

細線化処理を用いた手書き文書画像からの文字列領域分割法

1 D-6

足立雄介 吉川大弘 鶴岡信治
三重大学工学部

1. はじめに

OCRの開発が進む現在において、未だ手書き文書に対する文字認識では、満足する結果が得られていない。その要因として、個別文字の切り出しや、その前段階である文字列の切り出しが難しい事が挙げられる[1]。本稿では、文書の文字列が白領域で分割されている事に注目し、文書画像の白領域を細線化処理により線図形（芯線）で表し、その芯線から不要な部分を除去することにより、文字列を囲むループを抽出する方法を提案する。

2. 文字列領域分割法

文書画像からの文字列領域分割法は、細線化領域の限定、白領域細線化、不要な芯線の除去から成る。各処理に用いるパラメータは、文書画像の横方向の度数分布をフーリエ変換し、多く含まれる周波数成分の周期で最も小さいものを、文字列の高さ h_1 と考へ、 h_1 から求めたものを用いている。

2.1 細線化領域の限定

画像全体の白領域を直接細線化しようとする、処理時間が長くなり、細線化後の芯線も、文字列の領域とはかなり異なる分離線となる。そこで本稿では、明らかに細線化を行う必要のない領域の除去を、次の手順で行う。

① 2値画像（図1(a)）に対して、左右方向に h_1 画素の膨張収縮処理を行う（図1(b)）。

② 膨張収縮処理を行った画像を、大きさ $h_1/4 \times h_1/4$ 画素の正方領域で分割する。その正方領域の中に黒画素が存在しなければ、その領域は細線化を行わないことにする（図1(c)）。

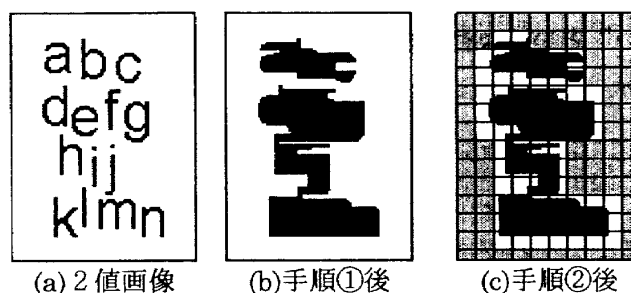


図1 細線化領域の限定

2.2 白領域の細線化

限定された白領域に対して細線化処理を行う[2]。ここで得られる芯線の画素の連なり t_1, \dots, t_n のうち、 t_1, t_n が共に特徴点（端点、分岐点、交差点）であり、他の画素 t_2, \dots, t_{n-1} が全て連結点であるものを“チェーン”と呼ぶ[3]。図2にチェーンの例を示す。

文字列の境界を表すループを得るためには、チェーンを単位とした取捨選択を考えればよいことになる。

2.3 チェインの評価

横書き文書の場合、文字列を分割するための線は、ほとんどが横線であると考えられる。そこで本稿では、不要なチェーンを発見するために、チェーンの傾き a 及び、高さ h_c を特徴量として用いる（図2）。

これらの特徴量により、次の①～⑤の各条件を満たすチェーンを段階的に除去していく。すなわち、まず①、②の条件により、傾き高さの大きなチェーンを除去し、次に③の条件で、高さの小さなチェーンを除去する。最後に④、⑤の条件により、傾きがそれほど大きくないチェーンを除去するのである。条件①→⑤の順で行うことで、より不要と思われるようなチェーンから順に除去していくことができる。

① $a \geq 4$, かつ $h_c \geq h_1/2$ ② $a \geq 2$, かつ $h_c \geq h_1/2$ ③ $a \geq 2$, かつ $h_c \geq h_1/8$

④ $a \geq 4/3$, かつ $h_c \geq h_1/4$

⑤ $a \geq 1$, かつ $h_c \geq h_1/4$

ただし、各条件によるチェーン除去終了後、芯線の中でループを構成していないチェーンは除去する。また、“1番外側に存在する線は枠なので消さない”という条件も加える。

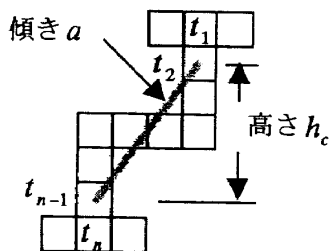


図2 チェインの例

2.4 過分割領域の結合

以上の処理では、どうしても過分割された領域が残ってしまう。そこで次の処理により、高さの小さい領域（高さが $h_1/2$ 未満の領域）を過分割領域とみなして、その領域を、本来含まれるはずの文字列の領域に結合させる。

過分割領域の左右を別々に見て、他の領域と接している部分の高さを求める。過分割領域は、その高さが最も大きな領域に含まれるものとする。図3では、 h_2 が最も大きいので領域②は領域③に含まれることになる。

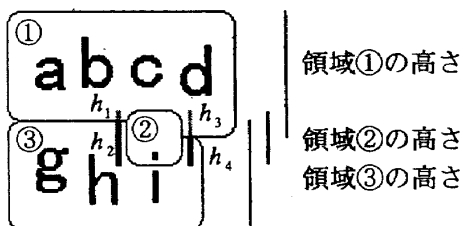
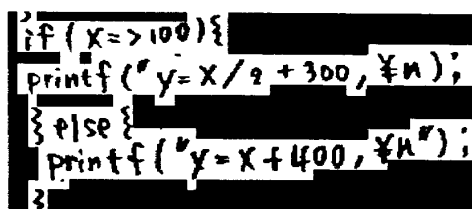


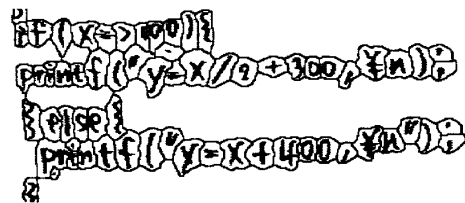
図3 過分割領域の結合

3. 実験結果

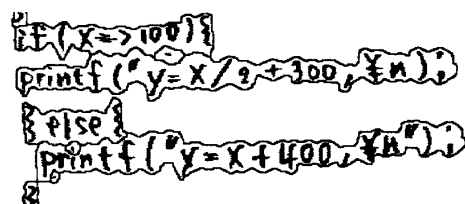
図4(a), (b), (c)及び(d)に、各処理後の画像の例を示す。白紙の紙に書かれた筆記者の異なる11画像に対し本手法を適用したところ、全215行の文字列のうち191行(89%)を適切に分割する事ができた。過分割は13行(6%)、過統合は15行(7%)に存在した。過統合は、主に文字列間の接触が原因であった。



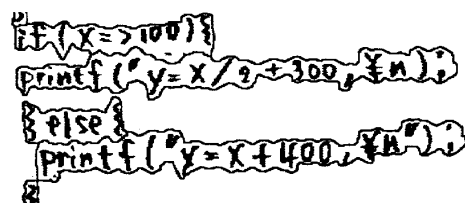
(a) 細線化領域の限定後の画像



(b) 白領域細線化後の画像



(c) チェイン除去後の画像



(d) 文字列領域分割結果

図4 各処理の結果例

4. おわりに

本稿では、領域を限定した細線化と、特徴量に傾きと高さを用いた、チェーンの段階的な除去により、手書きの文書領域から文字列領域を分割する手法を提案した。しかし、文字列間に接触がある場合、本手法ではそれらの文字列を分割することができないため、今後、それらに対する検討が必要である。

参考文献

[1] 上杉, 鶴岡, 三宅, 石田, “手書き文字認識のための文字列切り出し方法”, 電子情報通信学会パターン認識・メディア理解研究会, PRMU96-152, pp. 85-90, Jun. 1997.
 [2] 鶴岡, 木村, 吉村, 横井, 三宅, “デジタル図形の一細線化法と手書き文字認識システムへの応用”, 電子通信学会論文誌(D), Vol. J66-D, No. 5, pp. 525-532, 1983. 5.
 [3] 黄瀬浩一, 柳田修, “白領域の細線化を用いた文書画像の領域分割”, 電子情報通信学会論文誌(D-II), Vol. J80-D-II, No. 6, pp1608-1616, 1997. 6.