

部品イラスト図面からの文字列パターンの抽出

武田 晴夫 †

小野裕次郎 ††

† (株)日立製作所システム開発研究所

†† 日立京葉エンジニアリング(株)

部品イラスト図面を主対象として、画像入力した図面から文字および文字列パターンを抽出するアルゴリズムを提案する。本アルゴリズムは、画像を構成する要素に対して文字、図形、接触文字、分離文字などのラベルを付す過程と、この結果に基づいて要素を再抽出する過程を交互に繰り返すことによって、個別文字を抽出する。前者のラベル付けには文字列情報を利用するために、近接する他の要素のラベルが緩和過程を通じて利用される。本アルゴリズムによって、単独では抽出困難な図形や他の文字と接触した文字が抽出できることを示す。また文字に類似した図形の誤抽出を解消できることを示す。最後に本アルゴリズムの自動車の部品イラスト図面への適用結果について述べる。

Extraction of String Patterns from Illustration of Parts

Haruo Takeda †

Yuujiro Ono ††

† Systems Development Laboratory, HITACHI Ltd.

†† HITACHI Keiyo Engineering Ltd.

† 1099, Ohzenji, Asao-Ku, Kawasaki 215, Japan

An algorithm of extracting character and string patterns from illustration of parts is presented. The process of labeling to picture elements and the process of extracting elements by using the labels are repeated to extract characters. The relaxation method is used in the labeling process to use the string informations. This algorithm enables the extraction of characters touched by lines or other characters. It prohibits the extraction of parts similar to characters. The application to CD-ROM automobile parts catalogs are shown.

1. はじめに

光ディスク技術の進展などにより、紙面に書かれた図面・文書をイメージスキャナなどから直接読み込んで、データベース化するシステムが普及しつつある。このような分野では、蓄積されたデータの検索や利用を効率よく行うために、図面・文書上に記載された文字や記号を認識する技術が期待されている。

図面・文書上の文字を認識する場合、文字ピッチが必ずしも一定でないものから、文字間が接触するなどの状況下で如何に個別の文字パターンを抽出するかが重要な課題となる。従来、部分的な文字ピッチの乱れに関して組版知識を利用する方法[1]、非接触文字の切り出し情報を用いて接触文字を分離する方法[2]、文字の形状に関する知識を利用する方法[3]などが提案されている。これらは新聞・雑誌など、概ね規則的な行の配列から行の単位で抽出した文字列に対して、文字パターンの抽出方法を提案している。これに対して図形が文字領域に混在した一般の図面データについては、近接する線の密度[4]、一定の矩形単位での図形構造[5]などの特徴量を用いて文字パター

ンを直接抽出する方法が提案されている。これらの方法では、文字と類似した特徴を示す図形データへの配慮が必要である。

本報告では、画像を構成する各要素に対してラベルを付す過程と、この結果に基づいて要素を再抽出する過程を交互に繰り返すことにより、個別文字を抽出する方式を提案する。前者の過程で付すラベルは、当該要素が、文字パターン、図形パターン、複数の文字の接触したパターン、文字と図形の接触したパターン、分離した文字の部品パターンなどの何れであるかを示す。各要素ごとのラベルは確率で表現する。この確率の値は、初期値を各要素の画像処理によって求めた後、近接する他の要素のラベル、それとの位置関係などにより緩和法[6]を用いて調整する。後者の再抽出の過程では、他の文字・図形との接触パターンと判定された要素の分離、分離した文字部品パターンと判定された要素の合成などの処理を行なう。

このように図面の構成要素の編成を、要素の認識結果、および他の要素との関係を考慮して、動的に変更することにより、単独では抽出が困難な図形や他の文字と接触した文字、複数の部分に分

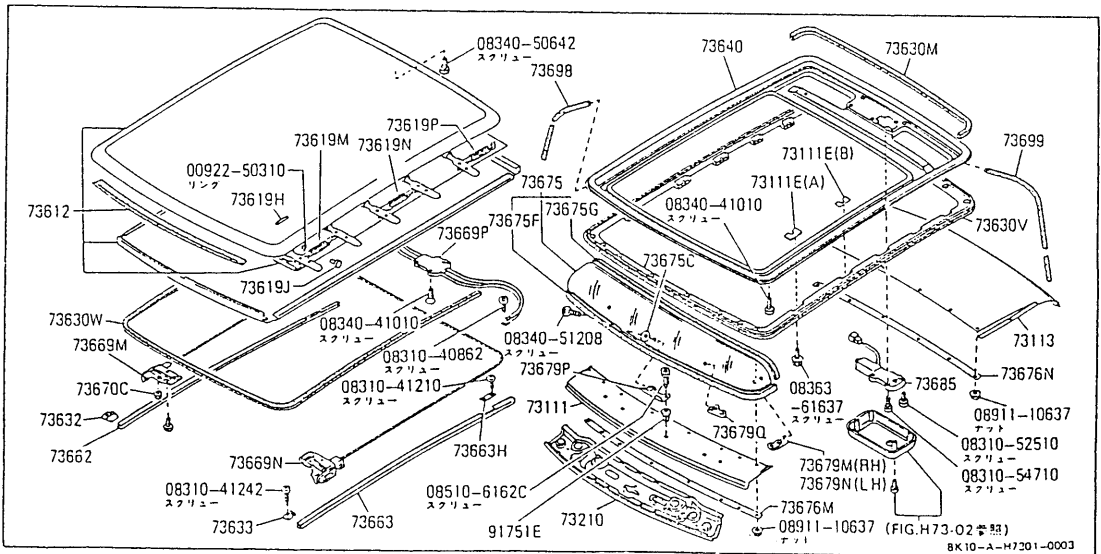


図1 部品イラスト図面の例

離した文字などを図面の中から抽出できることを示す。またこれによって、文字と類似した図形の誤抽出を抑制できることを示す。最後に本アルゴリズムの実験結果として、自動車の部品イラスト図面から部品コードを抽出・認識してCD-ROMによる電子カタログを作成するシステム[7]での適用例について述べる。

2. 前提条件

本報告で対象とする部品イラスト図面の一例として、自動車の補修用部品のカタログ図面を図1に示す。部品イラスト図面の一般的な特徴を以下に示す。

- a) 文字が文字列単位に、図形上の任意の位置に記入される
- b) 図形は任意の自由曲線であり、全体の大きさ、線分の長さ・幅、記号の種類などの限定はできない（文字と同等のサイズの部品図形も存在する）
- c) 文字と図形、文字と文字が接触するケースが存在する

このような図面から文字列情報を用いて文字パターンを抽出するために、以下文字列の印刷に関する前提条件を定義する。

通常このようなトレース図面における文字列は、写植、レーザプリンタ、ドットプリンタなどで印

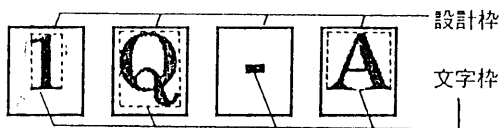


図2 文字枠と設計枠

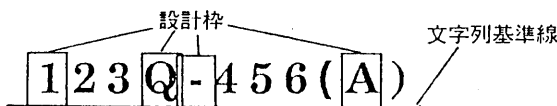


図3 文字列基準線

刷された紙片を、図形が描画された台紙上に文字列単位に貼付することによって作成される。ここでは印刷に用いる文字パターンの集合および組版パターンの集合が既知であると仮定して、印刷に関して次の3つの条件を設定する。（図2，3参照）

- (1) 1つの文字列を構成する文字の設計枠（文字枠周囲の余白を含む）の下辺は、1本の文字列基準線上に存在する。従って文字列基準線に対する文字枠（文字自身の外接矩形）の下辺の位置は、文字パターンに対して一意に決まる。
- (2) 文字の設計枠の高さは、同一の文字列内では互いに等しい。
- (3) 隣接する文字の間隔は、基準となるピッチと文字コードによって決まる（プロポーショナルピッチ）。ここでは文字枠間隔の設計枠の大きさに対する比の最大値が既知とする。

なお、通常同一書体、同一符号の文字パターンでサイズの異なるものは、光学的変換や輪郭曲線の座標変換などによって同一の字母から生成されるものも多い。このため本質的な条件ではないが、処理簡単化のために、以下このようなサイズの異なる同一の文字パターンは幾何学的に相似の関係にあるものとする。

3. 文字パターン抽出アルゴリズム

本アルゴリズムの基本手順を図4に示す。以下各ステップの詳細を説明する。

3.1 前処理

前処理では、図面全体の傾きを検出すると共に、処理対象範囲を限定する。ここでは傾きの基準とできる図面外枠が存在するものとして、外枠4辺の検出によって処理を行う。一般に文書の認識では傾きの補正が重要であるが、本アルゴリズムでは幾何学的変換に伴う局所的画質劣化を避けるために画像の補正自体は行わず、以下の処理において文字列の傾きに関するパラメータとして考慮する。

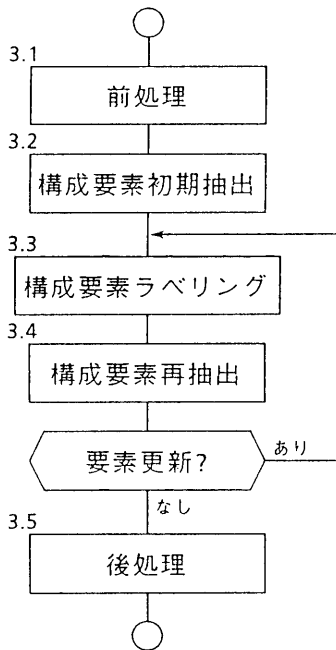


図4 文字パターン抽出アルゴリズムの基本手順

3.2 構成要素初期抽出処理

本処理では、黒画素の連結成分（8連結）を抽出し、各連結成分の、画像に対して正置した外接矩形の位置および大きさの情報を得る。ここでは図5に示す座標系について、第*i*外接矩形の位置は外接矩形左下隅座標(x_i, y_i)で表す。また第*i*外接矩形の大きさを、幅 w_i および高さ h_i で表す。本処理で求めた矩形情報を、当該図面の各構成要素の初期状態とし、以下 a_i ($i=1, 2, \dots$)と表す。

3.3 構成要素ラベリング処理

本処理では、まず本処理の開始以前にラベルが付されていない構成要素にラベルを付す。本ラベルは当該要素が、

- 1) 文字パターン (λ_1)
- 2) 図形パターン (λ_2)
- 3) 複数の文字の接触により生じた文字の合成パターン (λ_3)

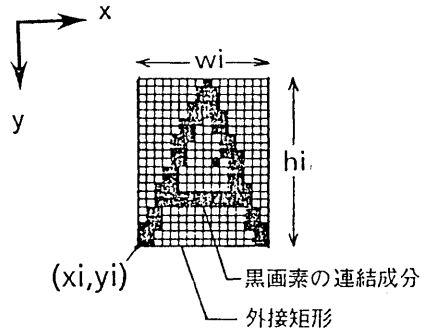


図5 黒画素の連結成分と外接矩形

- 4) 文字と図形の接触により生じた文字と図形の合成パターン (λ_4)
- 5) 分離した文字の一部である文字部品パターン (λ_5)

の何れであるかを示す。各要素 a_i に対するラベルは、それがラベル $\lambda_1, \dots, \lambda_5$ である確率の組 $p_i(\lambda)$ で表す。ここに

$$\sum_m p_i(\lambda_m) = 1$$

とする。

未ラベリングの構成要素 a_i へのラベル付加は、 a_i のラベルの初期値 $p_i^0(\lambda)$ を画像の認識処理により求めることによって行なう。画像認識では、構成要素の画像としての特徴量から、予め定義した方法に従って、上記ラベルの確率を求める。特徴量としては、ヒューリスティックな方法で選択した例えば

- a) 各要素 a_i の外接矩形の幅 w_i と高さ h_i
 - b) 上記の高さと幅の比率 h_i/w_i
 - c) 黒画素の占める面積の要素全体に対する割合
 - d) 内部の複雑度（黒白変化点の個数、空間周波数、単位面積当りのエッジの長さ等）
 - e) 標準文字パターンとのユークリッド距離
- などが考えられる。これらの特徴量からラベルの確率を求める処理は、一種のパターン認識処理と見做すことができる。ここでは対象図面の画素密度、使用する文字の大きさ・種類、線の太さ等に

適応させて精度を向上するために、標本図面の例示により計算方法を定義するものとする。

次に上記処理で求めたラベルおよび既にこれ以前に求めたラベルについて、近接する構成要素の間での整合をとることによって、その値を修正する。本処理は、 λ の整合性を表す関数 γ_{ij} を用いた $p_i^{(k)}(\lambda)$ に関する緩和法[6]によって求める。ここで関数 γ_{ij} は、要素 a_i と要素 a_j の各ラベルの組について、実数 $[-1, 1]$ を写像する関数であり、特に正の整合に対してその大きさに応じて実数 $(0, 1]$ 、負の整合に対してその大きさに応じて実数 $[-1, 0)$ 、整合なしの関係に対して値0を出力するものとする。

具体的には、各要素 a_i に対して、まずこれと近接する要素 a_j を求める。近接の条件は、要素 a_i の外接矩形と要素 a_j の外接矩形が所定のマージンについて共通部分を有するか否かとすることがで上記マージンの値は、文字が文字列を構成するための例えば次の条件から決定することができる。

$$-\varepsilon_1 < (Y_{i_2 j_2} - Y_{i_1 j_1}) / (X_{i_2 j_2} - X_{i_1 j_1}) - \tan(\theta - \theta') < \varepsilon_1 \quad (1)$$

$$-2\alpha' < H_{i_2 j_2} - H_{i_1 j_1} < 2\alpha' \quad (2)$$

$$-2\alpha < x_{i_2 j_2} - x_{i_1 j_1} - w_{i_1 j_1} < x w_{i_1 j_1} + 2\alpha \quad (3)$$

これらの条件は隣接する2つの文字候補が満たすべき条件であり、それぞれ第2章の条件(1)~(3)に対応する。ここで、 i_1, i_2 は前節で述べた矩形形状ラベルを示すインデックス、 j_1, j_2 は文字認識結果の候補を示すインデックス、 $(X_{i_1 j_1}, Y_{i_1 j_1}), (X_{i_2 j_2}, Y_{i_2 j_2})$ は当該ラベルおよび認識結果を前提としたときの文字設計枠の左下位置、 $w_{i_1 j_1}$ はその幅、 $H_{i_1 j_1}, H_{i_2 j_2}$ はその高さを表す。また α' は矩形サイズから上記設計枠を求めるために用いる比率に従って α を変換したもので、 ε_1 は $\varepsilon_1 = (1/w_{i_1 j_1} + 1/H_{i_1 j_1}) \cdot 2\alpha'$ と計算されるマージンである(付録参照)。 θ' は図面に対する文字列の傾きを示すパラメータで、ある文字候補

の組に対して共通の値が存在すれば(1)の条件が成立するものとする。 x は文字ピッチに関するパラメータで、条件(3)に対して、 θ' と同じ働きをする。

次に近接するすべての a_i, a_j の組に対して、それらの各ラベルの組合せの整合性 γ_{ij} を求める。本整合性は、次の情報を用いて決定することができる。

- ①各要素の外接矩形の包含関係
 - ②各要素の外接矩形の位置関係
 - ③2つの要素の最接近点間の距離
 - ④その他の図面に固有の知識
 - ⑤同一の黒画素に対する複数の要素の割当て
- 図面に固有の知識としては、例えば文字列に対する図形への引出線の存在、孤立した1文字が存在しない、ハイフンの連続は図形(破線)である、微小な点の集合は図形(ハッチング)である、などを考える。 γ_{ij} は上記①~⑤の値により分岐する決定木と、その値自身による関係の確からしきの情報を用いて、上記処理と同様にして標本図面から予め求める。

緩和処理の過程では、上記 γ_{ij} の値を用いて、次の漸化式を計算する。

$$p_i^{(k+1)}(\lambda) = (p_i^{(k)}(\lambda) (1 + q_i^{(k)})) / \sum_j (p_j^{(k)}(\lambda) (1 + q_j^{(k)}))$$

$$q_i^{(k)} = \sum_j (d_{ij} \sum_{\lambda'} \gamma_{ij}(\lambda, \lambda')) \cdot p_j^{(k)}(\lambda)$$

ここに

$$\sum_j d_{ij} = 1$$

とする。

3.4 構成要素再抽出処理

本処理では、上記ラベルのうち $\lambda_0 \sim \lambda_0$ である確率が高い要素について、構成要素の分離または統合処理を行なう。

本処理は、 λ_0 と λ_0 については、その構成要素と近接する要素のうちラベル λ_0 である確率が高

いものから、文字列を生成し、その延長上に文字候補を探索することによって新たな構成要素を抽出する。 λ_1 については、その構成要素と近接する要素のうちラベル λ_1 である確率が高いものの統合を行ない新たな構成要素を合成する。

次に本抽出が、図形や文字の一部を切り出す誤抽出である可能性を求める。誤抽出の可能性は、対象とする図形パターン・文字パターンの一部がある文字パターンとなっているケースについて

a) 認識結果が1, -, /の何れかである場合の、

抽出によって図形中の線分が切断される程度

b) 認識結果が3である場合の、8, B等の文字の一部である可能性の強さ

などの量から求める。これらの量は画像を直接参照することによって求める。この結果により図形である可能性が高いものはラベル λ_2 である確率が高いラベル値を、上記3.3における漸化式の初期値とする。

本再抽出処理において、ラベルが $\lambda_1 \sim \lambda_2$ であるにもかかわらず新たな構成要素を生成できないものについては、当該要素のラベル値をラベル λ_2 である確率が高いラベル値に変更する。新たな構成要素を生じた場合には、上記構成要素ラベリング処理を再度繰り返し、生じない場合には、本文抽出処理を終了する。

3.5 後処理

後処理では、文字コードの最終確認・訂正、抽出した文字パターン位置・大きさの正規化、認識結果に基づく画像データの修正等を行う。

4. 実験結果

本アルゴリズムの実験結果として、イラスト的に描かれた自動車部品のカatalog図面から部品コードを抽出・認識してCD-ROMを作成するシステムでの本アルゴリズムの適用例について述べる。本システムは、自動車の修理工場などで、整備・修理に使う部品の番号を検索するためのCD-ROMカタ

ログ[6]を作成するものである。CD-ROM検索時に図面上の部品番号をマウス等で指示するだけで該当部品を自動的に発注できるようにするために、CD-ROM作成時に図面上の部品コードなどの抽出・認識が必要となる。

4.1 処理例

この図面の認識における代表的な処理例について説明する。以下処理対象として図6のパターンを例に説明する。

(a) 端文字での図形との接触

左端の「7」および右端の「C」は図形に接触している(原画像上は分離していても、画像入力の変換誤差により接触するもの、以下も同様の意味)。まず構成要素ラベリング処理により中間の「3675」が抽出・認識される。これらの要素が文字である確率は周囲の文字により互いに強められる。図形に接触した文字については、同じラベリング処理の中でラベル λ_4 である確率が高められる。この結果次の構成要素再抽出処理において、上記中間部分で求めた文字列のピッチ、傾きのパラメータを用いて文字の探索が行われ、図形と接触した「7」および「C」が認識・抽出される。

(b) 中間文字での図形との接触

中間の「3」が図形(引出し線)に接触している。構成要素ラベリング処理で抽出・認識された「736」および「0M」の情報を用いることによって、上記(a)と全く同様に、これを抽出点認識することができる。

(c) 他の文字列との接触

上段の文字列の左端「7」と下段の文字列の右端「P」が接触している。2つの文字列でそれぞれ上記(a)と全く同様の処理が行われ、「7」および「P」がそれぞれの文字列の一部として抽出される。

(d) 同一文字列内での文字接触

中間の4文字「-000」が接触している。第1回の構成要素ラベリング処理で「BK10-A-H7301」と単独の文字からなる「3」が抽出されるが、これを基に第1回の再抽出処理で左端の「-」および右端の

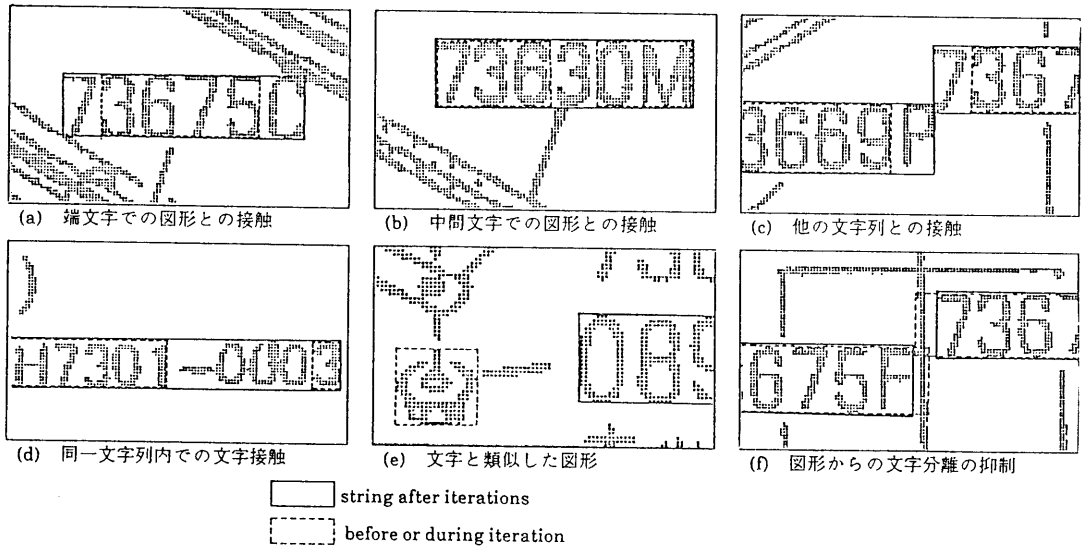


図6 特徴的なパターン

「0」が抽出される。更に第2回のラベリング処理で「BK10-A-H7301-」と「03」が抽出され、さらに第2回の再抽出処理で残りの「00」が抽出される。

(e) 文字と類似した図形

左端の「ナット」は、幅と高さの比率および単純類似度によれば、単独では「w」と認識される。但しこの例では右側の文字列「08911…」とはy方向の位置が異なるので近接する文字による文字である確率の強化が行なわれず、逆に図形中の単独文字は、図形である確率が高いとするルールにより、最終的には図形と判定される。同様に水平破線は3.3で述べた文字列の意味に関するルールで、垂直破線もここでは「単独の1は文字ではない」という意味に関するルールで図形として認識される。

(f) 図形からの文字分離の抑制

仮想文字列として「73675F」かつ右端が閉じていないものが得られると、文字抽出処理の前半では右端部分に「1」を認識する。ただしこれを抽出す

ると連続する図形線分を2つに分割するため、この抽出は上記3.4で述べた抑制ルールによって行われない。

4.2 実験結果

画像入力の特徴化密度が240dpi (dpi: 1インチ当たりの画素数)としたとき、 $\alpha = 2$ に設定したときの実験結果を以下に示す。対象図面をある1車種の全図面としたとき、文字抽出率は約99.5%、抽出された文字に対する文字認識率は約99.7%であった。これらの値は、特に3.3および3.4で述べた、処理対象図面に固有の知識を導入することによって高めることができた。

処理時間は10MHzの68020の場合、A4サイズ当り平均約40秒である。

5. まとめ

次の特徴をもつ部品イラスト図面を対象として、文字および文字列パターンを抽出する方法について述べた。

- 1) 文字が文字列単位に、図形上の任意の位置に記入される
- 2) 図形は任意の自由曲線であり、文字と同等のサイズの部品図形も存在する
- 3) 文字と図形、文字と文字が接触するケースが存在する

画像を構成する各要素に対してラベルを付す過程と、この結果に基づいて要素を再抽出する過程を交互に繰り返すことによって、個別文字を抽出する方式を提案した。図面の構成要素の編成を、要素の認識結果、および他の要素との関係を考慮して、動的に変更することにより、単独では抽出が困難な以下のようなパターンを抽出できることを示した。

- a) 図形と接触した文字
- b) 他の文字列と接触した文字
- c) 文字列内で接触・分離した文字

さらに、文字と類似した部品、図形の一部が文字と類似しているパターンが抽出されるのを抑制できることを示した。イラスト的に描かれた自動車部品のカタログ図面から部品コードを抽出・認識してCD-ROMを作成するシステムでの適用例では、文字抽出率は約99.5%であった。

本認識に用いるパラメータ値、条件などの知識を図面から獲得する方法について、別途報告する予定である。

謝辞 図面を提供して頂き、また実用評価にご協力頂いた日産自動車(株)部品事業部、佐久間輝雄主担、同、高相友三郎主査、(株)日産コーエー技術情報営業部、荒井勇三朗課長、(株)日立製作所大森ソフトウェア工場、鎌田修一技師、同、玉手健久氏に感謝致します。

参考文献

[1] 豊田, 野口, 西村: 日本語印刷文書における文字切り出し—新聞自動読み取りへの応用—, 情報処論, 24, 4, pp.481-487(1983).

[2] 秋山, 内藤, 増田: 非接触文字優先切出しによる印刷物からの文字切出し法, 信学論(D), J67-D, 10, pp.1194-1201(1984).

[3] 中村, 鈴木, 南: 横書き日本語文書における個別文字の抽出, 信学論(D), J68-D, 11, pp.1899-1909(1985).

[4] 岩城, 久保田, 荒川: 近接線密度法による文字・図形分離抽出, 信学論(D), J68-D, 4, pp.821-828(1985).

[5] 長田, 井上, 吉田: 論理回路図の自動入力処理, 信学論(D), J68-D, 4, pp.837-844(1985).

[6] S.W.Zucker, R.A.Hummel and A.Rosenfeld: An Application of Relaxation Labeling to Line and Curve Enhancement, IEEE trans. Comput., c-26, 4, pp.394-403 (1977).

[7] 佐久間, 神林, 鎌田, 他2名: パーソナルコンピュータとCD-ROMによる自動車部品画像検索システム, 日立評論, 71, 4, pp.47-52(1989).

付録 文字枠の高さと幅の比率に含まれる誤差

同一の文字を印刷、画像入力した場合の文字枠の幅、高さの観測値の範囲は、印刷する文字の大きさによらない定数 α に対して、

$w_0 - \alpha < w < w_0 + \alpha$, $h_0 - \alpha < h < h_0 + \alpha$
とできるとき、

$$\begin{aligned} (h_0 - \alpha)/(w_0 + \alpha) &< h/w \\ &< (h_0 + \alpha)/(w_0 - \alpha) \end{aligned}$$

となる。よって

$$\begin{aligned} (h_0 - \alpha)/(w_0 + \alpha) &= h_0/w_0(1 - \alpha/h_0)(1 + \alpha/w_0)^{-1} \\ &\doteq h_0/w_0(1 - \alpha/h_0)(1 - \alpha/w_0) \\ &\doteq h_0/w_0 - (1/h_0 + 1/w_0)\alpha \end{aligned}$$

などにより、 h/w の誤差は

$$\begin{aligned} -(1/h_0 + 1/w_0)\alpha &< w/h - h_0/w_0 \\ &< (1/h_0 + 1/w_0)\alpha \end{aligned}$$

となる。