_l = \lambda_l \phi_l \quad \dots\dots\dots ⑤$$

なる固有値問題を解く。 λ_l の大きい方から L 個に対応する ϕ_l ($l = 1 \dots L$) と M_i を辞書とする。平均ベクトル M_i を原点とし、 ϕ_l を各軸とする超平面を構成し⑥式によって距離を算出する。

$$D_p^2(X, i) = \|X - M_i\|^2 - \sum_{l=1}^L \{ \phi_l^t (X - M_i) \}^2 \quad \dots\dots\dots ⑥$$

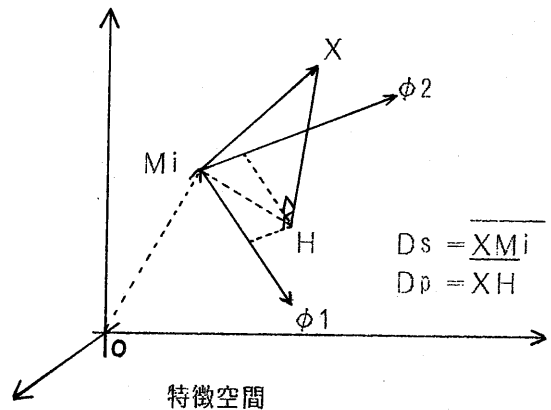


fig.1 投影類似度

3. 特徴抽出

手書き文字認識を行なうためには、対象となるすべての文字カテゴリーの種々の変形を吸収できるような特徴を選択する必要がある。今回はその一例として、単純類似度によってもかなりの認識率が報告されているSDF特徴・LSD特徴 [1] を用いた。(fig.2)

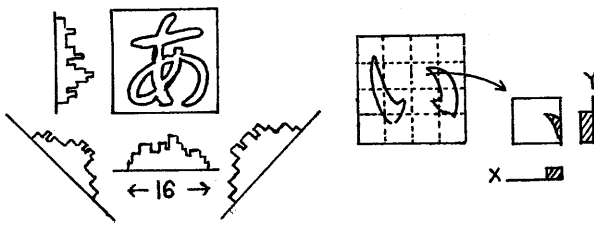


fig2-1. SDF特徴 fig2-2. LSD特徴

4. 認識実験

4-1. 使用データ

電総研手書教育漢字データベースETL-8よりひらがな71文字（濁点、半濁点のある文字を含み、小文字を除く）を用いて行なった。学習サンプルは1カテゴリー当り 100サンプルを用いた。なおSDF特徴は16×4次元、LSD特徴は16×2次元で行なった。

4-2. 結果

(1) 累積認識率（類似度算出手法による相違）

単純類似度、重み付け類似度、投影類似度（ $L=4$ ）について累積認識率をfig.3に示す。

(2) 正読率（超平面近似度 L による変化）

カテゴリーを表現する超平面を構成する主成分の個数 L を変えて正読率（第一候補が正解となる確率）を調べた。（fig.4） $L=0$ のものが単純類似度に対応する。

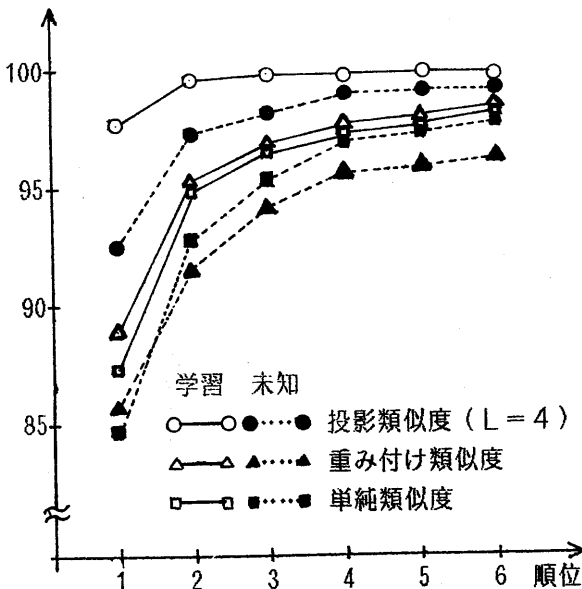


fig.3 累積認識率

5. 検討

fig.3 からわかるように、第一候補文字に着目すると S_s （単純類似度）、 S_v （重み付け類似度）、 S_p （投影類似度）の順に改善されており、 S_p 法の有効性が確かめられた。第二候補以下の未知サンプルにおいて S_v 法が S_s 法よりも劣っているのは S_v 法での変形に対するシャープさが原因と考えられ、 S_v, S_s 両者は優劣がつけがたい。

fig.4 より学習サンプルの正読率は L を大にするにつれて飽和するのに対して未知サンプルでは $L=3$ 前後をピークとして以降は低下する。これは L が大きくなると学習サンプルをより精度良く近似しようとするためにサンプル自身の偏りが辞書に反映されてしまうためであると考えられる。未知サンプルの正読率を高めるのが目的であれば $L=2\sim4$ 程度が適当であろう。

6. むすび

従来の類似度算出方式の拡張として、投影類似度を提案し、SDF・LSD特徴に関してその有効性を示した。今後は、他の特徴に対する投影類似度の適合性についても検討しつつ、より認識が困難とされている手書漢字等も、扱ってゆく予定である。

最後に本認識実験に使用した電総研手書教育漢字データベースETL-8の利用を許可して下さったことに対し、パターン情報部図形処理研究室各位に感謝の意を表す。

<文 献>

[1]. 萩田 他, “大局的・局所的線密度の併用による手書き漢字の大分類”, 信学技報PRL80-23

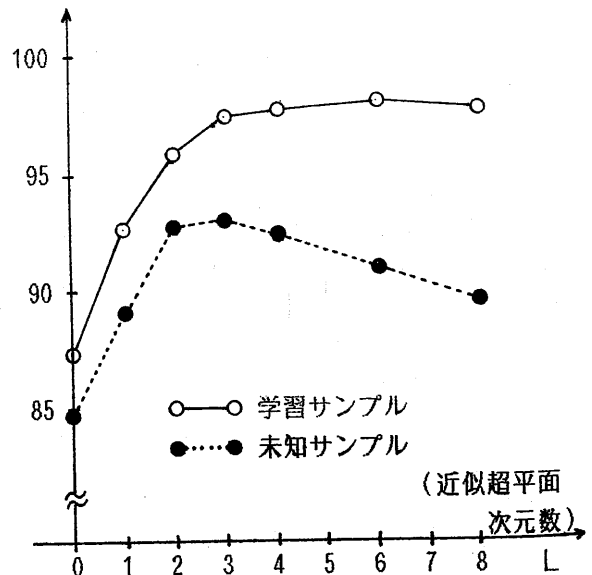


fig.4 投影類似度法の正読率