

つぶれ文字を認識するための一方法

電子制御工学科 志久 修

ご存知のようにOCR（光学式文字読取り）とは、紙にかかっている文字を認識し、その文字に対応する電子テキストを出力する装置です。読み取る対象に応じ、活字OCR、手書き文字OCR、帳票OCRなどがあります（ちなみに本校の授業評価アンケートの読み取りに使っているのは帳票OCRです）。これらのなかで活字OCRは速度・認識性能とも非常にすばらしく、さらに低価格で十分実用レベルにあると思います。私が使っている活字OCRソフトはスキャナのおまけとしてついてきたものですが、それでも文字サイズが10ポイント以上で、文書のレイアウトが単純で、プリンタから直接出力された高品質な文書なら100%近い精度で電子テキストに変換してくれます。

しかし、活字OCRにもまだまだ弱点があります。例えば、複雑なレイアウトの文書や図表が混在した文書がうまく読めないし、コピーやファックスで品質が落ちた文字も正しく読むことができません。特に品質が落ちた文字が読めないのは致命的な弱点となっています。というのも、活字OCRには電子テキストが残っていない古い文書の読み取りに大きな期待が寄せられていますが、このような文書にはコピーやファックスされたものも多く含まれており、さらには紙自体の劣化により、文字の品質が悪くなっているからです。これらの品質が悪い文字を高精度に読み取ることができるようにすることが、活字OCRの重要な研究目標となっています。

今回は品質が悪い文字の例として、つぶれた文字を対象とします。つぶれ文字対策として、つぶれた文字を認識のための標準文字パターンとして登録する方法、文字のつぶれた部分は認識に使わず、きれいな部分だけを使う方法、つぶれの影響を受けにくい特徴を用いる方法、などがあります。ここで述べるのは、のアプローチにもとづくものです。

文字認識方法としてこれまでたくさんの方法が開発されていますが、今回は文字を線図形として表現する方法を採用します。この理由は、郵政研究所が行った文字認識方法の大規模な比較実験により、文字を線図形として表現する方法の有効性が明らかにされているためです。

文字を表す線図形の代表例を図1に示します。図(a)はもとの文字画像で、図(b)は文字線の中心線（方法A）、図(c)は輪郭線（方法B）です。中心線は文字線の太さに依存せず、また輪郭線は文字の外形が表現できる利点それぞれあり、文字認識によく用いられています。

しかし、図2(a)のようなつぶれ文字に対しては、中心線が正しく求められない(図2(b))、内部の輪郭線が求められない(図2(c))、などの問題があります。その対策として、図2(d)のようにつぶれていない部分は中心線を使い、つぶれた部分はその輪郭線を用いる方法（方法C）、図2(e)のように文字の輪郭線とつぶれ部分の輪郭線を用いる方法（方法D）が考えられます。方法Cは東北大から提案された方法で、方法Dは今回比較する方法です。



図1 中心線（方法A）と輪郭線（方法B）

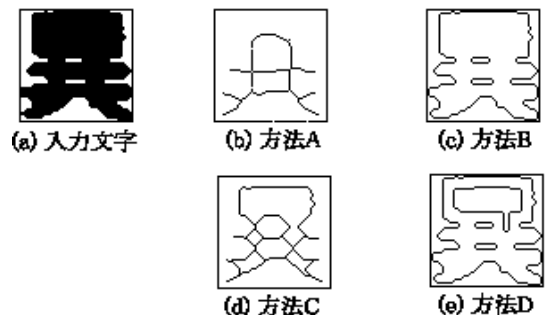


図2 つぶれ文字に対する4つの方法

以上の4つ方法を用いて実験を行いました。実験では図3のような品質を変えた文字画像を用意し、品質毎に文字認識率を調べていきました。詳細な実験条件は表1のとおりです。

図4に実験結果を示します。図4では品質を表す尺度として、文字画像を64×64画素で表したときの文字の平均線幅を使っています。線幅が太くなるほど文字のつぶれが生じ、逆に線幅が細くなるとかすれが生じていることを表しています。図4より線幅8程度で認識率がピークとなり、つぶれやかすれがひどくなるに従い認識性能は低下しています。つぶれ文字(線幅が太い場合)に注目すると、4つの方法とも認識率は低下していきませんが、方法Dが他の方法に比べ若干低下の割合が小さくなっています。さらにかすれ文字(線幅が細い場合)に対しては方法BとDは同じ認識率であることがわかります。以上の結果より、方法Dはつぶれ文字の認識に若干効果があり、さらにかすれ文字に対しては悪影響を与えないことがわかります。

ここで述べた簡単な改良で若干の性能向上が得られますが、さらに文字品質に対して頑健な認識方法の開発が必要だと思えます。

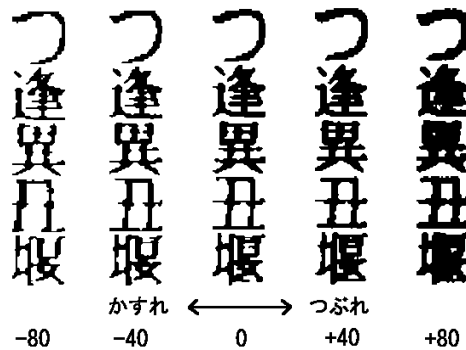


図3 文字データの例

表1 実験条件

文字データ	文字種 3,169 字種 (ひらがな・カタカナ・数字・アルファベット・漢字) 学習パターン数 30 パターン / 1 文字 (プリントアウトした文字) 評価パターン数: 6 パターン / 1 文字 (コピーした文字)
特徴量	方向線素特徴量 (19 6次元)
識別方法	部分空間法

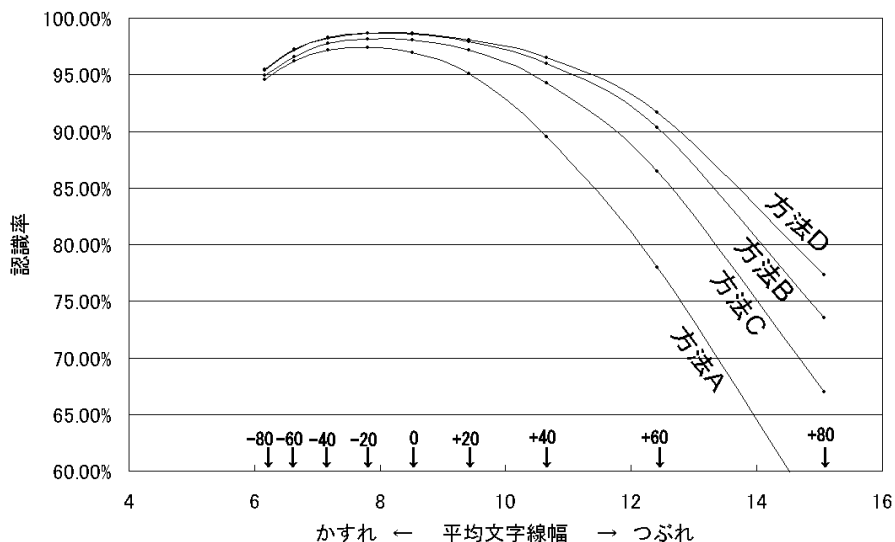


図4 実験結果

