

PC上で動く印刷英数字OCR

5P-3

中村洋治, 除村健俊, 豊川和治, 北山友
 (日本アイ・ビー・エム株式会社 大和研究所 画像システム開発)
 高橋弘晏
 (日本アイ・ビー・エム株式会社 サイエンス・インスティテュート)

1. まえがき

本論文では、特定化されていないフォントが混合されている文書を読むことができ、かつ処理が簡単な文字認識アルゴリズムについて述べる。このアルゴリズムは、読み込まれた文字から特徴を抽出し、階層的に2段階の大分類を行なった後、識別を行なう。さらに、類似文字の1対1の比較によるチェックを行なって認識率を高めている。

2. 認識アルゴリズム

文字の認識は図1に示すような手順で行なわれる。

- (1) イメージの入力: 文字を含むイメージ情報が入力され、記憶装置上にたくわえられる。
 (2) 文字の切り出し: 一文字ごとにそのイメージが切り出され、 32×32 画素の領域に格納される。
 (3) 特徴抽出1: 文字の各画素ごとに 2×2 のマスクをかけ、その中の白黒のパターンを調べることにより、局所輪郭特徴を抽出する。特徴抽出の様子を図2に示す。抽出された特徴量は図3に示すような領域ごとに集計され、4(領域の分割方向) \times 6(分割された区画の数) \times 4(局所輪郭の方向) = 96次元の特徴量とされる。以後、このように分割された領域をルック領域、分割の方向をルック方向と呼ぶ。この特徴は識別のために使用されるが、大分類1のためにはルック領域に関係なく4つの局所輪郭方向別に合計し、4次元の特徴量を求める。
 (4) 大分類1: (3)で求めた入力文字の4次元特徴量について、予め用意してある全てのカテゴリーに対応する標準パターンとの距離を次式で計算し、カテゴリーごとに定められた閾値を越えるカテゴリーを候補からはずす。

$$4\text{次元距離} = \sum_i |f_i - g_i|$$

ただし、 f_i : 標準パターンの4次元特徴量; g_i : 入力パターンの4次元特徴量; i : 局所輪郭の方向

- (5) 特徴抽出2: 図4のように、 32×32 中 30×30 の領域で4辺から中心に向かって進み、イメージ領域の半分に達するか、文字にぶつかるまでの距離を計る。4辺をそれぞれ6区画に分割し、その区画内で距離の合計を求めることにより、 $4 \times 6 = 24$ 次元の特徴量を抽出できる。

- (6) 大分類2: (5)で求めた入力文字の24次元の特徴量について、大分類1を通過したカテゴリーに対応する標準パターンとの距離を次の式で計算し、カテゴリーごとに定められた閾値を越えるものを候補からはずす。ここで、文字の大きさの変化による影響を打ち消すために、標準パターン、入力パターンの特徴量の平均 m_f , m_g を用いている。

$$24\text{次元距離} = \sum_{i=1}^4 \sum_{j=1}^6 | (f_{ij} - m_f) - (g_{ij} - m_g) | = \sum_{i=1}^4 \sum_{j=1}^6 | f_{ij} - g_{ij} + (m_g - m_f) |$$

ただし、 f_{ij} : 標準パターンの24次元特徴量 i : 辺の番号
 g_{ij} : 入力パターンの24次元特徴量 j : 区画の番号

$$m_f = \left(\sum_{i=1}^4 \sum_{j=1}^6 f_{ij} \right) / 24 \quad m_g = \left(\sum_{i=1}^4 \sum_{j=1}^6 g_{ij} \right) / 24 \quad (m_f - m_g \text{は1度だけ計算する。})$$

- (7) 識別: 次式により、大分類1と2を通過したカテゴリーの標準パターンと入力パターンとの距離を求め、各候補を距離の近い順に並べる。 K は24次元距離と96次元距離との間の大きさの違いを補正するための定数である。

$$\text{距離} = 24\text{次元距離} + K \times 96\text{次元距離}$$

96次元距離は次の式で求められる。

$$96\text{次元距離} = \sum_{i=1}^4 \sum_{j=1}^6 \sum_{k=1}^4 w_{ijk} | f_{ijk} - g_{ijk} |$$

ただし、 f_{ijk} : 標準パターンの96次元特徴量; i : 局所輪郭の方向; k : ルック領域
 g_{ijk} : 入力パターンの96次元特徴量; j : ルック方向; w_{ijk} : 重み定数

ここで、 w_{ijk} の値は次の式で決定される。

$$w_{ijk} = w_i \times w_j$$

ただし、 w_i : 局所輪郭の方向*i*による重み w_j : *i*とルック方向*j*との間の角度による重み
 w_i は局所輪郭のマスキング中の長さも考慮したものである。その値は、例えば次のようになる。

$$w_i = 3 \quad (i=1, 3 \text{ のとき}), \text{ 又は } 2 \quad (i=2, 4 \text{ のとき})$$

w_j は文字のずれによる影響を考慮したものである。一般に、ルック領域と輪郭線が垂直な場合に比べ、平行な場合の方が文字のずれが特徴量に及ぼす影響が大きい。そのため、 w_j の値は、局所輪郭方向とルック方向のなす角度が直角に近いほど大きく、平行に近いほど小さくしている。 w_j の値の例を次に示す。

Software OCR for Typed English Characters

Y. Nakamura, T. Yokemura, K. Toyokawa, Y. Kitayama, H. Takahashi

IBM Japan

$w_i = 1$ (角度: 0° のとき)、 3 (角度: 45° のとき)、又は 4 (角度: 90° のとき)

(8) 類似文字識別 : 一般に、類似している文字は全体のパターンの相違より局所的なパターンの相違により区別しやすいので、局所特徴を用いて認識率を向上させる。そのため、(7)で並べられた候補の中から次の2つの条件を共に満たす C_i について、局所パターンを用いた詳細チェックを行う。ただし、候補を距離の近い順に C_1, C_2, \dots とする。

- ◇ 入力パターンと C_i との距離 \leq 入力パターンと C_1 との距離 $\times M$
- ◇ $i \leq N$

ただし、 M, N は経験的に求めた定数であり、例えば英数字に関しては、 1.25 と 5 をそれぞれ使っている。こうして残った候補をあらたに C_1, C_2, \dots, C_n とすると、類似文字識別の手順は次のようになる。

【1】2候補 C_n, C_{n-1} について、次の(i), (ii)を行う。

- (i) 予め定義された類似文字識別用特徴テーブルを見て、24次元と96次元の特徴の中からこの特定の2候補を区別するのに最も適した局所特徴を数個選ぶ。ここでは2個にしている。このテーブルの作成方法については後述する。
- (ii) 入力パターンと2候補との距離を次の式によって夫々求め、近い方を勝者とする。ただし、(ii)で選ばれた2つの特徴を i, j とする。

$$\text{距離} = K_i |f_i - g_i| + K_j |f_j - g_j|$$

ただし、 f_i, f_j : 標準パターンの特徴量 g_i, g_j : 入力パターンの特徴量

K_i, K_j : 24次元距離と96次元距離との間の大きさの違いを補正するための重み定数

($1 \leq i \leq 24$ ならば f_i は24次元の第 i 番目の特徴、 $25 \leq i \leq 120$ ならば96次元の第 $(i-24)$ 番目の特徴を表わす。)

【2】 C_n, C_{n-1} 間の勝者と C_{n-2} について、(i), (ii)を行う。

【3】こうして勝ち抜き制で比較をしていき、最後に C_1 とその前の勝者とを比較して距離の近かった方を最終結果とする。この方法では、すでに抽出された特徴のうちからいくつかを選んで使っているため、新たに特徴抽出を行う手間が不要であり、また一時に2つずつの候補を取り出してきて比較することにより、簡単できめ細かいチェックが可能である。

3. 類似文字識別のための局所特徴テーブルの作成

類似文字識別において2つの候補を比較するために使われる特徴は、次の式の値を最大及び次大にするように120次元の特徴の中から予め選ばれ、局所特徴テーブルに納められたものである。

$$(M_{ix} - M_{jx})^2 / (V_{ix} + V_{jx}) \quad (1 \leq x \leq 120)$$

ただし、 M_{ix} : カテゴリー C_i の x 番目の特徴量の平均 V_{ix} : カテゴリー C_i の x 番目の特徴量の分散

M_{jx} : カテゴリー C_j の x 番目の特徴量の平均 V_{jx} : カテゴリー C_j の x 番目の特徴量の分散

ここでいう平均、分散とは標準パターン作成時に一つのカテゴリーに対し使用する複数のトレーニング・データから求められるものである。

4. むすび

本論文で述べたアルゴリズムにより、簡単な手法でマルチフォント印刷文字の分類、識別、類似文字のチェックを行なうことができる。実験では、IBM-PC/AT (80286, 6MHz) を使い、2回コピーされた文書に対して、99%以上の認識率と約7.5文字/秒の認識速度が得られた。

5. 参考文献

- (1) 高橋 : "細線連結素方向による簡易手書漢字認識", 信学技報PRL82-8 (1982)

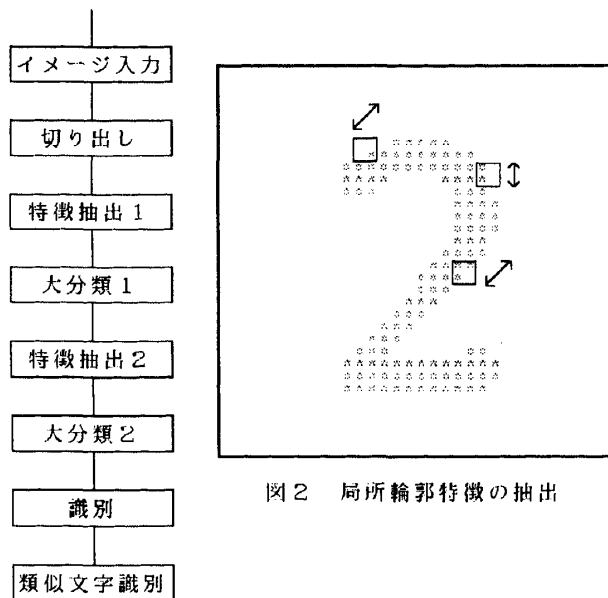


図2 局所輪郭特徴の抽出

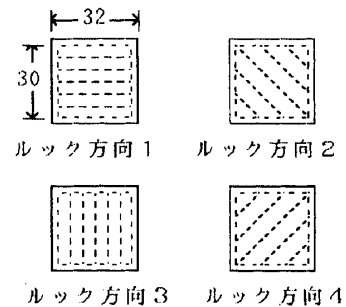


図3 特徴抽出1における領域の分け方

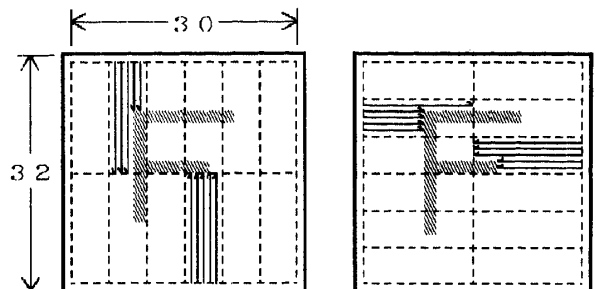


図4 外形特徴の抽出

図1 全体の処理の流れ