

特徴選択による文字認識高速化の一手法

5P-8

大田 裕、西村 康、富本 哲雄

(松下電器株式会社 無線研究所)

1、まえがき

漢字OCRにおいては、多次元の特徴を使用する特徴マッチングが広く使用されている。しかし特徴数を増していくつれて、マッチングの計算量が増し、認識時間が長くなる。

そこで我々は、文字ごとに識別に有効な特徴を選択し、この特徴を多段階に構成した識別に使用して認識の高速化を図る手法について検討を進めた。

ここでは、本手法の内容と特徴選択の有効性を確認する予備実験および多段階構成の認識についての実験結果について報告する。

2、分離率の定義

特徴の識別力を数値化するため、ある特徴 k によってカテゴリ i [平均 m_i , 標準偏差 σ_i] よりカテゴリ j [平均 m_j , 標準偏差 σ_j] が分離される確率である分離率 $B(k)_{ij}$ を図1の斜線部分の面積になるように定義する。

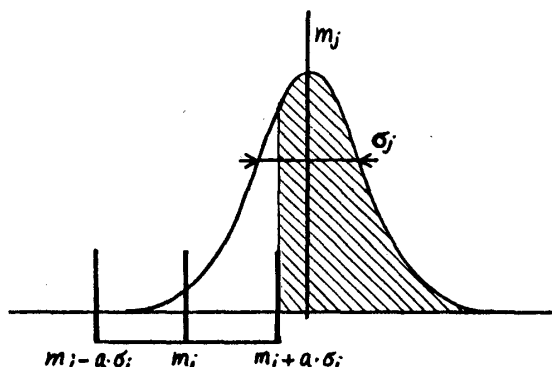


図1

3、逐次特徴選択の方法

各特徴の分離率の総和 U を式1で算出し、 U が最大の特徴を逐次的に選択する。[カテゴリ数を M とする]

$$U(k)_i = \sum_{j=1}^M B(k)_{ij} \quad \dots \text{式1}$$

ただし第2番目の特徴選択からは、それまでに選択した N 個の特徴 [$B(m)_{ij}, m=1..N$] と識別力が重複するのを避けるために分離率 B の代わりに式2で算出した累積分離率 B' を使用して総和 U を求める。

$$B'(k)_{ij} = 1 - (1 - B(k)_{ij}) \cdot \prod_{m=1}^N (1 - B(m)_{ij}) \quad \dots \text{式2}$$

4、特徴選択の有効性を確認する予備実験

前述の方法で各カテゴリごとに特徴を選択して認識する。選択する特徴数の上限を順次変えて、認識率・マッチング時間の変化をみた。実験の方法は次のとおりである。

| | | |
|--------|--------|------|
| 文字数 | 316 | カテゴリ |
| 全特徴次元数 | 370 | 次元 |
| 使用データ | ETL8B2 | |
| 学習サンプル | 40 | セット |
| 未知サンプル | 40 | セット |

[実験結果] 図2に示す。

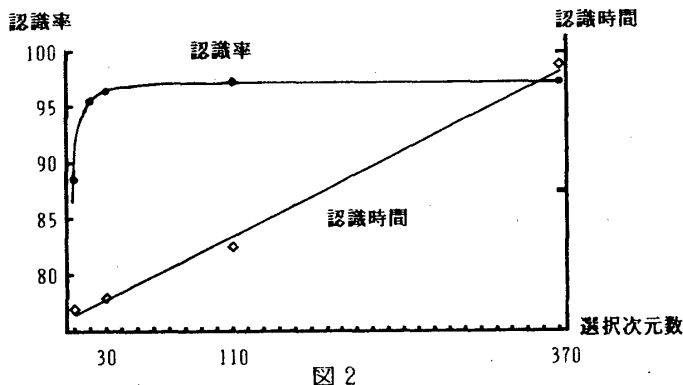


図2

High speed character recognition method using feature selection
 Hiroshi Ohta, Yasushi Nishimura, Tetsuo Tomimoto
 Matsushita Electric Industrial Co., Ltd.

- ・ 特徴30次元で、ほぼ370次元の特徴の認識性能を得ており前述の特徴選択の方法が有効であることがわかった。
- ・ 特徴選択が認識時間の短縮につながるということがわかった。

5、多段階識別の構成手法

直列に並べた3段階の識別で構成した。第1段と第2段は、前述の方法でカテゴリごとに370次元の特徴から選択した30次元で識別する。第3段は370次元の全特徴を使って識別する。各段階の処理は次のとおりで、従来の予め定められた数に候補カテゴリ数をしぼり込む大分類とは異なる構成である。

【 第1段 】

30次元の特徴〔平均 m_k , 標準偏差 σ_k , $k=1 \sim 30$ 〕について入力文字の特徴値 x_k が、

$$| (x_k - m_k) / \sigma_k | < c_1$$

をみたくカテゴリを候補とし、候補カテゴリ数が1ならばそのカテゴリを出力し、第1段で認識処理を終る。〔 c_1 は定数〕

【 第2段 】

30次元の特徴〔平均 m_k , 標準偏差 σ_k , $k=1 \sim 30$ 〕について入力文字の特徴値 x_k が、

$$\sum_{k=1}^{30} | (x_k - m_k) / \sigma_k | < c_2 \cdot 30$$

をみたくカテゴリを候補とし、候補カテゴリ数が1ならばそのカテゴリを出力し、第2段で認識処理を終る。〔 c_2 は定数〕

【 第3段 】

370次元の特徴〔平均 m_k , 標準偏差 σ_k , $k=1 \sim 370$ 〕で入力文字〔特徴値 x_k 〕と各カテゴリの辞書との距離 D が最小のカテゴリを出力する。

$$D = \sum_{k=1}^{370} | (x_k - m_k) / \sigma_k |$$

6、多段階識別による認識実験

前述の方法で第1段、第2段の定数 c_1 、 c_2 を順次変えて、認識率・マッチング時間の変化をみた。実験の方法は次のとおりである。

| | | |
|--------|------------|------|
| 文字数 | 952 | カテゴリ |
| 特徴次元数 | 第1、2段 30次元 | 次元 |
| | 第3段 370次元 | 次元 |
| 使用データ | ETL8B2 | |
| 学習サンプル | 80セット | |
| | (奇数データセット) | |
| 未知サンプル | 80セット | |
| | (偶数データセット) | |

〔実験結果〕 図3に示す。

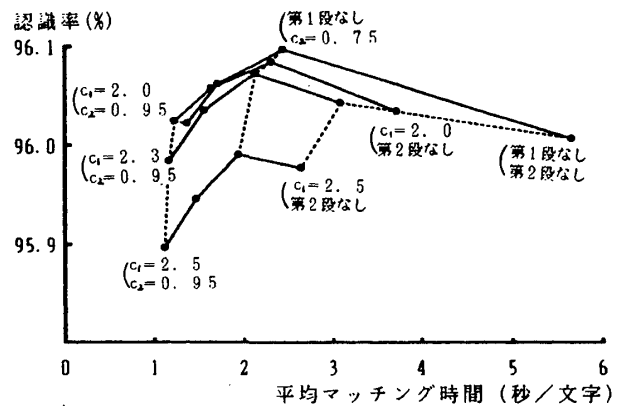


図3

- ・ $c = 2.0, c = 0.95$ で、認識率を下げることなく平均マッチング時間を約1/5に減少させることができたことがわかった。

6、まとめ

各カテゴリの平均とマッチング度の高い入力文字について計算量を削減することによって、平均マッチング時間を減少させた。そのための特徴選択として分離率による特徴選択を提案した。

今後は、多段階識別の構成方法についてさらに検討を進めていきたい。