

5N-4

ランデータ画像処理を利用した
高速文字枠抽出方式

伊勢 広敏*、宗政 成大*、町田 哲夫*、道野 正雄**

* (株)日立製作所システム開発研究所

** (株)日立製作所小田原工場

1. まえがき

OCRはデータ入力手段として、事務処理の分野を中心に広く利用されている。OCRを使用する場合、帳票上の文字枠の位置、大きさなどを計測し、フォーマットデータを定義する必要がある。この定義作業を効率化する方法として、帳票画像から文字枠を自動抽出する方式を検討している。

本稿では、上記文字枠抽出の高速化を目的として、イメージリダにより入力した帳票画像をランデータとして扱う方式を開発した。帳票の場合、ラン数は、画素数に比べて、1/10程度であるため、処理時間の高速化が可能となる。ここでは、帳票画像から文字枠抽出する処理、および、この処理で利用しているランデータの画像処理方式について報告する。

2. 文字枠抽出の概要

ランデータは、画素塊の水平成分を保存しているため、文字枠等の矩形を中心とした図形処理に適している。本文字枠抽出方式では、まず、黒画素塊抽出により、イメージリダから入力した帳票画像をランデータに変換する。次に、図1に示すように、縦横変換、候補選別、文字枠統合、パターンマッチングの各処理を経て、ランデータから文字枠候補を抽出する。

(1) 黒画素塊抽出

帳票画像から黒ランデータを生成する。イメージリダで入力した帳票画像は、転送効率を向上するため、画像圧縮(MMR符号)されている。ここでは、MMR符号をラベリングした黒ランデータに変換し、同一ラベルを持つ黒ランデータを黒画素塊として処理する。

(2) 縦横変換

イメージリダのセンササイズに関する制約から、読み取った方向により、帳票画像の縦横変換を必要とする場合がある。この場合は、黒ランデータを90度回転することにより、帳票画像の縦横変換を実現する。

(3) 候補選別

ランデータから文字枠候補を選別する。具体的には、かすれによる画質劣化を補正しながら、黒画素塊および白画素塊を利用して文字枠らしい図形を抽出し、抽出した図形から文字枠候補を選別する。

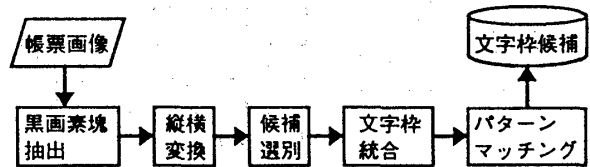


図1. 文字枠候補抽出

a) 結線処理

かすれ等により分離して抽出された黒画素塊から黒枠線を復元するため、一定距離内にある複数の黒ランを結線して、1つの黒ランに変換する。

b) 白画素塊抽出

黒枠線で囲まれた余白領域を取り出すため、白画素塊を抽出する。なお、白画素塊は、同一のラベルを持つ白ランデータである。

c) 文字枠選別

抽出した外接矩形から文字枠候補を選別する。選別の基準は、外接矩形の大きさ、形状等により、文字枠の確からしさをチェックする。

(4) 文字枠統合

類似する位置に複数の文字枠候補が、重畳して抽出された場合、最適なものを識別する。ここでは、原画像に近い状態で抽出された候補を最適としているので、識別基準は、結線幅が小さいこと、文字枠候補サイズが周辺の文字枠に近いことなどである。

(5) パターンマッチング

既抽出文字枠の位置、大きさなどを利用して未抽出の文字枠を補足する。ここでは、書式ルールを利用して、文字枠が存在する確率が高い位置を推定し、推定した位置に相応するサイズ、形状の文字枠を生成する。

3. ランデータをベースとした画像処理

2. で示した縦横変換、候補選別、パターンマッチングを高速に実現する画像処理方式を以下に示す。

3.1 ランデータの形式

本方式では、ランテーブルとyリストで、ランデータを構成する。図2にランデータの例を示す。x1は各黒ランの始点のx座標、x2は終点のx座標に1を加えた値、LRはランに付加されたラベルである。rは、黒ランの番号であり、ランテーブルを格納した領

域のアドレスにより判別できる。また、yリストは、各ラインの先頭にある黒ランの番号rの値を格納している。

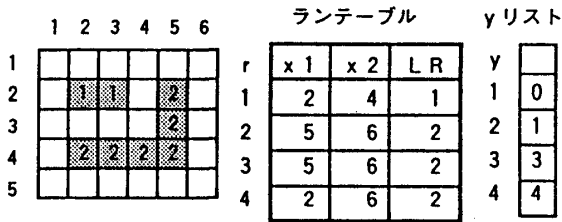


図2. ランデータ

3.2 論理演算

ランデータに対する論理演算は、縦横変換、パターンマッチング等に利用している。表1に、論理演算の一例として、ライン間の排他的論理和 (EOR) 演算に関する処理内容を示す。ここでは、上下に隣接するランの各状態 (10通り) に対する、ラン設定、ラン更新、次の状態を示している。ラン設定では、上ランR_iおよび下ランR_jの始点x₁、終点x₂の位置関係により、演算後に生成されるランの始点、終点の位置を表している。ラン更新における+1は、次のランへの更新を意味する。

表1 ライン間のEOR演算

状態	ランイメージ	ラン設定		ラン更新	次の状態
		始点	終点		
A	0	R _i .x ₁			B
	1				C
	2	R _j .x ₁			D
B	3		R _i .x ₂	i+1	A
	4		R _j .x ₁		C
C	5		R _i .x ₂	i+1	D
	6			i+1, j+1	A
	7		R _j .x ₂	j+1	B
D	8		R _j .x ₂	j+1	A
	9		R _i .x ₁		C

3.3 90度回転

ランデータの90度回転は、図3に示すように、水平ランを垂直ランに変換することにより実現できる。変換手順は、まず、下のラインとEOR演算することにより変化点を検出し、次に、変化点の状態から垂直ランの始点、終点を算出する。

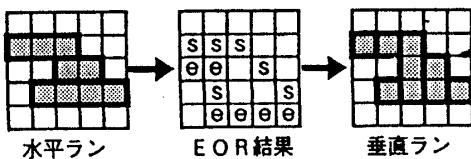


図3. 縦横変換

3.4 ラベリング処理

結線処理、縦横変換等により、ランの結合、分離が発生した場合、ランデータに付加されたラベルを付け直す。ラベリング処理では、ランテーブル以外に、ラベルの統合を示す連結テーブルを利用する。ここでは、処理対象のランに対して、以下の処理を実行する。

- (1) 参照ランのラベルを付加
- (2) 新しいラベルの生成
- (3) 連結テーブルの更新

3.5 外接矩形生成

断線したランを復元するため、一定距離内にあるランを結線しながら外接矩形を生成する。

(1) 水平結線

同一ラインにある黒ランを結線する場合には、一定距離内にある黒ランを統合する。したがって、結線後、黒ランデータ量は増加することはない。

(2) 垂直結線

黒ランを垂直方向に結線する場合、結線幅分の黒ランが新たに生成される。したがって、原画像に比較して、結線後には黒ランデータ量が増大し、処理速度が劣化する場合がある。これを防ぐために、ここでは、一定距離内にある黒ランデータのラベルだけを変更し、同一ラベルを持つランデータを利用して外接矩形を抽出している。

4. 処理時間評価

帳票入力 (MMR符号) してから文字枠候補を抽出するまでの処理時間を評価した。評価対象とした帳票は、A4サイズであり、200dpiの線密度で取り込んだものである。入力画像を構成する黒ラン数、ラベル数および黒画素数は、130127、22961、653219である。

帳票の処理時間を各処理ステップごとに算出した。黒画素から画素塊を識別し、文字枠候補を抽出する方式と比較して、結果を表2に示す。

表2 評価結果

	黒画素塊抽出	画素変換	候補識別	文字枠結合	マッチング	合計
ラン	2.3	2.9	52.1	0.6	32.3	90.2
画素	11.5	10.1	137.7	0.6	46.3	206.2

(単位:秒)

5. あとがき

帳票から文字枠候補を抽出することを目的として、ランデータをベースとした文字枠抽出方式を開発した。本方式により、黒画素から文字枠を抽出する方式に比較して、数倍程度、抽出速度を向上した。

参考文献

- 1) 伊勢他: 枠線認識を特徴とする帳票記入欄抽出方式の提案, 情報処理学会第43回全国大会, pp.2-385~386 (平成3後期)