

文書画像における矩形領域の抽出とその分類

伊藤昭治

(日本アイ・ビー・エム 東京サイエンティフィック・センター)

1. はじめに

半導体技術等の進歩により、コンピュータの主記憶装置および演算処理装置の性能は著しく向上している。その反面、コンピュータへの「情報」の入力は依然キイーン方式がほとんどであり、通常、逐次的に1文字づつ入力されるが、大量のデータ、あるいは、毎日発生するデータを人手をかけて入力することは、現実的に不可能な場合がある。また、これからのコンピュータがエンドユーザーの仕事場で使われる傾向が強くなり、ユーザーは最小限の労力で自動的に「情報」をコンピュータに入力できることが望ましい。

特に、日本語文書は英語文書に比較して、キイーン作業に労力を必要とする。また、一般の文書は、情報としてコード化可能な文字データばかりでなく、図および写真も含んでいる。これらの情報も簡単にコンピュータに入力できることが望ましい。

一方、日本においては、ファクシミリ装置が急速に普及し、コンピュータに接続することにより、文字データばかりでなく、図や写真を入力することが可能になってきた。

このような環境のもとでは、コンピュータに文書画像を自動入力する上で矩形領域の抽出とその分類が重要な技術となる。

上記の観点から、文書画像処理に矩形領域の概念を導入した。文書はいくつかの矩形領域より成り立ち、矩形領域間の距離を判断することにより、上位概念の特定の矩形領域を形成する。矩形領域を抽出した後、「写真」、「文章」、「文字」、「表」、「グラフ」および「線画」に分類する。

本報告では、文書画像のグラフ表現方法、矩形領域の抽出アルゴリズムおよび矩形領域の分類のための判別分析アルゴリズムとそれらの実験結果について述べる。

2. 文書構造の表現

処理対象をどのように表現するかは正しい解析を行うために重要である。また、解析の目的によって適切な表現方法をとらなければならない。2次元に広がった文書画像は、特定の機能を有する有限個の部分空間より構成され、部分と部分との関係をフローとして表現できる。すなわち、全体に対する部分と部分の結合構造や結合特性を解析するための文書構造の表現方法について述べる。

2.1 矩形領域の概念

文書画像を構成する情報の単位として、矩形領域の概念を導入する。また、その構造は、矩形領域とその関係によって組み立てられ、階層性を有している。矩形領域は走査の方向を一つの辺としているので、処理の単純化のうえで有利であるとともに、階層構造の観点から、上位概念の領域は、それが集合として含んでいる矩形領域を囲む最小の矩形領域として表現できる。矩形領域間の関係は、フローの概念を導入して、矩形領域間の位置の関数として定量化できる。従って、矩形領域は視覚的独立の条件を必要条件とする領域であって、情報の独立性および関連性については考慮しない。

以上の考え方に基いて、文書構造を視覚的に独立な矩形領域とその位置の関数によって表現する。

2.2 基本矩形領域の抽出

文書画像の2値画像行列を同一の大きさの単位局所領域 (ϵ - δ 領域) に分割し圧縮した画像行列を求める。この行列に対して図1のごとく十字のマスク論理演算子を用い、繰り返し処理して矩形化する。画像の欠け、つぶれなどの誤差に対して安定であり、その後の処理も高速化される。

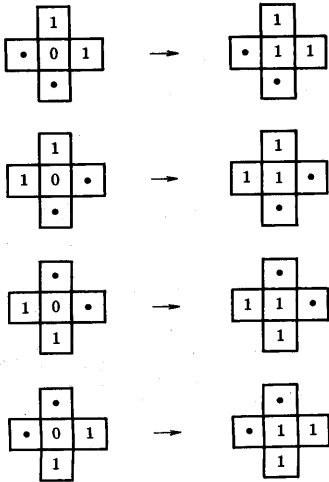


図1. 矩形化論理

2.3 矩形領域の関連

一般の文書(英文の場合)は下記のルールに従って構成されている。

- 1) ϵ - δ 矩形領域が、文書画像の構造を解析するための単位である。
- 2) 文書画像の構造は矩形領域間の関係で表現される。
- 3) 矩形領域間の関係は文書上の情報の流れと矩形領域の隣接によって表現される。
- 4) 情報の流れは上から下へ、左から右へ向う(英文の場合)。
- 5) 情報の流れの方向に隣接する矩形領域は情報も隣接する。

各矩形領域 v_i は文書画像中の位置で標識づけられる。

$$V = \{v_i; i = 1 \sim p\}$$

v_i の位置は矩形領域の左上の座標 (x_i, y_i) と右下の座標 (x'_i, y'_i) で表現される。文書画像は p 個の矩形領域として定義される。

v_i と v_j の y 方向の隣接関係は下記のごとく定義される。

$$\text{If}(a) \wedge \{ \text{If}(b) \vee \text{If}(c) \vee \text{If}(d) \} = 0$$

$$\left\{ \begin{array}{l} \text{If(条件が真)} = 1 \\ \text{If(条件が偽)} = 0 \\ a: \text{条件 } y_i \leq y_j \leq y'_i \leq y'_j \\ b: \text{条件 } x_i \leq x_j \leq x'_i \\ c: \text{条件 } x_i \leq x'_j \leq x'_i \\ d: \text{条件 } x_i \leq x \leq x' \leq x'_i \\ x_i \leq x_j \\ X = \text{Max}\{x_i, x_j\} \\ X' = \text{Min}\{x'_i, x'_j\} \end{array} \right.$$

かすすべての $k \{k; v_k \in V, k \neq i, k \neq j\}$ で満足する時、 v_i と v_j は y 方向に隣接する(図2)。

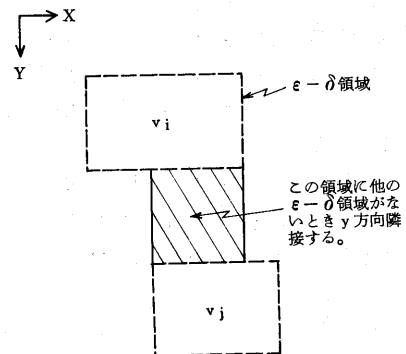


図2. Y方向隣接

X方向隣接についても同様に定義できる。

2.4 文書構造のグラフ表現

文書構造の基本単位となる ε - δ 領域とそれらの関係をグラフ理論を応用した隣接関係を文書構造を表現する方法を提案する。

- グラフ理論を用いる理由としては、
- 1) ε - δ 領域を点 (Point または Node) 隣接関係を線 (Edge または Line) と考えれば、文書構造をそのままグラフとして表現できる。
 - 2) グラフ表現は行列を主体に隣接関係を分析できるので、コンピュータの利用が容易である。

が上げられる。

ε - δ 領域を点、隣接を線とするグラフ G を定義する。

$$\text{文書} \cong G(V, E)$$

ただし

$$V = \{v_i; i = 1 \sim p\}$$

p 個の矩形領域を示す点の集合

$$E = \{e_k; k = 1 \sim q\}$$

q 本の隣接関係を示す線の集合

e_k は v_i, v_j の 2 対で表わされる

矩形領域の位置, 大きさ, 面積, 中心点などの点 v_i の特性 および 隣接距離, 中心点距離, 線の重みなどの線 e_k の特性を表 1 に示す。

3. 文書構造の解析

文書画像上の対象領域の抽出はグラフの操作の問題に置き換えられる。すなわち、サブグラフ G' の生成および連結グラフの抽出により対象矩形領域を抽出できる。

3.1 部分グラフの生成

グラフの点 v_i の次数を 4 方向の隣接に対して最大 1 とする条件を満足するようにグラフを正規化する。任意の点 v_i が同一隣接方向に対して複数の線 $\{e_i; 1 \leq i \leq n\}$ を持っている時、 $\{e_i\}$ より唯一選択される線 e_i は、次のうちのいずれかの条件を満足し、上から順に優先する。

- 1) X 方向隣接の場合

$$\text{Min}_{1 \leq i \leq n} BX_i$$

表 1. グラフの特性値

ε - δ 領域である点 v_i の特性		隣接を表わす線 ($e_k = v_i v_j$) の特性	
特性	式	特性	式
位置	左上 (x_i, y_i)	隣接距離 (Adjacent Distance)	x 方向隣接 $A_k = X' - X $
	右下 (x_i, y_i)		y 方向隣接 $k = Y' - Y $
大きさ	$LX_i = x'_i - y_i $		$X = \text{min}(x'_i, x_j), X' = \text{max}(x'_i, x_j)$
	$LY_i = y'_i - y_i $		$Y = \text{min}(y'_i, y_j), Y' = \text{max}(y'_i, y_j)$
面積	$SR_i = LX_i \times LY_i$	中心点距離	$B_k = \sqrt{BX_k^2 + BY_k^2}$
中心点	$CX_i = \frac{x_i + x'_i}{2}$		$BX_k = CX_i - CX_j $
	$CY_i = \frac{y_i + y'_i}{2}$	面積を考慮した線の重み	$BY_k = CY_i - CY_j $
		$SA_k = \frac{SR_i \times SR_j}{A_k^2}$	
		$SB_k = \frac{SR_i \times SR_j}{B_k^2}$	

4方向隣接の場合

$$\text{Min } B_i \quad 1 \leq i \leq n$$

- 2) $\text{Min } A_i \quad 1 \leq i \leq n$
- 3) $\text{Min } B_i \quad 1 \leq i \leq n$
- 4) $\text{Max } SA_i \quad 1 \leq i \leq n$
- 5) $\text{Max } SB_i \quad 1 \leq i \leq n$

隣接方向は、 x 方向、 y 方向あるいは x や y 両方向から選択できる。この処理により情報の流れをより良く表現することができる。

次に、部分グラフを生成することによって、目的とする矩形領域を抽出できる。

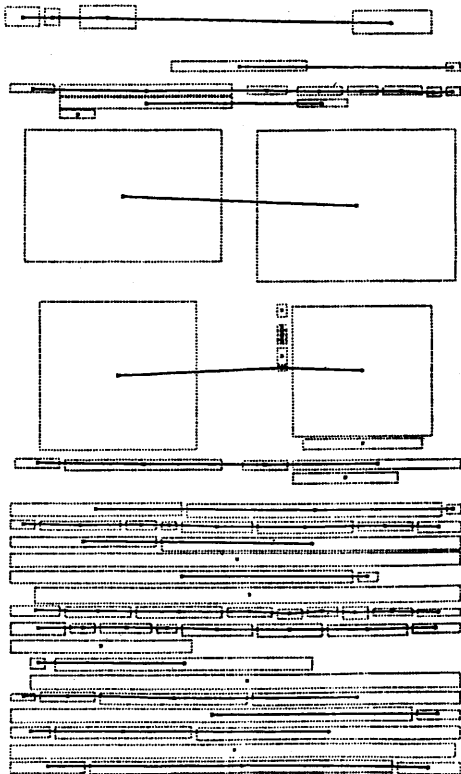


図3. x 方向の正規化グラフ

文書画像中の文章の行成分は x 方向のみの隣接によってその構造を示すことができる（縦書きの場合は y 方向隣接）。すなわち、文章のグラフから y 方向隣接を除去した部分グラフを求めると、情報の x 方向のみの流れを保持し、文章の構造を良く示すことができる。このように、特定の方向の隣接関係を除去することによって、縦あるいは横につながりのある矩形領域を抽出することができる。

3.2 矩形領域の抽出

一般に文書のグラフ G は連結グラフとなる。隣接している点を一つの集合としてグループ化し、新たな矩形領域を抽出することができる。得られる矩形領域は2点の座標で表現される。

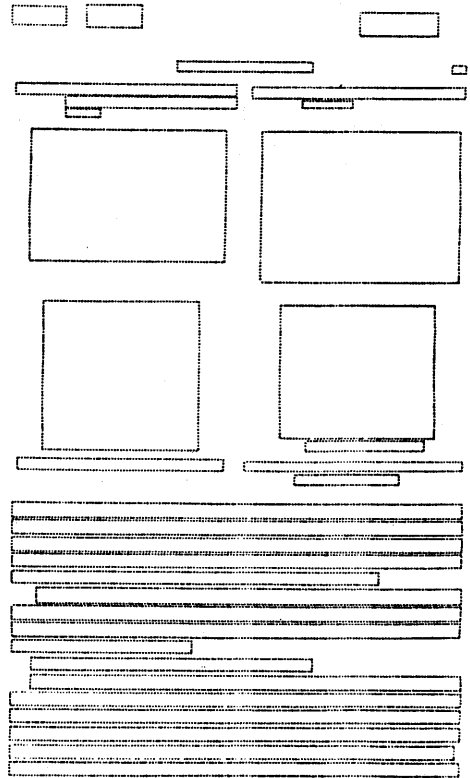


図4. グループ化した矩形領域

そこで、領域の抽出とは「グラフ G よりその部分グラフ G' を生成し、連結成分を抽出する」処理に他ならない。

連結グラフから任意の部分グラフを分離させるために除去すべき線の最小限の集合を Cut Set と呼ぶ。また、点 v が連線か非連結かによって、点 v と線長との連結係数を求め、これを n 行 m 列においた行列を連結行列と呼ぶ。

これらの手法を使用して解析した結果を図3～図6に示す。原文書は英語で書かれた論文で、手書きで記入された文書識記号、文章領域、文字領域、表領域、図領域を含み、メカニカルスキャナーによって走査して文書画像データを得た。文書画像のサイズは、縦 230mm 、横 160mm であり、走査線密度は 10本/mm で、サンプリングピッチ

は $100\ \mu\text{m}$ である。

図3は x 方向の隣接を正規化した後の隣接関係のサブグラフを示す。図4は図3のグラフから、さらに、 x 方向の一定以上の隣接距離にある線を除去し、求めたグラフの連結線をグループ化して、行毎の文章領域、独立した文字領域、表領域および図領域を抽出した結果を示す。図5は図4で求めた新たな矩形領域に対して y 方向の隣接関係をグラフ表現した後、一定以上の隣接距離にあるものを除去した結果を示す。図6は図5に基づいて連結成分をグループ化して抽出した矩形領域を示す。文章領域全体が一つの矩形領域として抽出されている。

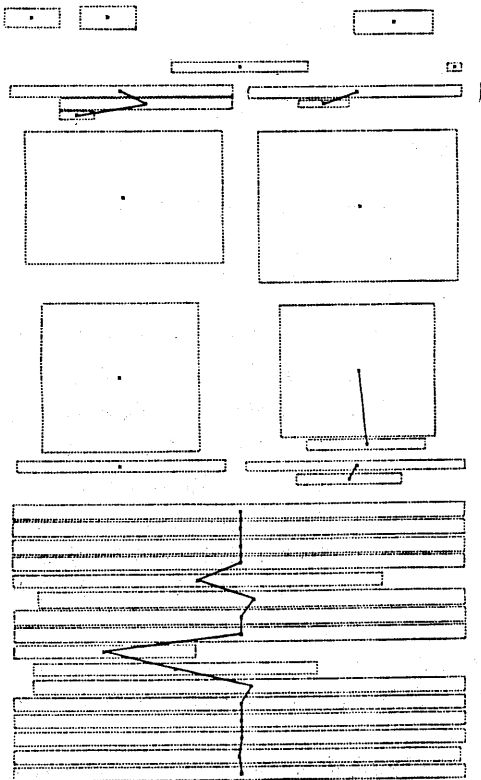


図5. y 方向隣接を除去したグラフ

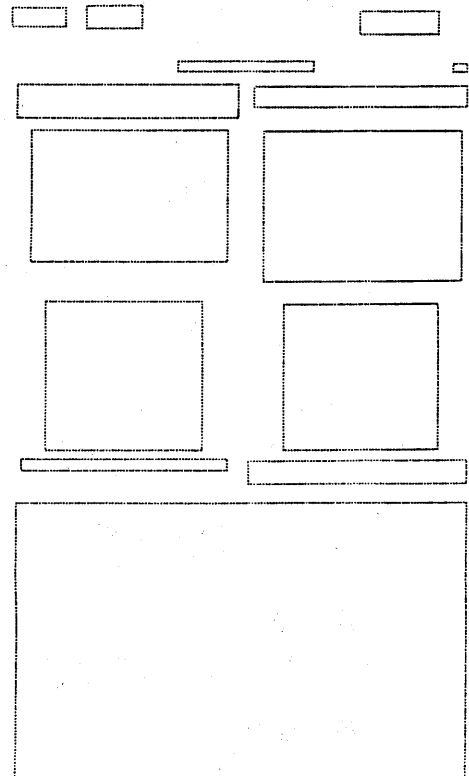


図6. 文章領域を抽出したものの

4. 矩形領域の分類

矩形領域を次の6種類のカテゴリーに分類する。

- 1) 文章領域 (Text)
- 2) 文字領域 (Character)
- 3) 写真領域 (Picture)
- 4) 表領域 (Table)
- 5) グラフ領域 (Graph)
- 6) 線画領域 (Line Drawings)

2) の文字領域は文字を主体として構成される2行以下の領域を示す。

4.1 矩形領域の特徴抽出

分類を効果的に行なうためには、各カテゴリーの情報かどの程度含まれているかという点から特徴を選択しなければならない。矩形領域から下記の9種の特徴を抽出した。

1) 黒画素の濃度

- 黒画素の発生確率

$$P_B = \frac{N_B}{N}$$

- 黒画素から黒画素への遷移確率

$$P_{BB} = \frac{N_{BB}}{N}$$

2) 黒画素のランレングス度数分布

- 平均

$$\bar{x} = \frac{1}{N} \sum_{i=1}^k f_i x_i$$

- 分散

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^k f_i (x_i - \bar{x})^2$$

- ヒズミ

$$\beta = \frac{\frac{1}{N} \sum_{i=1}^k (x_i - \bar{x})^3 f_i}{\sigma^3}$$

- 最頻値

$$m_0 = x_j, \quad f_j = \text{Max}_{1 \leq i \leq k} f_i$$

- 最大度数

$$m_f = \text{Max}_{1 \leq i \leq k} f_i$$

3) 周辺分布

周期関数の自己共分散関数は、周期的な増減、ピークを示す。そこで、X軸方向の周辺分布の自己共分散関数と自己共分散関数のパワースペクトルを求め、下記の特徴を抽出する。

- 自己共分散関数のピーク値
- 自己共分散関数のパワースペクトルの最大値と平均値の比

周辺分布, 自己共分散関数, パワースペクトルは次式で定義される。

X軸方向の周辺分布

$$h(y) = \int_a^b F(x, y) dx$$

ただし、対象領域は $a \leq x \leq b$ である。

関数 $x = x(t)$ の自己共分散関数

$$C_{xx}(t) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T (x(t) - \mu_x)(x(t+\tau) - \mu_x) d\tau$$

ただし、 $\mu_x = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T x(t) dt$

関数 $x = x(t)$ のパワースペクトル

$$\hat{x}(f) = \int_{-\infty}^{+\infty} x(t) e^{-2\pi i f t} dt$$

ここで、ピークとは次のように定義したものである。

- 1) タイムラグ $\tau = 0$ のときこれを1番目のピーク P_1 とする。
- 2) $C_{xx}(\tau_{ik})$ がピーク P_k を示すとは
 - $\tau = \tau_{ik}$ において関数 $C_{xx}(\tau)$ が極値をとる。
 - τ_0 に最も近いピーク P_{k-1} における関数値の符号と $C_{xx}(\tau_{ik})$ の符号が異符号である。

また、最大タイムラグ τ_m は、領域の大きさを考慮して、領域のX軸方向の長さの1/2とした。

処理対象数59領域 について特徴を抽出した。表2はカテゴリー別の領域数を示す。表3は59領域に対する特性等

表2. カテゴリ別処理対象領域数

カテゴリー	カテゴリコード	カテゴリナンバー	データ数(59)
文章	KT	1	14
文字	G1	2	22
表	G2	3	9
グラフ	G3	4	8
線画	G4	5	4
写真	KP	6	2

表3. 特性要因に対する反応

NO.	CAT.	Pa	P4	PaM	VARIANCE	SAGANESS	PAIJAAN	MODE	RAZIO	PSA
1307	1	0.0790	0.0670	0.0008	0.0930	92.0020	0.6520	2.6200	17.6200	2.0220
1310	1	0.0620	0.0620	0.0008	0.0410	72.0000	0.7100	2.6000	22.0020	3.0220
1503	1	0.0470	0.0470	0.0008	0.0410	215.0000	0.6000	3.2000	18.2020	18.2020
1613	1	0.0370	0.0260	0.0002	0.0007	163.0000	0.7400	2.4000	61.2020	6.2020
1204	1	0.0220	0.0460	0.0008	0.0009	176.0000	0.6120	2.4000	77.2020	77.2020
1502	1	0.0710	0.0490	0.0002	0.0010	104.0020	0.3900	2.4000	55.1000	7.2020
1111	1	0.0400	0.0280	0.0005	0.0020	94.0000	0.4020	2.2020	29.0000	3.2020
1411	1	0.0220	0.0400	0.0005	0.0020	176.0000	0.6400	3.0000	64.0000	7.2020
2393	1	0.0420	0.0280	0.0005	0.0025	176.0000	0.6400	2.7000	53.7000	9.2020
2004	1	0.0420	0.0280	0.0005	0.0013	175.0000	0.6000	2.1000	50.5000	11.2020
2095	1	0.0490	0.0780	0.0003	0.0016	154.0000	0.4800	2.3020	44.0000	12.2020
2096	1	0.0490	0.0780	0.0003	0.0012	184.0000	0.5000	2.7000	57.2020	16.2020
2007	1	0.0490	0.0640	0.0003	0.0014	156.0000	0.4700	2.5000	50.4000	12.2020
2004	1	0.0790	0.0670	0.0002	0.0009	180.0000	0.5000	2.4000	74.0020	17.0020
1203	2	0.1490	0.2540	0.0012	0.0030	25.0000	0.6200	3.1000	5.1000	4.2000
1202	2	0.1620	0.2860	0.0010	0.0020	82.0000	0.6300	2.5000	7.5000	4.2000
1201	2	0.0750	0.0870	0.0010	0.0010	42.0000	0.5000	2.7000	3.1000	2.0000
1302	2	0.1730	0.2440	0.0008	0.0030	44.0000	0.5200	3.0000	6.0000	1.0000
1507	2	0.1610	0.2870	0.0009	0.0040	47.0000	0.5400	2.7000	5.2000	3.0000
1506	2	0.1940	0.2620	0.0005	0.0020	44.0000	0.4600	3.0000	4.0000	1.0000
1505	2	0.1140	0.0830	0.0009	0.0020	70.0000	0.6000	2.5000	6.2000	1.0000
1501	2	0.1220	0.0810	0.0009	0.0020	14.0000	0.5000	2.4000	3.1000	2.0000
1612	2	0.1160	0.0930	0.0010	0.0020	42.0000	0.4000	3.2000	9.0000	3.0000
1810	2	0.0750	0.0800	0.0010	0.0010	44.0000	0.3700	2.4000	14.0000	2.0000
1402	2	0.0980	0.0800	0.0010	0.0010	52.0000	0.4200	1.9000	7.0000	4.0000
1404	2	0.0700	0.0920	0.0010	0.0010	85.0000	0.5000	2.5000	12.0000	2.0000
1403	2	0.0930	0.0910	0.0010	0.0010	54.0000	0.4000	2.4000	12.0000	2.0000
1401	2	0.1100	0.0870	0.0010	0.0010	10.0000	0.5000	2.4000	5.1000	3.0000
1304	2	0.0920	0.0830	0.0009	0.0010	44.0000	0.4000	3.0000	7.0000	1.0000
1310	2	0.1010	0.0630	0.0007	0.0010	27.0000	0.3700	3.0000	7.0000	3.0000
1308	2	0.1040	0.0600	0.0009	0.0010	22.0000	0.4000	3.0000	7.0000	3.0000
1305	2	0.1210	0.0630	0.0009	0.0010	44.0000	0.4000	2.7000	5.0000	1.0000
1206	2	0.0700	0.0700	0.0009	0.0010	42.0000	0.4000	1.9000	12.0000	1.0000
1207	2	0.1100	0.0620	0.0009	0.0010	42.0000	0.4000	1.9000	12.0000	1.0000
1302	2	0.1170	0.0630	0.0010	0.0010	7.0000	0.5000	3.0000	3.0000	3.0000
2032	2	0.0400	0.0400	0.0011	0.0011	16.0000	0.5000	3.0000	2.0000	1.0000
1306	3	0.0380	0.0490	0.0022	0.0022	557.0000	0.9000	19.2000	16.0000	30.2000
1405	3	0.0400	0.0400	0.0009	0.0009	16.0000	0.5000	3.0000	2.0000	1.0000
1405	3	0.0310	0.0430	0.0009	0.0009	205.0000	0.9000	19.2000	12.0000	20.2000
2101	3	0.0360	0.0330	0.0008	0.0008	500.0000	0.9000	17.2000	11.5000	32.0000
2103	3	0.0430	0.0430	0.0008	0.0008	491.0000	0.9000	17.2000	4.0000	31.2000
2104	3	0.0450	0.0380	0.0021	0.0021	263.0000	0.9000	16.2000	12.0000	31.2000
2105	3	0.0400	0.0450	0.0008	0.0008	491.0000	0.9000	17.2000	16.0000	17.0000
2110	3	0.0700	0.0900	0.0077	0.0077	245.0000	0.9000	13.7000	48.0000	9.2000
2111	3	0.0700	0.0900	0.0008	0.0008	492.0000	0.9000	16.2000	23.0000	23.0000
1408	4	0.0400	0.0380	0.0009	0.0009	492.0000	0.9000	16.2000	12.0000	39.2000
1407	4	0.0420	0.0400	0.0009	0.0009	492.0000	0.9000	11.7000	11.0000	41.2000
2107	4	0.0400	0.0400	0.0007	0.0007	492.0000	0.9000	8.0000	15.0000	31.2000
2112	4	0.0470	0.0370	0.0017	0.0017	22.0000	0.9000	15.1000	5.0000	59.2000
2106	4	0.0400	0.0400	0.0008	0.0008	216.0000	0.9000	21.0000	3.0000	79.2000
2104	4	0.0400	0.0400	0.0008	0.0008	271.0000	0.9000	8.0000	11.1000	46.2000
2105	4	0.0280	0.0400	0.0009	0.0009	397.0000	0.9000	40.0000	3.0000	60.2000
1205	5	0.0450	0.0450	0.0009	0.0009	492.0000	0.9000	7.7000	8.1000	35.2000
1205	5	0.0480	0.0480	0.0020	0.0020	247.0000	0.9000	10.1000	15.0000	55.2000
1502	5	0.0400	0.0400	0.0009	0.0009	492.0000	0.9000	9.2000	18.1000	79.2000
1503	5	0.0470	0.0470	0.0009	0.0009	492.0000	0.9000	9.2000	18.1000	79.2000
1502	5	0.0470	0.0470	0.0010	0.0010	247.0000	0.9000	9.0000	22.2000	55.2000
1503	5	0.0400	0.0400	0.0009	0.0009	492.0000	0.9000	5.1000	17.2000	55.2000
202	6	0.0320	0.0610	0.0001	0.0011	226.0000	0.6100	10.2000	165.0000	31.2000
302	6	0.0300	0.0650	0.0009	0.0009	49.0000	0.2000	6.1000	62.0000	41.0000

因の反応を示す。図7は表領域のランレングス度数分布を示し、図8は周辺分布の特性を示す。

4.2 判別分析

特性要因に対するデータの統計的解法として判別分析を用いる。6カテゴリのどれに属するかあらかじめ分っている領域の9特性要因に対する反応を特性変量データ(X₁, X₂, ..., X₉)とし判別分析を行い領域の分類を行なう。

この9次元表現を、各特性成分を有したままで1次元表現に変換する。すなわち、各特性変量にそれぞれ適当な大きさの重みW_iが1~9をつけて、次式のごとく、線型結合することにより、新しい特性変量Yを合成し判別関数を得る。

$$Y = \sum_{i=1}^9 W_i X_i$$

重みW_iは相関W_iが最大になるように決定する。

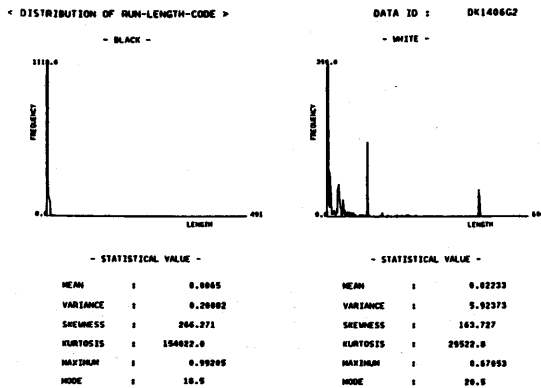


図7. ランレングス度数分布

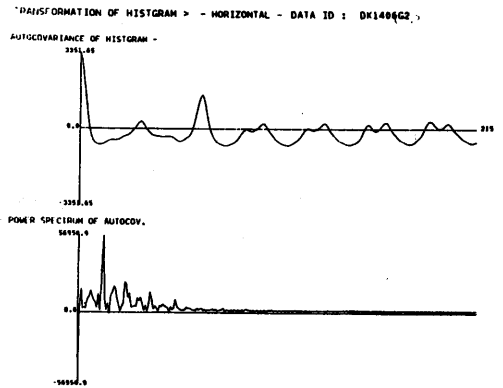


図8. 周辺分布の特性

相関比 η^2 は次のように定義される。

$$\eta^2 = \frac{\sum_{t=1}^k \pi_t \{E(Y) - E(Y^{(t)})\}^2}{D^2(Y)}$$

ただし、

π_t : t 番目のカテゴリーに Y が属する確率

$E(Y)$: 全体の Y の平均値

$D^2(Y)$: 全体の Y の分散

$E(Y^{(t)})$: t 番目のカテゴリー内での Y の平均値

全サンプル領域の特性変量データは (X_1, X_2, \dots, X_9) の行列で表現される。行列の各行が各領域の特性変量データである。このとき、最大相関比 λ と w_j は次の固有方程式の最大固有値とそれに対応する固有ベクトルとして求まる。

固有方程式 $AW = \lambda BW$

ただし、

$$A = (a_{jka})$$

$$a_{jka} = \sum_{t=1}^k \{ \pi_t E(X_j^{(t)}) E(X_k^{(t)}) \} - E(X_j) E(X_k)$$

$$B = (b_{jka})$$

$$b_{jka} = E(X_j X_k) - E(X_j) E(X_k)$$

$$W = (w_j)$$

$$j, k = 1, \dots, 9$$

$$a_{jka} = a_{kja}$$

$$b_{jka} = b_{kja}$$

この固有方程式をエバーライン法で解を求めた。求まる重みから判別関数を決定し、正規化してカテゴリー別頻度分布から判別点をミニマックス法で求め的中率を計算した。図9は前記の対象領域に対する分類結果を示す。

5. まとめ

文書画像の矩形領域の抽出と分類の汎用的アルゴリズムを報告した。分類した矩形領域の解析・認識については別に報告がある。また、これらの実験結果の応用が今後の課題である。

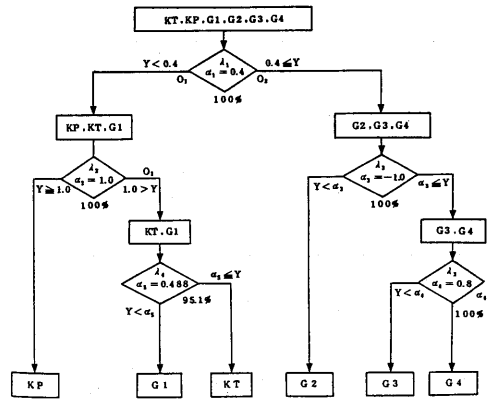


図9. 判別分析結果

参考文献

- 1) 村尾, 坂井: "文書画像における構造情報の抽出", 情報学会全国大会, 1980.
- 2) 秋山, 増田: "印刷物の記事領域における文字の切り出し", 信学技報, PRL 80-70, 1981.
- 3) 伊藤, 他: "文書構造の表現と解析", TSCレポート, N:G318-1574, 1982.
- 4) S. ITO, et al.: "Field Segmentation and Classification", 6th ICPR, 1982.
- 5) 伊藤, 他: "文書領域の特徴抽出と判別分析", TSCレポート, N:G318-1573, 1982.
- 6) 長谷, 星野: "2次元フーリエ変換を用いた文書画像の領域判別法", 情報学会, コンピュータビジョン研 20, 1982.
- 7) S. ITO, "Automatic Input of Flow Chart in Document Image", 6th ICSE, 1982.
- 8) 伊藤, 他: "文書画像における表領域の抽出とその認識", 第13回画像工学コンファレンス, 1982.
- 9) 伊藤, 他: "ファクシミリ文書の矩形フィールドと手書き識別番号の認識", 第12回画像工学コンファレンス, 1981.