

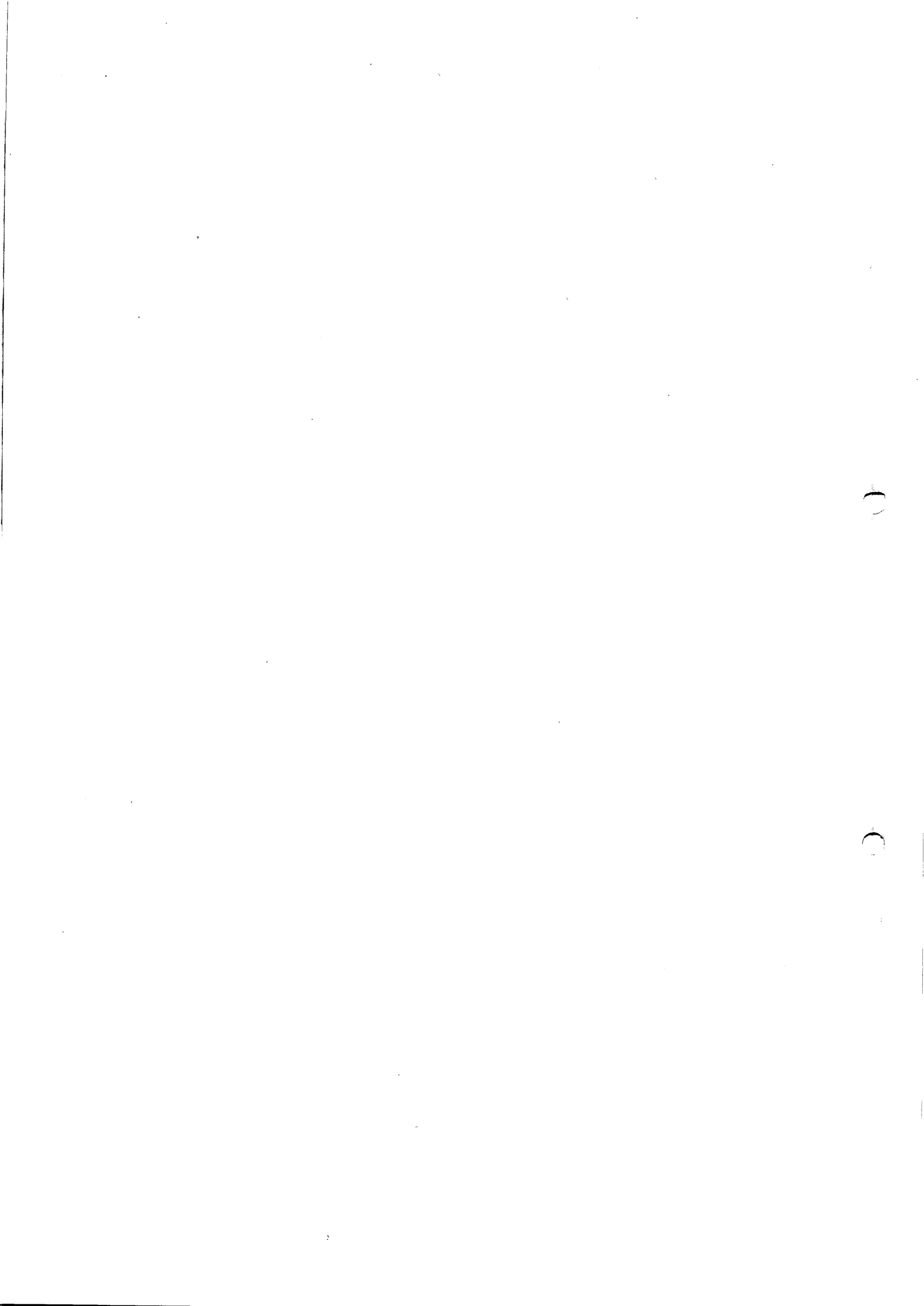
修士学位論文本審査資料

# 多フォント文書認識に関する基礎的研究

東北大学大学院工学研究科

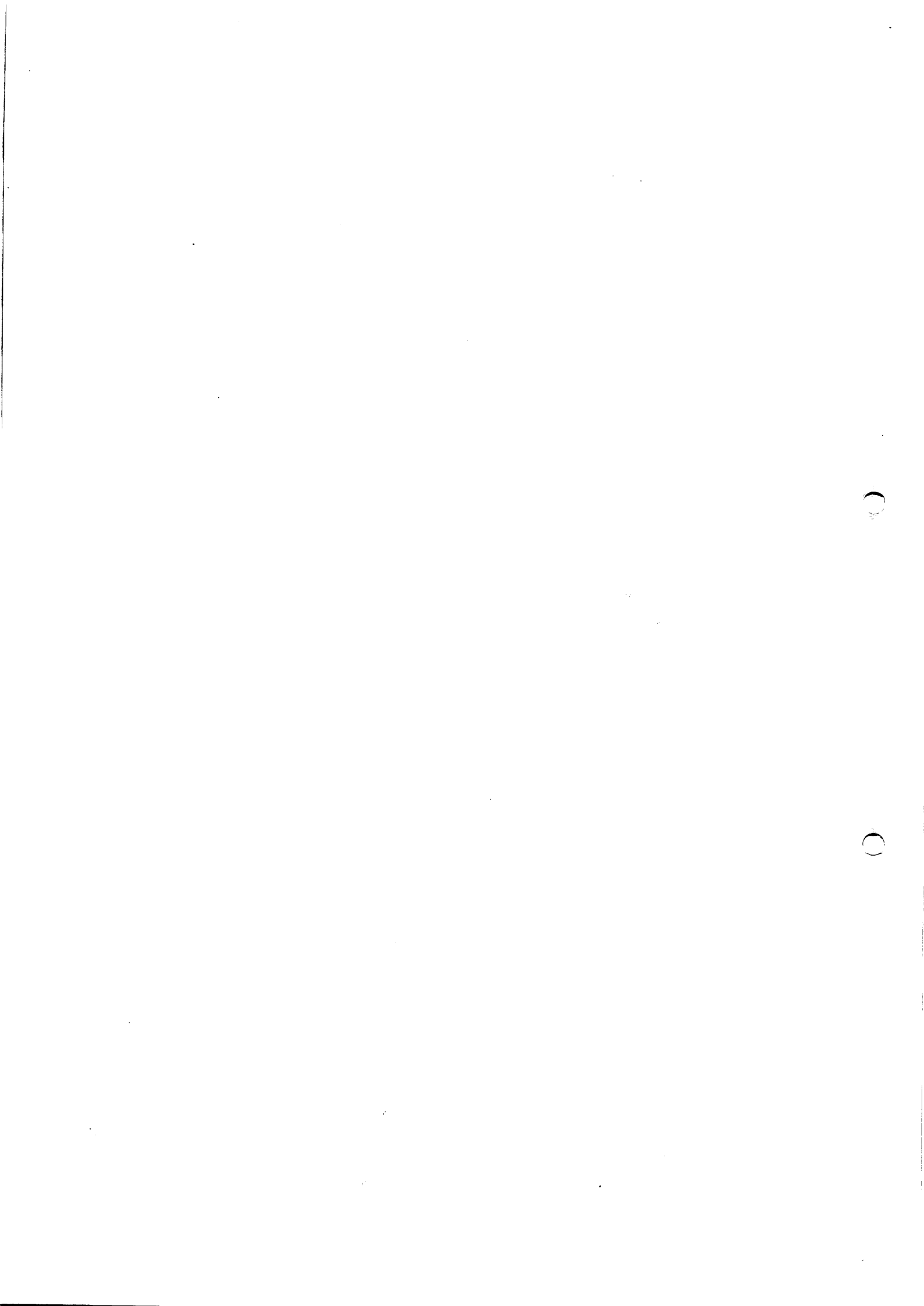
情報工学専攻

池田 啓明

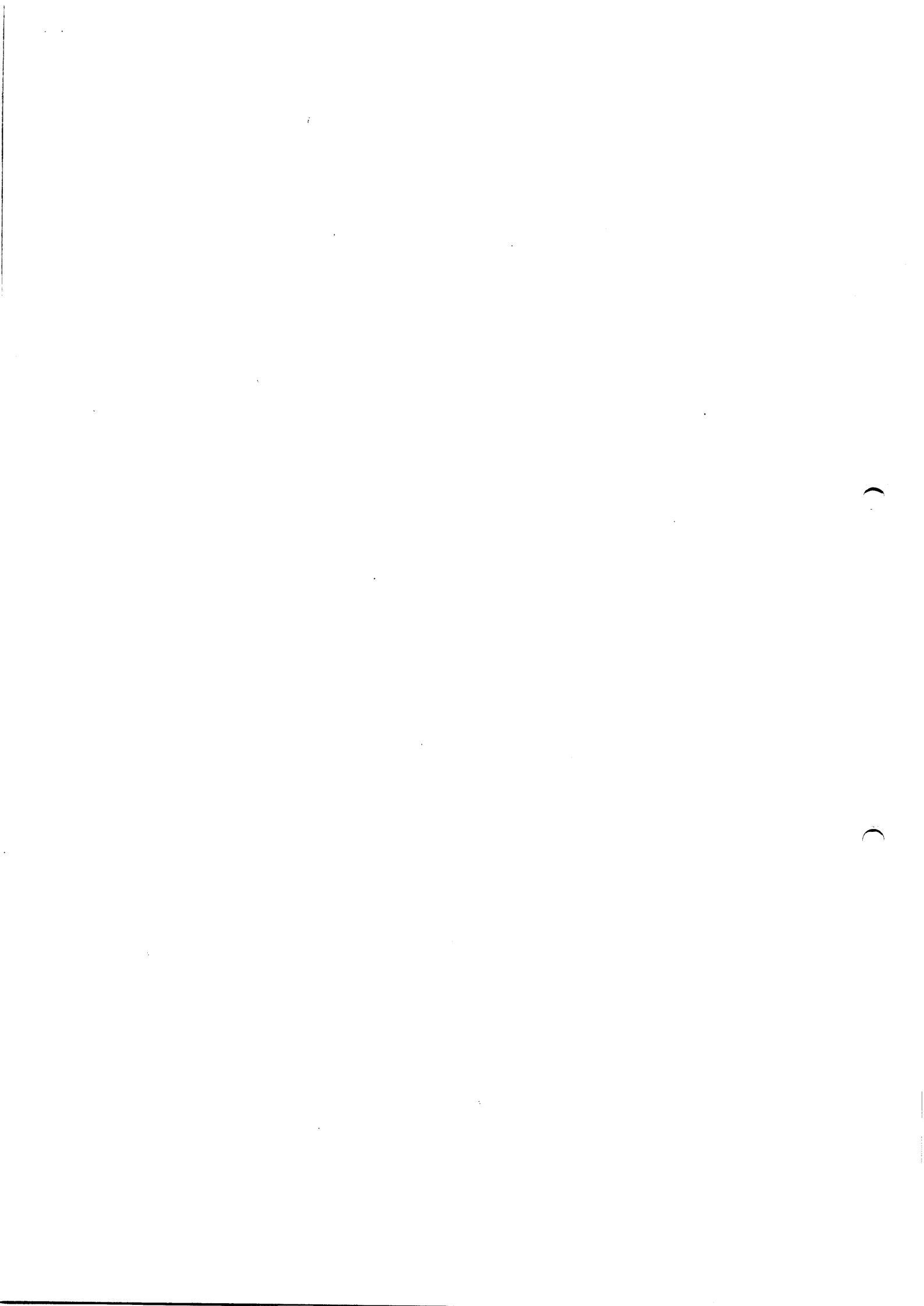


# 目次

<b>1</b>	<b>序論</b>	<b>1</b>
1.1	本研究の背景及び目的	1
1.2	本論文の構成	2
<b>2</b>	<b>文書認識システム</b>	<b>3</b>
2.1	まえがき	3
2.2	文書認識システムの構成	3
2.3	文字認識の方法	4
2.3.1	文字入力・前処理	4
2.3.2	特徴抽出	5
2.3.3	認識	6
2.3.4	候補出力	7
<b>3</b>	<b>フォントの抽出</b>	<b>12</b>
3.1	前書き	12
3.2	フォントの抽出	13
3.2.1	認識・結果出力	13
3.2.2	リストの作成・候補集合の作成	14
3.2.3	除去候補文字の処理	14
3.2.3.1	詳細識別	16
3.2.4	細分類	16
3.2.5	除去候補文字の学習	17
3.3	抽出実験	19
3.3.1	実験方法	19
3.3.2	実験結果	19

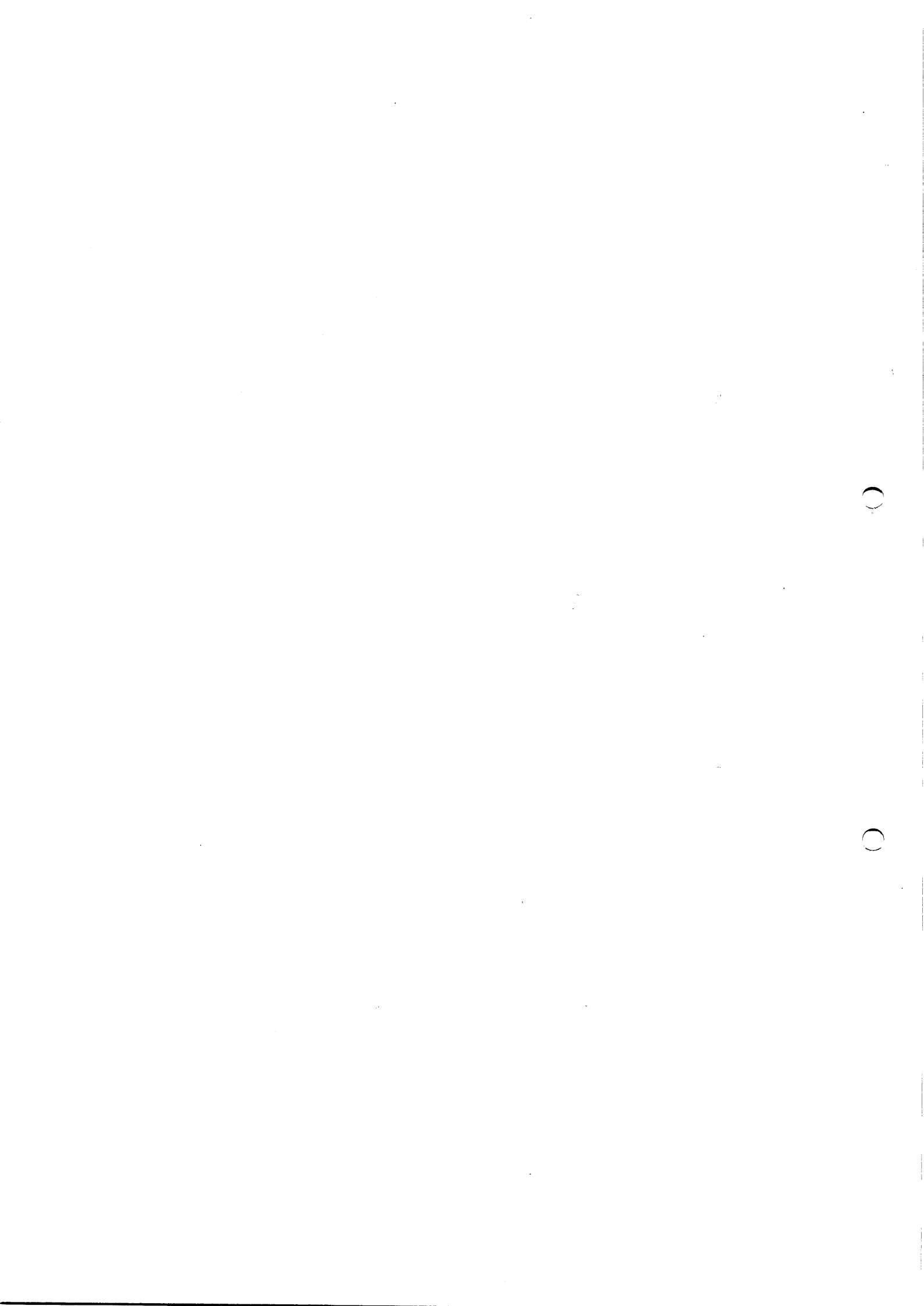


目次	ii
3.4 考察	19
4 認識実験	41
4.1 前書き	41
4.2 実験方法	41
4.3 新辞書の作成方法	41
4.3.1 マルチテンプレート法	42
4.3.2 置換法	42
4.4 クローズ実験	42
4.5 オープン実験	42
4.6 更新した辞書による認識実験	43
4.6.1 辞書の更新方法	43
4.6.2 実験結果	44
4.7 考察	45
5 結論	68
謝辞	70
参考文献	71
A 社説の構成字種	付録1



# 目 次

2.1	文字認識アルゴリズム	8
2.2	文書認識システムのアルゴリズム	9
2.3	各処理後のイメージ	10
2.4	方向線素特徴量	11
3.1	重心を用いたフォントの抽出法	21
3.2	未知文書からのフォント抽出アルゴリズム	22
3.3	フォント抽出の対象とした文字	23
3.4	辞書中の仮名文字	23
3.5	新聞社説の例	24
3.6	認識結果の例	25
3.7	コードとフォントの対応	25
3.8	候補集合の例	26
3.9	除去候補文字の学習例 1	27
3.10	除去候補文字の学習例 2-1	28
3.11	除去候補文字の学習例 2-2	29
3.12	詳細識別文字	33
3.13	識別領域の例	33
3.14	識別領域	34
3.15	候補数と出力数	37
3.16	誤抽出の例	39
4.1	辞書の更新の例 1	48
4.2	辞書の更新の例 2	49
4.3	候補数と認識率(ひらがな 65 文字)	54
4.4	候補数と認識率(その他)	55





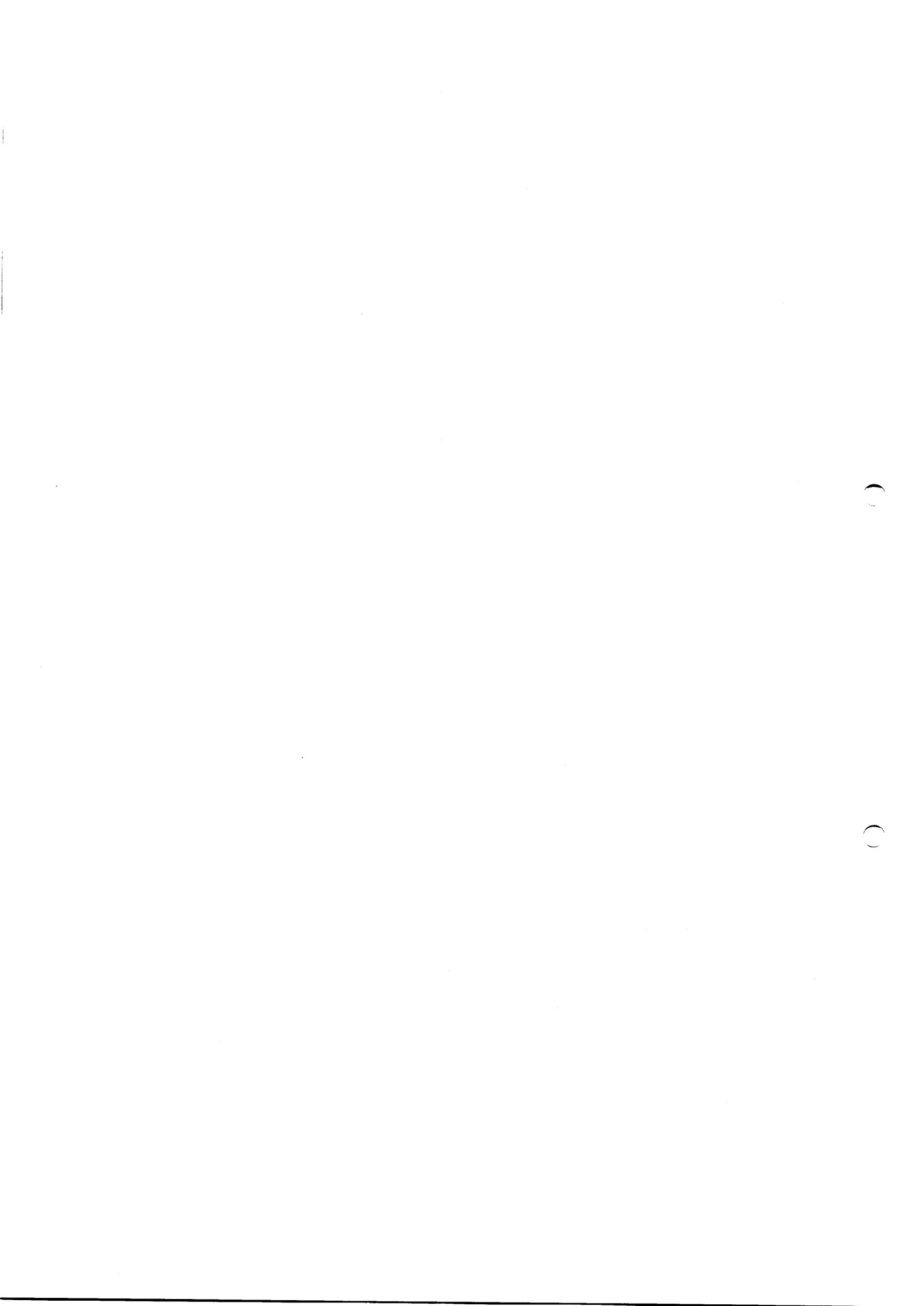
## 目 次

4.5	認識文書毎の1位認識率 . . . . .	61
4.6	学習文書数と抽出辞書中の文字 . . . . .	62
4.7	学習文書数毎の1位認識率 . . . . .	67



# 表 目 次

3.1	除去候補文字 1	30
3.2	除去候補文字 2	31
3.3	除去候補文字数と評価値	32
3.4	候補数と候補集合	35
3.5	候補数と出力数	36
3.6	誤抽出された文字	38
3.7	抽出できなかった文字	40
4.1	従来の辞書の認識率 (社説 11~20)	50
4.2	1 位認識率 (クローズ実験)	51
4.3	改善の結果 (マルチテンプレート法)	52
4.4	改善の結果 (置換法)	53
4.5	認識文書中に存在しないひらがな 65 種の文字	56
4.6	マルチテンプレート法の 1 位認識率 (オープン実験)	57
4.7	置換法の 1 位認識率 (オープン実験)	58
4.8	抽出フォント毎の 1 位認識率 (マルチテンプレート法)	59
4.9	抽出フォント毎の 1 位認識率 (置換法)	60
4.10	抽出辞書にない文字	62
4.11	従来の辞書の認識率 (社説 21~30)	63
4.12	1 位認識率 (辞書の更新による実験)	64
4.13	学習文書数毎の 1 位認識率 (マルチテンプレート法)	65
4.14	学習文書数毎の 1 位認識率 (置換法)	66



# 第 1 章

## 序論

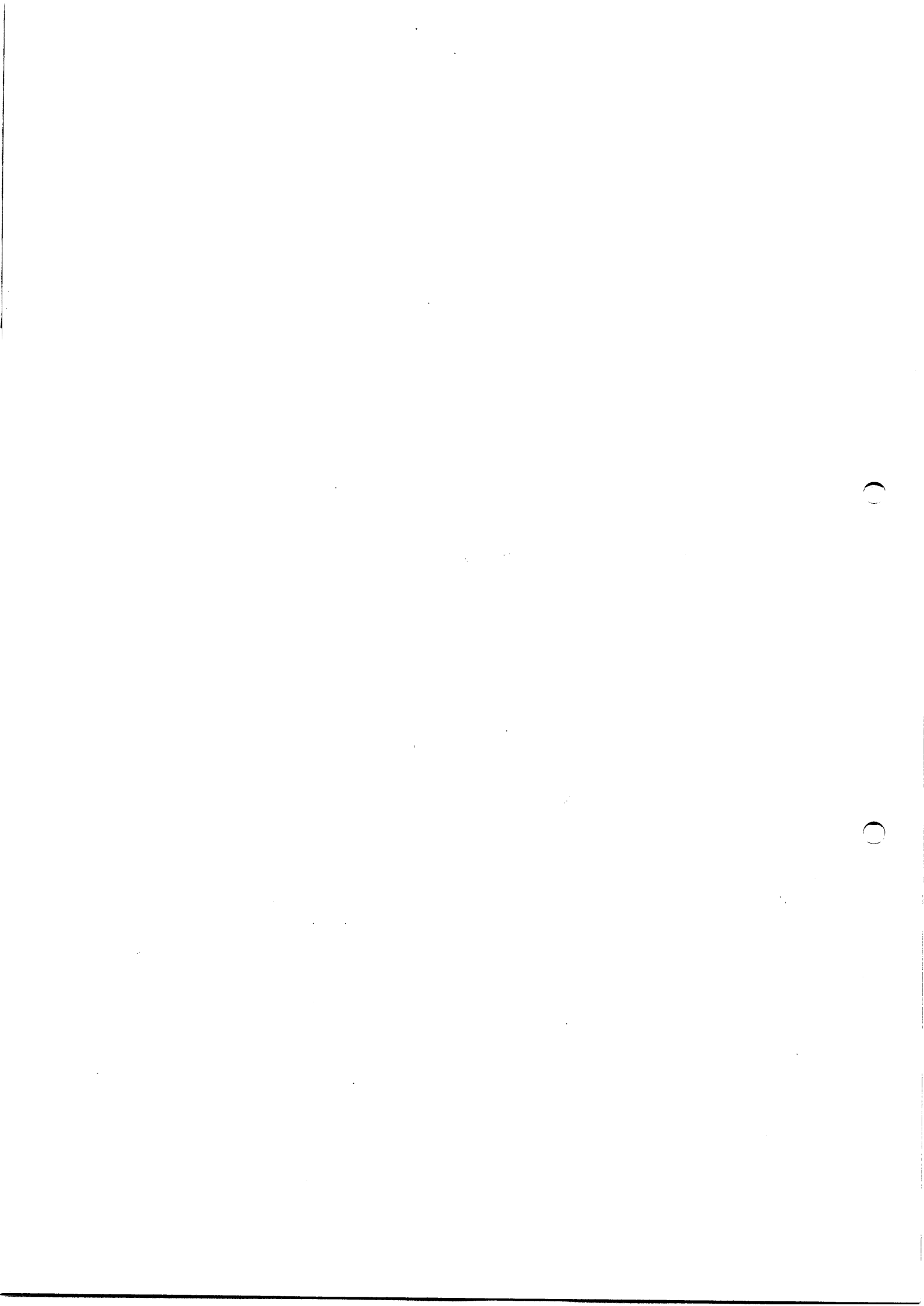
### 1.1 本研究の背景及び目的

各種の文字認識システムの研究は 30 年余り昔から始まり、十数年におよぶ基礎研究が積み重ねられた後に、文字読取り装置 OCR(Optical character Reader) が開発された。開発当初は、読取り可能な字種は数字 10 種程度が限度で、郵便番号の自動読取り区分などにしか実用化できなかった。また OCR の価格も非常に高価で普及には至らなかった。しかし半導体技術の発達、とりわけ LSI 技術を基盤としたマイクロプロセッサの発明で、廉価な OCR が登場、以来 OCR は急速に普及していき、それにともない、認識技術も高度化し、現在では数千種にのぼる手書き漢字をも認識の対象とできるようにまで到達した。

さらに、近年では OA 化により、ワープロやファックス、コピー機等の事務機器が産業から家庭にまで広がっており、大量の文書が出回っている。このような大量の文書を保存するために文書そのものではなく、データとして保存するというデータベースは、現在でも百科事典などで実用化されているが、さらに幅広い分野にわたって普及と発達が見込まれている。紙面上の文字として印字してある情報を、この様なデータベースに登録するにはコンピュータに入力しなければならない。この入力の際に、人間がキーボードから入力するのではなく、OCR を用いて自動的に入力させることできれば作業の効率化をはかれる。今後の情報化社会において OCR を用いた文書入力、人間とコンピュータとを結ぶマン・マシンインターフェースとして、今まで以上に重要な役割を担うことが期待されている。

[1][2][3]

しかし、現状では認識率 100% の文書認識システムは達成されておらず、どうしても誤認識を生じてしまう。実際の文書認識においては誤認識があるかどうかは、人間が確かめなければならない。さらに誤認識がある場合には、修正に大変な労力をさかななければならない。



このため、文書認識の自動化が大変に期待されている。

自動化を実現する際、現在の文書認識で問題になるのは、その文書のフォントが、既に登録されている辞書のパターンからかけはなれている場合である。認識方法として一般的に用いられているパターンマッチング法では、辞書(標準パターン・代表ベクトル)に依存する部分が多い。このため、このような場合には単純にパターンマッチング法だけでは正解を選出することができず、認識結果の上位の数候補に関して詳細に識別しなければ正確な認識を行うことができない。また、詳細識別を行うことでプロセス等は複雑になり、単純パターンマッチング法のみで自動化するよりも難しくなる。そこで、認識方法で高精度化をはかるのではなく、その文書のフォントを抽出し、辞書を作成するという辞書の高精度が考えられる。辞書を新しく作成する方法は同じフォントの文書を大量に認識する、というような場合に有効であると考えられる。<sup>[8]</sup>

本研究では既存辞書を基に、認識対象となる文書からフォントを抽出し、新辞書を作成すること、その辞書を用いてフォントを抽出した文書を再認識させ誤りを自動的に修正することなど、抽出したフォントを用いて文書認識の高精度化をはかる方法について実験し検討する。

## 1.2 本論文の構成

本論文の構成は次の通りである。

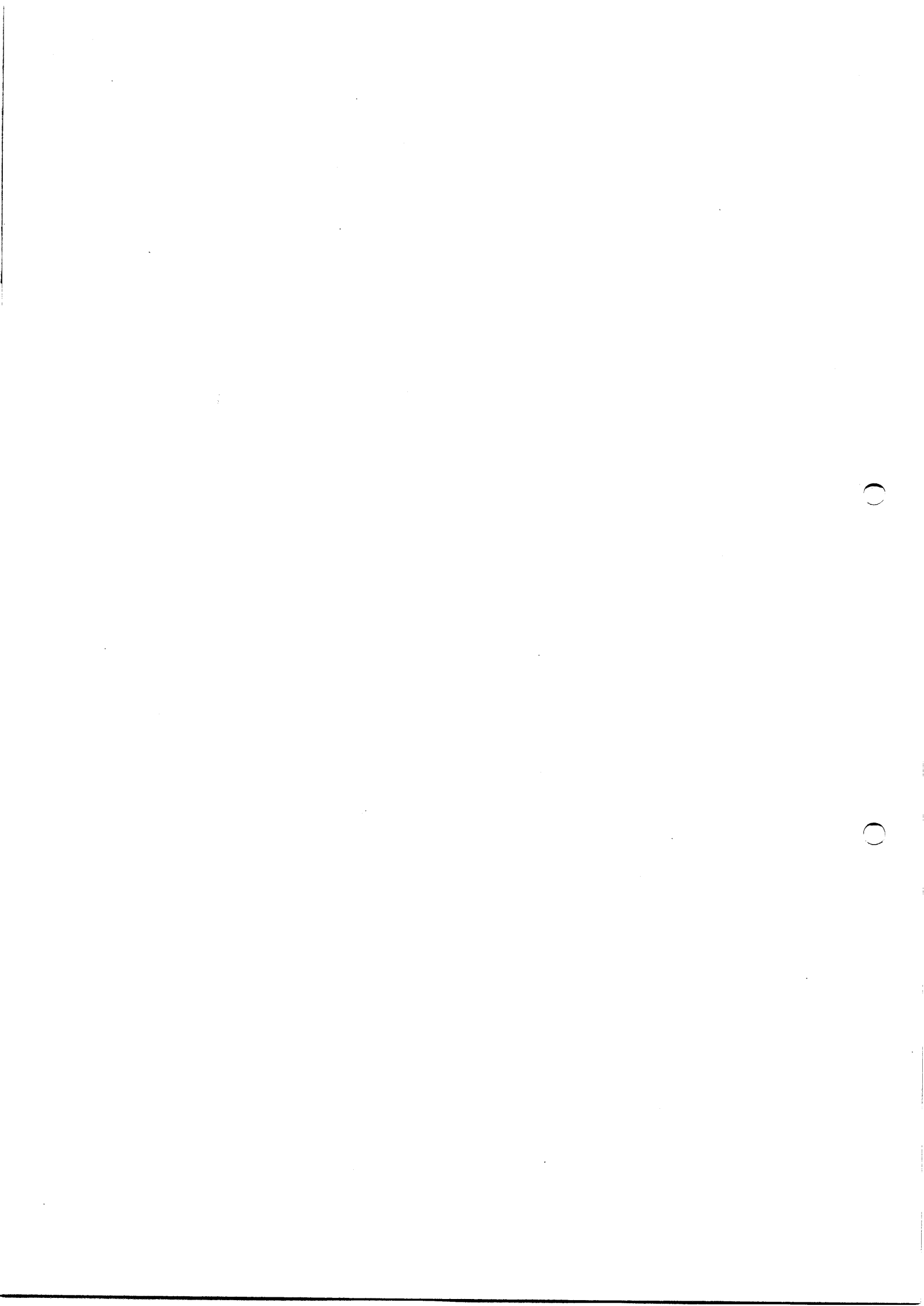
第1章 序論であり、本研究の背景と目的を述べる。

第2章 本研究で用いる特徴量や、文書認識システムについて述べる。

第3章 文書中から文字のフォントを抽出するアルゴリズムについて説明し、フォントの抽出実験を行い、結果について検討する。

第4章 3章で抽出したフォントより辞書を作成して認識実験を行い、再認識での誤りの自動修正、新辞書の有効性などについて検討する。

第5章 本研究の結論・今後の課題について述べる。





## 第 2 章

# 文書認識システム

### 2.1 まえがき

本章では本研究で提案する文書認識システムの概略、システムで用いる文字認識の手法について説明する。

### 2.2 文書認識システムの構成

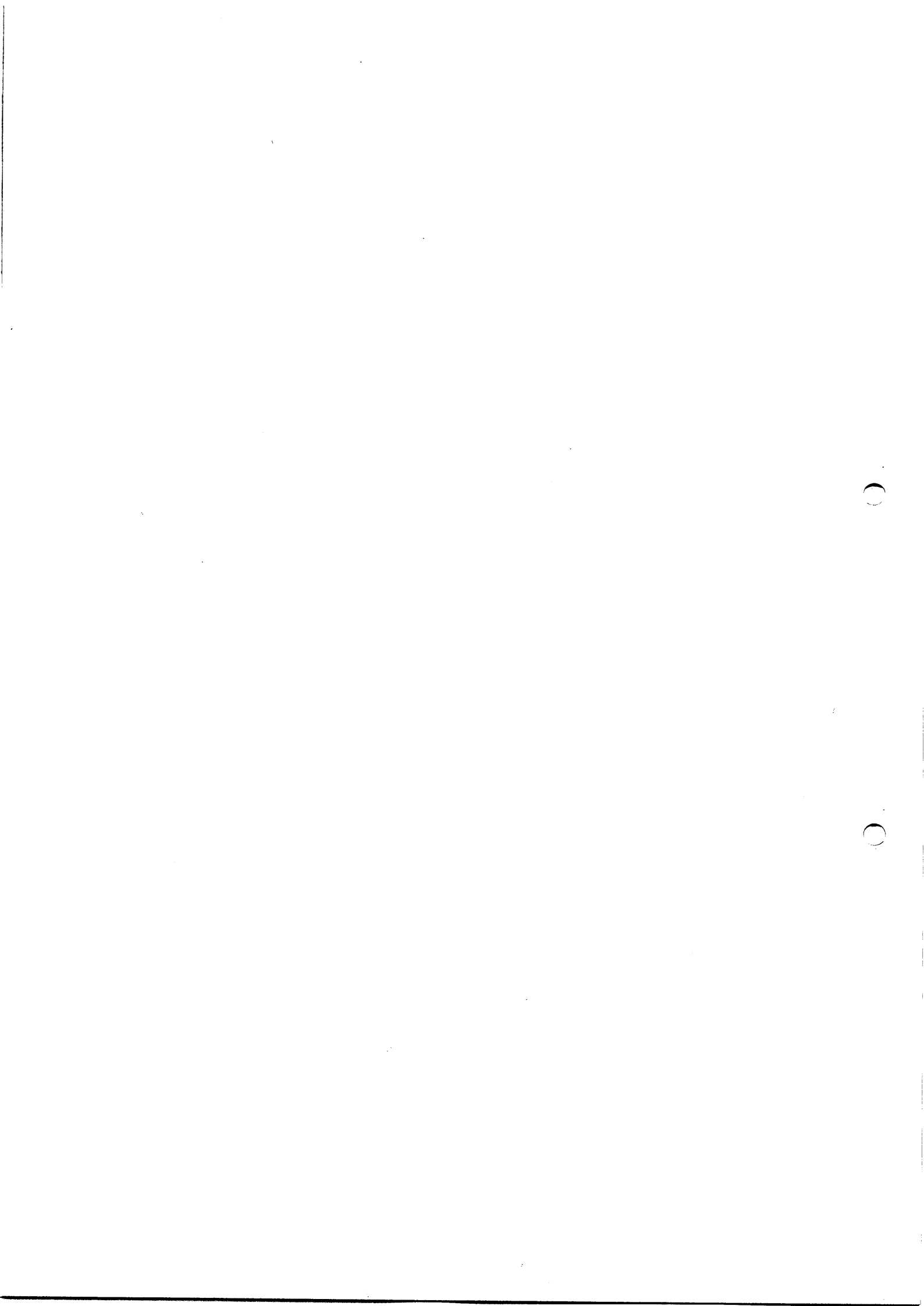
認識率を向上させるためには、パターンマッチング法だけでなく、他の認識方法と併用するなど、認識手法を複雑にすることは有効である。しかし、認識対象からフォントを抽出することが可能で、そのフォントを基に標準パターンを新しく作成できれば、1つの認識手法のみ(パターンマッチング法)を用いても良い認識結果が得られるはずである。つまり、認識手法を改善・改良するのではなく、辞書を改善することによってパターンマッチング法だけを用いても認識結果は向上できるのではないかと、ということである。図 2.2に示す本研究で提案する文書認識システムでは、この点に着目し対象文書からのフォントの抽出を中心に据えている。この認識システムでは辞書を抽出した文書を再認識し、誤認識を自動修正することを目的にしているが、この抽出した辞書を、他の文書の認識に用いることで、パターンマッチング法のみで従来法よりも認識精度が向上することも期待している。

- 認識

認識では予め用意しておいた辞書を用いて、認識対象文書を認識させる。このとき認識結果としては1位の候補だけでなく、2位・3位…と複数の結果を出力する。

- フォント抽出

文書中の全ての文字について認識結果を基に、フォントを抽出する。フォント抽出の



アルゴリズムは3章で詳しく説明する。

- 新辞書作成

抽出されたフォントから、新たに標準パターンを作成する。本研究では、この新標準パターンを既存辞書に加える方法のマルチテンプレート法と、もともとの標準パターンと新標準パターンとを置換する置換法の2つの方法で新しい辞書を作成している。

- 再認識

再認識では新辞書を用いる。

## 2.3 文字認識の方法

文字認識のアルゴリズムには、大きく分けて次の二つの方法がある。

- パターンマッチング法

- 構造解析法

パターンマッチング法は、パターン同士の重なりぐあいで評価し、認識を行う。このため、文字の多少の変形やノイズには強いが、類似文字の識別は難しい。また、計算機上で容易に高速で実現できる等の特徴をもつ。

これに対し後者の構造解析法は、線分の接続関係や位置関係などの文字構造に着目し、構造の類似性で認識を行う。この手法は類似文字の多い漢字の認識や、個人のくせが顕著である手書き文字の認識に有効であるとされているが、特徴量の定義や抽出が難しく、また計算機上での処理に時間がかかる等の問題も多い。

以後、本研究で用いるパターンマッチング法について説明する。図 2.1に文字認識のアルゴリズムを示す。

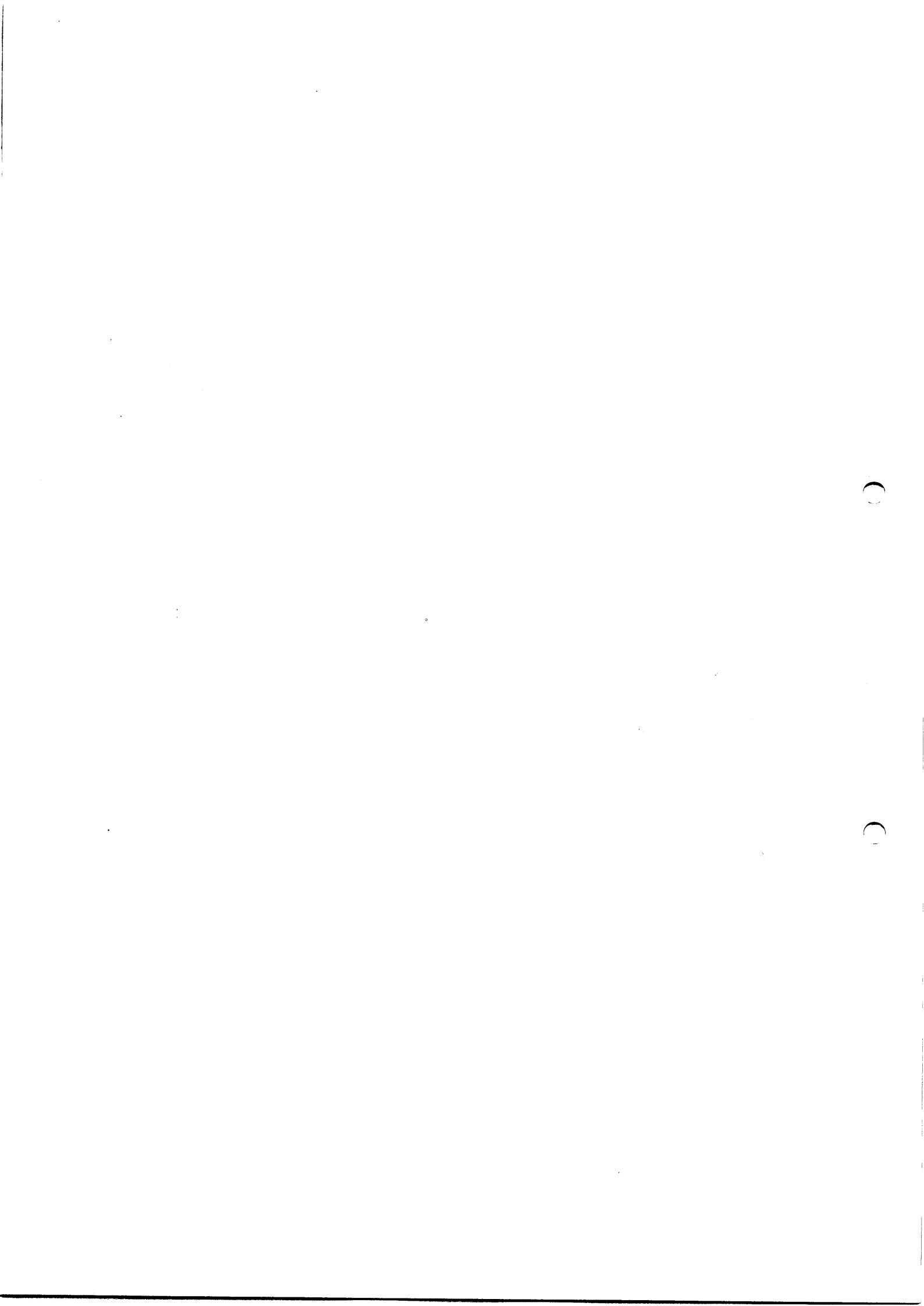
### 2.3.1 文字入力・前処理

- 文字入力

認識される文字データは文書の形でイメージスキャナによって入力される。出力されたイメージデータは、切り出し処理により、各文字毎に切り出される。

- 前処理

切り出された文字データには、実際にマッチングに使われる特徴量を求めるために前



処理が施される。前処理は、ノイズ除去・スムージング、正規化、細線化、線素化の 4 つの処理によって構成されている。

- ノイズ除去・スムージング

入力画像 (切り出されたデータ) は印刷時のプリンタの状態やイメージスキャナなどの性能によって、ノイズや線分上の凹凸が少なからず発生している。これらは認識精度に悪い影響を及ぼすため、まず入力画像に対してノイズ除去を行う。本研究では  $2 \times 2$  ドット以下の孤立点をノイズと見なし除去することにする。スムージングは、周囲 8 近傍を参照し、中心のドットを決定する。実際には  $3 \times 3$  ドットのマスクを用いて行う。

- 正規化

正規化は、もとの入力イメージ画像を一定の大きさ (本研究では  $64 \times 64$  ドット) に拡大又は縮小する処理であり、これによって文字の大きさや領域の位置ずれによる影響を吸収する。正規化の手法には大きく分けると、線形正規化と非線形正規化の二つの手法に分けられるが、本研究では線形正規化を用いている。

- 細線化

細線化は、通常数ドットの太さを持った文字の線分を 1 ドットの幅に変換する処理である。この処理によって、文字の線分幅の違いを吸収する。本研究では Hilditch の細線化アルゴリズム<sup>[5]</sup> を用いる。この方法は、各画素においてその周りの  $3 \times 3$  のマスクを用いて、着目するドットの近傍の 8 ドットの連結状態を見ながら細線化していく方法である。

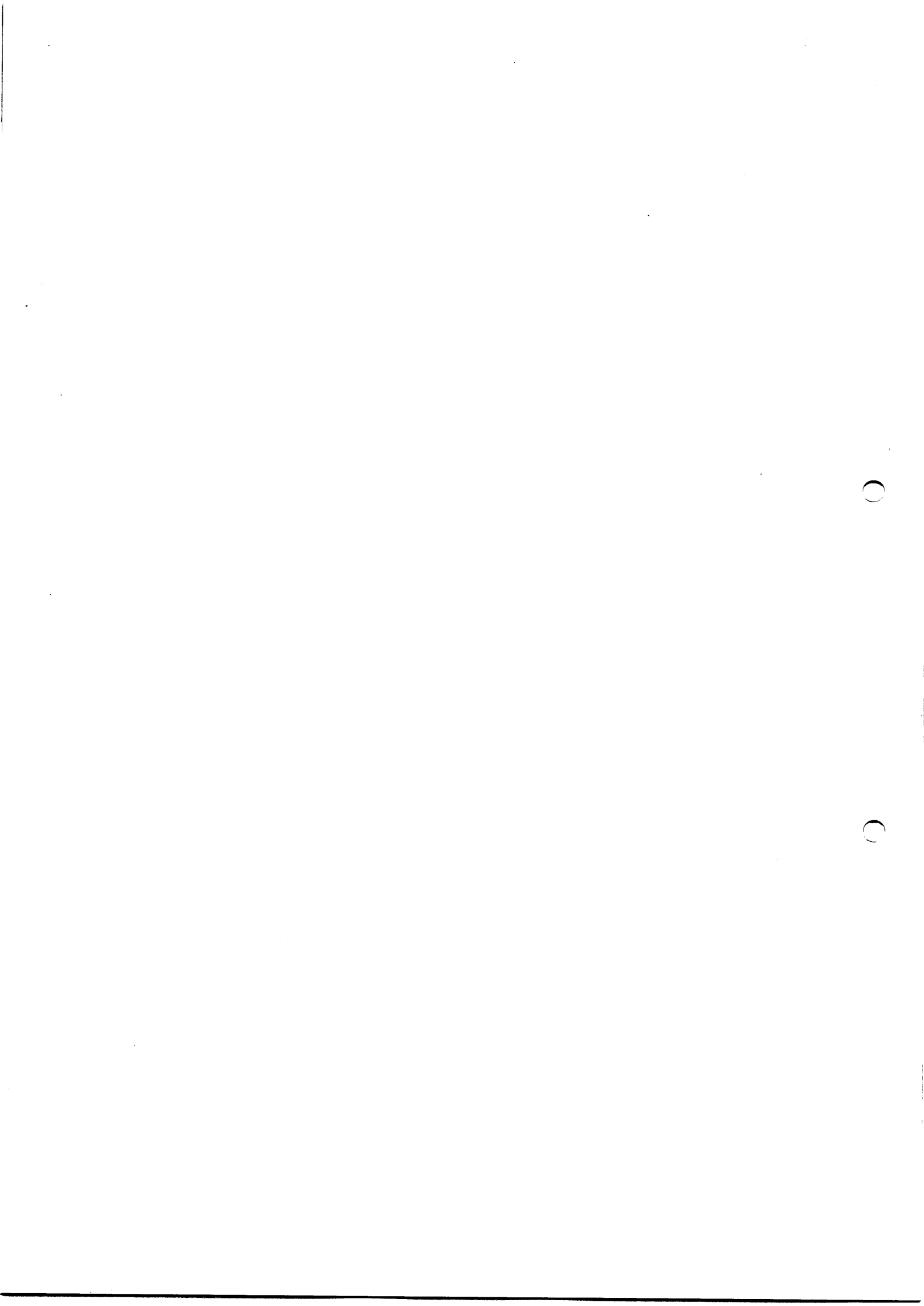
- 線素化

線素化は、細線化された図形の各画素について、黒画素であればその周囲  $3 \times 3$  の小領域を参照し最も自然な方向と考えられる 4 つの縦「|」、横「—」、斜め  $45^\circ$ 「/」、斜め  $135^\circ$ 「\」、のうちの 1 つの線素に対応させる。

前処理の各処理による効果を図 2.3 に示す。

### 2.3.2 特徴抽出

パターンマッチング法では、処理の高速化、パターン分離の効率化などのために、文字パターンを特徴量という数値ベクトルに変換する。この課程を特徴抽出という。以降は本研究で使用する方向線素特徴量<sup>[4]</sup>、について説明しておく。



- 方向線素特徴量

方向線素特徴量<sup>[4]</sup>の抽出法は図 2.4 のように、まず、 $64 \times 64$  ドットの線素化文字を 8 ドット間隔に縦横を分割する。次に、左上から  $16 \times 16$  ドットを半分ずつ重複させて 49 領域(左上を 0 とし、左から右、上から下へ順に並ぶ)ができる。1 領域のベクトルは 4 方向のカウンタからなる 4 次元ベクトルでできてあり、1 領域内の重みは図 2.4 の下図のようになっている。よって 1 文字当りの次元数は  $N = 196(49 \times 7)$  次元のベクトルとなる。

### 2.3.3 認識

ここでは最も一般的な手法である全数整合法を説明する。全数整合法は、各文字の標準パターンと特徴抽出で得られた未知入力文字の特徴量で評価値(距離)を求め、近いものから順に識別する方法である。アルゴリズムの簡潔性やある程度の高い認識率を得られることから、一般的に用いられている。認識で用いられる辞書・評価値について以下で説明する。

- 辞書の作成方法

パターンマッチング法では、入力されたパターンと標準パターンとを比較し、一番似ている(距離的に近い)標準パターンのコードを結果として出力する。この標準パターンは辞書と呼ばれ、予め全ての登録文字に対して用意しておかなければならない。辞書作成法としては、

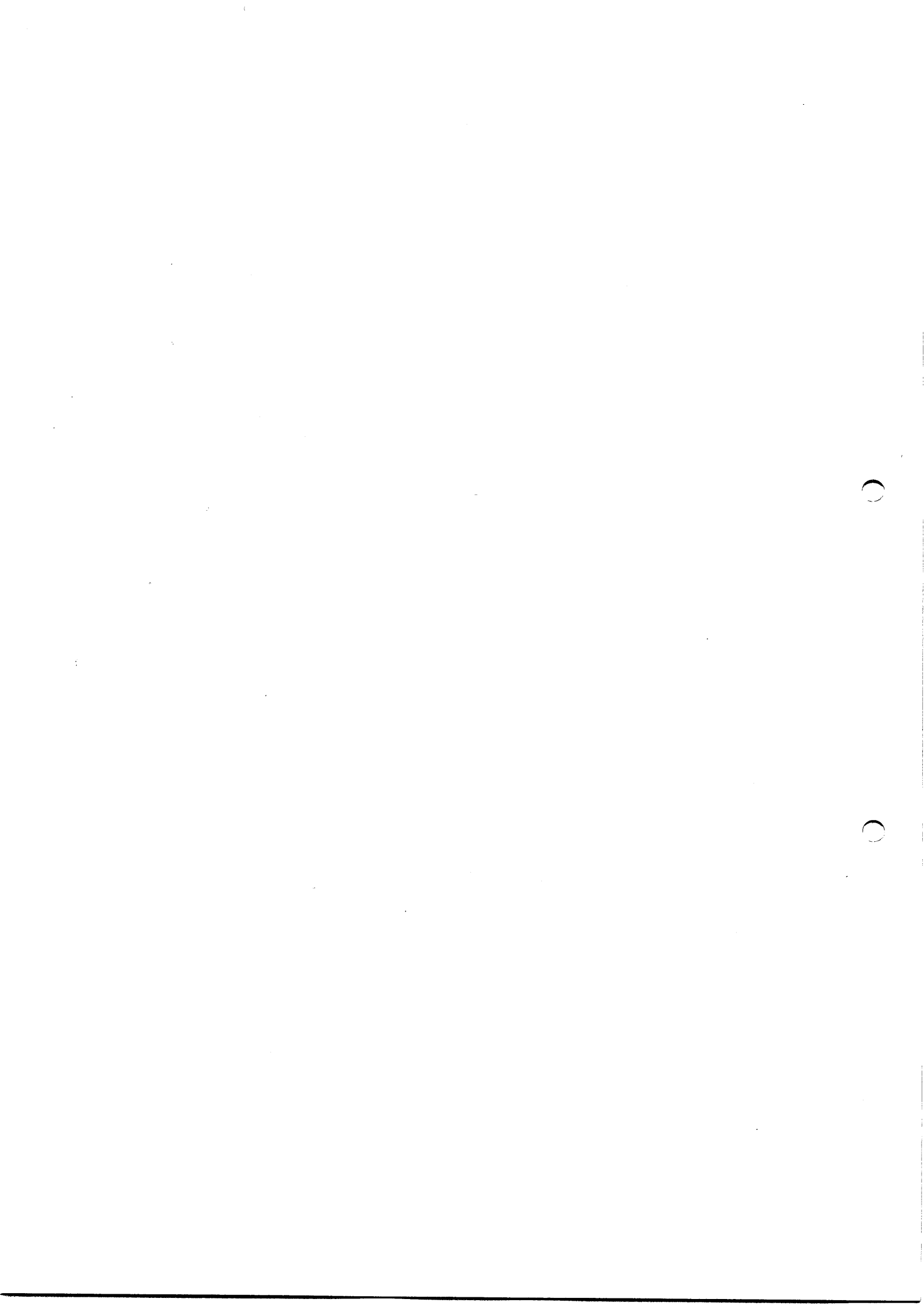
- 文字  $k$  のサンプルパターンの集合  $S(k)$  の重心ベクトル

が一般的であり、本研究ではサンプルパターンの集合として 42 種のプリンター出力を用いて作成した重心ベクトルを辞書(既存辞書)として使用している。

- 評価値

未知入力文字の特徴量と標準パターンとの類似性や整合性の評価は、評価値により行われる。評価値の例としてよく使用されるものとして、ユークリッド距離、重み付きユークリッド距離<sup>[6]</sup>がある。ここで、 $x$ 、 $u$  をそれぞれ未知入力ベクトルと、ある文字の標準パターンベクトルとすると、それぞれ次のように定義される。

ユークリッド距離：





$$d_e(x, u) = \sum_{i=1}^n (x_i - u_i)^2 \quad (2.1)$$

重み付きユークリッド距離：

$$d_w(x, u) = \sum_{i=1}^n w_i (x_i - u_i)^2 \quad (2.2)$$

ここで  $w_i$  は重み付きユークリッド距離の重み<sup>[7]</sup>で、各次元ごとに次のように定義する。

$$w_i = A \frac{1/v_i}{\sum_{j=1}^n 1/v_j} \quad (2.3)$$

但し、 $A$  は定数、 $v_i$  は標準パターン作成用データセットにおける各次元の分散、 $n$  は次元数である。本研究ではこの 2 つを使用する。この 2 つの他にシテイブロック距離

$$d_c(x, u) = \sum_{i=1}^n |(x_i - u_i)| \quad (2.4)$$

があり、計算時間が速い等の特徴をもつ。

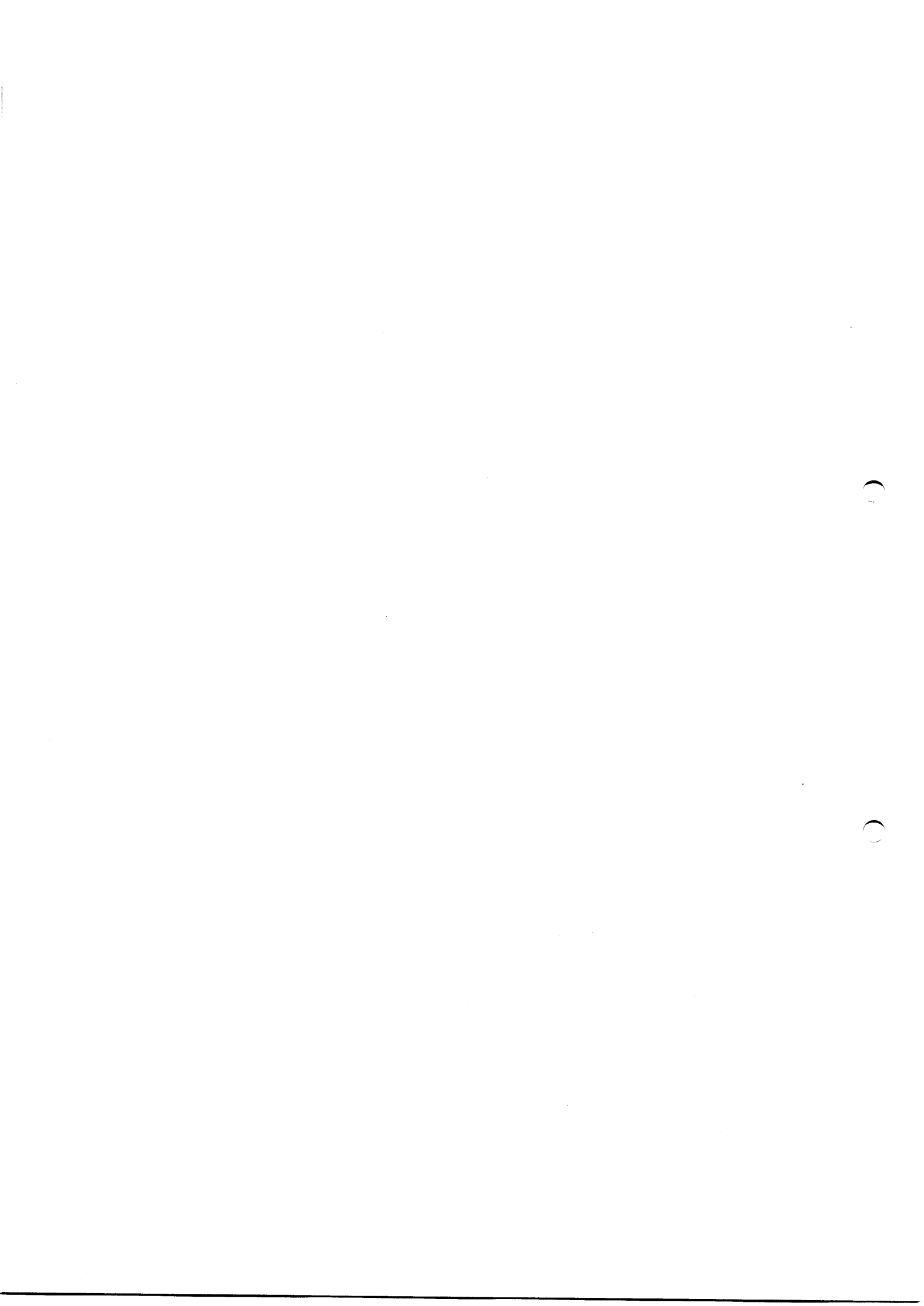
距離以外の評価値としては、2 つの特徴量の内積から定まる類似度<sup>[2]</sup>が使われている。

$$E(u, x) = \frac{u \cdot x}{|u||x|} = \cos \theta \quad (2.5)$$

である。

### 2.3.4 候補出力

評価値計算をした後評価値の小さい順に、第 1 位候補、第 2 位候補、…、第  $N$  候補を出力する。



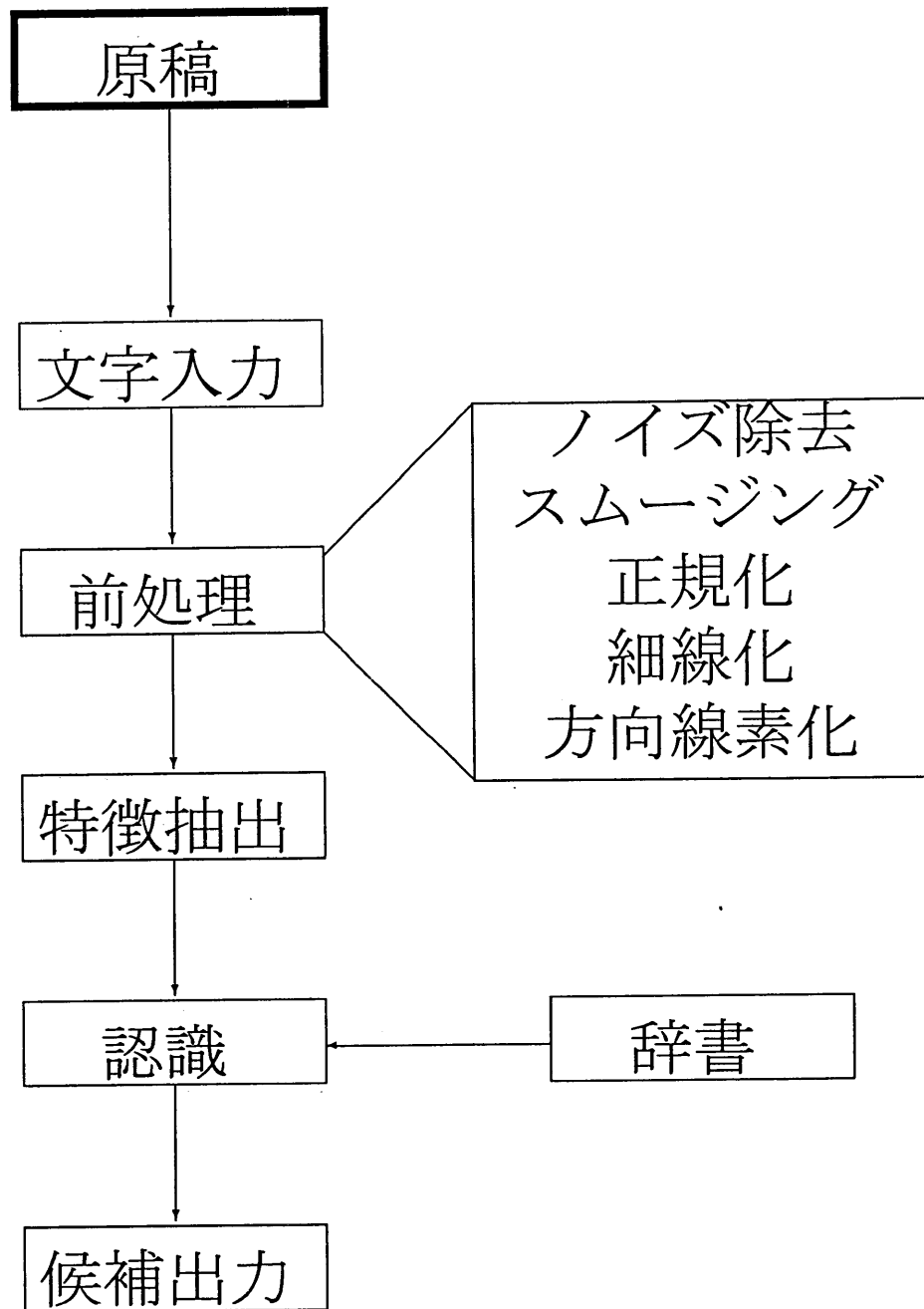
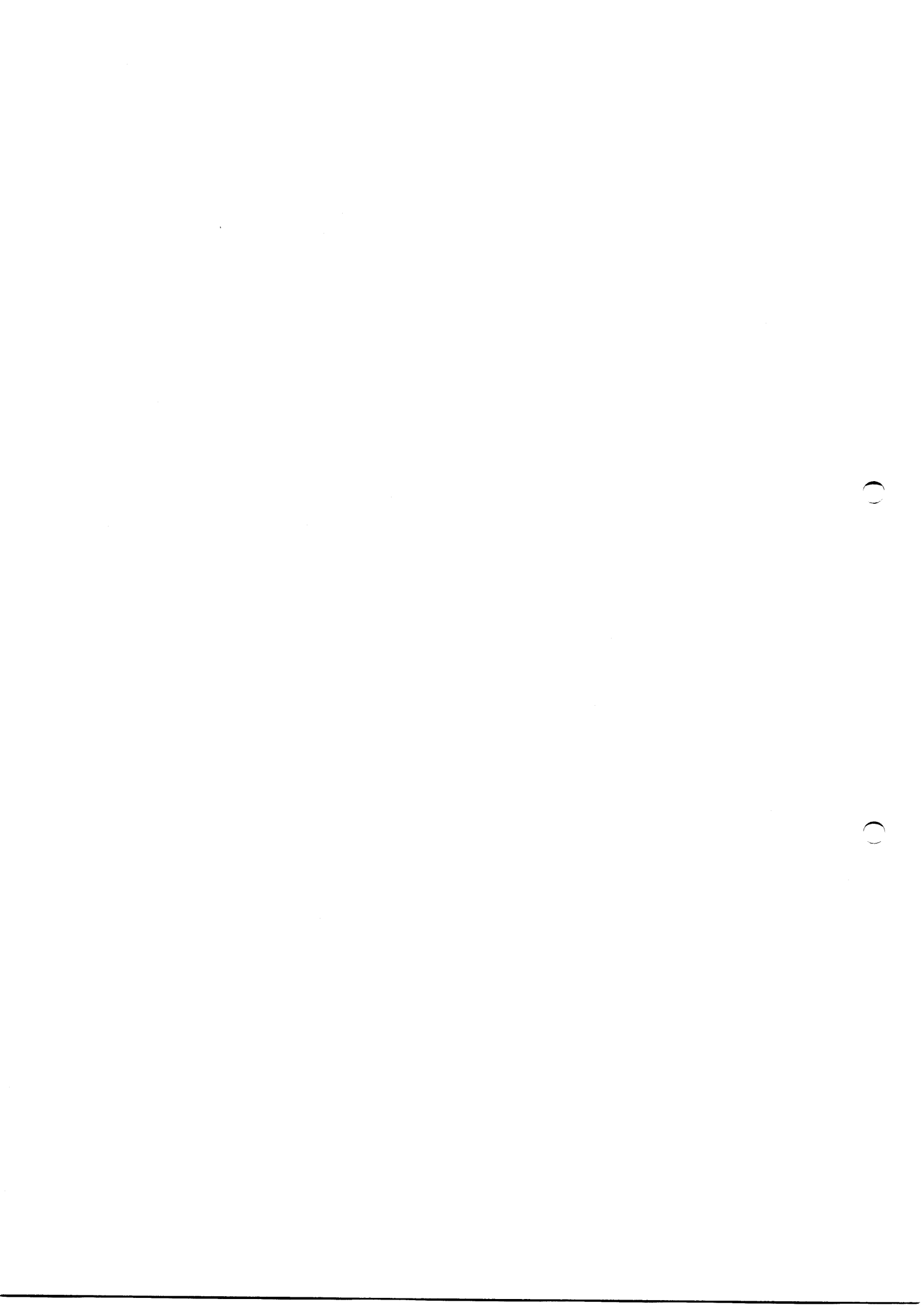


図 2.1: 文字認識アルゴリズム



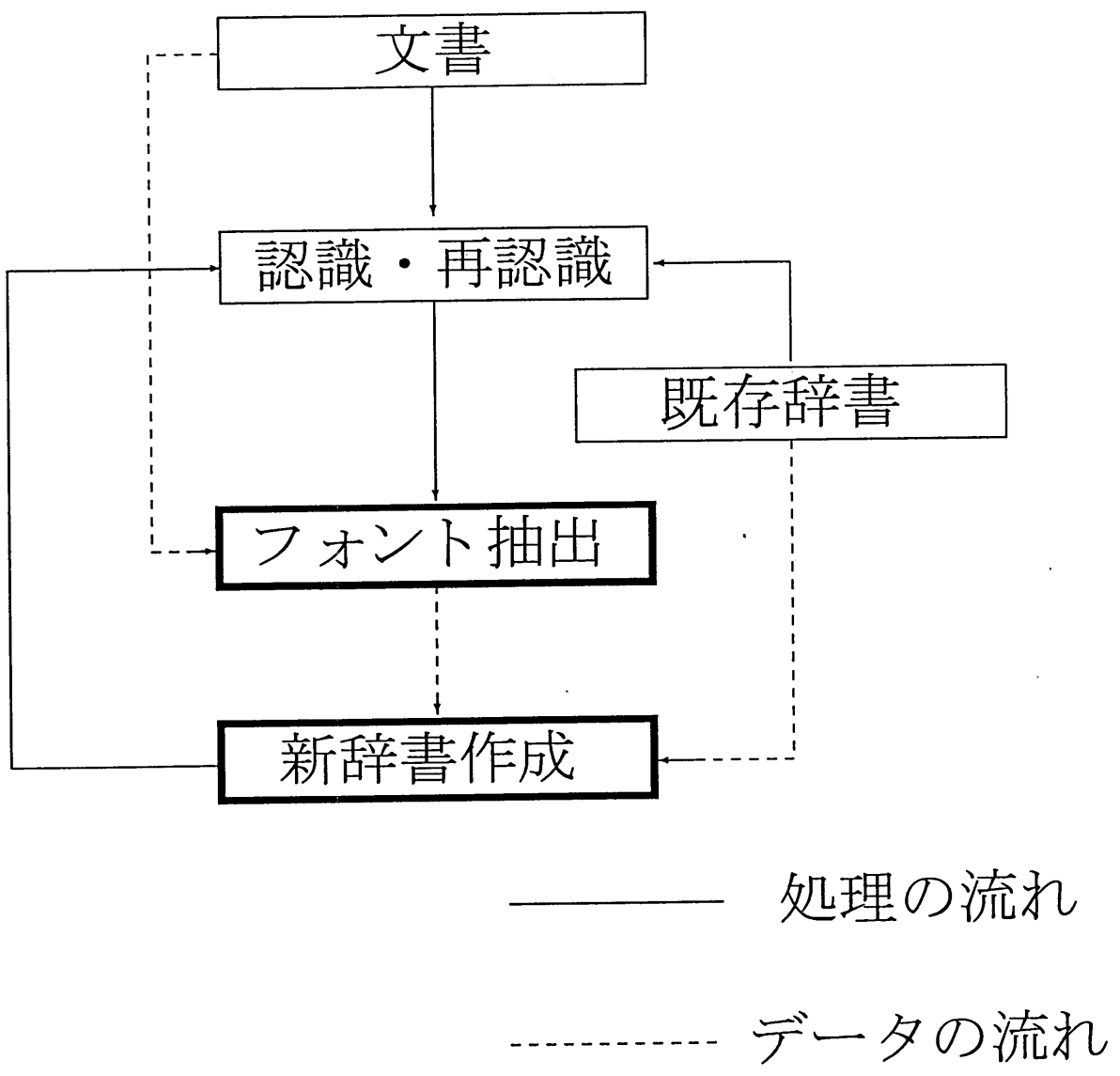
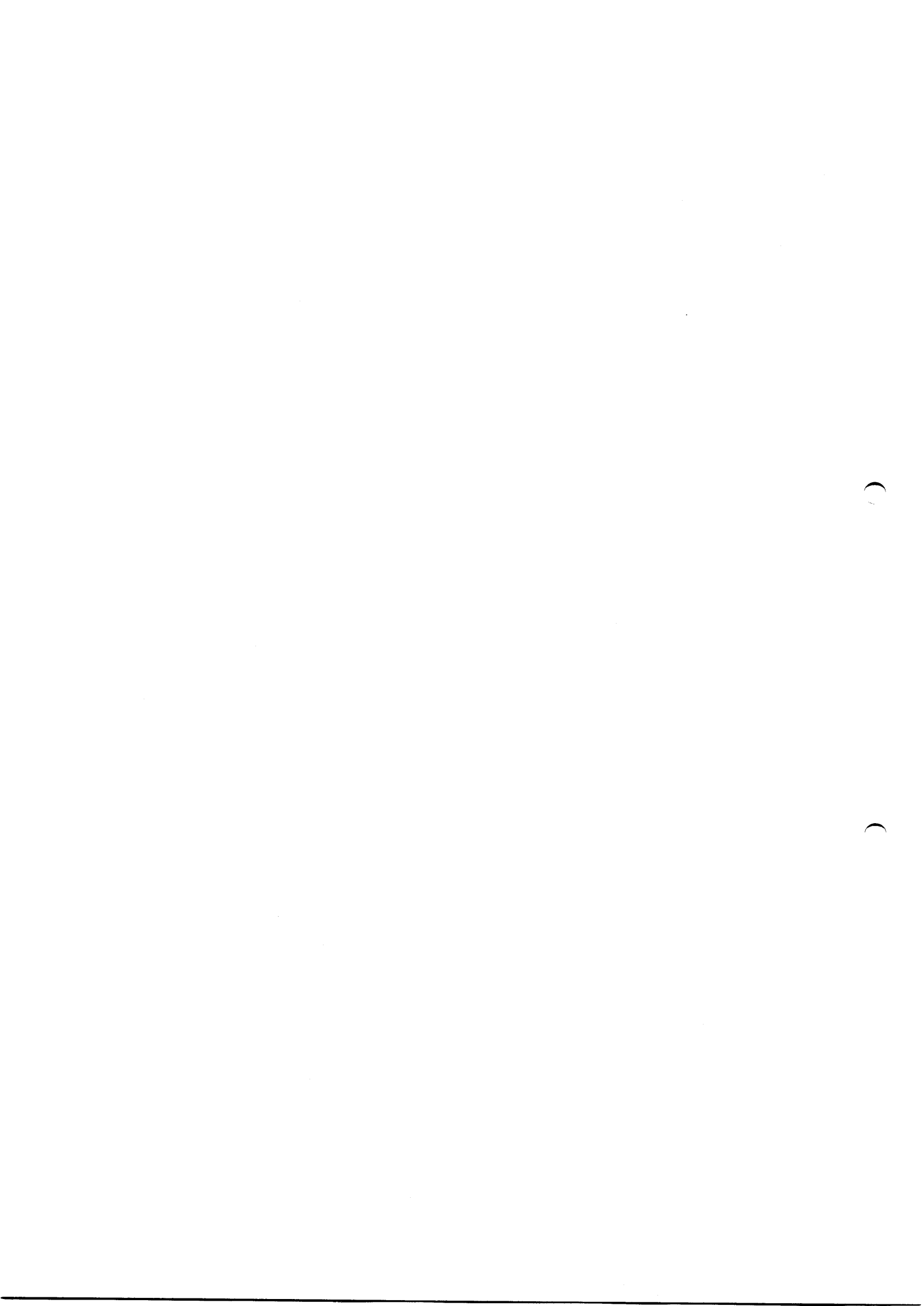
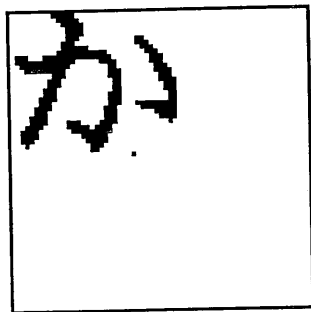
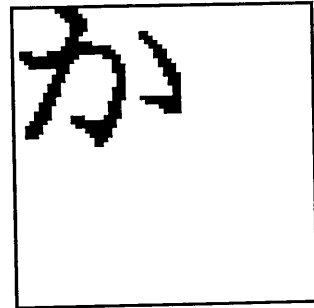


図 2.2: 文書認識システムのアルゴリズム

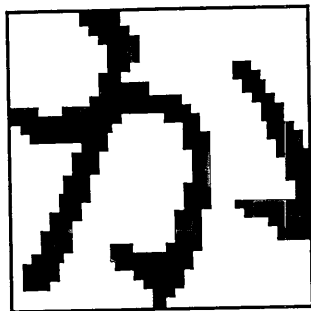




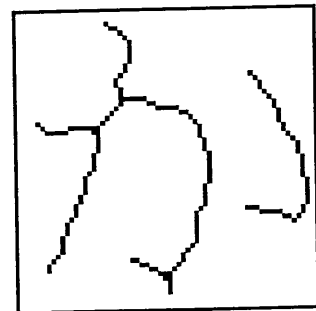
入力画像



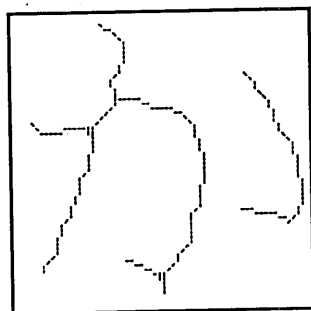
ノイズ除去・スムージング



正規化



細線化

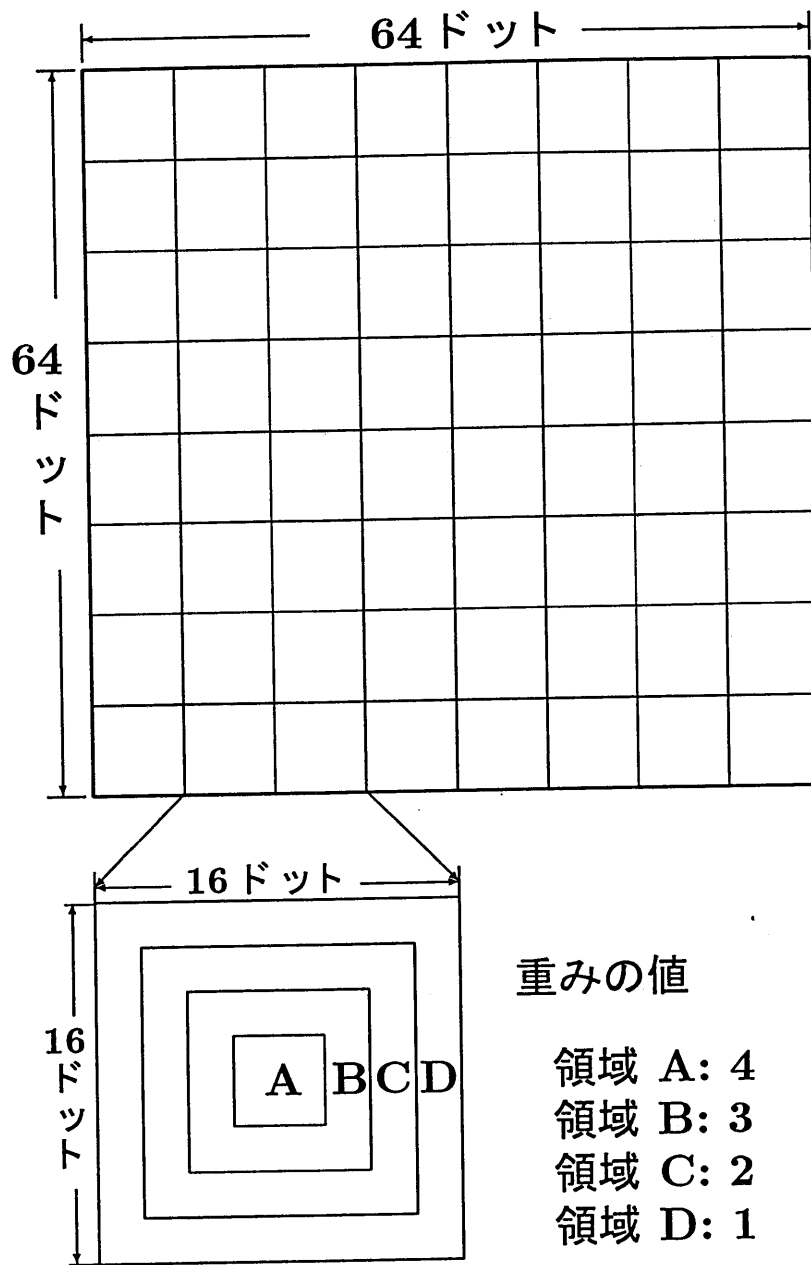


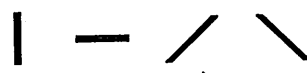
線素化

図 2.3: 各処理後のイメージ





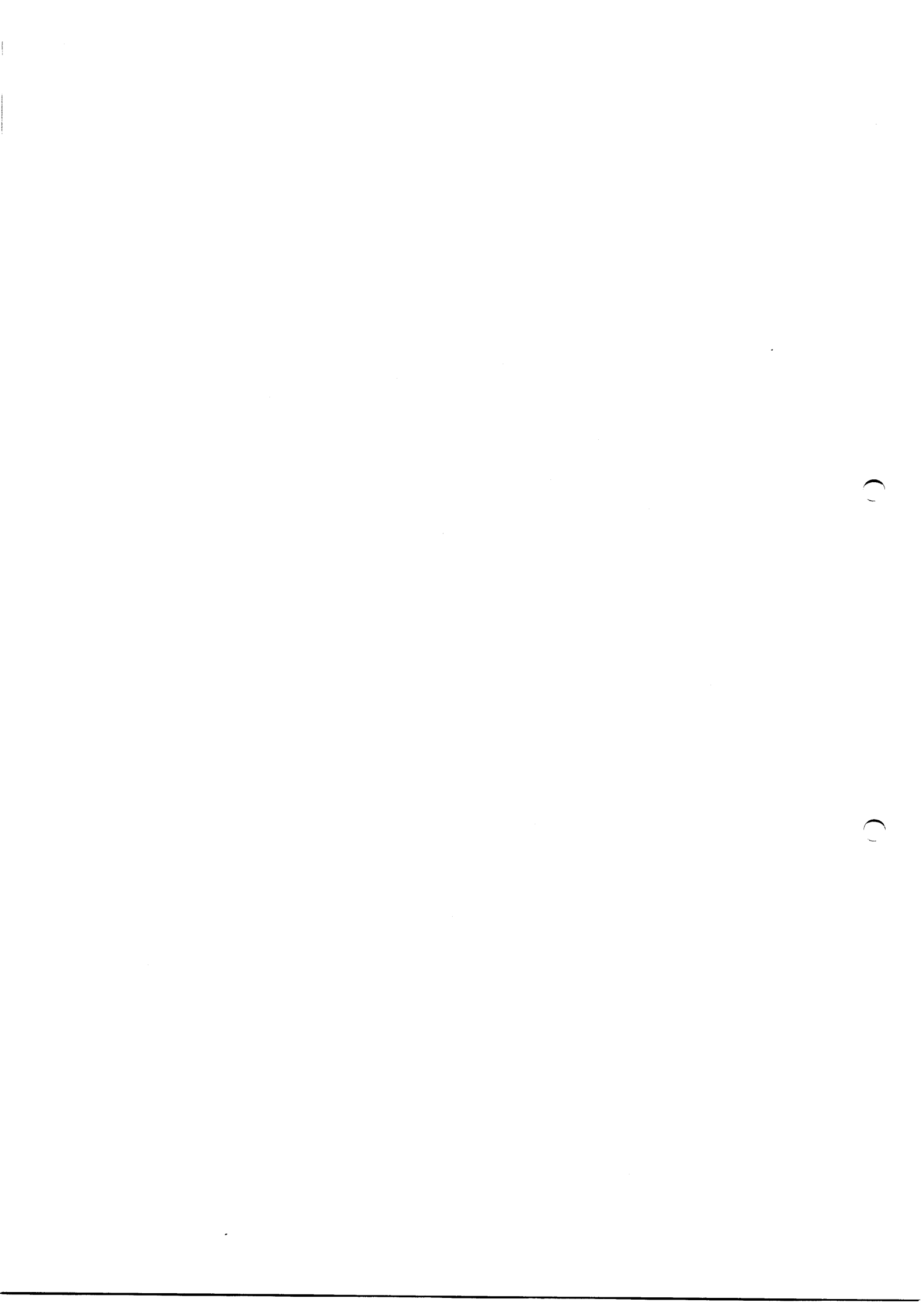


方向線素: 

方向線素特徴量:  $V_m = (V_{m_1}, V_{m_2}, V_{m_3}, V_{m_4})$   
 where  $m = 1 \sim 49$

小領域の数:  $49(7 \times 7)$   
 方向線素特徴量の次元数:  $196(4 \times 49)$

図 2.4: 方向線素特徴量



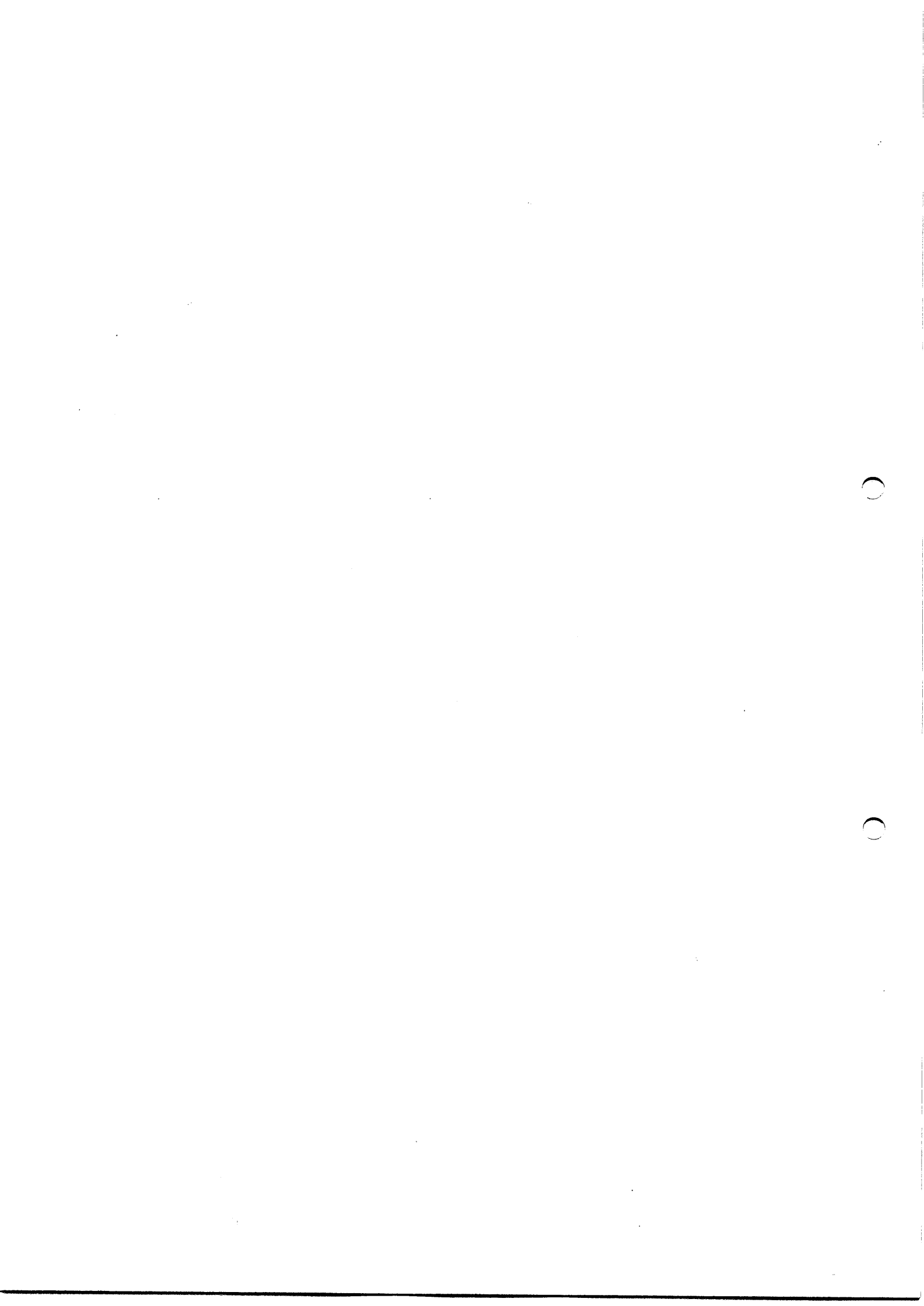
## 第 3 章

# フォントの抽出

### 3.1 前書き

この章以降で、抽出すべきコードと同じコードのフォントのことを正解のフォントと表し、抽出すべきコードとは違うコードのフォントを不正解のフォントと表すことにする。本文中から未知のフォントを正確に抽出することは容易ではない<sup>[9]</sup>。特に次に挙げる 2 つの点が問題になる。

1. 正解 (抽出すべきコード) のフォントを集める。
  2. 正解のフォントだけを抽出する。
- 前者の場合は、認識結果を用いて正解のフォントを集める際に、認識結果が 1 位のコードに対応しているフォントだけを集めた場合には、現在誤認識しているものからは、コードとフォントが正しく対応していないフォントが正解のフォントとして集められてしまい、2 位以降にある真の正解を集めることができない。このため入力されたフォントに、2 位以降の複数の候補のコードと対応させたからもフォントを集めてこなければ、全ての正解のフォントを集めることができない。
  - 次に上記の様に 1 つの未知フォントに複数のコードを割当てて、フォントを集めた場合、集まったフォントは正解のフォントだけとなることはほとんどない。正解のフォントは含まれていても、その他には不正解のフォントも含まれているはずである。正解のフォントがその集合の大多数を占める場合であれば、集めた全てのフォントのベクトルの重心をとり、重心から離れているようなベクトルを持つフォントを取り除くことで、正解のフォントのみを抽出することは可能である。(図 3.1 の成功例) しかし、



出現頻度が低い文字や、類似している文字が多い文字の場合は、正解のフォントがごく少数の場合や、正解のフォントと、類似している文字のコードのフォントが同数ほどの場合には、この方法では正解のフォントが全て取り除かれたり、不正解のフォントも同時に抽出されたりする可能性がある。図 3.1での失敗の例1では、抽出したいフォントは「か」であるが、その母集団の中に「か」に類似している「が」のフォントも多く含まれるため、「か」と「が」が抽出されてしまっている。また、失敗例の2では「お」が母集団の中の少数なため、誤って「も」が抽出されてしまう場合がある。この様にフォントを抽出する母集団の構成により、抽出結果が変わってしまうため、この重心を用いた手法だけではフォントを正確には抽出できない。

そこで、本研究では図 3.2で示すアルゴリズムのフォント抽出法を提案する。

また、今回フォントを抽出する対象としては、事前に新聞社説(1~20)の構成する字種について調査し(A章参考)、100文字以下で、社説中の全ての文字の50%弱を占める、ひらがなとし、その中から3.2.5節の除去候補文字の学習用に用いた新聞社説(1~10)の10部中に出現した65種とした。(表 3.3) 実験で用いた新聞社説は図 3.5に例示する。

## 3.2 フォントの抽出

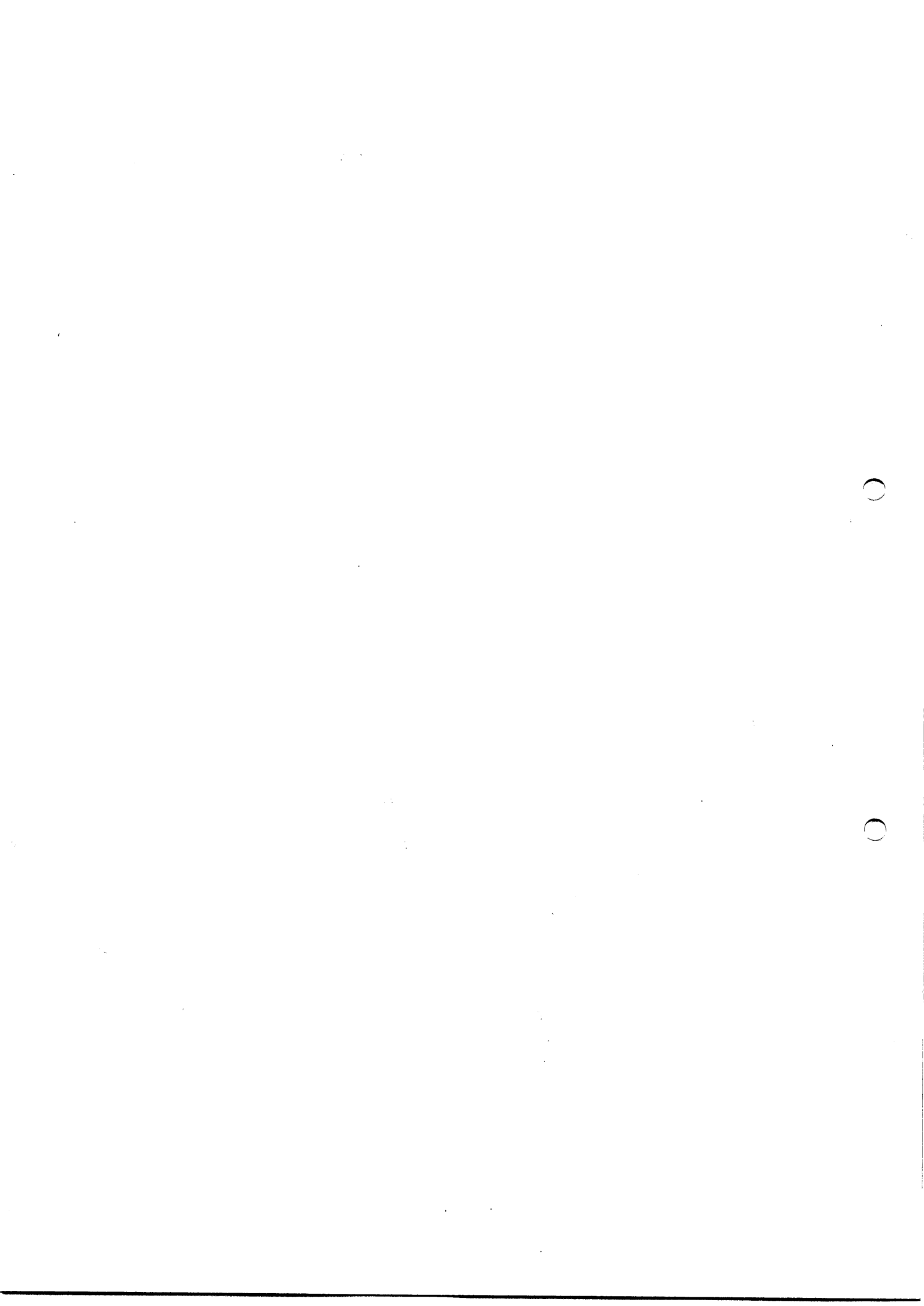
図 3.2の各処理について以下に詳しく説明する。

### 3.2.1 認識・結果出力

認識の処理では、予め準備しておいた既存辞書を用いて、未知文書の各文字を全数整合法により認識させる。誤認識する文字を救済するため1位の候補だけでなく下位候補まで複数の候補を正解として出力する。ここで何位まで正解として出力するかを候補数と定義する。通常の認識と同様に、1位だけを正解とする場合は候補数が1であり、5位候補まで正解とするならば候補数は5となる。

認識に用いた辞書の文字数は、JIS第一水準の漢字2965文字と、ひらがな・カタカナのうち小文字「っ」「ぁ」等を取り除き、さらに人間でも字体からだけでは、ひらがな・カタカナか判断がつかない「へ」「べ」「ぺ」の3文字に関しては、同一のものとして扱い、ひらがな73文字、カタカナ71文字の合わせて3109文字である。辞書に含まれるひらがなとカタカナを表 3.4に示す。

この処理を文書中の全ての文字に行い、各文字に対して、候補数の複数の正解を与える。通常の認識では入力されたフォント(文字、パターン)と、その入力フォントに一



番近い辞書中の標準パターンとのコードとを1対1で対応させているが、ここでは入力フォント1つに対して、複数のコードを対応付けることになる。実際に図3.6と図3.7(a)の例で説明する。まず、図3.6中の「A」、「B」、「C」はそれぞれ字種の異なる文字の入力フォントを表し、「A」と「A'」は同じ字種の文字の異なる入力フォントを表している。右側の1位~4位までの欄は、そのときの上位の認識結果である。通常の認識では入力フォント(「A」)に対し、1位認識結果である「は」というコードを1対1で対応させるが、今回用いる手法では入力フォント(「A」)に対し、「は」だけでなく、「け」、「ほ」等複数のコードを対応付ける。(図3.7(a))この対応付けた複数のコードの中に、フォントのコードと同じコードがある(フォントとコードが正しく対応している)と考える。

### 3.2.2 リストの作成・候補集合の作成

認識・結果出力は文書中の各文字毎に処理を行ったが、以降の処理では、結果出力で対応付けした各コード毎に行う。リストの作成では、まず出力したコードそれぞれについて候補数までの認識順位で、そのコードを出力した入力フォント(文字)と対応付ける。次に出力したコード毎にクラスタリングし、リストを作成する。図3.6から対応関係を表したのが図3.7(b)になる。「は」というコードを認識時に出力したのは、入力フォント「A」と「C」、「A'」であるので、これをコード「は」と対応させる。同様に、コード「ば」に対しては、入力フォント「B」と「C」を対応させる。

候補集合の作成では、リストに従って各コード毎に実際の入力フォントを集めてくる。この各コードに集まったデータを候補集合と定義する。この候補集合が正解のフォントを抽出するための母集団になり、認識する際の候補数が、この候補集合の大きさを決定する。図3.8に、実際に候補集合として集められた入力フォントを示す。これは、コード「け」に対して候補数を3にしたときの例である。この候補集合中には、「は」、「け」、「ゆ」の3文字のフォントが集められているが、不正解のコードである「は」、「ゆ」のフォントを除去するためには次の除去候補文字の処理を行う。

### 3.2.3 除去候補文字の処理

認識において複数の正解を与えるため、候補集合中には抽出すべき正解のフォントのみではなく、不正解のフォントも含まれている。特に文書中で出現頻度が低いコードの場合、候補集合中のデータのほとんどは不正解で、正解のフォントは少数しか含まれていない場合もある。この様に他種のフォントが混在している中から正解のフォントのみを抽出するには、候補集合中に存在する不正解のフォントを除去しなければならない。この候補集合

0

0



中に存在する除去しなければならないフォントのコードを除去候補文字と定義する。除去候補文字は各コード毎に固有であり、その数もコードによって異なる。

除去候補文字の大部分は、フォントによって変わるのではなく、その文字の字体に依存しているため統計的に決定できる。しかし、一部はフォントの違いで大きく異なるため事前に、そのフォントでの除去候補文字を学習しておくことが必要になる。除去候補文字の学習については3.2.5節で述べる。

この除去候補文字には、「け」と「ゆ」の様に明らかに異なる文字とわかるものと、「か」と「が」の様に類似度が大きく識別しがたい文字との2つがある。除去候補文字の処理もこの2つに合せて2段階で行う。

まず、明らかに異なる文字の場合は、パターンマッチング法で誤認識した場合の、正解と誤認識との評価値の性質を用いて削除することができる。

いま、「A」というコードの候補集合中に「A」、「B」の2種のフォントが集められたとする。ここで「A」、「B」のフォントの入力ベクトルと、それぞれの標準パターンとの評価値の関係は次式の様になると考えられる。

$$E(u_A, v_A) < E(u_B, v_A) \quad (3.1)$$

$$E(u_A, v_B) > E(u_B, v_B) \quad (3.2)$$

ただし、 $v_A$ 、 $v_B$ 、 $u_A$ 、 $u_B$ はそれぞれ、A、Bの入力ベクトル、標準パターンを表し、 $E(u_A, v_A)$ は、「A」の入力ベクトルと「A」の標準パターンとの評価値を表すとする。

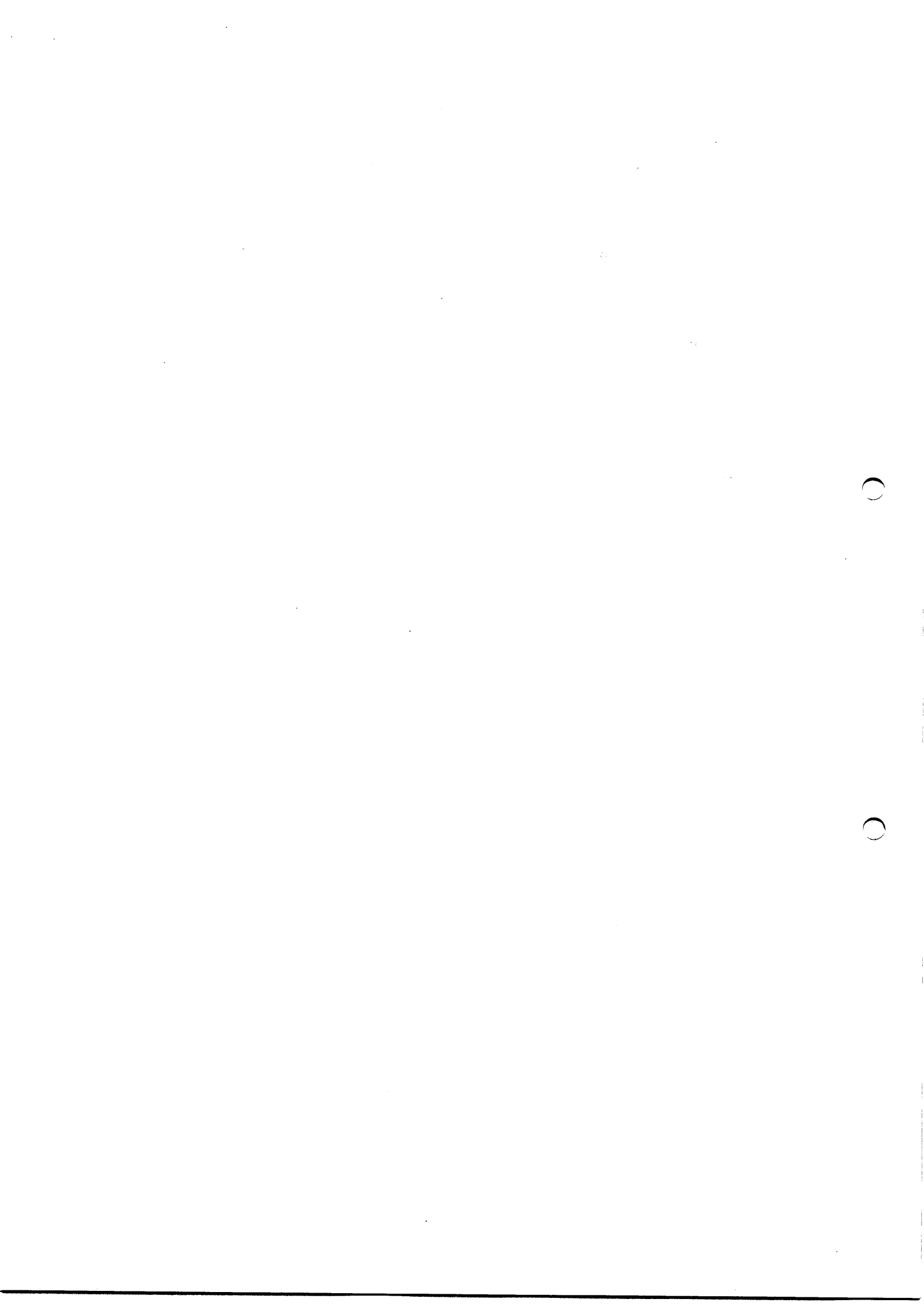
ここで、不正解である「B」の入力ベクトルは、「A」の標準パターンより、「B」の標準パターンに距離が近いと考えられる。この性質を用いて、予め各コード毎に誤って候補集合中に選出されるような文字(除去候補文字)を学習させておき、除去候補文字のコードの標準パターンと、正解のコードの標準パターンとを、候補集合中の未知フォント(ベクトル)とマッチングさせ、未知フォントが、正解のコードの標準パターンよりも、除去候補文字のコードの標準パターンに近ければ、その未知フォントは、正解のコードのフォントではないとみなして、除去を行う。

正解のコードを「T」、除去候補文字のコードを「 $x_n$ 」、候補集合中のあるフォントを $Y_m$ とすると、

$$E(u_T, v_{Y_m}) > E(v_{x_n}, v_{Y_m}) \quad (3.3)$$

のとき、フォント $Y_m$ を候補集合から除去する。

この処理で候補集合中の不正解のフォントを除去し、候補集合中のフォント数を、それに含まれる正解のフォント数に近づける。



### 3.2.3.1 詳細識別

次に類似度が大きく、識別しがたい文字の処理について説明する。「こ」と「ご」のように清音と濁音がある文字や、「し」と「り」のように全体の字形がベクトル的に似ている文字の識別の際には、文字間の距離は、通常の異なる2文字間のものよりも差異が小さいため、上記の方法では、正解かどうか識別が難しくなる。このような文字は詳細識別文字として新しく定義する。詳細識別文字と正解のコードとの識別は文字の画像領域全体ではなく、差異の大きい部分だけを用いて行う。

除去候補文字の処理、詳細識別については、3.2.5節の除去候補文字の学習で具体例を示しながら説明する。

### 3.2.4 細分類

候補集合中に、各コードの除去候補文字以外の除去すべきコードのフォントが存在していた場合には、上記の除去候補文字の処理だけでは除去できない場合が考えられる。そのようなフォントを除去するために次の示す細分類の処理を行う。

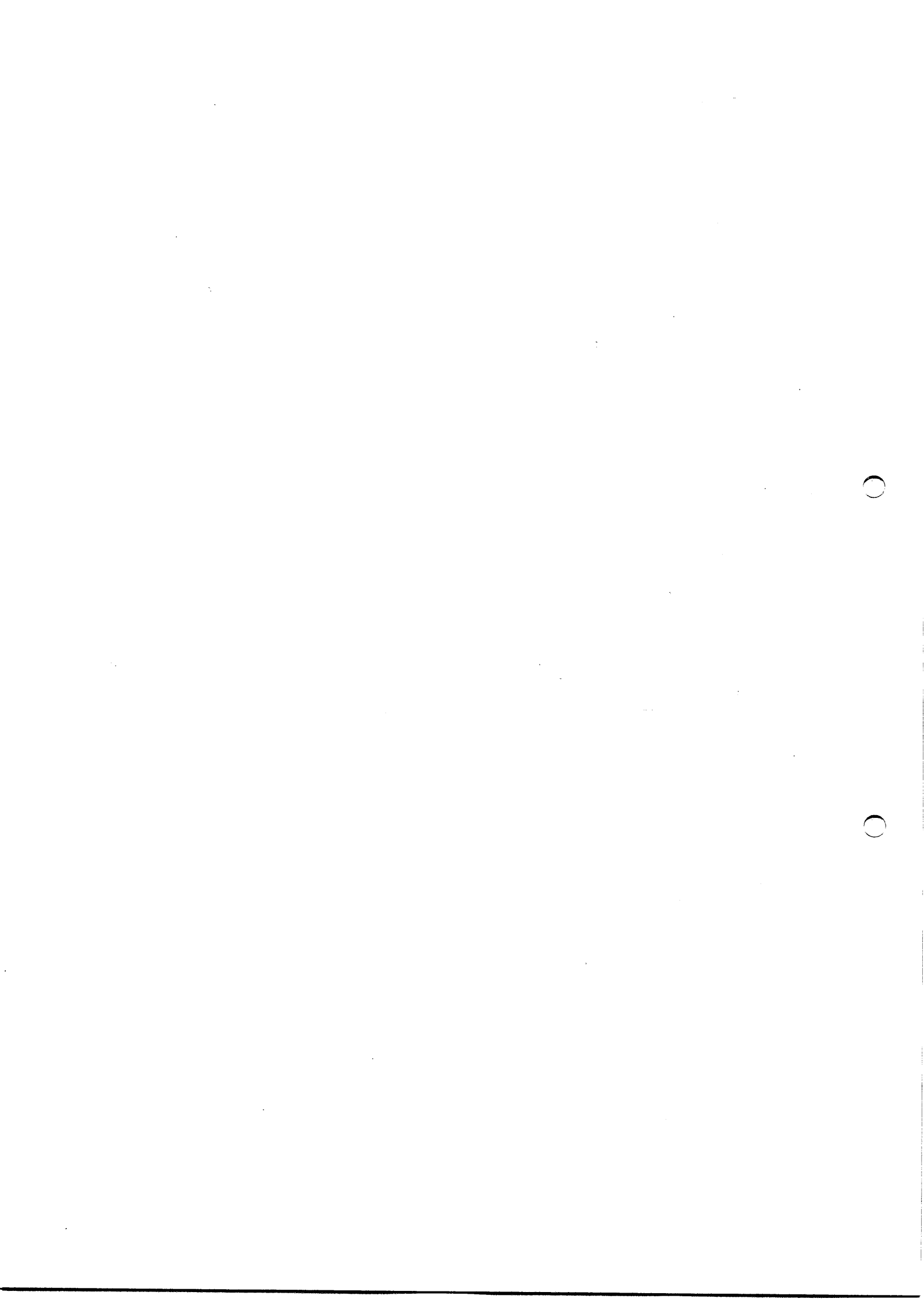
1. 候補集合中の全てのフォント (ベクトル) から重心ベクトルを作成する。
2. 全てのフォントベクトルとこの重心ベクトルとマッチングを行い、ユークリッド距離が、個々のコード毎に予備実験で定めた閾値よりも大きくなれば、そのフォントを候補集合から削除する。
3. 候補集合中のフォントで再び重心ベクトルを作成し、マッチングを行う。
4. これを収束するまで繰り返す。

なお、各コード毎の閾値は次のように定めた。あるコードの標準パターンを  $U$ 、その学習サンプルパターンを  $P$  とし、 $U$  と  $P$  のユークリッド距離 ( $dis$ ) の平均 ( $x$ ) と分散 ( $\sigma^2$ ) を求める。このときの平均と偏差の和をそのコードの閾値 ( $threshold$ ) とした。

$$U = \frac{1}{N} \sum_{j=1}^N P_j \quad (3.4)$$

$$dis_n = \frac{1}{N} \sum_{j=1}^N (P_j - U)^2 \quad (3.5)$$

$$x = \frac{1}{N} \sum_{j=1}^N dis_j \quad (3.6)$$



$$\sigma^2 = \frac{1}{N} \sum_{j=1}^N (dis_j - x)^2 \quad (3.7)$$

$$threshold = x + \sigma \quad (3.8)$$

この細分類の処理を終えて候補集合中に残っていたフォントは、そのコードのフォントとして抽出する。

### 3.2.5 除去候補文字の学習

本研究で用いた除去候補文字は表 3.1、表 3.2に示すが、これらは次のようにして学習させた。

- データとして新聞社説(1~10)の10部を使用し、候補数を10として今回提案したアルゴリズムを用いて、各コード毎に候補集合を作成する。図 3.9は「み」のコードの場合で、候補集合中には「ん」、「み」、「入」、「ふ」、「ぶ」の5字種が合せて112個集まった。
- この候補集合中に存在した正解以外のコードを、仮の除去候補文字として決定する。例では「ん」、「入」、「ふ」、「ぶ」の4字が仮の除去候補文字となる。この仮の除去候補文字を用いて、式(3.3)により候補集合中の不正解のフォントを取り除いていくが、仮の除去候補文字  $x_n$  が以下の式(3.9)、式(3.10)の条件を満たすときは、その  $x_n$  を除去候補文字から削除し、除去候補文字の削減・最適化を行った。式(3.3)中の評価値は、各コード毎にユークリッド距離(式(2.1))、重み付きユークリッド距離(式(2.2))のどちらかを使用した。(図 3.3)

$$\begin{cases} T(x_1, x_2, \dots, x_{n-1}) = T(x_1, x_2, \dots, x_n) = T_{in} \\ D(x_1, x_2, \dots, x_{n-1}) \geq D(x_1, x_2, \dots, x_n) \end{cases} \quad (3.9)$$

$$\begin{cases} T(x_1, x_2, \dots, x_{n-1}) = T_{in} \\ T(x_1, x_2, \dots, x_n) < T_{in} \end{cases} \quad (3.10)$$

ただし、 $T_{in}$ は最初の候補集合中の正解のフォント数、 $x_j$ は除去候補文字とし、 $T(x_1, x_2, \dots, x_{n-1})$   $D(x_1, x_2, \dots, x_n)$ はそれぞれ除去候補文字が  $x_1 \sim x_n$ のときの除去候補文字の処理を行った後の候補集合中の正解フォント数と、候補集合中のフォント数を表すものとする。式(3.9)は  $x_n$ を除去候補文字から外しても、残りの  $n-1$ 字の除去候補文字により不正解のフォントを除去できる場合である。図 3.9の場合は仮の除去候補文字の



うち「ぶ」が、除去候補文字では外れている。これは、「ぶ」を外しても残っている「ふ」によって不正解のコードである「ふ」と「ぶ」のフォントが除去できることを示している。

式(3.10)は $x_n$ を除去候補文字としたことで、不正解のフォントだけでなく、抽出すべき正解のフォントまで除去している場合を示している。図3.10では、仮の除去候補文字の「ざ」を用いて除去候補文字の処理をした場合に正解である「ぎ」も除去されてしまったことを示す。このような場合には、「ざ」を除去候補文字から外す。

ここで最適化された除去候補文字を $x_1 \sim x_N$ とすると、

$$T(x_1, x_2, \dots, x_N) = T_{in} \quad (3.11)$$

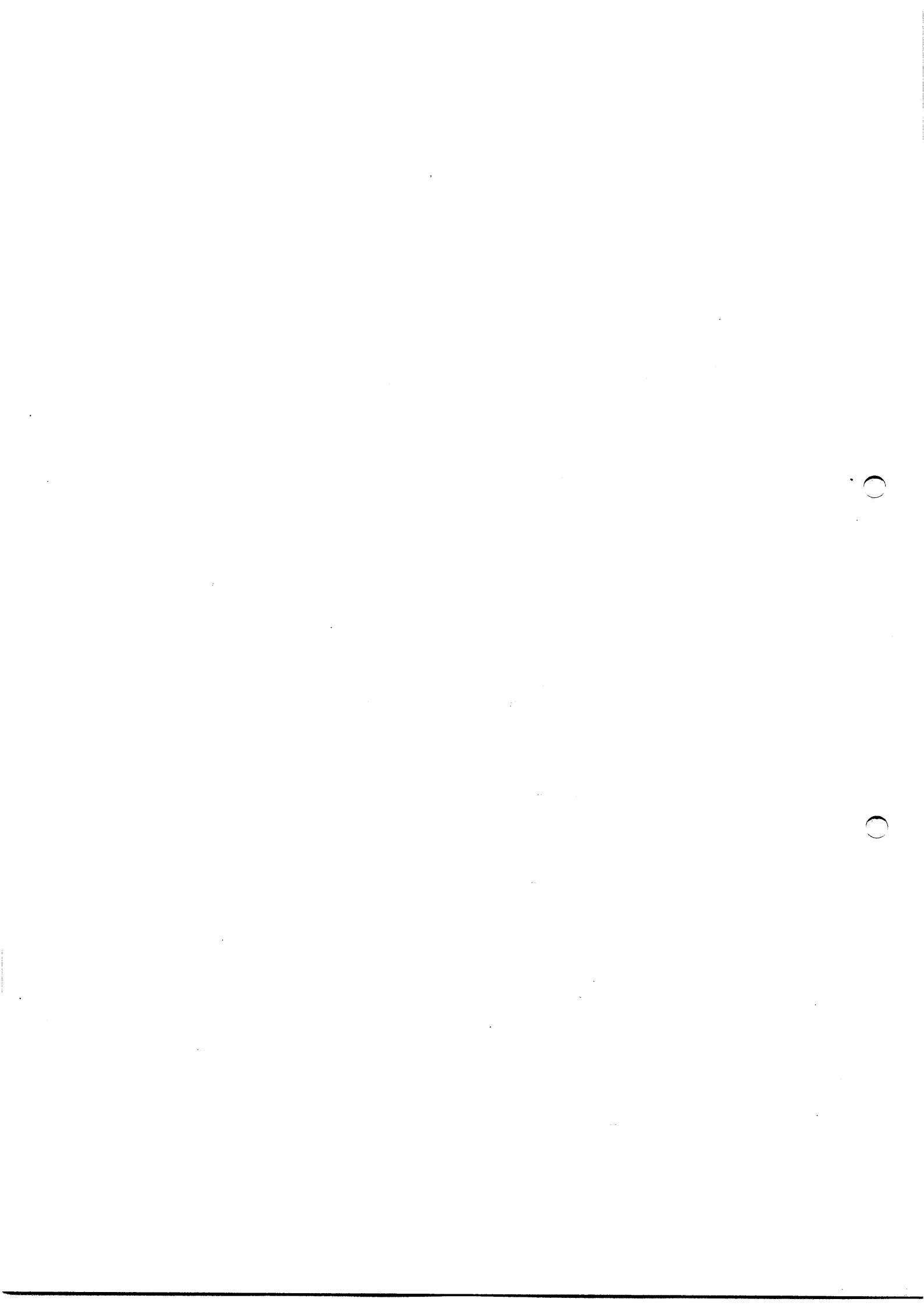
となることがフォントを正確に抽出するための条件となるが、式(3.10)の条件で除去候補文字から外された $x_n$ のフォントが候補集合に残った場合や、清音と濁音の文字の様な類似文字では、不正解のフォントを全て除去できない場合がある。そこで、この段階で候補集合中に不正解のフォントがある場合は、その不正解のフォントのコードを詳細識別文字として新たに詳細識別を行うものとする。(図3.10の「ざ」)

- 詳細識別文字と正解のコードとの識別・除去でも、式(3.3)を用いるが、マッチングを行う領域は図3.14で示される $64 \times 64$ ドットの文字の画像領域全体ではなく、差異が大きい部分のみでマッチングを行う。この部分を識別領域とする。2.3.2節で説明した方向線素特徴量の画像領域での49個の $16 \times 16$ ドットの小領域(マス)のうちの2マスを識別領域とする。通常のマッチングでは49マス全てを用いるため、ベクトル(方向線素特徴量)は196次元であるが、詳細識別では $2(\text{マス}) \times 4(\text{方向})$ の8次元となる。各コード毎に、2マスを次のように決定した。

候補集合中に正解「T」のサンプルが $N$ 個、不正解「F」のものが $M$ 個あり、それぞれ $t_1 \sim t_N$ 、 $f_1 \sim f_M$ とする。このとき、

$$E_m(u_T, v_{t_i}) > (<) E_m(u_f, v_{f_j}) \quad (3.12)$$

ただし、 $E_m(u_T, v_{t_i})$ は $m(0 \leq m \leq 48)$ マスにおける「T」の標準パターンと、「 $t_i$ 」の入力フォントとのユークリッド距離を表す。式(3.12)が $(0 \leq i \leq N, 0 \leq j \leq M)$ の全ての $(i, j)$ で成り立つようなマス $m$ を識別領域とした。 $m$ が2以上あるときは、式(3.12)の右辺と左辺の差が大きい順に2つを識別領域とした。式3.12を満たす $m$ が1つしかない場合はこの1つを識別領域としている。図3.11では、 $N = 13$ 、 $M = 10$ であり、識別領域としては、マス(18, 19)となった。(図3.13)。今回用いた全ての詳細識別文字と識別領域は表3.12に示す。





- 以上の除去候補文字の処理と、詳細識別文字の除去を行い、候補集合中には抽出すべき正解のデータのみが残るようになり、除去候補文字の学習を終了した。(学習データ新聞社説(1~10)の10部において)

### 3.3 抽出実験

本研究で提案するフォント抽出のアルゴリズムによってフォントが正確に抽出できるかどうかを確認するためにフォントの抽出実験を行った。

#### 3.3.1 実験方法

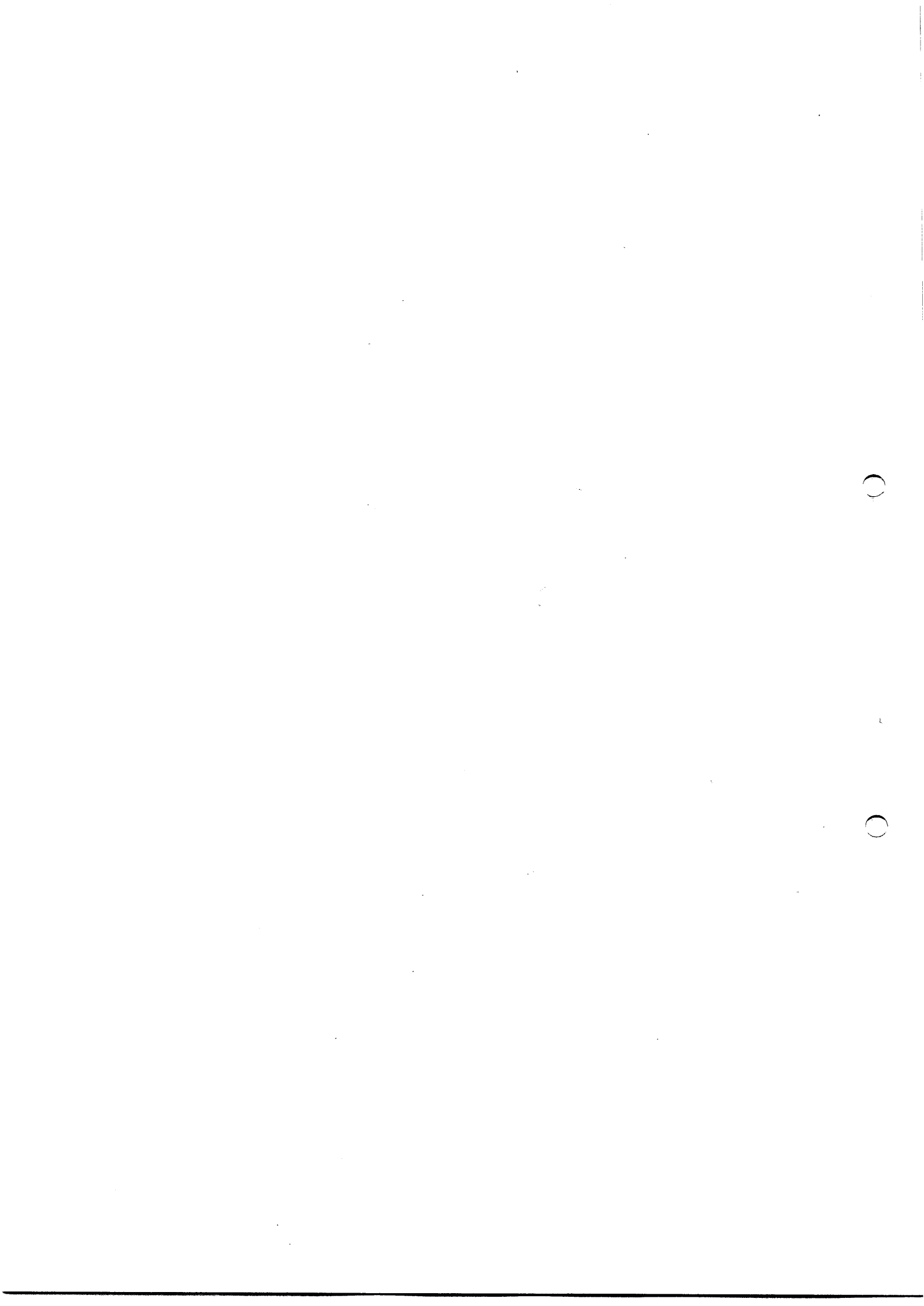
- フォントを抽出する対象文書として、除去候補文字の学習に用いなかった新聞社説(11~20)の10部を対象文書とし、1部ずつに対して実験を行う。また抽出する文字は表 3.3で示したひらがな 65 文字とする。
- フォント抽出のための認識で使用する既存辞書は各種プリンタのフォントセット 42 種から作成した重心ベクトルを用いる。辞書の文字数は 3109 字である。
- 図 3.2で示したアルゴリズムでフォントの抽出を行う。ここで、候補数をどれだけとれば全ての正解を候補集合に出力できるかを調べるため、候補数を 1 から 10 まで変えて実験を行った。

#### 3.3.2 実験結果

候補数と候補集合に集められる正解について表 3.4に示し、抽出の結果については表 3.5に示す。また、候補集合中の正解の割合と、真の正解数までに不足している正解数をグラフにし、図 3.15に示す。実際に抽出できなかった文字は表 3.7に示す。

### 3.4 考察

表 3.4の A は、除去候補文字の処理を一切行わなかった作成直後の候補集合のデータ数である。実際の実験では社説 1 部毎に行ったが、ここでは 10 部の結果を和で示してある。B は同様に候補集合に取り入れられた正解数を表している。C は、 $B / 10328$  で、D は、 $B/A$  で候補集合中の正解数の割合、つまり母集団中の有効なデータの割合を示している。B の欄に注目すると候補数 6 で全ての正解を候補集合に取り込むことができおり、この



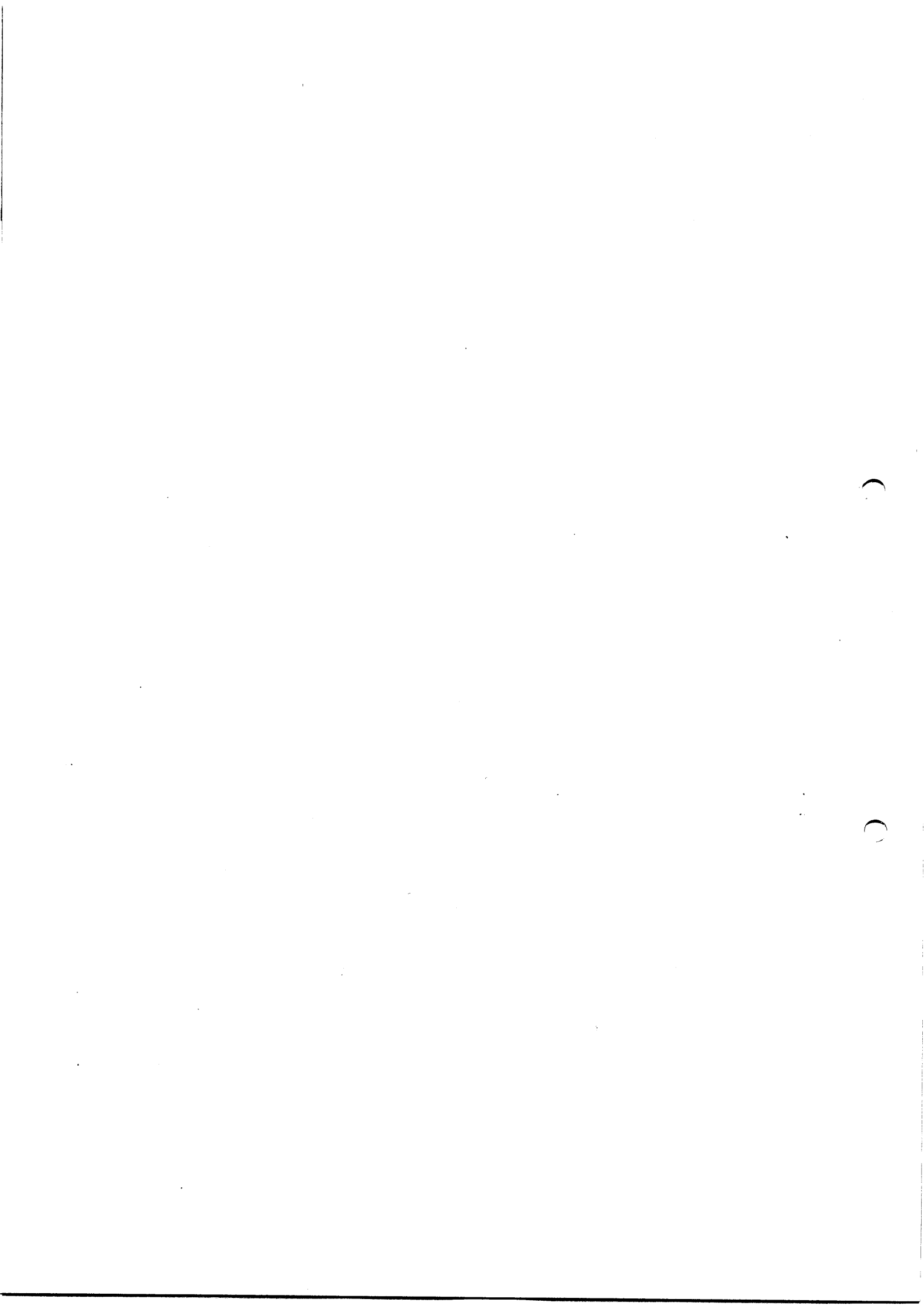
結果、3.1節で示した条件の1つである、文書中の全ての正解を集めることは達成できた。候補数7以降では新しく候補集合に取り込まれるデータは全て誤ったものである。候補数6以下で抽出できなかった文字の詳細は表3.7に示す。次にD欄に注目すると、有効なデータの割合は候補数6では32.12%なのが候補数10では23.80%と10ポイント近くさがっており、候補数7以上のデータは全体の処理の実行時間を無駄に使っていることになり、候補数の上限として6が適当であると考えられる。

表3.5で、真の出力数575は10部中に存在したひらがな65文字の字種数の和であり、フォントの抽出を期待される数になる。A欄の正解の出力数であるが、このカウントは社説毎に、1つでも抽出のあった文字を1として数えている。このため、1つの社説中に同じ抽出すべき文字が複数存在していた場合でも、その全てを抽出しなくても、1文字を抽出すれば抽出としている。A欄の()の内に数字は、存在した抽出すべき文字を全て抽出した文字数をカウントした。この正解全てを抽出した文字数は全体の8割であり、候補集合中に集められた全ての正解の抽出には成功しなかった。次に、正解以外を抽出したBとCについて述べる。Bは正解と違う誤った文字を出力した場合で、Cは複数字種のフォントを1つの字種として出力した場合であるが、今回の実験ではCの出力は見られなかった。しかし、候補数2~10でBの誤った抽出が確認された。表3.6に実際に誤って抽出された文字を示すが、誤抽出が生じると抽出結果から作成した辞書を用いる際に、誤抽出したコード(例「液」)と、実際に抽出すべきコード(例「ゆ」)の両方のコードを持つ文字の認識に影響があると考えられる。実際の影響については、4章での実験で確認する。

誤って抽出された文字は実質「液」、「点」、「ザ」の3文字であるが、これらは全て除去候補文字の学習に用いた社説10部中には出現しなかった文字である。3文字とも出力数が1であり、候補集合中から除去候補文字を除去する際に、1文字だけが除去されなかったため、細分類を行っても除去できなかったものである。

また、「ザ」を除く2文字は候補数10にしたときのみに抽出されている。図3.16では、「液」の場合について示している。候補数が9のときには、候補集合に「液」のフォントは選出されてないため、除去候補文字の除去の処理後に候補集合中にデータは残らない。しかし、候補数を10にすると「液」のフォントが、最初の候補集合に選出され、除去候補文字の除去後も候補集合中にあり、そのまま誤抽出されてしまった。

このことより学習したデータにない文字が候補集合中にある場合でも、候補数を抑えることによって誤抽出を避けることも可能であるが、「ザ」の様に候補数が2から抽出されてしまう場合には除去候補文字を更新する等の処置を施さなければ、誤抽出を改めることは難しい。



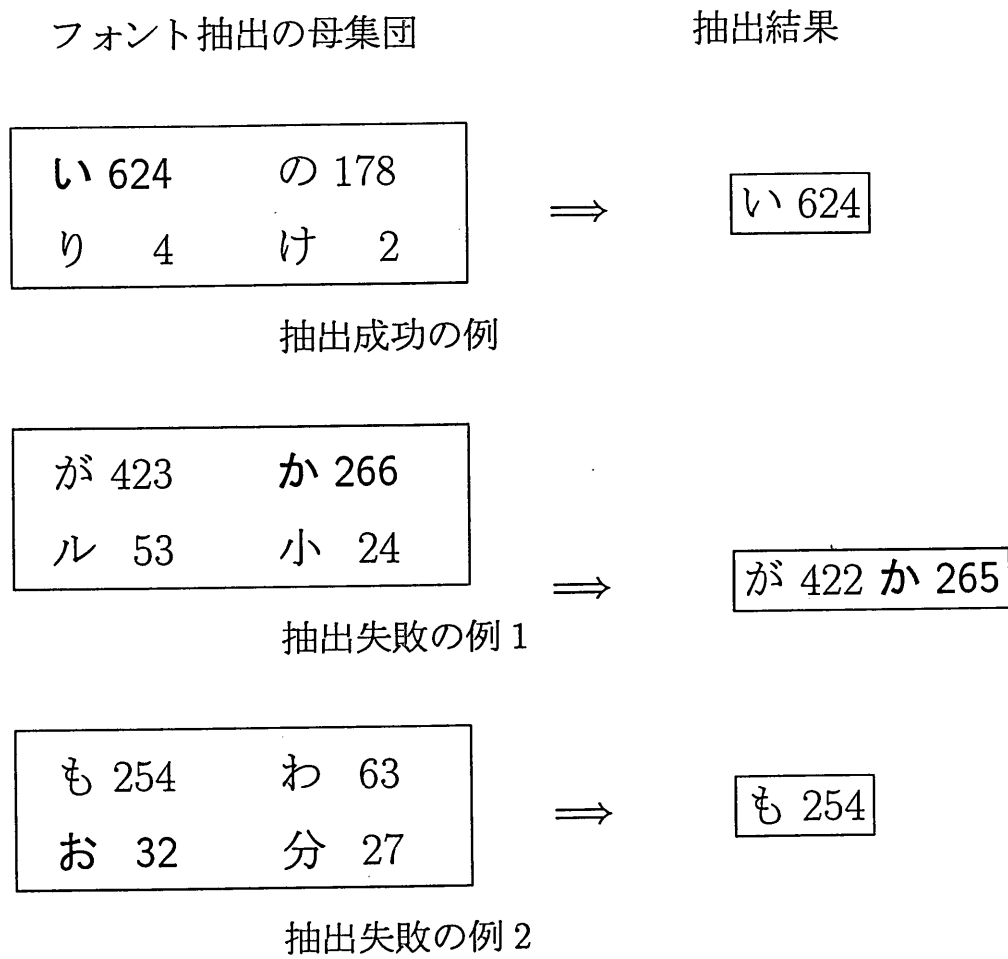
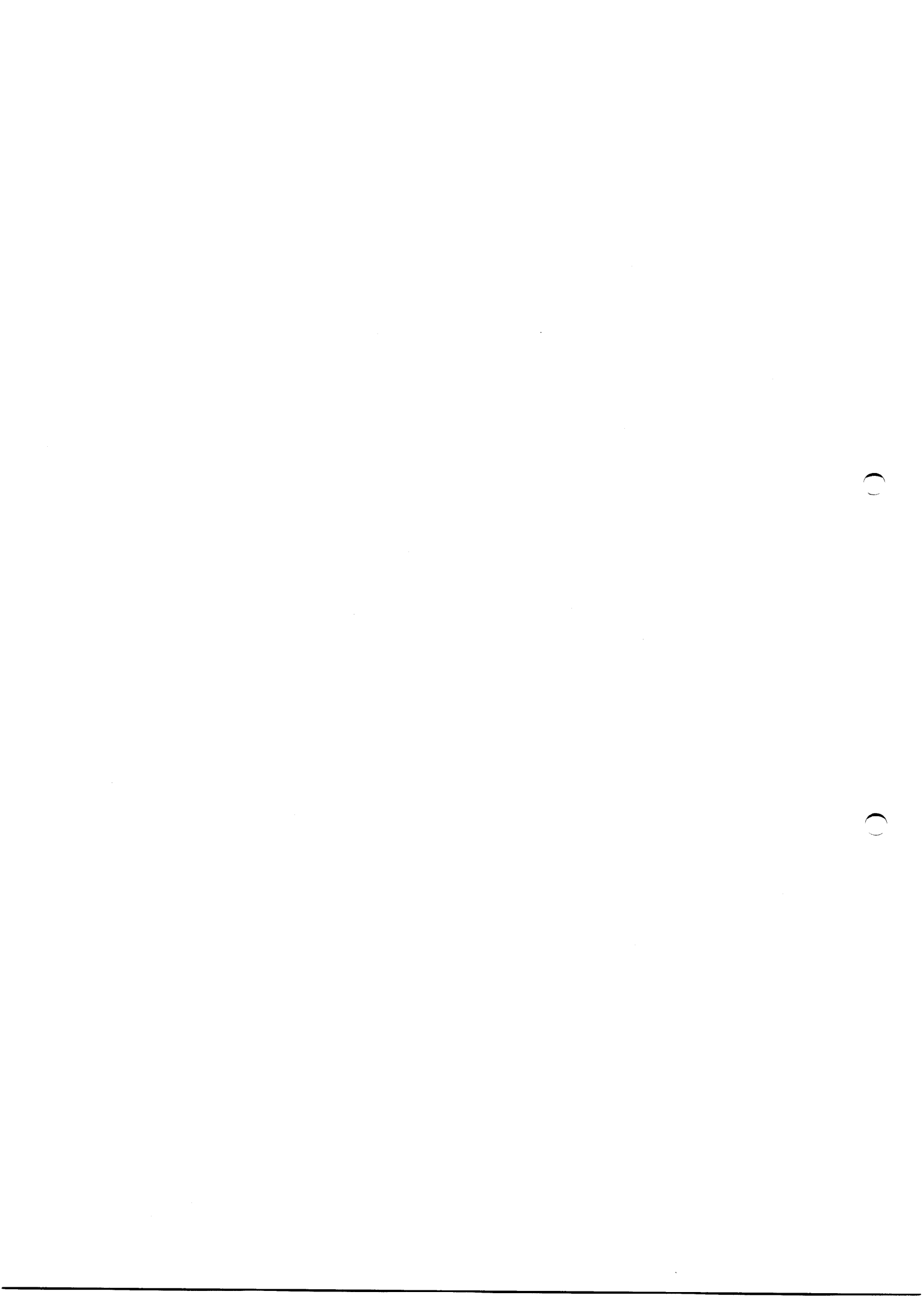


図 3.1: 重心を用いたフォントの抽出法



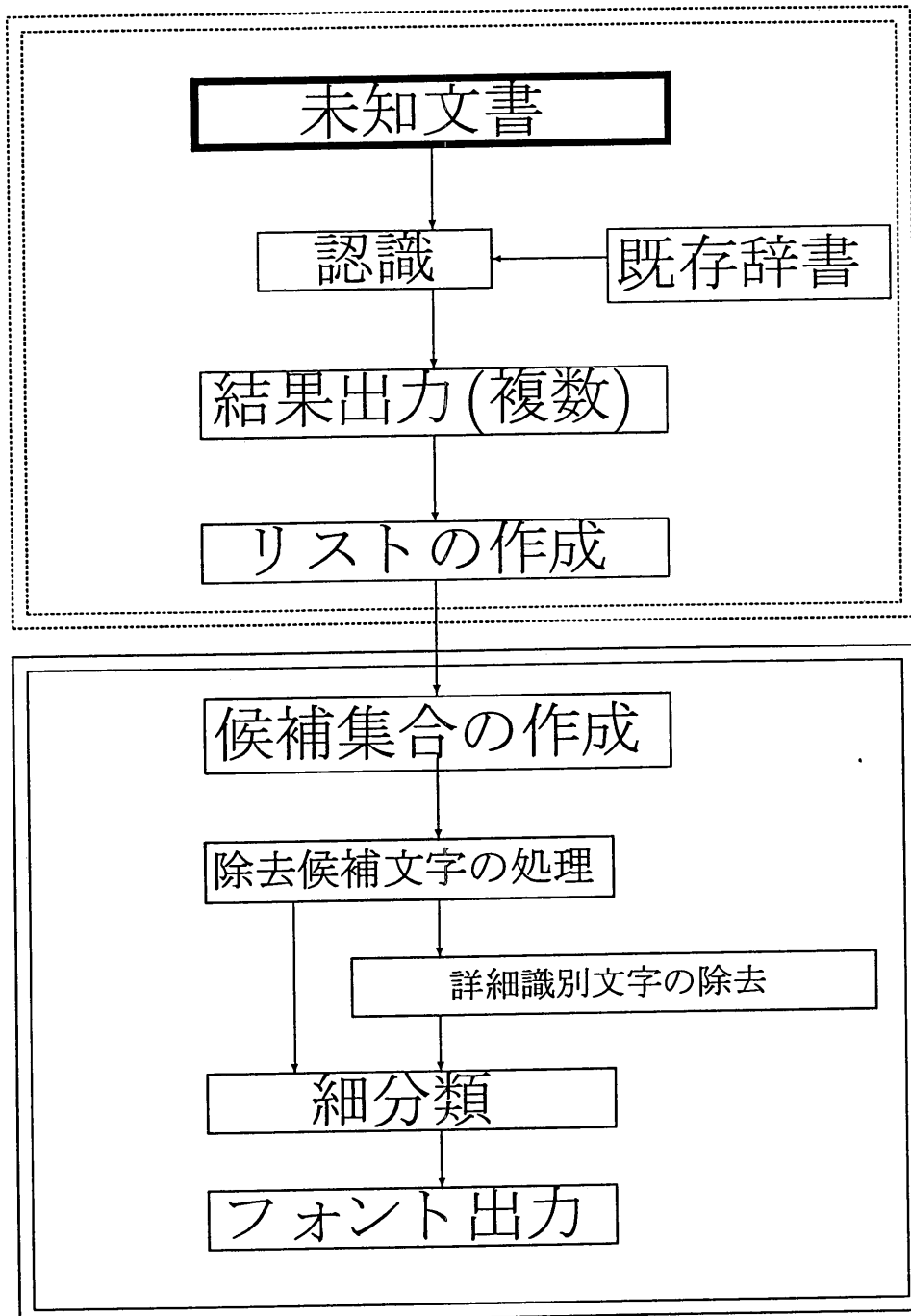


図 3.2: 未知文書からのフォント抽出アルゴリズム





あ	い	う	え	お	か	が	き	ぎ	く	ぐ	け	げ	こ	ご
さ	ざ	し	じ	す	ず	せ	ぜ	そ	ぞ	た	だ	ち	つ	づ
て	で	と	ど	な	に	ぬ	ね	の	は	ば	ひ	び	ふ	ぶ
へ	べ	ほ	ぼ	ま	み	む	め	も	や	ゆ	よ	ら	り	る
れ	ろ	わ	を	ん										

図 3.3: フォント抽出の対象とした文字

あ	い	う	え	お	か	が	き	ぎ	く	ぐ	け	げ	こ	ご
	ざ	し	じ	す	ず	せ	ぜ	そ	ぞ	た	だ	ち	ち	つ
づ	て	で	と	ど	な	に	ぬ	ね	の	は	ば	ば	ひ	び
び	ふ	ぶ	ぶ	へ	べ	ぺ	ほ	ぼ	ぼ	ま	み	む	め	も
や	ゆ	よ	ら	り	る	れ	ろ	わ	ゐ	ゑ	を	ん		
ア	イ	ウ	エ	オ	カ	ガ	キ	ギ	ク	グ	ケ	ゲ	コ	ゴ
サ	ザ	シ	ジ	ス	ズ	セ	ゼ	ソ	ゾ	タ	ダ	チ	チ	ツ
ヅ	テ	デ	ト	ド	ナ	ニ	ヌ	ネ	ノ	ハ	バ	パ	ヒ	ビ
ピ	フ	ブ	プ	ホ	ボ	ポ	マ	ミ	ム	メ	モ	ヤ	ユ	ヨ
ラ	リ	ル	レ	ロ	ワ	ヰ	エ	ヲ	ン	ヴ				

図 3.4: 辞書中の仮名文字



# 社説



## なぜいま四年制高校なのか

中央教育審議会が、五年半ぶりに再開された。高校教育の改革と生涯学習の基盤整備を中心課題としている。

諮問された内容を一読してまず気づかされるのは、高校の修業年限に踏み込んだ学制改革案が含まれていることだ。

それに沿って、極めて具体的な提案を盛り込み、それぞれの当否を問う諮問の仕方にもなっている。よく抽象的な課題を投げかけるだけでなく、またこれまでの諮問方式に比し、大きく異なる点があった。

いまの三年制高校のほかは四年制の高校

についても検討を求めている。

これらはいずれも、画一的な学校教育に一石を投じる要素を持っている。個々に見れば悪くないと思われるものもある。

だが、多くの国民の受け止め方は、なぜいま四年制高校なのか、というものではないだろうか。職業高校だとせよ、普通科にも広げるとせよ、どうしてもいま必要なものなのかどうか。このあたりは、いからの検討が待たなければならない。

ただ、そうした違和感を覚えるのは、今回の制度改革案が、高校と、大学の一部だ

そうした環境の中

対象に検討を進め、

教育を含めた全面改

としているように思

だが、学校の区切

度全体をにらんだう

筋というものだろう

今回の諮問がわか

た大きな背景を持ち

限のみに手をつける

も知れない。それに

に乗り出す必要性に

という問題もあるう

とほ言っても、高

元  
**秘書**

リクルート事件で

図 3.5: 新聞社説の例



フォント	1位	2位	3位	4位	...
「A」	「は」	「け」	「ほ」	「吋」	...
「B」	「ぽ」	「ぽ」	「ば」	「ぼ」	...
「C」	「は」	「ほ」	「ば」	「け」	...
「A'」	「け」	「は」	「ほ」	「げ」	...

図 3.6: 認識結果の例

フォント	コード	コード	フォント
「A」	←	「は」	「は」 ← 「A」 「C」 「A'」
「A」	←	「け」	「け」 ← 「A」 「C」 「A'」
「A」	←	「ほ」	「ほ」 ← 「A」 「C」 「A'」
⋮	⋮	⋮	「ば」 ← 「B」 「C」

(a)

(b)

図 3.7: コードとフォントの対応



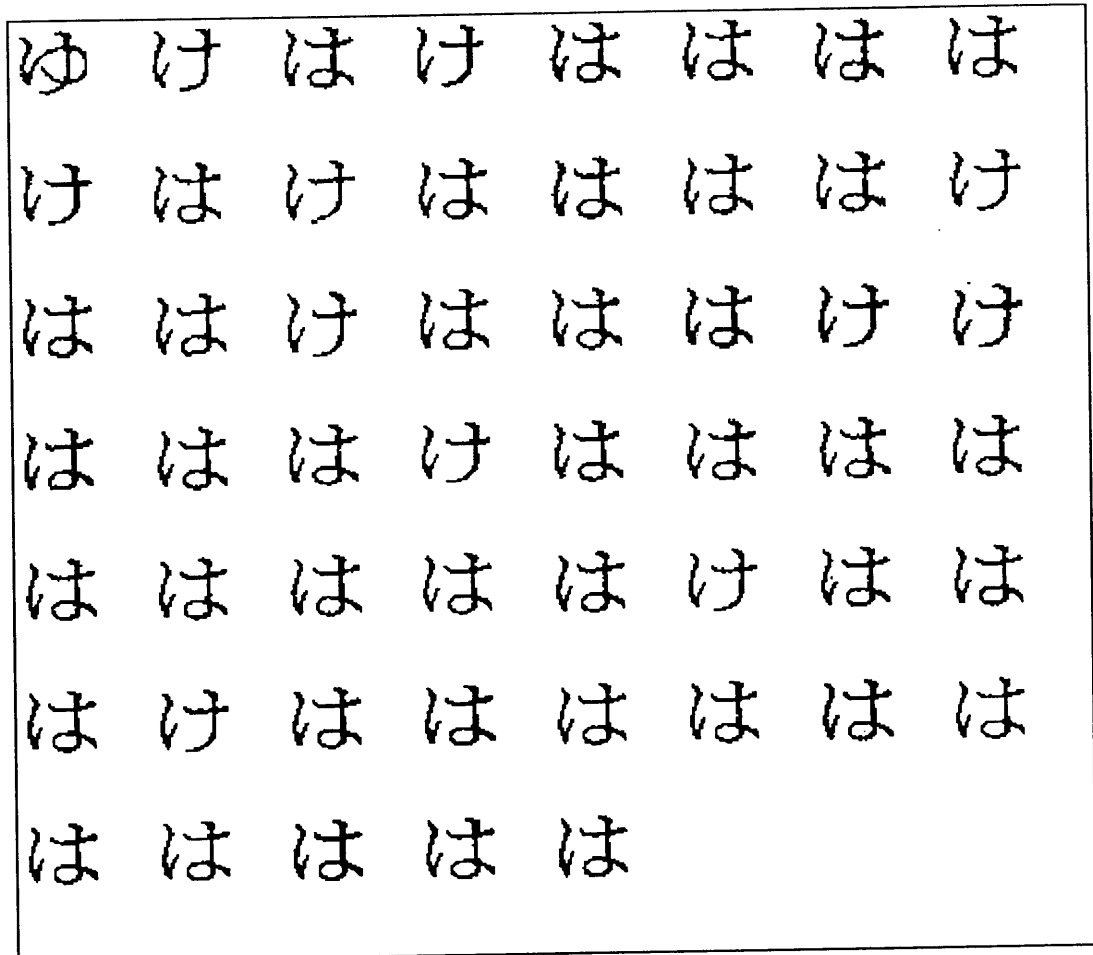


図 3.8: 候補集合の例





「み」の候補集合中のフォント

(5 字種 112 文字)

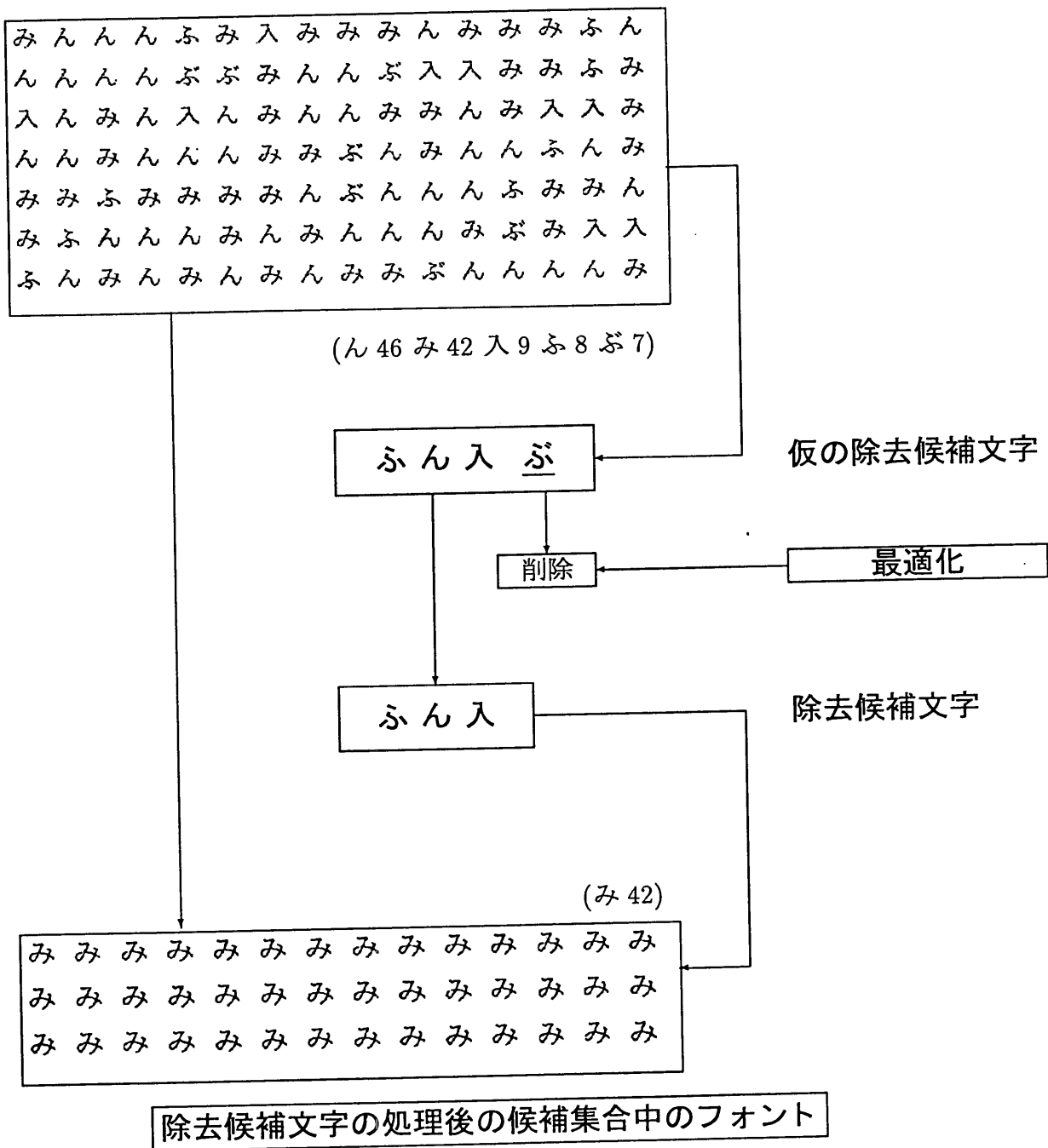
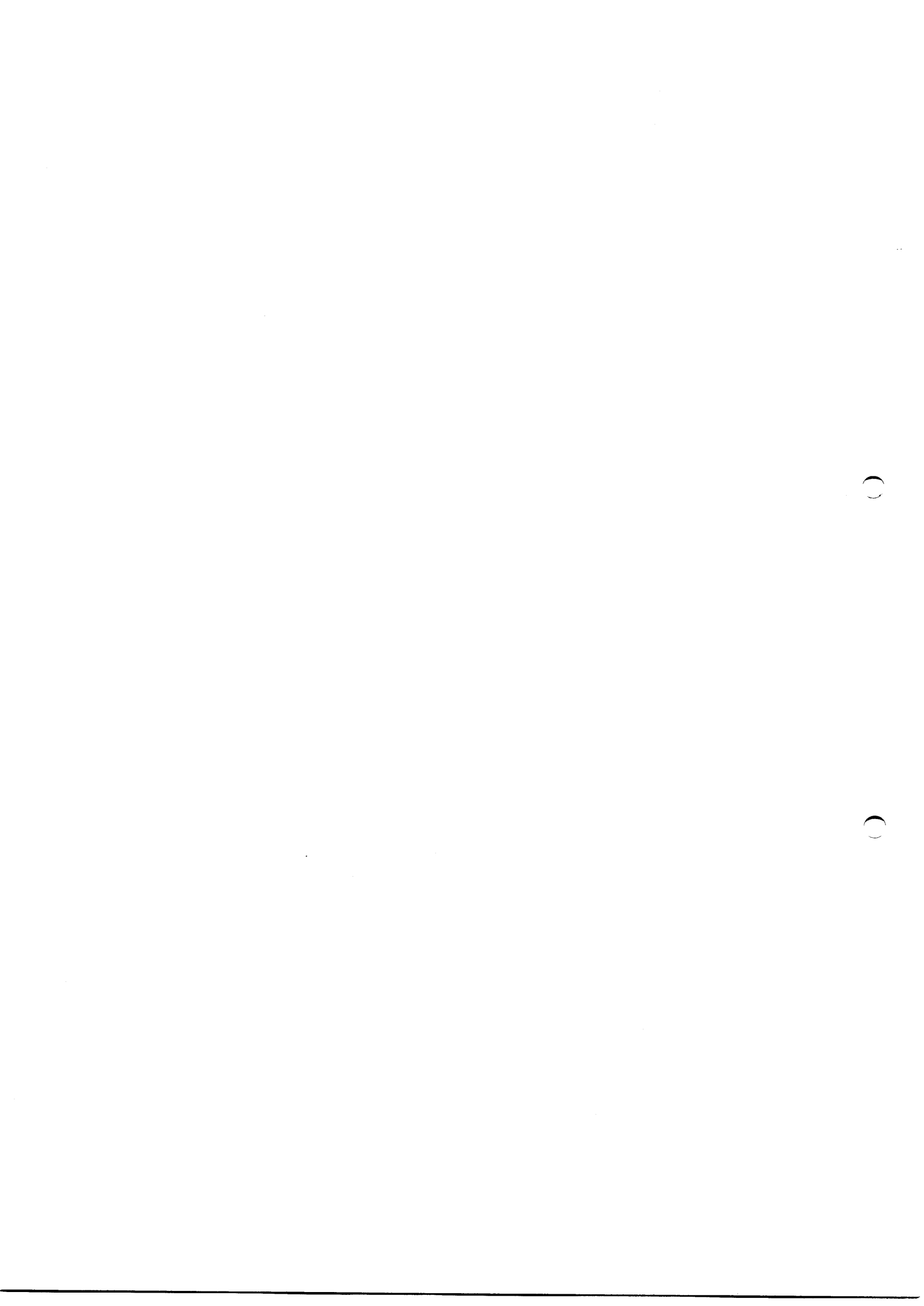
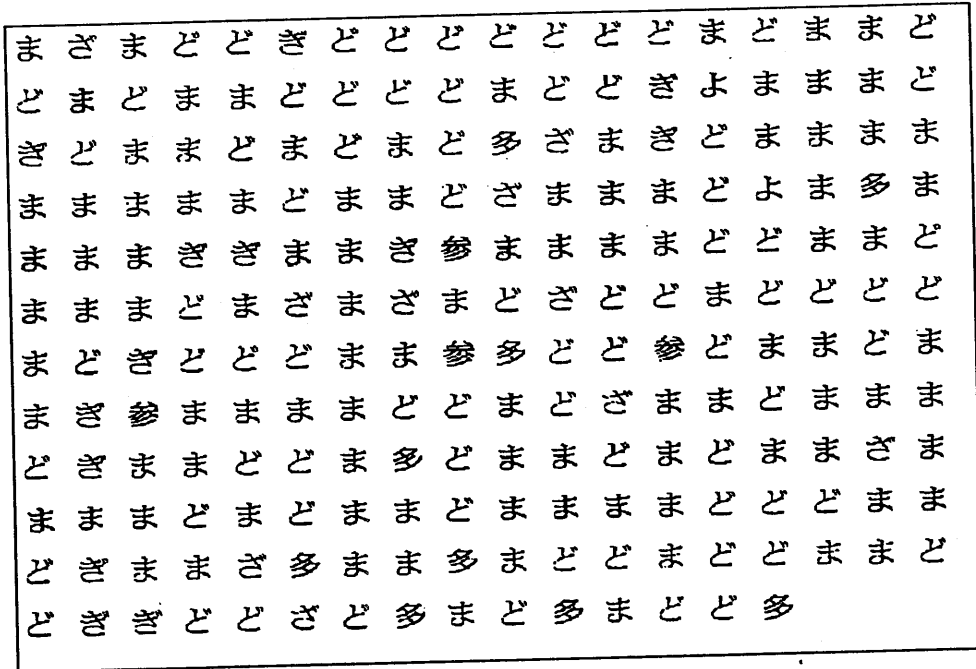


図 3.9: 除去候補文字の学習例 1



「ぎ」の候補集合中のフォント

(7 字種 213 文字)



(ま 99 ど 76 ぎ 13 ま 10 多 9 参 4 よ 2)

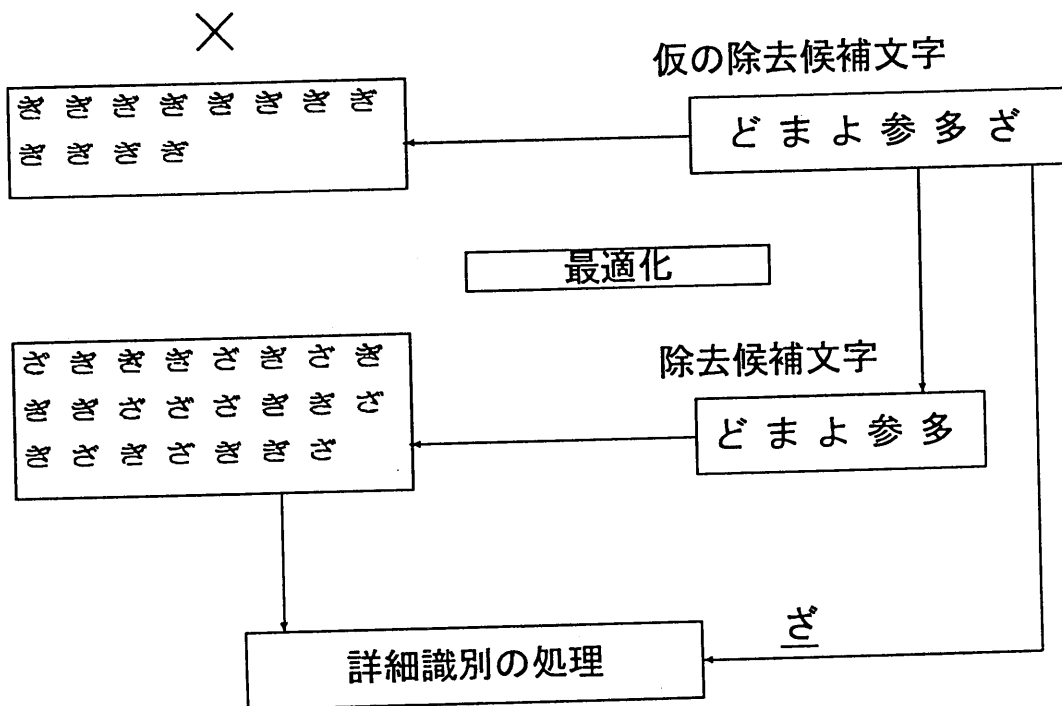
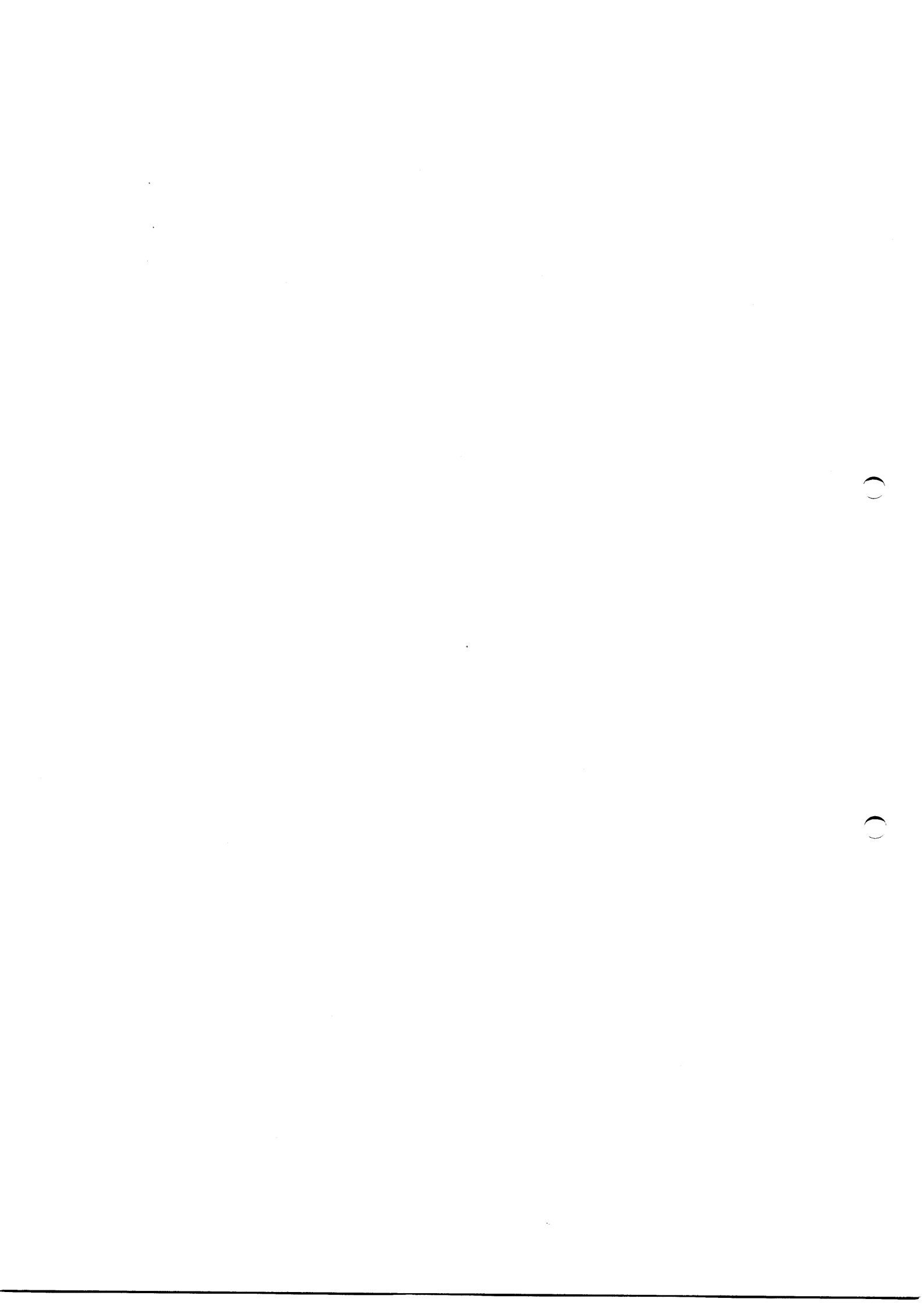


図 3.10: 除去候補文字の学習例 2-1



除去候補文字の処理後の候補集合中のフォント (2 字種 23 文字)

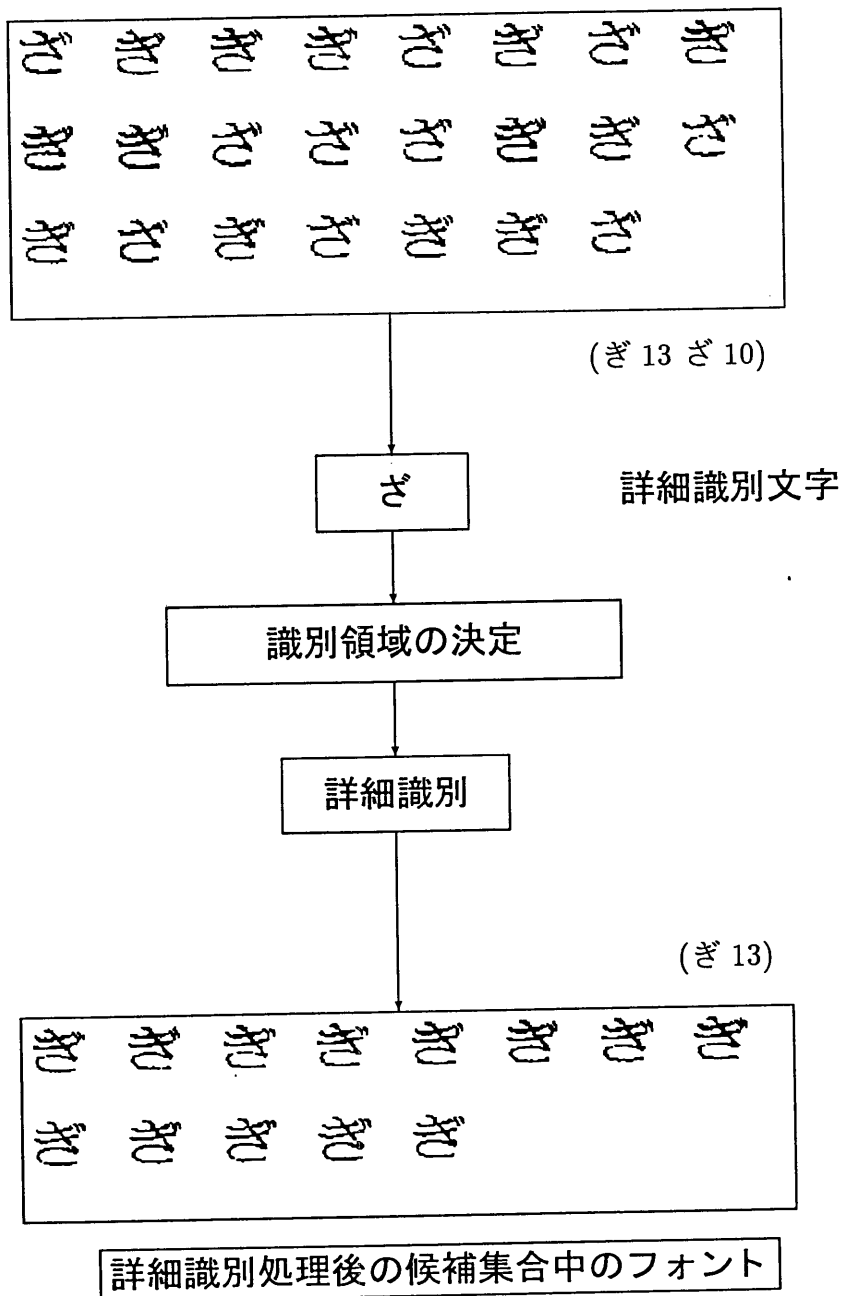


図 3.11: 除去候補文字の学習例 2-2

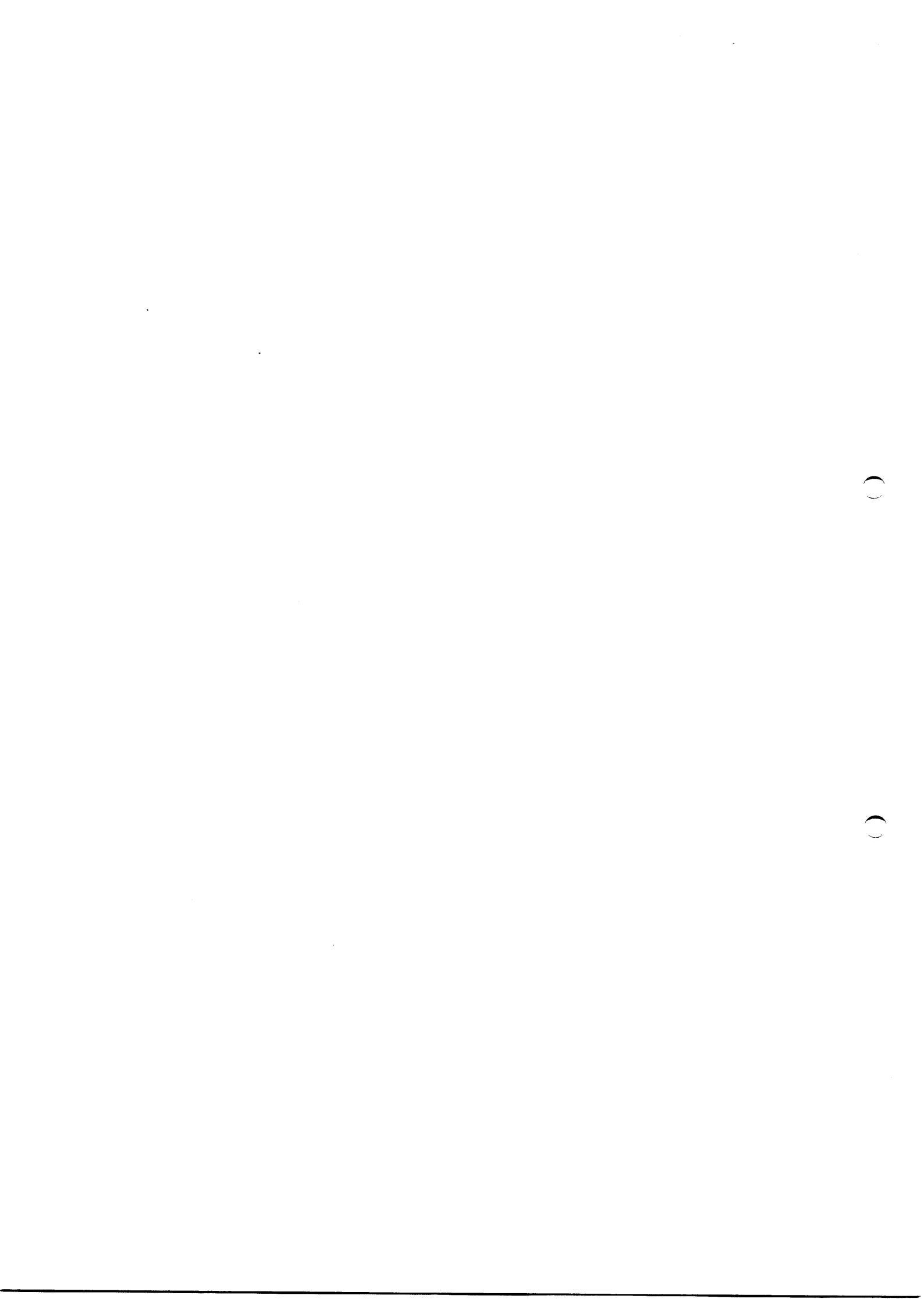


表 3.1: 除去候補文字 1

コード	除去候補文字	コード	除去候補文字
あ	おぬみもん	し	じにらニレ
	占方命力	じ	けしにらり
い	けのり		ニビリ亡
う	つらる	す	ずオキケナ
え	くそたてミ		求水本六
お	あねも分	ず	すオナブボ
か	がなやケル		求歩
が	かだべやケ	せ	ぜサ
	トボ	ぜ	せだどぜ
き	さまを吉参	そ	ぞるろを
	表	ぞ	ぐそでみろ
ぎ	どまよ参多		ぜデ守分
く	えぐ今	た	えだなにわ
ぐ	くへイウジ		セ
	ー	だ	がぜたづて
け	げにびぼゆ		なにもわぜ
げ	はび庁	ち	おたもらを
こ	ごとになるエ	っ	うぐづでの
	ニ		コ
こ	ざじだどゴ	づ	うがぐっウ
	ビユンニ		ガゴゼ分
さ	うきざとも	て	くぞっエマ
	を台		下
ざ	ぎすどぶを	で	ぞつてのや
	ゴ		フ及





表 3.2: 除去候補文字 2

コード	除去候補文字	コード	除去候補文字
と	どら公上占	ぼ	けにぼほ序
ど	じとらビ	ま	す安拳主全
な	すたをウ女 珍立		歩
に	いけした広 仁	み	ふん入
ぬ	ねめ応山内 約	む	だなもれわ 以位心打
ね	おかむめも れ付料	め	あいぬのみ わか小心
		も	ちむら合
の	いつみめり 山	や	みネヤ一
		ゆ	がのや中仲
は	なにぼほ	よ	まエ少上歩
ぼ	げなはびぼ ゆビ体	ら	しちふもる 一各
		り	いうしつの むゆウリロ
ひ	いけしなび ユレ以	る	ちろミ各石
び	がげひゆブ	れ	かなん介九
ふ	ぶよんムン	ろ	あうそちみ
ぶ	ふよんム小		る
へ	なべんレ外 人	わ	あおかみも れん心分
べ	ぐなへんイ ドバ	を	あさちとら る去合
ほ	にはぼ行低	ん	たふみル

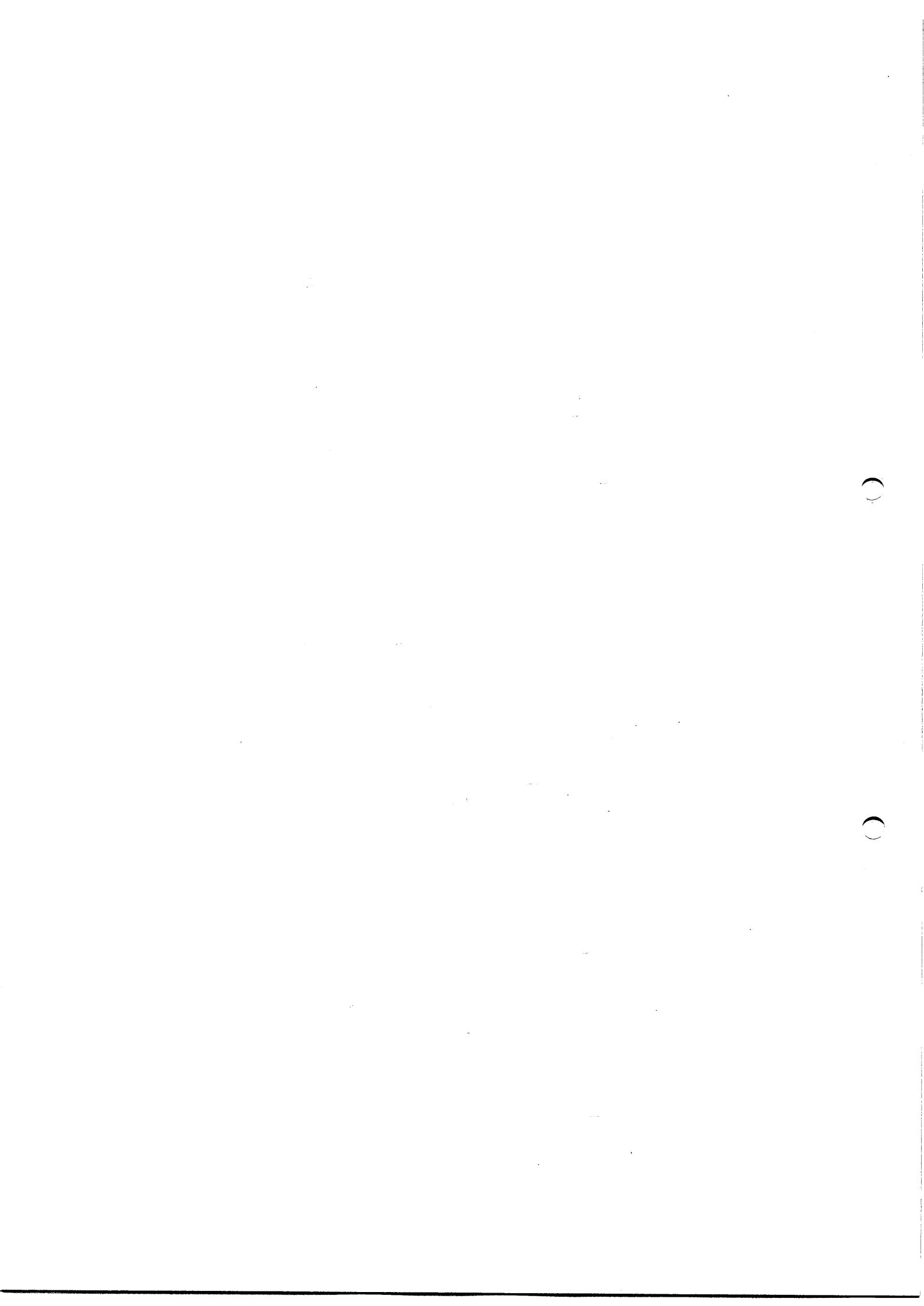
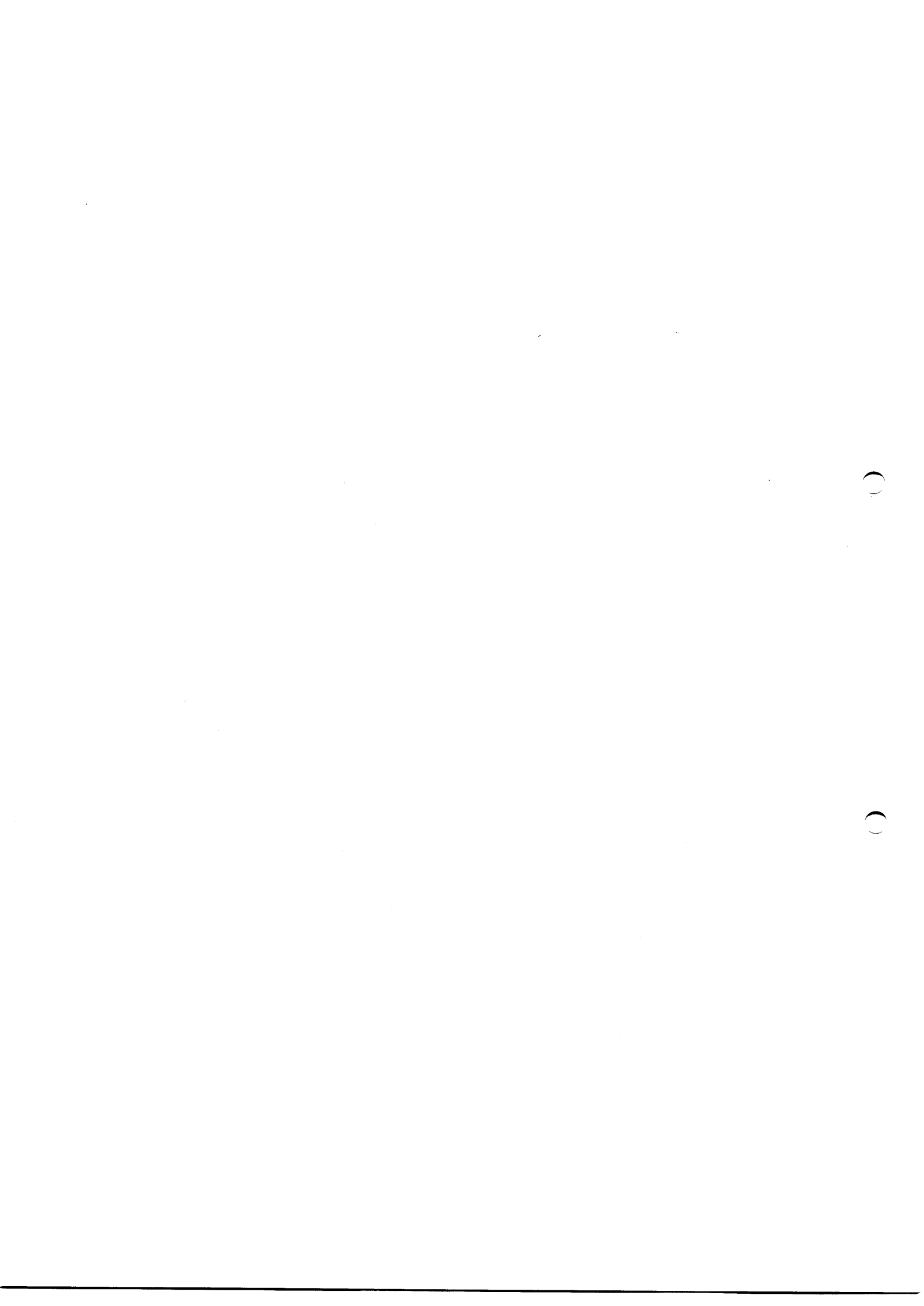


表 3.3: 除去候補文字数と評価値

文字数	用いる評価値	
	ユークリッド距離	重み付き
2	せ	
3	いくみ	うげ
4	ぜはもやん	おそど
5	えかぎしちふぶ ほぼゆよれ	けとびる
6	ぐざたつにぬの へま	きこてろ
7	がずなべら	さで
8	ばを	ねひ
9	あじすぞづむめ	ごわ
10	だり	



識別文字	識別領域	識別文字	識別領域
う ら	6 13	う ラ	0 0
ぎ ざ	18 19	こ ご	13 48
ざ さ	13 48	し り	13 20
そ も	2 3	だ た	18 18
て で	19 20	ね れ	39 46
ね わ	38 39	は け	39 40
は ゆ	27 34	ば は	6 6
ひ か	34 41	ぶ ふ	6 13
ほ ぼ	6 13	わ れ	46 47

図 3.12: 詳細識別文字

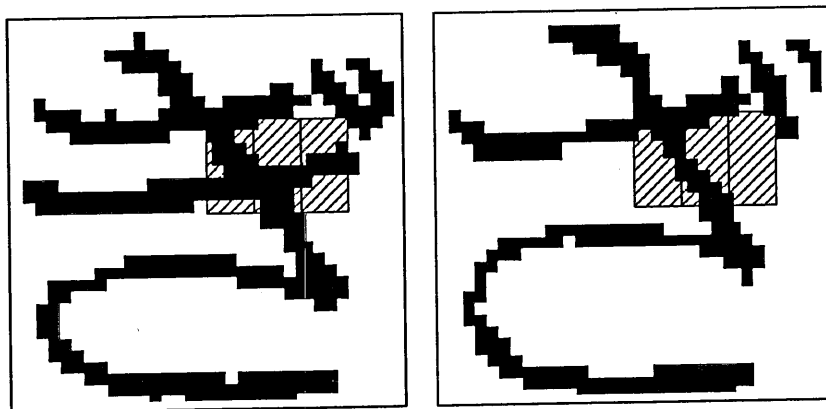
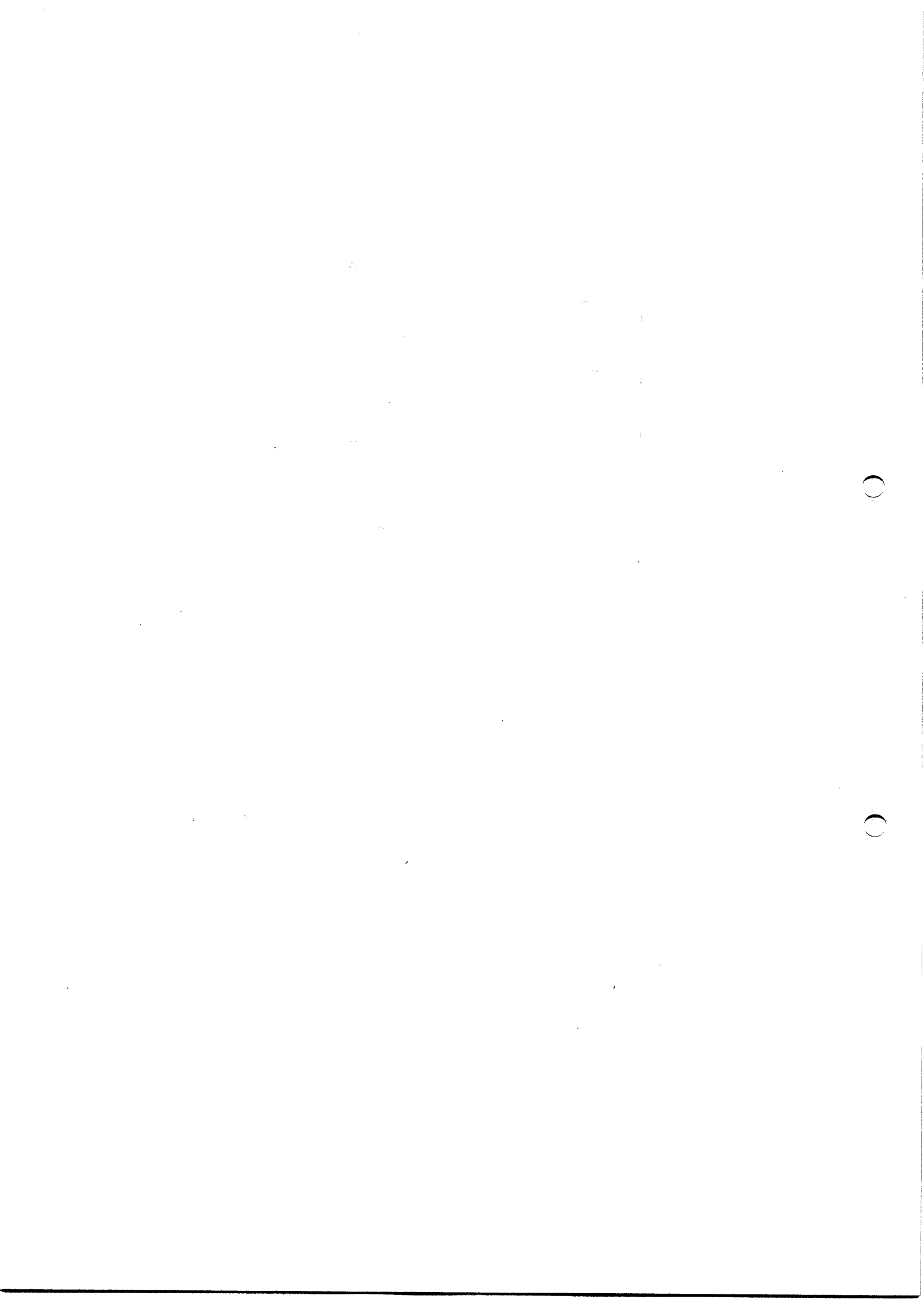


図 3.13: 識別領域の例



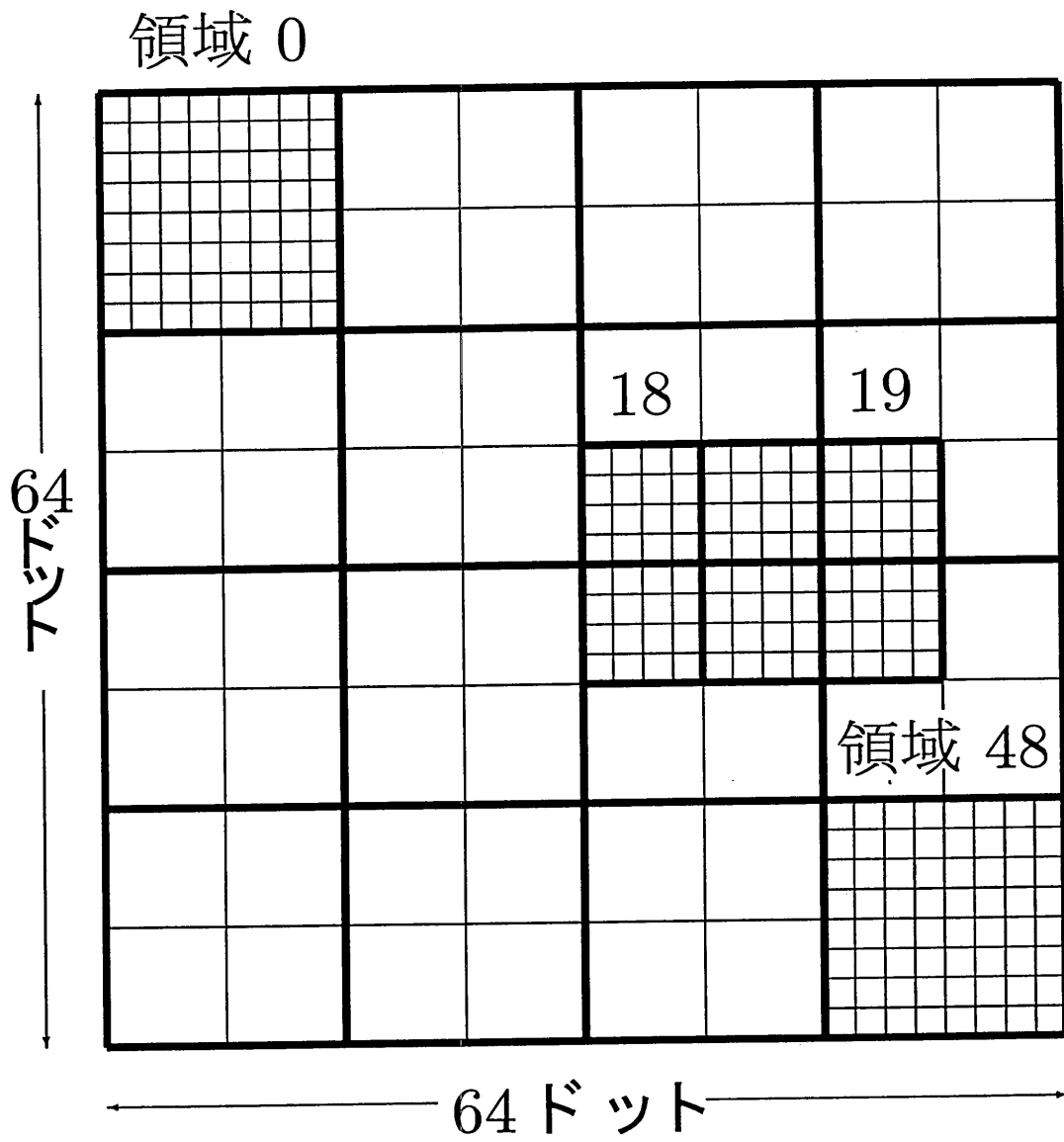


図 3.14: 識別領域

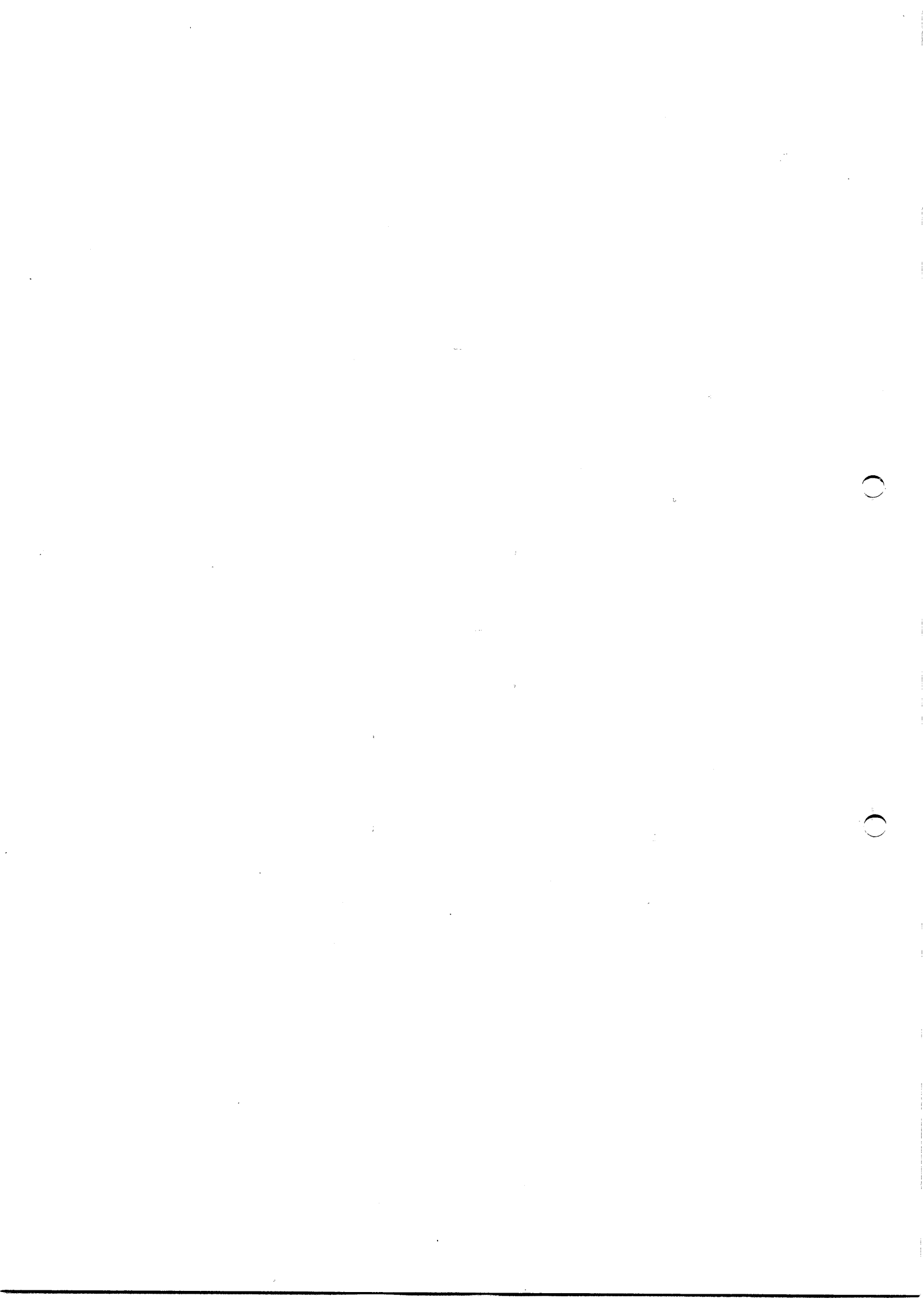




表 3.4: 候補数と候補集合

候補数	A	B	C(%)	D(%)
1	10131	9958	96.42	98.29
2	16391	10264	99.38	62.62
3	21898	10316	99.88	47.11
4	25904	10323	99.95	39.85
5	29296	10325	99.97	35.24
6	32150	10328	100.00	32.12
7	34984	10328	100.00	29.52
8	38006	10328	100.00	27.17
9	40905	10328	100.00	25.25
10	43390	10328	100.00	23.80

A:全候補集合のデータ数 B:正解数

C: B/全正解数 (10328) D: B/A

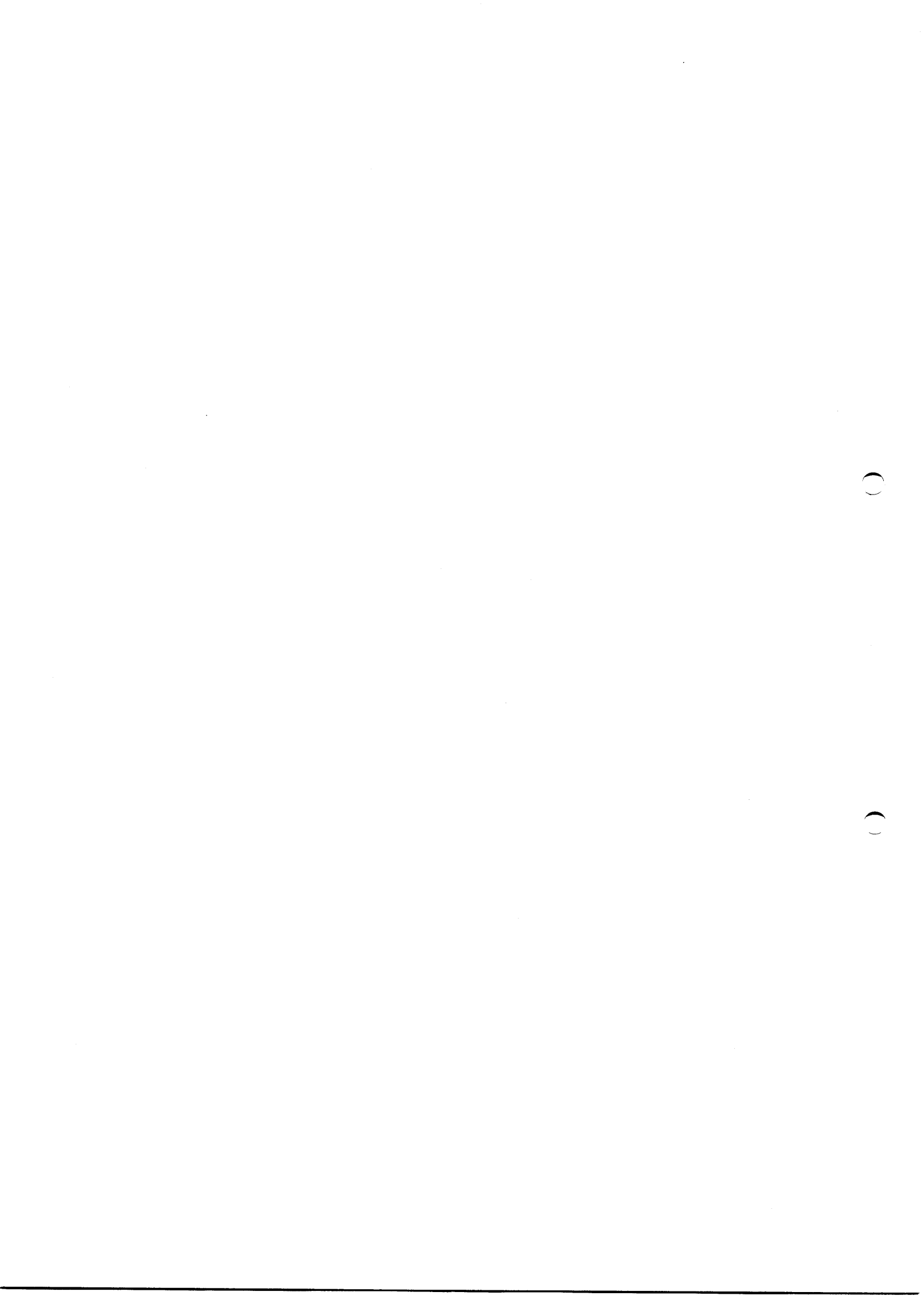


表 3.5: 候補数と出力数

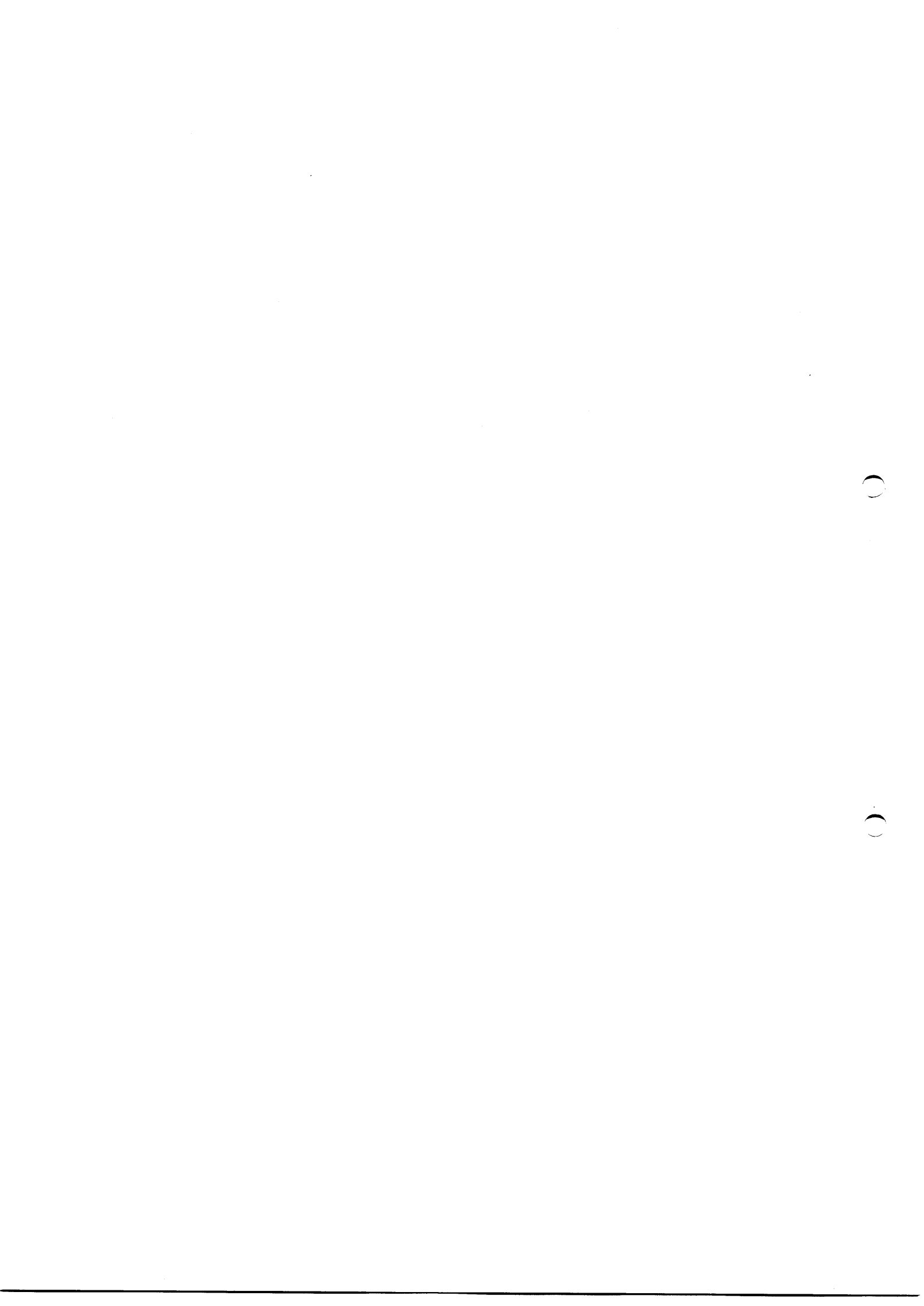
候補数	出力数			
	A	B	C	合計
1	533(389)	0	0	533
2	568(449)	1	0	569
3	572(465)	1	0	573
4	574(469)	1	0	575
5	574(472)	1	0	575
6	575(473)	1	0	576
7	575(473)	1	0	576
8	575(471)	1	0	576
9	575(471)	1	0	576
10	575(471)	3	0	578

A: 全て正解が出力 B: 全て不正解を出力

C: 正解と不正解の混在を出力

() は文書中の全正解を出力した数

※期待される真の出力数は 575



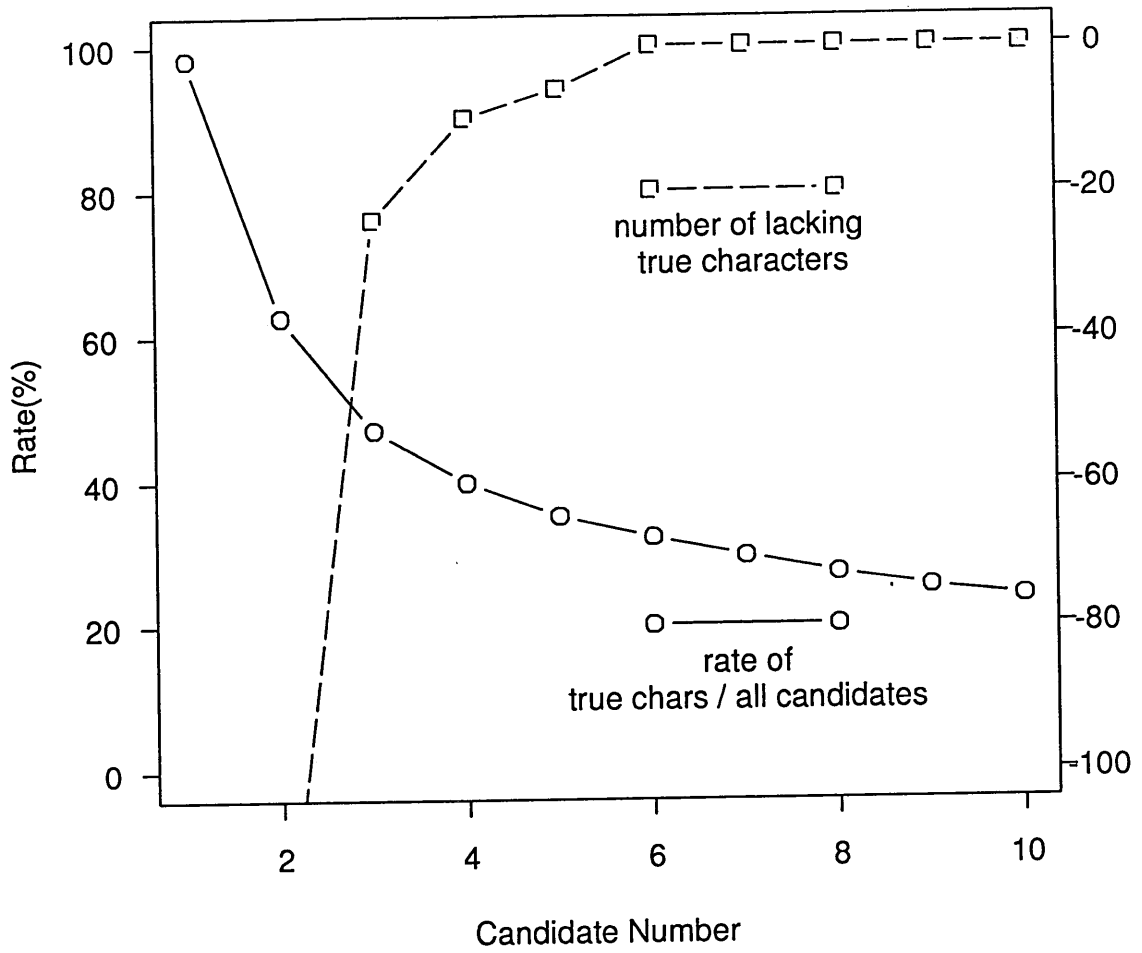


図 3.15: 候補数と出力数

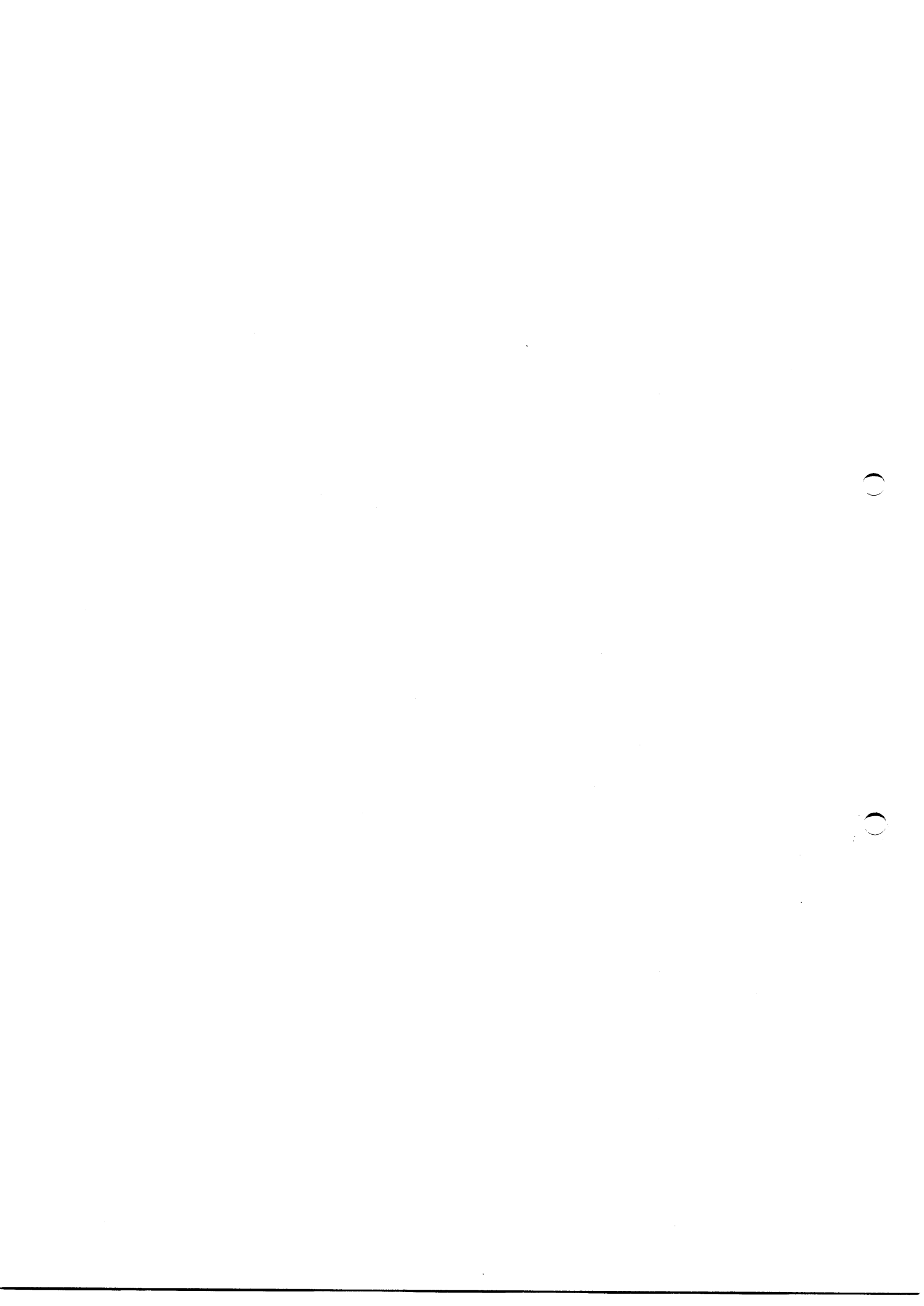
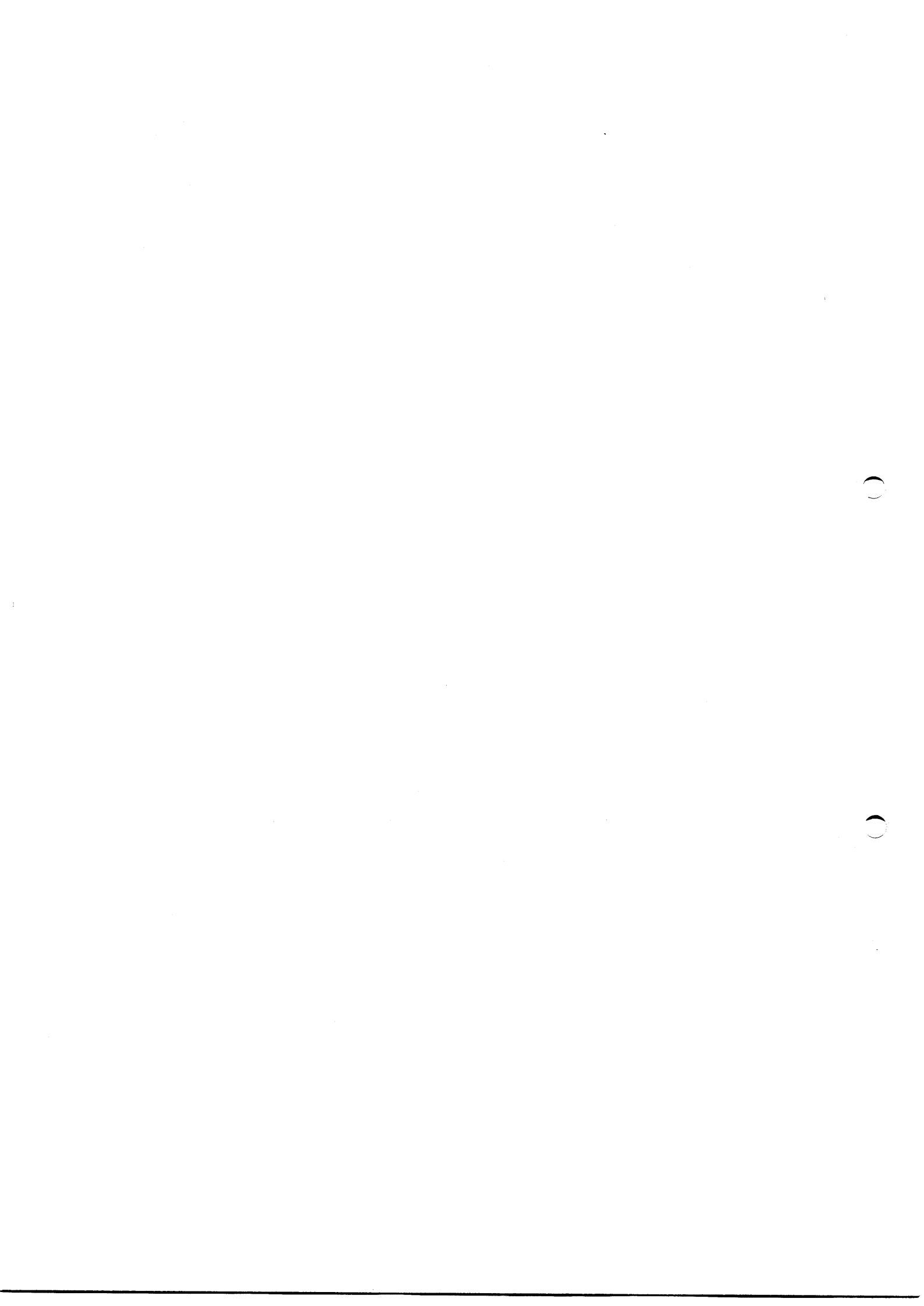


表 3.6: 誤抽出された文字

文書	候補数	抽出すべき文字	出力数	誤抽出した文字
11	10	ゆ	1	液 (1)
13	10	ぬ	1	点 (1)
14	2	ぜ	1	ザ (1)
14	3	ぜ	1	ザ (1)
14	4	ぜ	1	ザ (1)
14	5	ぜ	1	ザ (1)
14	6	ぜ	1	ザ (1)
14	7	ぜ	1	ザ (1)
14	8	ぜ	1	ザ (1)
14	9	ぜ	1	ザ (1)
14	10	ぜ	1	ザ (1)





候補集合

除去候補文字の処理後

中	中	の
の	の	

⇒

データなし
-------

抽出成功の例 (候補数 9)

中	中	の
の	液	の
の	の	の

⇒

液
---

誤抽出の例 (候補数 10)

図 3.16: 誤抽出の例



表 3.7: 抽出できなかった文字

文書	候補数	抽出できなかった文字
11	1	ぎ ば び わ
12	1	ぎ ば び ゆ わ
12	2	ゆ
12	3	ゆ
13	1	さ ば び わ
14	1	ぎ ご ば び わ
14	2	ば
15	1	ぎ ば わ
16	1	ぎ ご さ ば び わ
16	2	ぎ ば
17	1	ぎ ざ わ
18	1	ぎ ご ぼ ぼ わ
18	2	ご
18	3	ご
19	1	ば ゆ わ
19	2	ゆ
19	3	ゆ
19	4	ゆ
19	5	ゆ
20	1	ぎ ば び わ
20	2	ぎ



## 第 4 章

# 認識実験

### 4.1 前書き

抽出されたフォントを用いて認識実験を行う。この認識実験の結果により、本研究で提案した文書認識システムの有効性を確認する。

### 4.2 実験方法

3章で抽出されたフォントの重心ベクトルと既存辞書から新辞書を以下に示す2つの方法で作成し、認識実験を行う。実験は次の3種類について行った。

1. クローズ実験
2. オープン実験
3. 更新した辞書による認識実験

また比較のために3章でフォントの抽出に用いた辞書(既存辞書)を用いても認識実験を行う。認識方法は評価値としてユークリッド距離を用いた全数整合法によって行う。

### 4.3 新辞書の作成方法

社説からフォントを抽出し作成した辞書を以降は抽出辞書と呼ぶことにする。新辞書は、抽出辞書と既存辞書から作成されることになる。



### 4.3.1 マルチテンプレート法

既存辞書に、抽出辞書を付加して新辞書を作成する。フォントが抽出されなかった文字に関しては、全て0の値を持つベクトルを代用する。この新辞書は、最大でひらがな65文字に対しテンプレート数2のマルチテンプレートの辞書となる。辞書の字数は3109+65で3174文字になる。

### 4.3.2 置換法

フォントを抽出したコードの抽出辞書を既存辞書のものと置換して新辞書を作成する。抽出されなかった文字は既存辞書のベクトルをそのまま用いる。辞書の字数は3109文字のままである。

## 4.4 クローズ実験

フォントの抽出に用いた新聞社説(11~20)の10部について、抽出したフォントから作成した新辞書を用いて認識実験を行う。フォントの抽出の際には候補数を1~10まで変えたものを用いる。

フォントを抽出した文書を、その抽出したフォントから作成した新辞書を用いて認識するので、再認識することになり、誤認識の自動修正の有効性を確認する実験となる。

表4.1に既存辞書を用いた認識結果を、表4.2にクローズ実験の1位の認識結果を示す。また、表4.3、表4.4に従来法と比較して、クローズ実験でどれだけ認識結果が改善されたかを示す。マルチテンプレート法、置換法ともに候補数6以上ではフォント抽出の対象であるひらがなの65種に対して、従来法では誤認識してしまった370文字を全て正しく認識している。また、表4.2を見ると候補数6でひらがな65種に対する認識では、マルチテンプレート法で100%、置換法でも99.99%の1位認識率が得られている。3章の結果とも合わせて、以降の実験では候補数を6に固定して実験を行うことにした。

## 4.5 オープン実験

フォントの抽出に用いた新聞社説(11~20)の10部について、認識する文書以外の9つの文書からそれぞれ作成した9つの新辞書を用いて認識実験を行う。

この実験で1部の文書からフォントを抽出して作成した新辞書に、他の文書の認識に対する有効性を確認する。この有効性が確認できれば、認識する文書全てにフォント抽出・





再認識(クローズ実験)を行わなくても、1部の文書でフォントの抽出を行い、それから作成した新辞書を用いてパターンマッチング法のみで、クローズ実験程度の認識結果が得られることになる。

以降の実験では、抽出辞書を抽出する文書と、認識を行う文書が異なる。そこで表4.5に本研究で認識の対象とした社説(11~30)の20部の文書に存在しないひらがな65種を示しておく。

表4.6、表4.7にマルチテンプレート法と置換法による1位認識結果を示す。表中のAの欄は前節のクローズ実験の結果である。B欄では、9つの新辞書を用いた認識結果のうち一番悪かったものと一番良かったものを、C欄では9つの新辞書の平均認識率を示した。また表4.8、表4.9には、フォントを抽出した文書毎の1位認識率を示した。

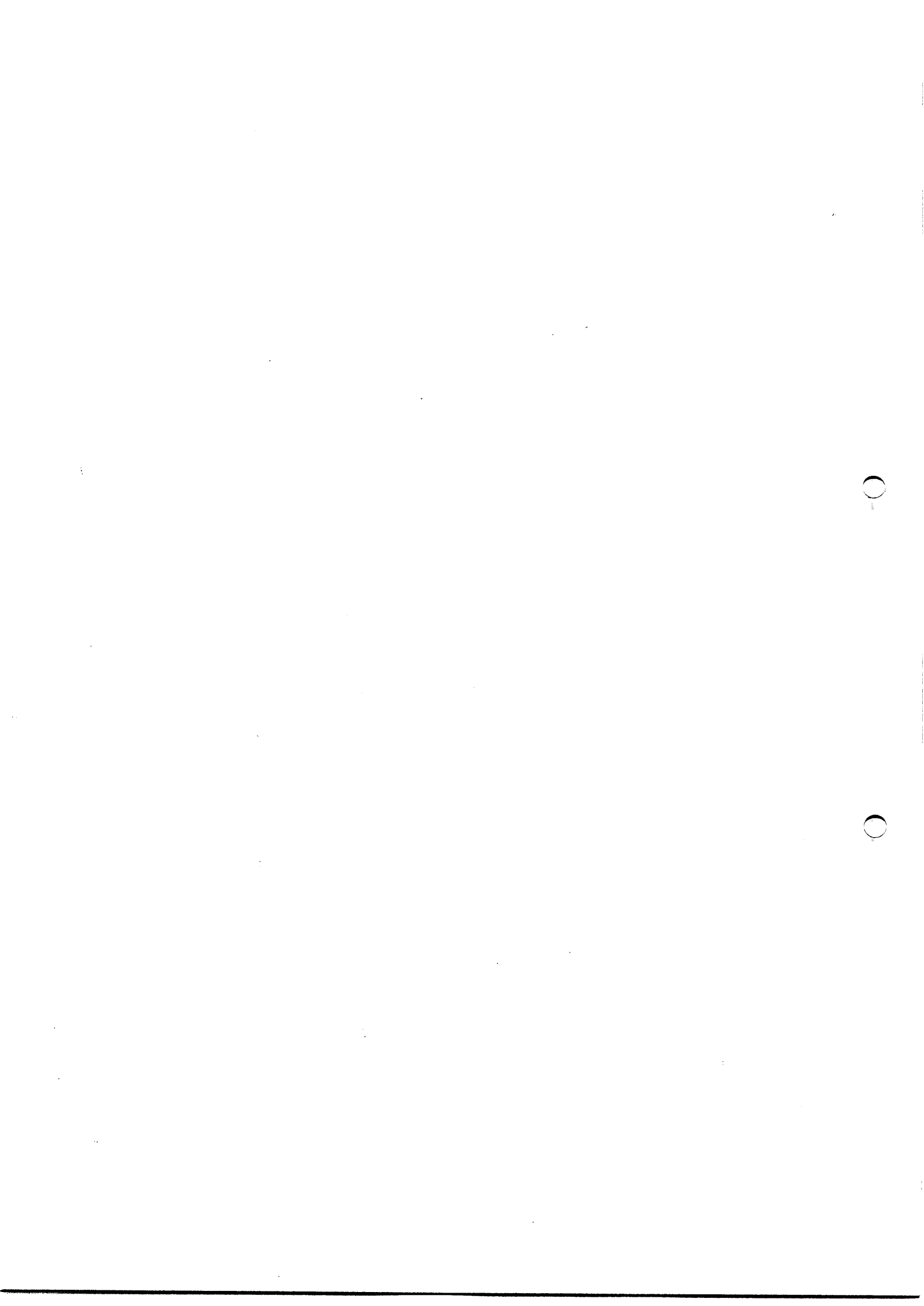
表4.1と表4.6、表4.7とを比較すると、従来法の3位認識率と、オープン実験の平均の1位認識率がほぼ同じになる。一番悪いもので比較すると従来法の2位認識率よりも劣るものもある。これは認識文書中の対象文字と、辞書を抽出した文書とが一致していないため、クローズ実験と同程度の認識結果が得られたのはいずれも、認識文書中の対象文字を全て含んだ抽出辞書を用いた場合であった。この様に単に社説1部から作成した新辞書を用いてパターンマッチング法のみで他の文書を認識した際には、従来法よりは良い結果が得られるが、クローズ実験にまでは及ばなかった。

## 4.6 更新した辞書による認識実験

4.5節の結果より、社説1部だけから作成した新辞書を用いた場合、その認識結果は、抽出に用いた社説中の文字の内容に依存してしまい、他の文書の認識に用いた場合にクローズ実験と同程度の結果が必ずしも得られないことがわかった。そこで、フォントの抽出用の社説(11~20)の10部を用いて抽出辞書を更新し、抽出辞書を十分に学習させて、他の社説(21~30)の10部に対して認識実験を行った。次に辞書の更新法を説明する。

### 4.6.1 辞書の更新方法

抽出辞書を更新させる場合には、学習させる順番が問題になるが、今回は社説の内容は考慮せず、社説の番号が若い順に学習させ更新することにした。更新に用いた社説数を学習文書数とする。例えば、社説11と社説12で抽出辞書を更新・作成した場合は学習文書数2とする。社説11~20まで全て用いた場合は学習文書数10となる。学習文書数と、抽出された文字数との関係は、図4.6、表4.10に示す。これらから、今回用いた更新法では、



学習文書数5で、ひらがな65種の全ての抽出辞書を獲得できたことがわかる。

- 社説(11~20)からフォントを抽出し、各社説毎に重心ベクトルを作成する。
- 社説の番号の若い順に、フォントが抽出された場合、その重心ベクトルを抽出辞書とする。
- 以後の社説でフォントが抽出された場合、抽出数が現在の抽出辞書のものより多い場合はその社説から抽出した重心ベクトルを抽出辞書とする。

$か_n$ が社説 $n$ から抽出されたフォントの重心ベクトルを表すものとして図4.1で説明する。最初は社説11から抽出された19個のフォントの重心ベクトルが抽出辞書になるが、次の社説12では34個が抽出されたので抽出辞書が $か_{12}$ に更新される。

ただし、現在の抽出辞書の抽出数が1のときに、次の社説の抽出数が1の場合は、現在の抽出辞書と、抽出されたベクトルの重心ベクトルを抽出辞書とする。これは3章での結果より、抽出数が1の場合、誤抽出している可能性がある。このために抽出数が少ない場合には抽出辞書に用いない方がより正確な抽出辞書を作成できるが、出現頻度の低い文字が正確に抽出されている場合も抽出辞書が作成できなくなるという弊害が生じる。そこで、抽出数が1のときは上記の処理を行う。ただし、この後に2以上の抽出数が得られた場合には、通常通りにその重心ベクトルが抽出辞書になる。図4.2中の $G(A, B)$ は、 $A$ と $B$ の重心ベクトルを示している。文書12でフォントの抽出があり、抽出数が1なのでこのベクトル( $ぼ_{12}$ )がそのまま抽出辞書になる。次に抽出があるのは文書16だが、抽出数が1であるので、この $ぼ_{16}$ と現在の抽出辞書の $ぼ_{12}$ の重心ベクトル $G(ぼ_{16}, ぼ_{12})$ が新しい抽出辞書になる。

#### 4.6.2 実験結果

1位認識率は表4.12に、学習文書数毎の1位認識率を表4.13、表4.14に示す。比較のために従来法による社説(21~30)までの認識結果を表4.11に示す。

マルチテンプレート法は学習文書数6で99.99%、置換法は学習文書数4以上で100%の1位認識率を達成できた。これはクローズ実験と同程度の結果であり、この辞書の更新による学習法が有効であると考えられる。



## 4.7 考察

初めに、3つの実験で正解を得られなかった文字について詳しく解析してみる。その理由は大きく次の3つに大別できる。

### 1. 正解のフォントが抽出されない場合

社説のフォントが抽出できないため、従来法で誤認識した文字をクローズ実験でも誤認識する例や、フォントを誤抽出し、誤って抽出したフォントのコードに誤認識する例がクローズ実験の候補数が少ない場合にみられた。

### 2. 正解のフォントが全て抽出していない場合

クローズ実験で誤認識した場合を例として挙げる。誤認識した文字は「ば」で、その社説中には2つの「ば」が存在していて、1つは正しく認識できたが他方は誤認識した。便宜上誤認識した方のベクトルを  $ば_f$ 、認識できた方のを  $ば_t$  で表す。 $ば_f$  は候補数1で候補集合に取り入れられたが詳細識別で除去されている。一方の  $ば_t$  は候補数3で候補集合に取り入れられ、そのまま  $ば_t$  だけが抽出されており、候補数2までは抽出がなく、候補数3以降では  $ば_t$  が抽出辞書として用いられることになる。 $ば_f$  との評価値は

$$E(ば_f, ば_f) < E(ば_f, ば_s) < E(ば_f, ば_t) < E(ば_f, ば_t) \quad (4.1)$$

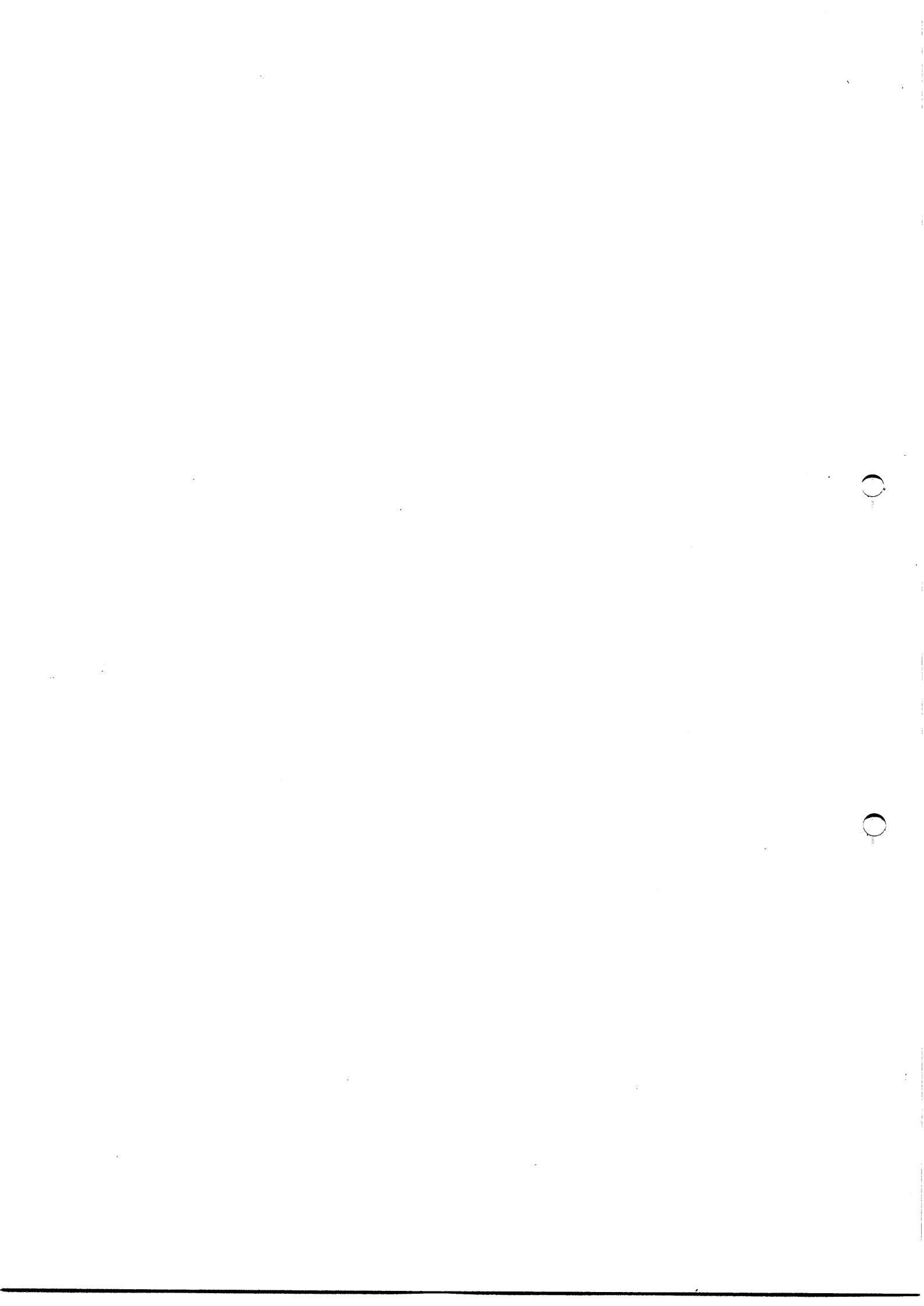
の順となる。このため候補数2まででは従来の既存の辞書のベクトルで正解が得られているが、候補数3以降での置換法では  $ば_s$  が  $ば_t$  に更新されてしまい、誤認識している。また辞書の更新実験で不正解になった文字の場合は、逆にマルチテンプレート法でだけ誤認識している。これは「し」で、抽出辞書を  $し_n$  と表すと、

$$E(し_f, し_s) < E(し_f, し_s) < E(し_f, し_n) < E(し_f, し_n) \quad (4.2)$$

の順になるため誤認識した例である。この場合でも、抽出辞書には文書中の全ての「し」のフォントが抽出されないため、正しく認識できなかったものである。これらを正しく認識させるには、類似文字と詳細識別類似文字の学習、さらに細分類での閾値の値を十分に検討しなければならない。

### 3. 類似文字も対象の文字の場合

「わ」が「ね」に誤認識する場合の例である。「わ」と「ね」は相互も他方の詳細識別類似文字になっている。このため「わ」が抽出されず「ね」だけが抽出されている場合に評価値は



$$E(w_f, n_s) < E(w_f, w_s) < E(w_f, n_n) \quad (4.3)$$

の順になり、置換法の場合は  $n_s$  が  $n_n$  に更新されているので正解を得るが、マルチテンプレート法では不正解となってしまう。

次に、誤って抽出された「液」、「点」、「ザ」からそれぞれ作成された「ゆ」、「ぬ」、「ぜ」の抽出辞書を用いた認識実験での影響について述べる。「ザ」を除く2つの場合は候補数10のときに誤抽出しているため、クローズ実験でしか評価できず、それぞれの社説中には「ゆ」、「ぬ」が存在していないため、従来法では正しく認識できた「液」、「点」をマルチテンプレート法、置換法とも「ゆ」、「ぬ」に誤認識している。一方、「ザ」の場合は社説(11~20)中に「ザ」は誤抽出されたものだけで、「ぜ」は9文字が存在していた。この9文字に対する認識結果ではマルチテンプレート法は全てを正解し、置換法では9文字中4文字を「せ」と誤認識していた。この様に、マルチテンプレート法では誤抽出した場合でも既存辞書のベクトルを認識に使用できるので救済できるが、置換法では既存辞書のベクトルを削除するため、実際は「ザ」であるものを「ぜ」の標準パターンとして用いなければならなくなり、誤認識が生じる。しかし、9文字の「ぜ」を全て誤認識するのではなく、5文字がこの条件で正解を得たということは、同じフォントの類似文字(除去候補文字)が認識に及ぼす影響が大きいことを示している。

続いて、フォント抽出の対象ではない文字に対する影響について考察する。今回の実験では、フォント抽出の対象文字(ひらがな65種)の認識精度を従来法と比較して向上させることが目的であった。実際、1位認識率で比較すれば従来法よりも認識精度は上がっているが、対象文字以外の認識結果が従来法よりも大きく下回れば問題になる。従来法で正解だったもののうち、今回の実験で不正解になったものは、マルチテンプレート法、置換法ともクローズ実験では4字(社説10部中)、オープン実験では25字(社説10部×9セット中)であった。辞書の更新による実験では、どちらも従来法と同じ認識結果になった。誤認識した理由は、フォントの抽出の失敗によるものが一部で、ほとんどはひらがな65種の除去候補文字になっているものが、そのひらがなと誤認識してしまっている。カタカナの「リ」がひらがなの「り」に誤認識するのがこの例になる。この結果より、抽出対象とする文字の拡張を検討しなければならない。特に、対象文字の除去候補文字・詳細識別文字となっている文字を抽出の対象とすることで、対象文字以外への影響を緩和できると思われる。

今迄述べてきたように誤認識した文字はフォントの抽出に何らかの問題があると考えられる。今回用いたフォントの抽出法は3の考察でも触れたが、全てのフォントを正確に抽出するには至っていない。候補集合に選出された全ての正解のフォントを、全て抽出して





いないものが全体の2割あり、これらから作成した抽出辞書で誤認識したものが数例確認された。特に「し」と「じ」の様な類似文字の認識の際に、この影響が顕著であり、このままでは抽出辞書を用いてもパターンマッチング法のみでは完全に類似文字の識別を行うことは難しい。正確にフォントの抽出を行うために、除去候補文字・詳細識別文字、細分類での閾値を検討しなければならない。

また、辞書の作成法に関しては、マルチテンプレート法と、置換法で対象文字に関して認識結果に差が出たのもフォントの抽出の影響であると考えられる。上記したフォントの抽出法の問題点が解決し、実際に文書認識を行う際には、誤認識の影響を考え、マルチテンプレート法が有効であると考えられる。

0

0

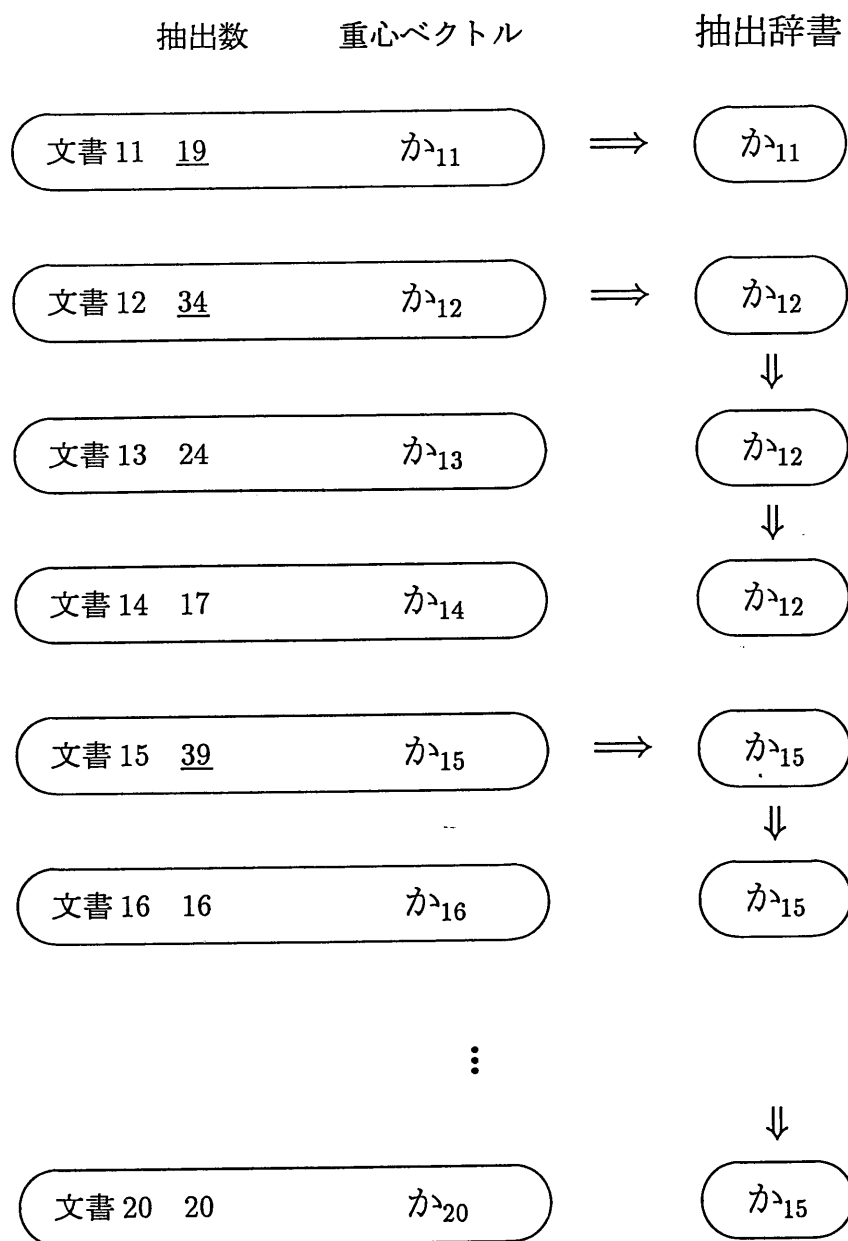
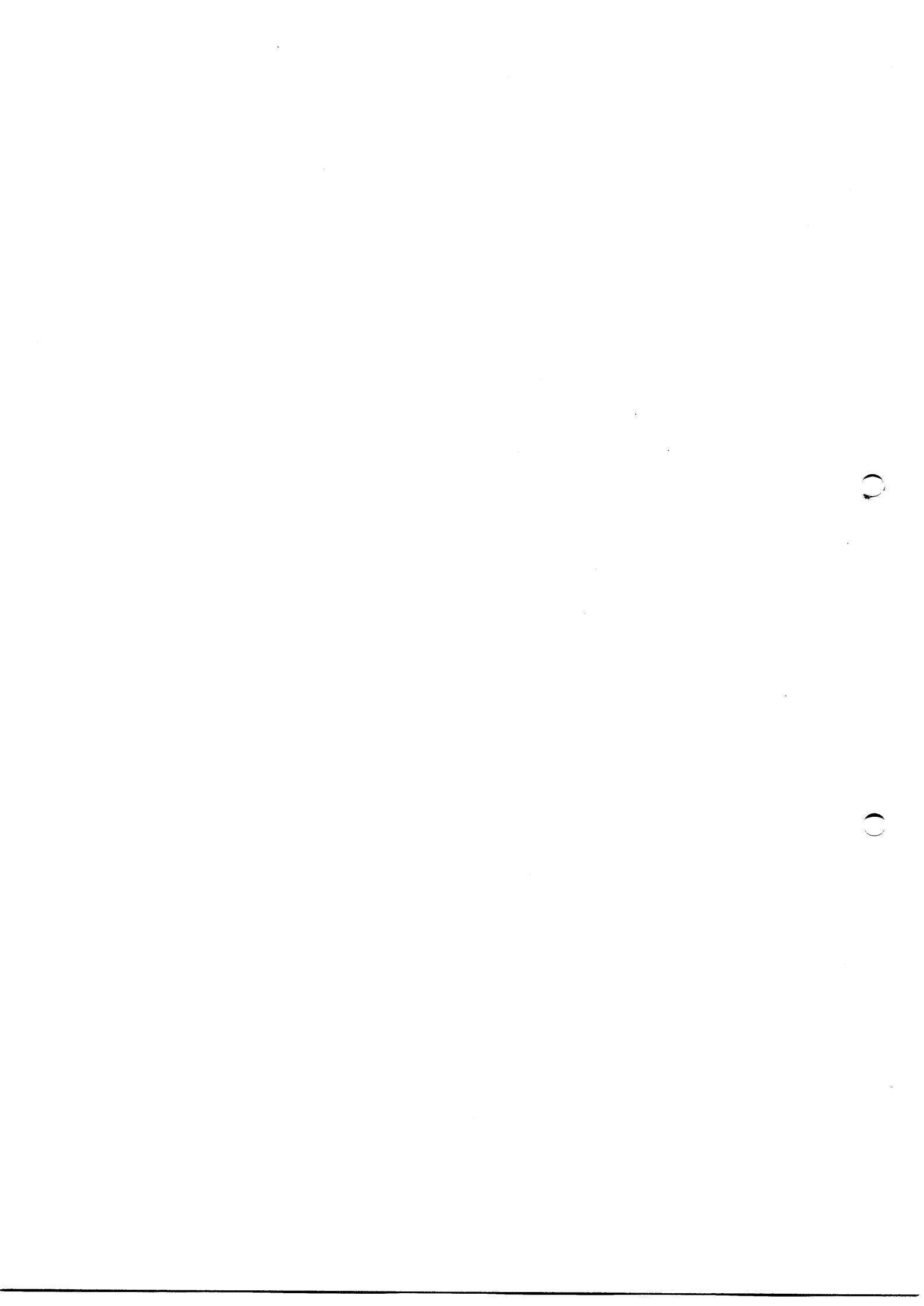


図 4.1: 辞書の更新の例 1



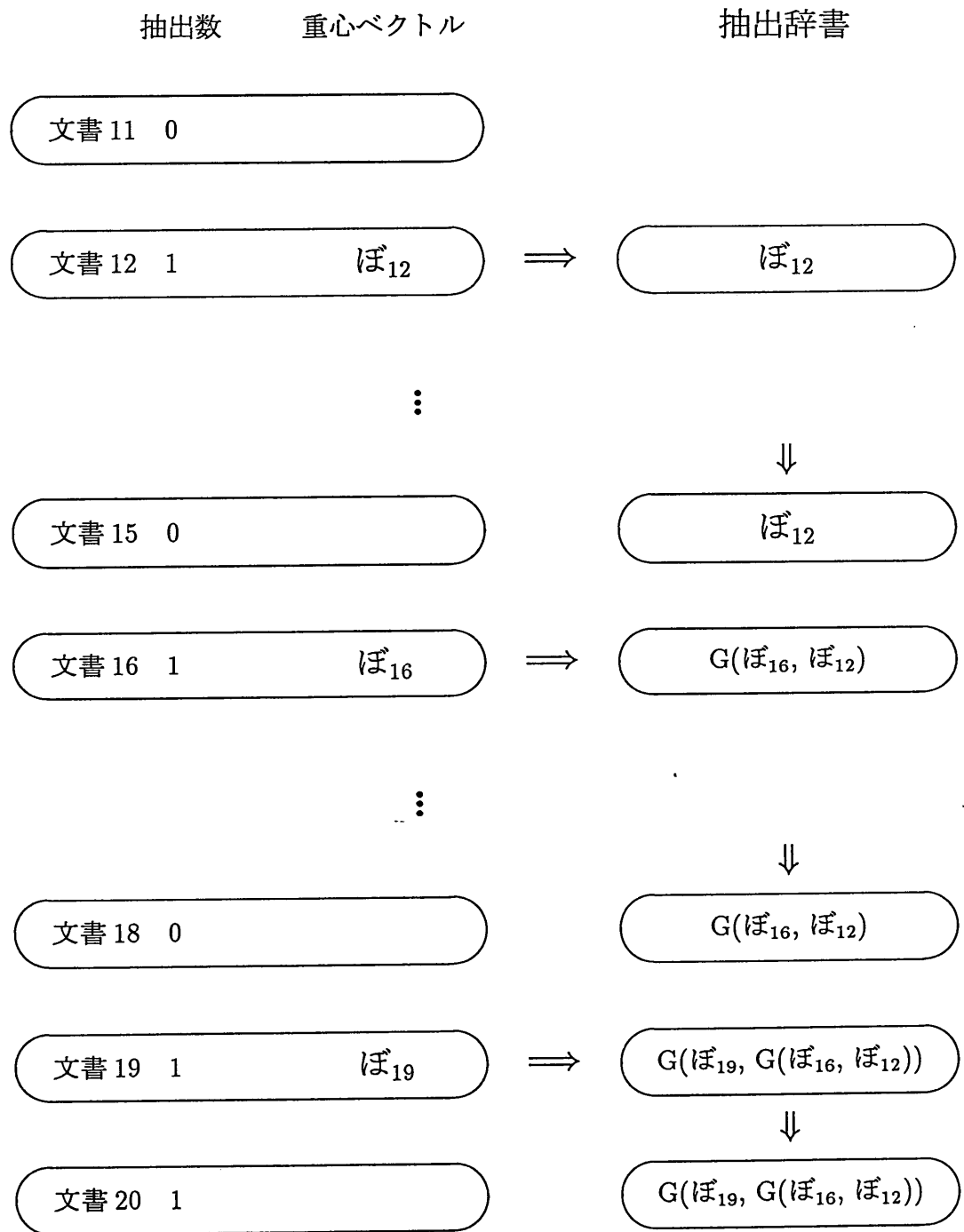


図 4.2: 辞書の更新の例 2

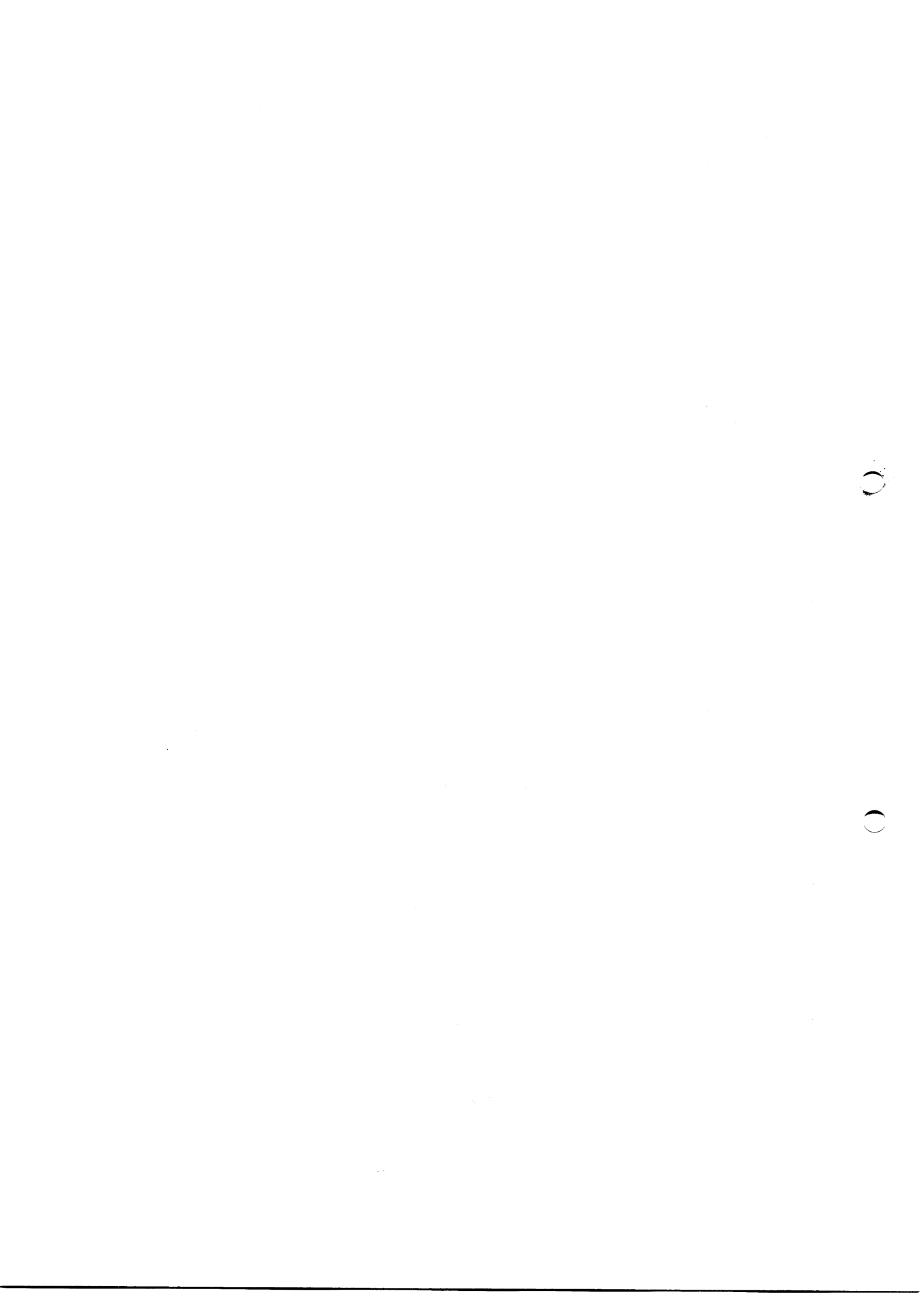


表 4.1: 従来の辞書の認識率 (社説 11~20)

社説	ひらがな 65 種		
	1 位認識率 (%)	2 位認識率 (%)	3 位認識率 (%)
平均	96.34 (37)	99.37 (6)	99.88 (1)
11	96.77 (34)	99.71 (3)	99.90 (1)
12	95.90 (43)	99.33 (7)	99.81 (2)
13	96.18 (34)	99.10 (8)	99.89 (1)
14	95.80 (41)	99.18 (8)	99.90 (1)
15	96.55 (39)	99.29 (8)	99.91 (1)
16	96.50 (35)	99.50 (5)	100.00 (0)
17	95.85 (45)	99.45 (6)	100.00 (0)
18	97.07 (32)	99.63 (4)	99.82 (2)
19	96.46 (36)	99.51 (5)	99.90 (1)
20	96.44 (31)	98.85 (10)	99.66 (3)

社説	その他		
	1 位認識率 (%)	2 位認識率 (%)	3 位認識率 (%)
平均	93.11 (70)	97.97 (21)	98.86 (12)
11	91.83 (83)	97.93 (21)	98.82 (12)
12	93.42 (68)	97.87 (22)	98.55 (15)
13	93.20 (67)	98.17 (18)	98.78 (12)
14	94.42 (59)	97.73 (24)	99.05 (10)
15	92.23 (74)	97.69 (22)	98.53 (14)
16	92.03 (80)	97.81 (22)	99.30 (7)
17	93.97 (62)	98.93 (11)	99.61 (4)
18	93.62 (64)	97.81 (22)	98.50 (15)
19	92.71 (76)	97.70 (24)	98.66 (14)
20	93.56 (63)	98.06 (19)	98.77 (12)

() の内は誤認識数

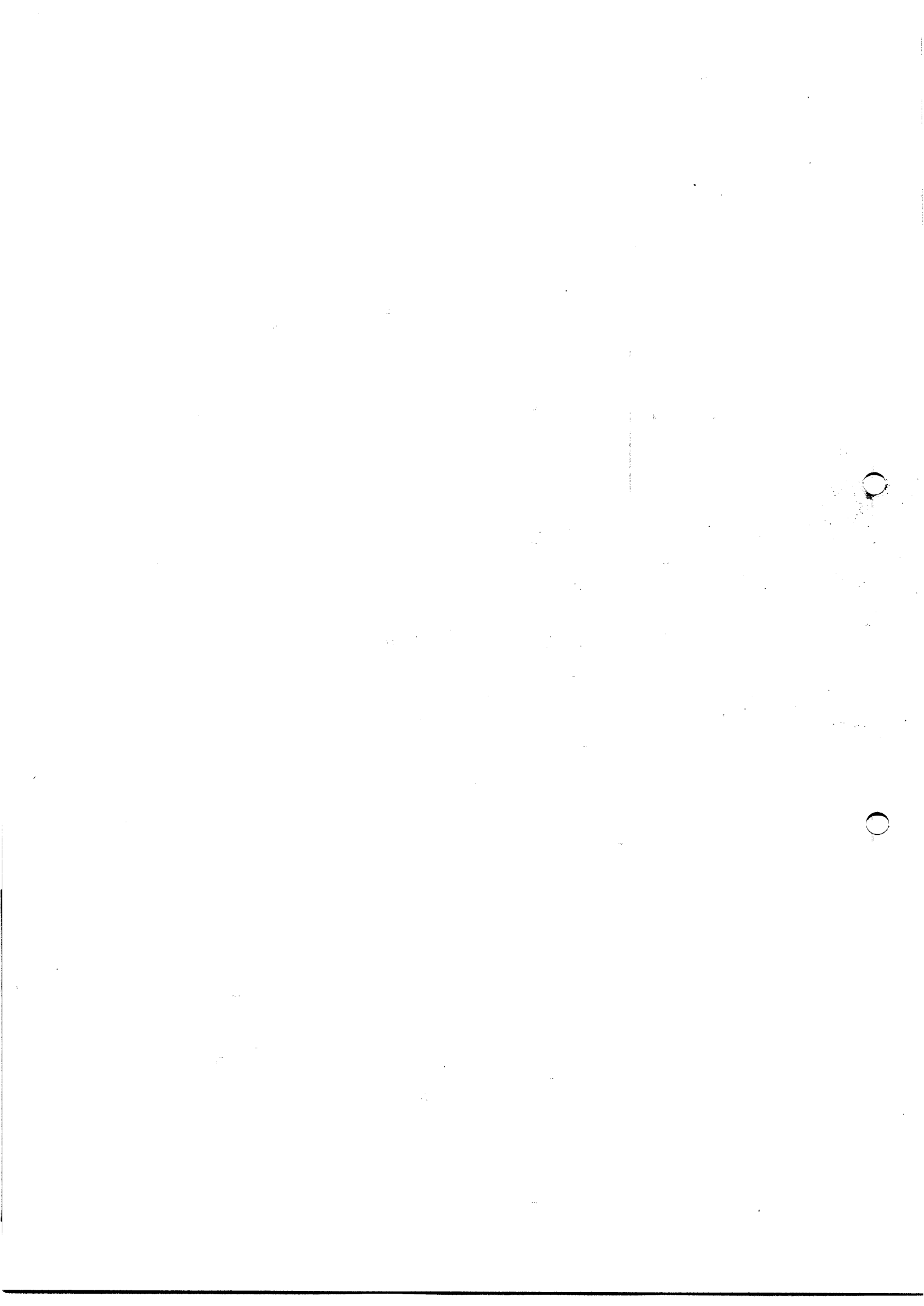




表 4.2: 1 位認識率 (クローズ実験)

候補数	マルチテンプレート法	
	ひらがな 65 種 (%)	その他 (%)
1	94.46 (157)	93.10 (697)
2	99.90 (10)	93.09 (698)
3	99.96 (6)	93.08 (699)
4	99.99 (1)	93.08 (699)
5	99.99 (1)	93.08 (699)
6	100.00 (0)	93.08 (699)
7	100.00 (0)	93.08 (699)
8	100.00 (0)	93.08 (699)
9	100.00 (0)	93.08 (699)
10	100.00 (0)	93.05 (702)

候補数	置換法	
	ひらがな 65 種 (%)	その他 (%)
1	99.09 (92)	93.09 (698)
2	99.90 (10)	93.08 (699)
3	99.95 (5)	93.07 (700)
4	99.98 (2)	93.07 (700)
5	99.98 (2)	93.07 (700)
6	99.99 (1)	93.07 (700)
7	99.99 (1)	93.07 (700)
8	99.99 (1)	93.07 (700)
9	99.99 (1)	93.07 (700)
10	99.99 (1)	93.07 (700)

() の内は誤認識数



表 4.3: 改善の結果 (マルチテンプレート法)

候補数	全体		ひらがな 65 種				その他			
	A	C	A	B	C	D	A	B	C	D
1	213	2	213	1	0	33	0	0	2	0
2	360	3	360	0	0	7	0	1	3	0
3	366	4	366	0	0	4	0	1	4	0
4	369	4	369	0	0	1	0	1	4	0
5	369	4	369	0	0	1	0	1	4	0
6	370	4	370	0	0	0	0	1	4	0
7	370	4	370	0	0	0	0	1	4	0
8	370	4	370	0	0	0	0	1	4	0
9	370	4	370	0	0	0	0	1	4	0
10	370	7	370	0	0	0	0	0	7	0

A: 誤認識が正しく認識された数

B: A 以外で認識順位が上がった数

C: 正解から誤認識へ変った数

D: C 以外で認識順位が下がった数

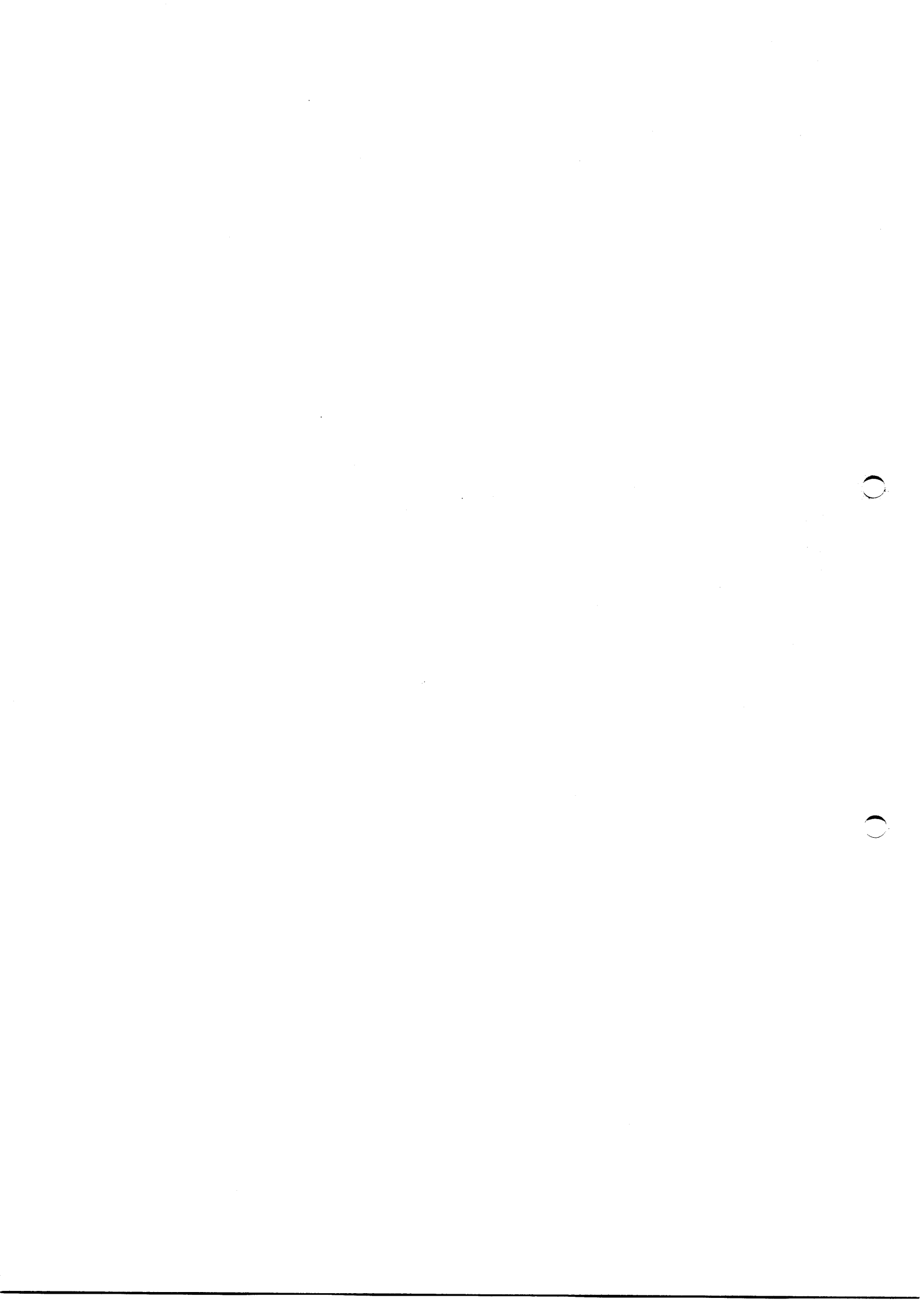


表 4.4: 改善の結果 (置換法)

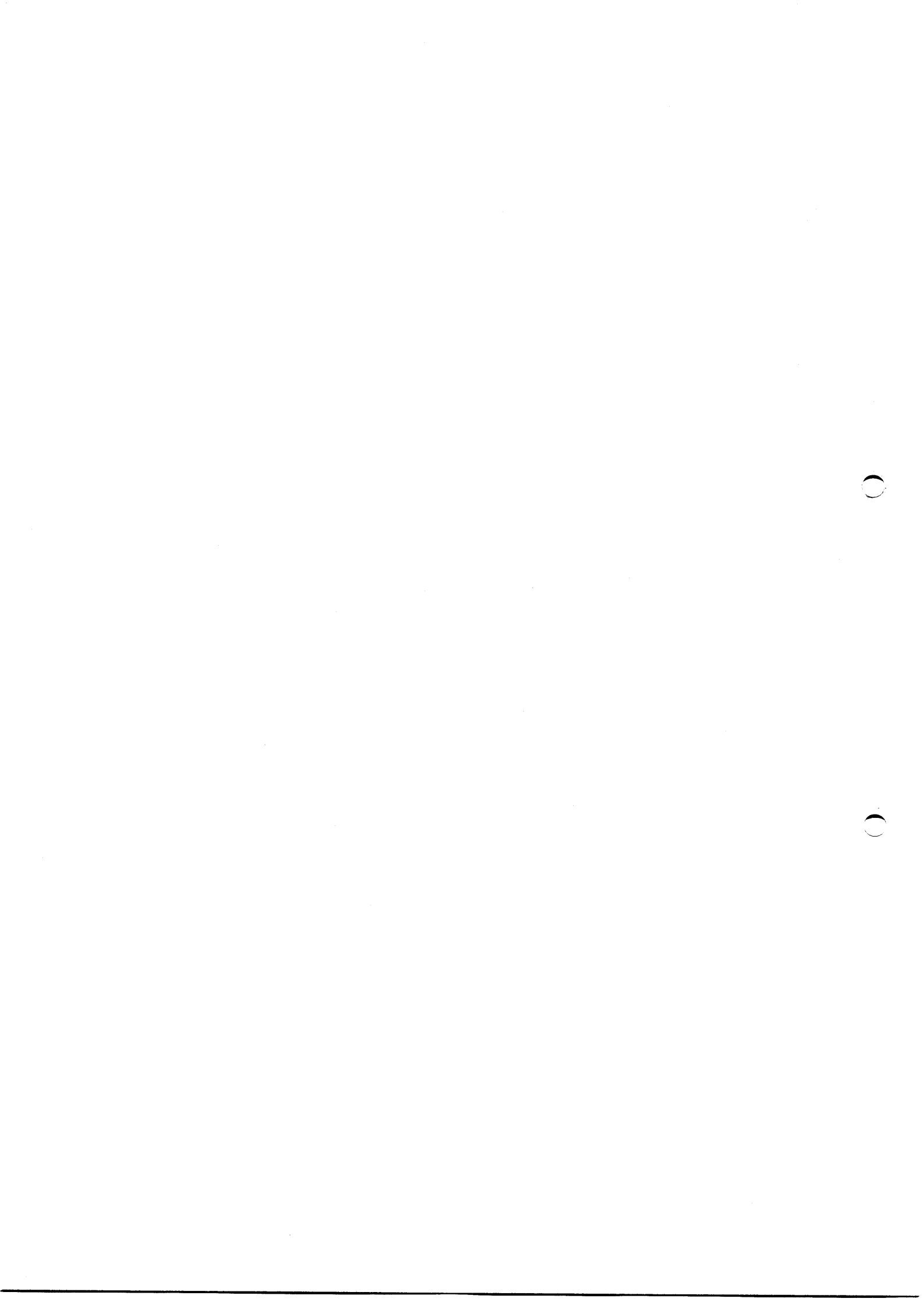
候補数	全体		ひらがな 65 種				その他			
	A	C	A	B	C	D	A	B	C	D
1	278	2	278	5	0	12	0	0	2	0
2	360	4	360	4	0	1	0	1	3	0
3	366	5	366	3	1	0	0	1	4	0
4	369	5	369	1	1	0	0	1	4	0
5	369	5	369	1	1	0	0	1	4	0
6	370	5	370	0	1	0	0	1	4	0
7	370	5	370	0	1	0	0	1	4	0
8	370	5	370	0	1	0	0	1	4	0
9	370	5	370	0	1	0	0	1	4	0
10	370	8	370	0	1	0	0	1	7	0

A: 誤認識が正しく認識された数

B: A 以外で認識順位が上がった数

C: 正解から誤認識へ変った数

D: C 以外で認識順位が下がった数



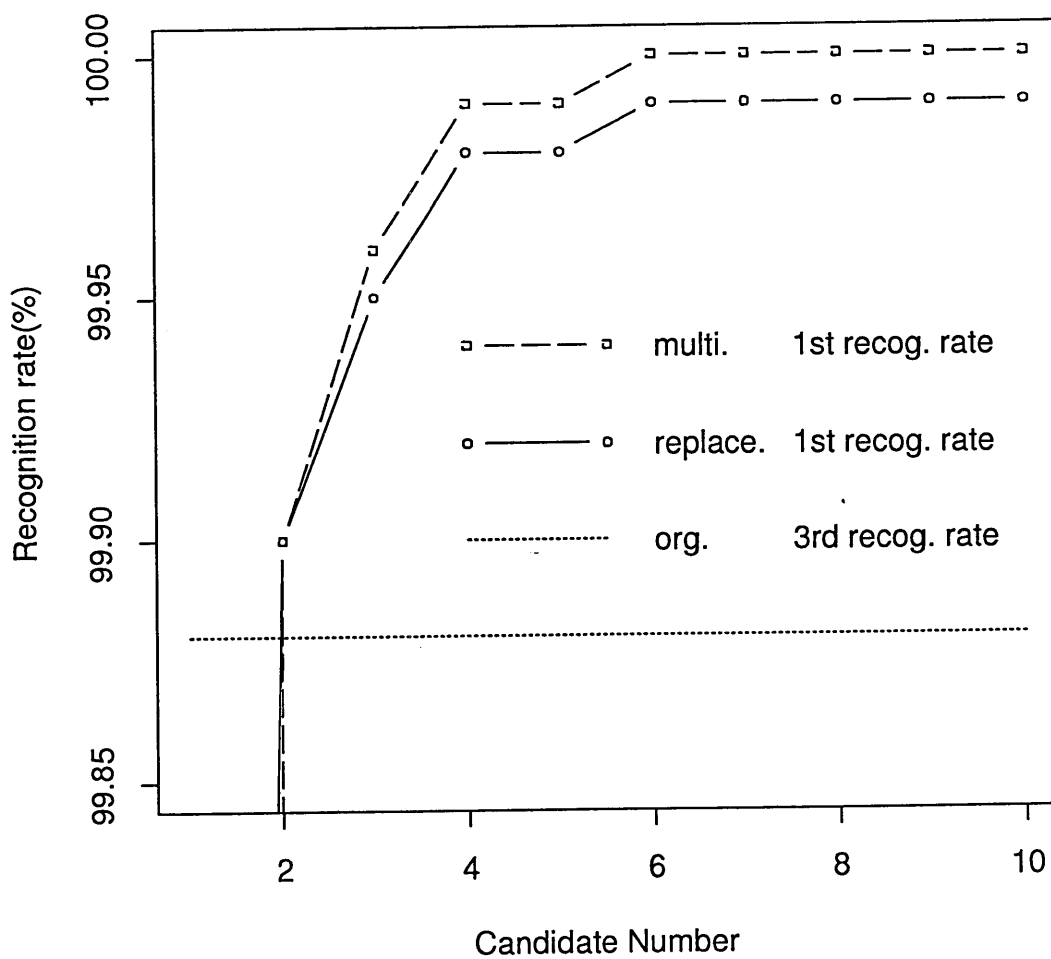
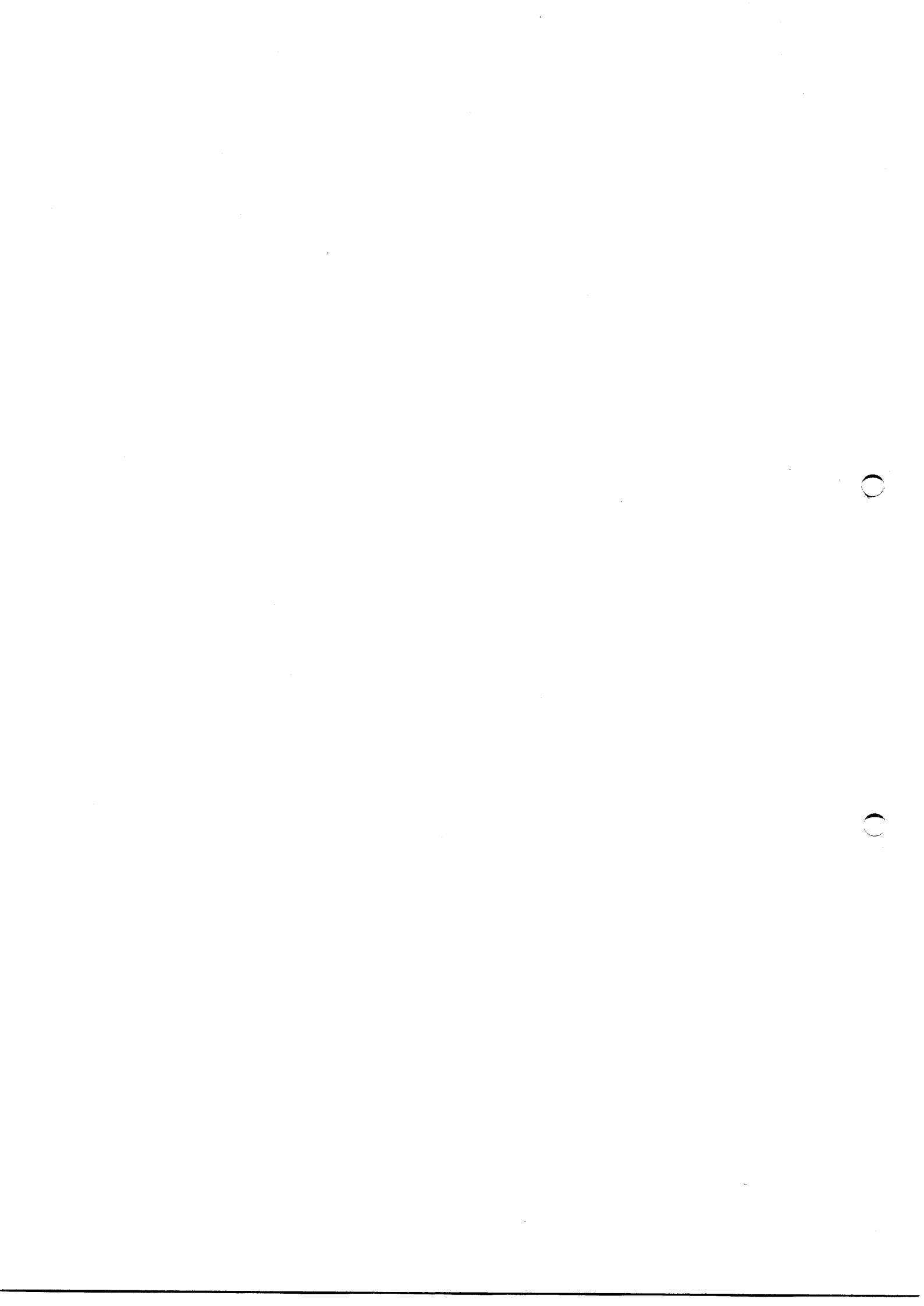


図 4.3: 候補数と認識率 (ひらがな 65 文字)





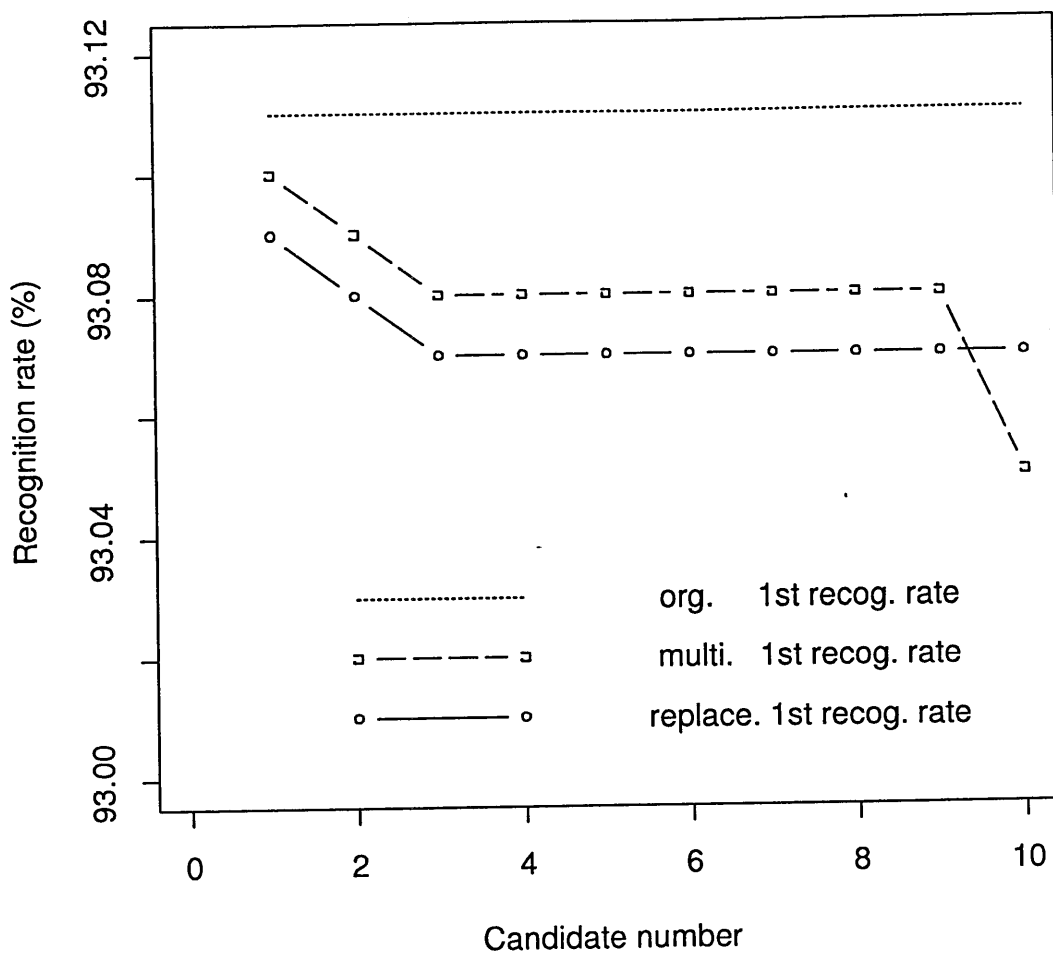


図 4.4: 候補数と認識率 (その他)



表 4.5: 認識文書中に存在しないひらがな 65 種の文字

社説	ひらがな 65 種で文書中に存在しない文字
11	ご ざ ぞ へ ぼ ゆ
12	ご ざ ぜ ぞ づ ね ひ ふ
13	お ぎ げ ご ぜ ぞ づ ぬ ね ふ ぼ ゆ
14	ぐ じ ぜ ぞ ぬ ふ ぶ ぼ
15	ぬ ふ ぼ
16	げ ぜ ぞ ぬ ゆ
17	げ ご ぞ づ び ふ ぶ ぼ ゆ
18	ざ ぜ ぞ ふ ほ ぼ ゆ
19	ぎ ご ぜ ぬ ひ ほ む ん
20	ご ぜ ぞ ぬ ね ほ ぼ む ゆ
21	ご ざ ぜ ぞ ぬ ね ひ ふ ぶ ぼ ゆ
22	ぐ ぞ ひ ふ ぶ ほ ぼ
23	お げ ご ぜ ぬ ひ ほ ゆ
24	ぜ ぞ ぬ ひ び ふ ほ
25	ざ ぜ ぬ ひ ふ ぶ む ろ
26	ぐ ご ぞ ぬ ね び ぶ む ゆ ん
27	ぞ ぬ ゆ
28	ざ せ づ ぬ ね ひ び ふ ぼ ゆ
29	ご ざ ぜ ぞ ね び む ゆ
30	お ぎ ご ぜ ぞ ぬ ひ へ



表 4.6: マルチテンプレート法の1位認識率(オープン実験)

社説	ひらがな 65 種			その他		
	A	B	C	A	B	C
11	100.00	99.43~100.00	99.79	91.83	91.83	91.83
12	100.00	99.24~100.00	99.69	93.42	93.32~93.42	93.41
13	100.00	99.21~100.00	99.79	93.20	93.20	93.20
14	100.00	99.49~100.00	99.68	94.32	94.32~94.42	94.41
15	100.00	99.47~99.91	99.71	92.12	92.12	92.12
16	100.00	99.60~100.00	99.79	92.03	91.93~92.03	92.02
17	100.00	99.63~100.00	99.89	93.97	93.97	93.97
18	100.00	99.54~100.00	99.81	93.72	93.62~93.72	93.66
19	100.00	99.51~100.00	99.74	92.71	92.71	92.71
20	100.00	99.77~100.00	99.87	93.36	93.36~93.46	93.42

A: 認識文書から辞書を抽出(クローズ)

B: 認識文書以外から辞書を抽出したときの min~max

C: B の平均

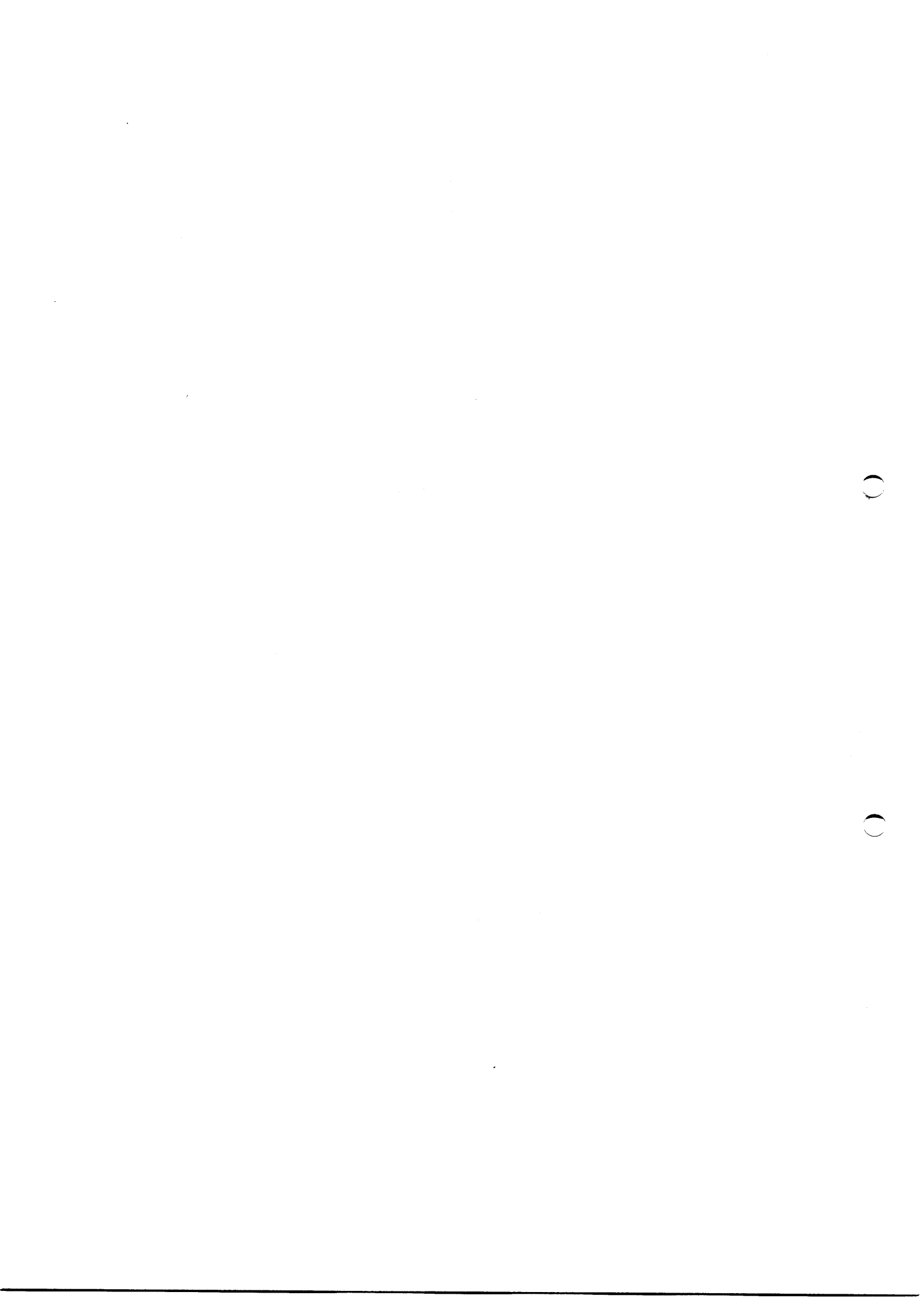


表 4.7: 置換法の1位認識率 (オープン実験)

社説	ひらがな 65 種			その他		
	A	B	C	A	B	C
11	100.00	99.24~100.00	99.76	91.83	91.83	91.83
12	100.00	99.24~100.00	99.67	93.42	93.32~93.42	93.41
13	100.00	99.21~100.00	99.79	93.20	93.20	93.20
14	99.90	99.49~99.90	99.61	94.32	94.32~94.42	94.41
15	100.00	99.20~99.91	99.68	92.12	92.12	92.12
16	100.00	99.65~100.00	99.79	92.03	91.93~92.03	92.02
17	100.00	99.54~100.00	99.87	93.97	93.97	93.97
18	100.00	99.54~100.00	99.82	93.62	93.62~93.72	93.67
19	100.00	99.51~100.00	99.72	92.71	92.71	92.71
20	100.00	99.77~100.00	99.87	93.36	93.36~93.46	93.42

A: 認識文書から辞書を抽出 (クローズ)

B: 認識文書以外から辞書を抽出したときの min~max

C: B の平均

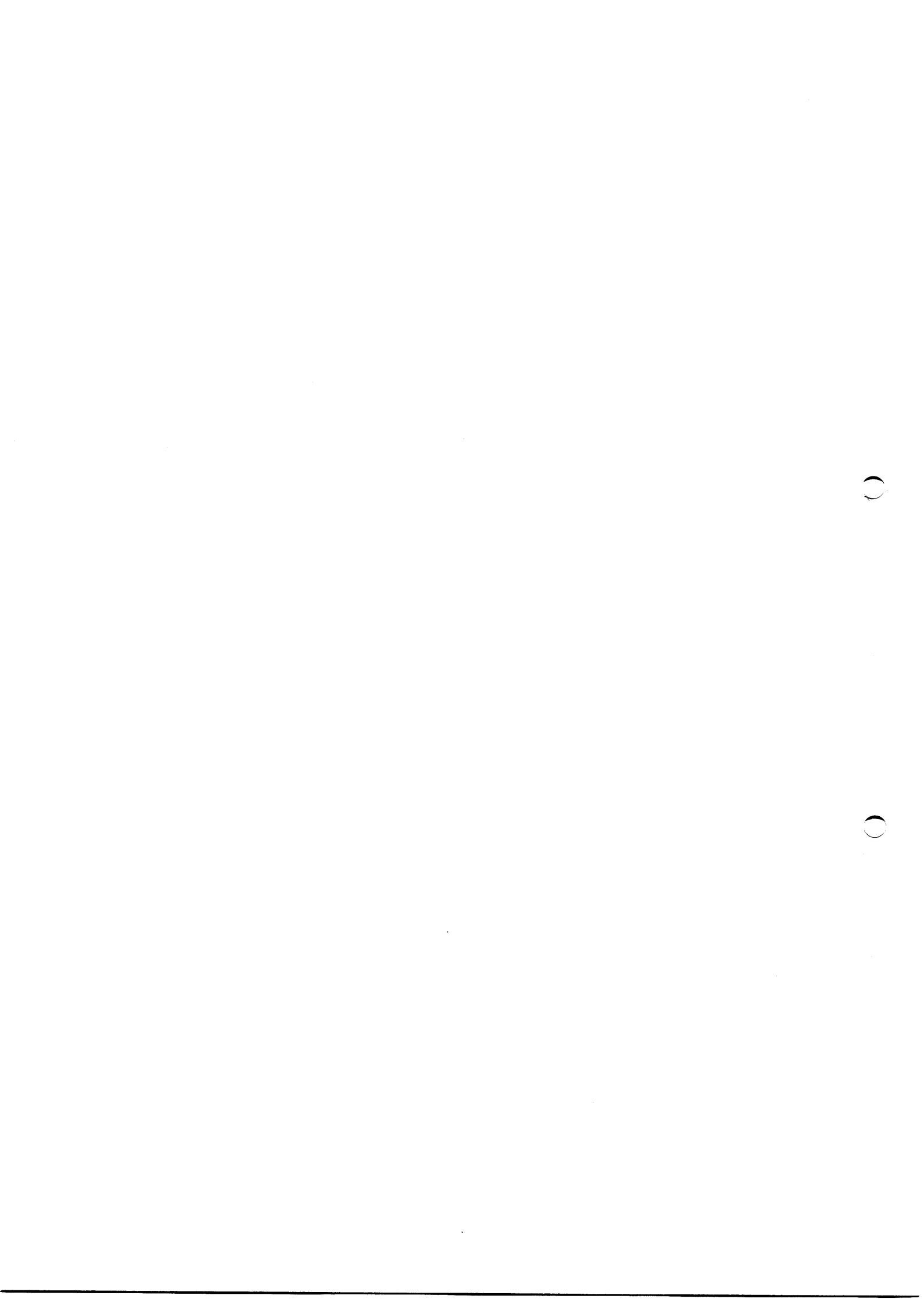




表 4.8: 抽出フォント毎の 1 位認識率 (マルチテンプレート法)

認識 社説	フォント抽出社説									
	11	12	13	14	15	16	17	18	19	20
11	100.00	100.00	99.62	99.43	100.00	99.90	99.81	99.71	99.71	99.90
	91.83	91.83	91.83	91.83	91.83	91.83	91.83	91.83	91.83	91.83
12	99.81	100.00	99.24	99.81	100.00	99.62	99.52	99.71	99.71	99.81
	93.42	93.42	93.32	93.42	93.42	93.42	93.42	93.42	93.42	93.42
13	99.89	99.78	100.00	99.21	99.89	100.00	99.89	99.66	99.78	100.00
	93.20	93.20	93.20	93.20	93.20	93.20	93.20	93.20	93.20	93.20
14	99.49	99.59	99.59	100.00	100.00	99.90	99.59	99.59	99.59	99.80
	94.32	94.42	94.42	94.32	94.42	94.42	94.42	94.42	94.42	94.42
15	99.73	99.82	99.56	99.47	100.00	99.82	99.47	99.73	99.91	99.82
	92.12	92.12	92.12	92.12	92.12	92.12	92.12	92.12	92.12	92.12
16	99.80	99.80	99.80	99.60	100.00	100.00	99.80	99.70	99.70	99.90
	92.03	92.03	92.03	92.03	92.03	92.03	92.03	92.03	91.93	92.03
17	99.91	99.91	99.91	99.63	100.00	100.00	100.00	99.72	99.91	100.00
	93.97	93.97	93.97	93.97	93.97	93.97	93.97	93.97	93.97	93.97
18	99.91	99.91	99.63	99.82	100.00	99.82	99.54	100.00	99.82	99.82
	93.62	93.72	93.72	93.62	93.72	93.62	93.62	93.72	93.62	93.72
19	99.51	99.61	99.70	99.70	100.00	99.80	99.70	99.70	100.00	99.90
	92.71	92.71	92.71	92.71	92.71	92.71	92.71	92.71	92.71	92.71
20	99.89	99.89	99.77	99.89	100.00	99.89	99.77	100.00	99.77	100.00
	93.36	93.46	93.36	93.46	93.36	93.46	93.36	93.46	93.46	93.36

上段 ひらがな 65 種の 1 位認識率 (%)

下段 ひらがな 65 種以外の 1 位認識率 (%)

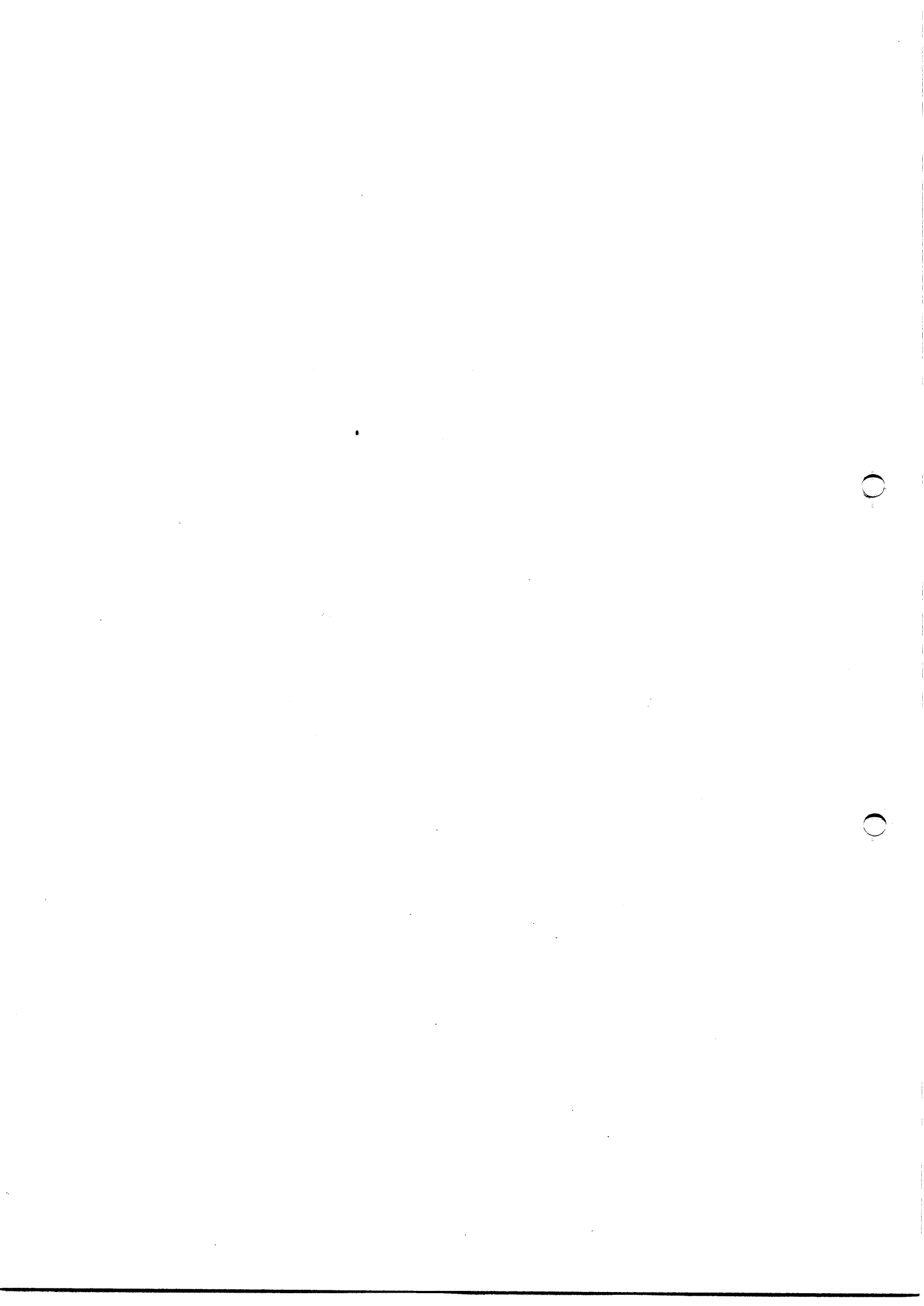
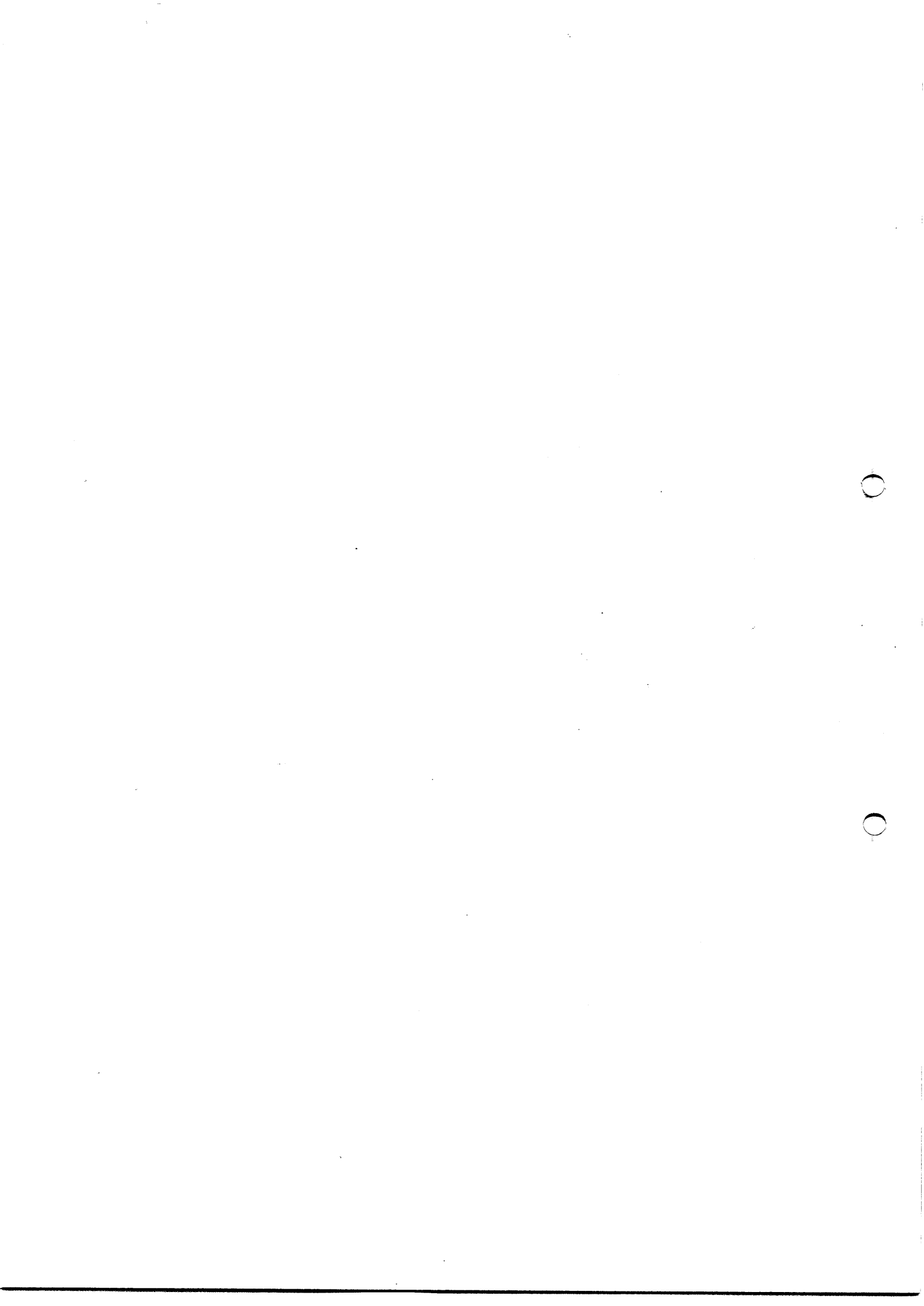


表 4.9: 抽出フォント毎の 1 位認識率 (置換法)

認識 社説	フォント抽出社説									
	11	12	13	14	15	16	17	18	19	20
11	100.00	100.00	99.62	99.24	100.00	99.81	99.81	99.71	99.71	99.90
	91.83	91.83	91.83	91.83	91.83	91.83	91.83	91.83	91.83	91.83
12	99.81	100.00	99.24	99.71	100.00	99.52	99.52	99.71	99.71	99.81
	93.42	93.42	93.32	93.42	93.42	93.42	93.42	93.42	93.42	93.42
13	99.89	99.78	100.00	99.21	99.89	100.00	99.89	99.66	99.78	100.00
	93.20	93.20	93.20	93.20	93.20	93.20	93.20	93.20	93.20	93.20
14	99.49	99.49	99.59	99.90	99.90	99.80	99.49	99.59	99.49	99.69
	94.32	94.42	94.42	94.32	94.42	94.42	94.42	94.42	94.42	94.42
15	99.73	99.82	99.56	99.20	100.00	99.82	99.47	99.73	99.91	99.82
	92.12	92.12	92.12	92.12	92.12	92.12	92.12	92.12	92.12	92.12
16	99.80	99.80	99.80	99.50	100.00	100.00	99.80	99.80	99.70	99.90
	92.03	92.03	92.03	92.03	92.03	92.03	92.03	92.03	91.93	92.03
17	99.91	99.91	99.91	99.54	100.00	100.00	100.00	99.63	99.91	100.00
	93.97	93.97	93.97	93.97	93.97	93.97	93.97	93.97	93.97	93.97
18	99.91	99.91	99.73	99.82	100.00	99.82	99.54	100.00	99.82	99.82
	93.72	93.62	93.62	93.72	93.62	93.72	93.62	93.62	93.72	93.72
19	99.51	99.61	99.70	99.51	100.00	99.80	99.70	99.70	100.00	99.90
	92.71	92.71	92.71	92.71	92.71	92.71	92.71	92.71	92.71	92.71
20	99.89	99.89	99.77	99.89	100.00	99.89	99.77	100.00	99.77	100.00
	93.36	93.46	93.36	93.46	93.36	93.46	93.36	93.46	93.46	93.36

上段 ひらがな 65 種の 1 位認識率 (%)

下段 ひらがな 65 種以外の 1 位認識率 (%)



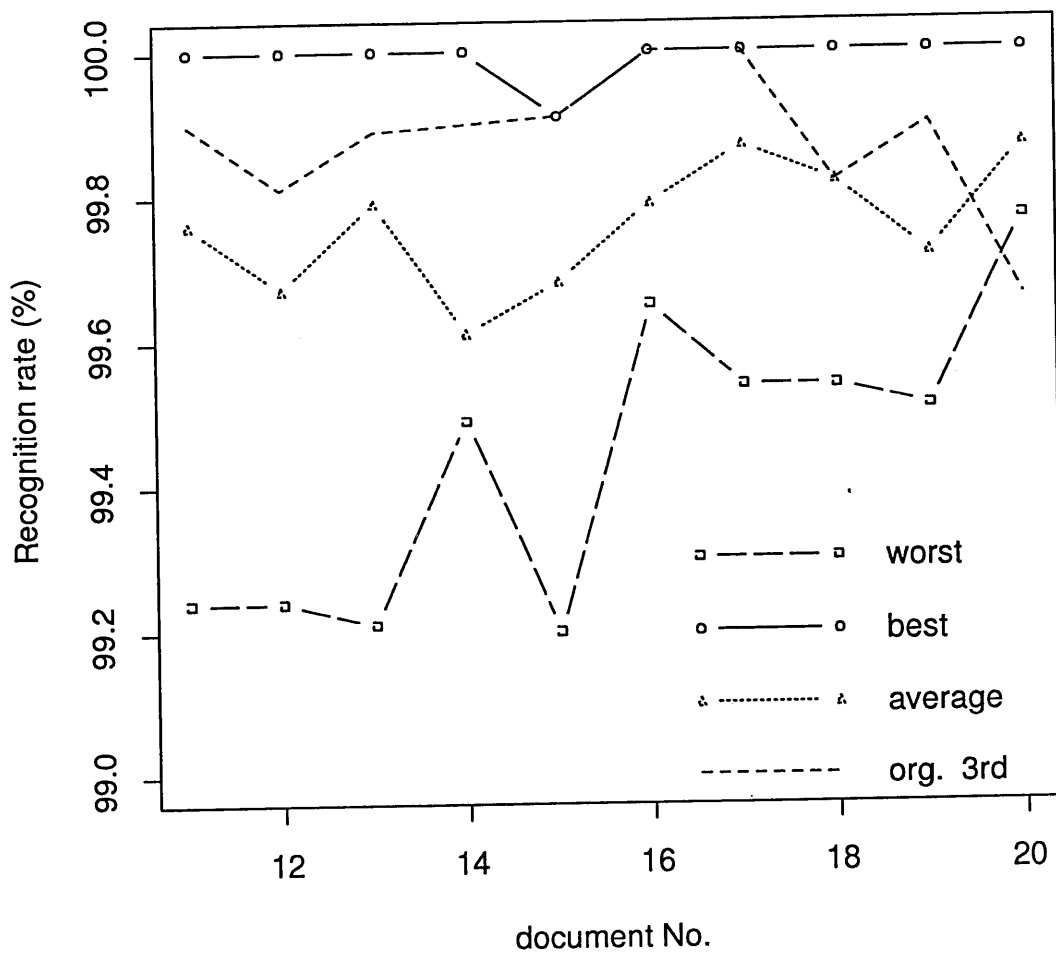


図 4.5: 認識文書毎の 1 位認識率



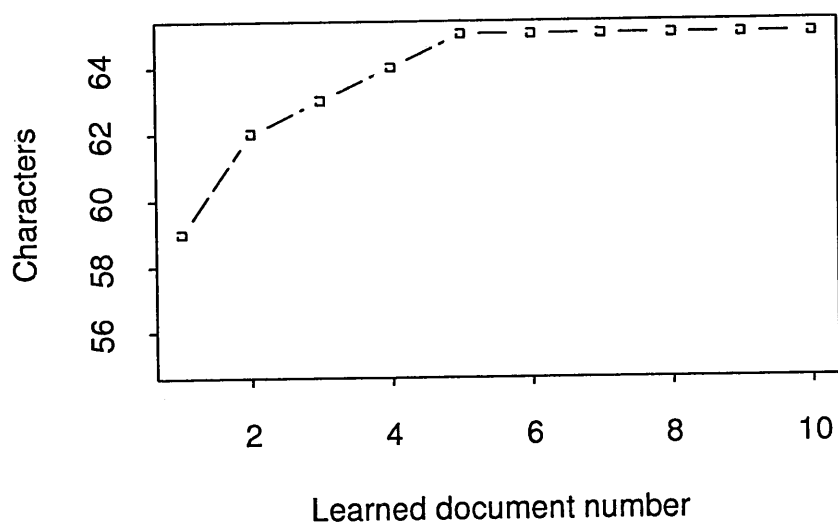


図 4.6: 学習文書数と抽出辞書中の文字

表 4.10: 抽出辞書にない文字

学習文書数	辞書にない文字
1	ご さ ぞ へ ぼ ゆ
2	ご さ ぞ
3	ご ぞ
4	ぞ

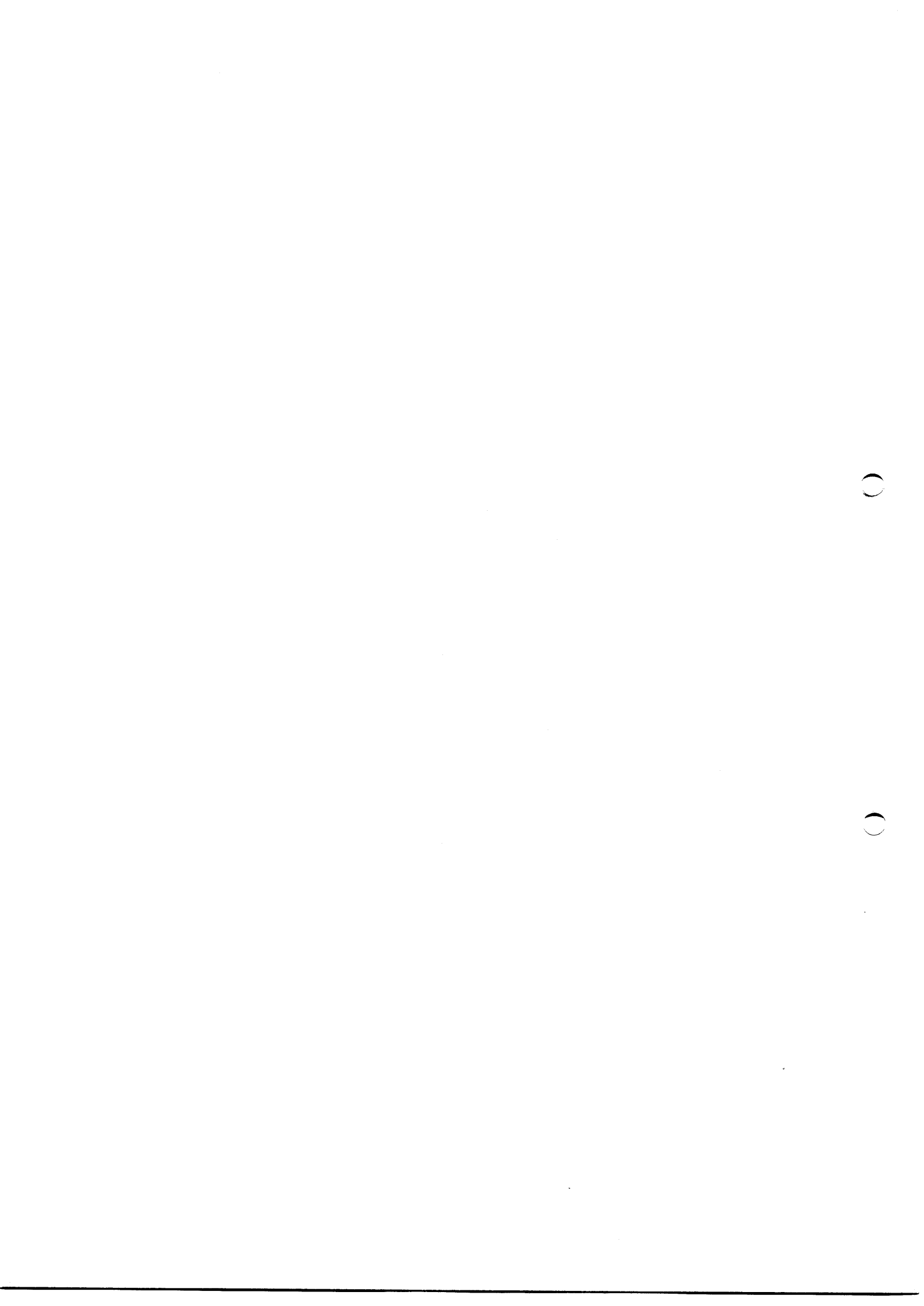




表 4.11: 従来の辞書の認識率 (社説 21~30)

社説	ひらがな 65 種		
	1 位認識率 (%)	2 位認識率 (%)	3 位認識率 (%)
平均	96.08 (41)	99.35 (7)	99.88 (3)
21	97.65 (26)	99.82 (2)	100.00 (0)
22	95.84 (39)	99.04 (9)	99.68 (3)
23	95.63 (53)	99.18 (10)	99.67 (4)
24	95.71 (44)	99.80 (2)	99.80 (2)
25	95.30 (47)	99.10 (9)	99.80 (2)
26	96.99 (33)	99.63 (4)	100.00 (0)
27	95.45 (52)	98.95 (12)	99.82 (2)
28	96.56 (35)	99.41 (6)	99.90 (1)
29	96.16 (40)	99.23 (8)	99.81 (2)
30	95.44 (42)	99.35 (6)	99.89 (1)

社説	その他		
	1 位認識率 (%)	2 位認識率 (%)	3 位認識率 (%)
平均	92.50 (78)	97.66 (24)	98.80 (13)
21	93.92 (60)	98.48 (15)	99.09 (9)
22	91.51 (98)	97.14 (33)	98.96 (12)
23	93.77 (53)	98.24 (15)	98.82 (10)
24	90.60 (99)	97.72 (24)	99.34 (19)
25	92.41 (81)	97.28 (29)	98.22 (19)
26	94.00 (63)	98.29 (18)	98.95 (11)
27	91.58 (81)	97.30 (26)	98.65 (13)
28	91.55 (84)	97.08 (29)	98.59 (14)
29	93.64 (71)	97.31 (30)	98.66 (15)
30	92.28 (106)	97.88 (25)	98.73 (15)

( ) の内は誤認識数



表 4.12: 1 位認識率 (辞書の更新による実験)

学 習 文 書 数	マルチテンプレート法	
	ひらがな 65 種 (%)	その他 (%)
1	99.66 (36)	92.50 (781)
2	99.76 (25)	92.50 (781)
3	99.90 (11)	92.50 (781)
4	99.98 (2)	92.50 (781)
5	99.98 (2)	92.50 (781)
6	99.99 (1)	92.50 (781)
7	99.98 (2)	92.50 (781)
8	99.98 (2)	92.50 (781)
9	99.98 (2)	92.50 (781)
10	99.98 (2)	92.50 (781)

学 習 文 書 数	置換法	
	ひらがな 65 種 (%)	その他 (%)
1	99.68 (34)	92.50 (781)
2	99.78 (23)	92.50 (781)
3	99.91 (9)	92.50 (781)
4	100.00 (0)	92.50 (781)
5	100.00 (0)	92.50 (781)
6	100.00 (0)	92.50 (781)
7	100.00 (0)	92.50 (781)
8	100.00 (0)	92.50 (781)
9	100.00 (0)	92.50 (781)
10	100.00 (0)	92.50 (781)

() の内は誤認識数



表 4.13: 学習文書数毎の1位認識率 (マルチテンプレート法)

社説	学習文書数										
	1	2	3	4	5	6	7	8	9	10	
21	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92
22	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51
23	99.75	99.75	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77
24	99.61	99.71	99.90	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60
25	99.70	99.80	99.80	99.90	99.90	100.00	99.90	99.90	99.90	99.90	99.90
	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41
26	99.63	99.82	99.91	99.91	99.91	99.91	99.91	99.91	99.91	99.91	99.91
	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00
27	99.30	99.30	99.82	100.00	100.00	100.00	100.00	99.72	99.91	100.00	100.00
	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58
28	99.80	99.80	99.80	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55
29	99.90	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64
30	99.57	99.89	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28

上段 ひらがな 65 種の 1 位認識率 (%)

下段 ひらがな 65 種以外の 1 位認識率 (%)

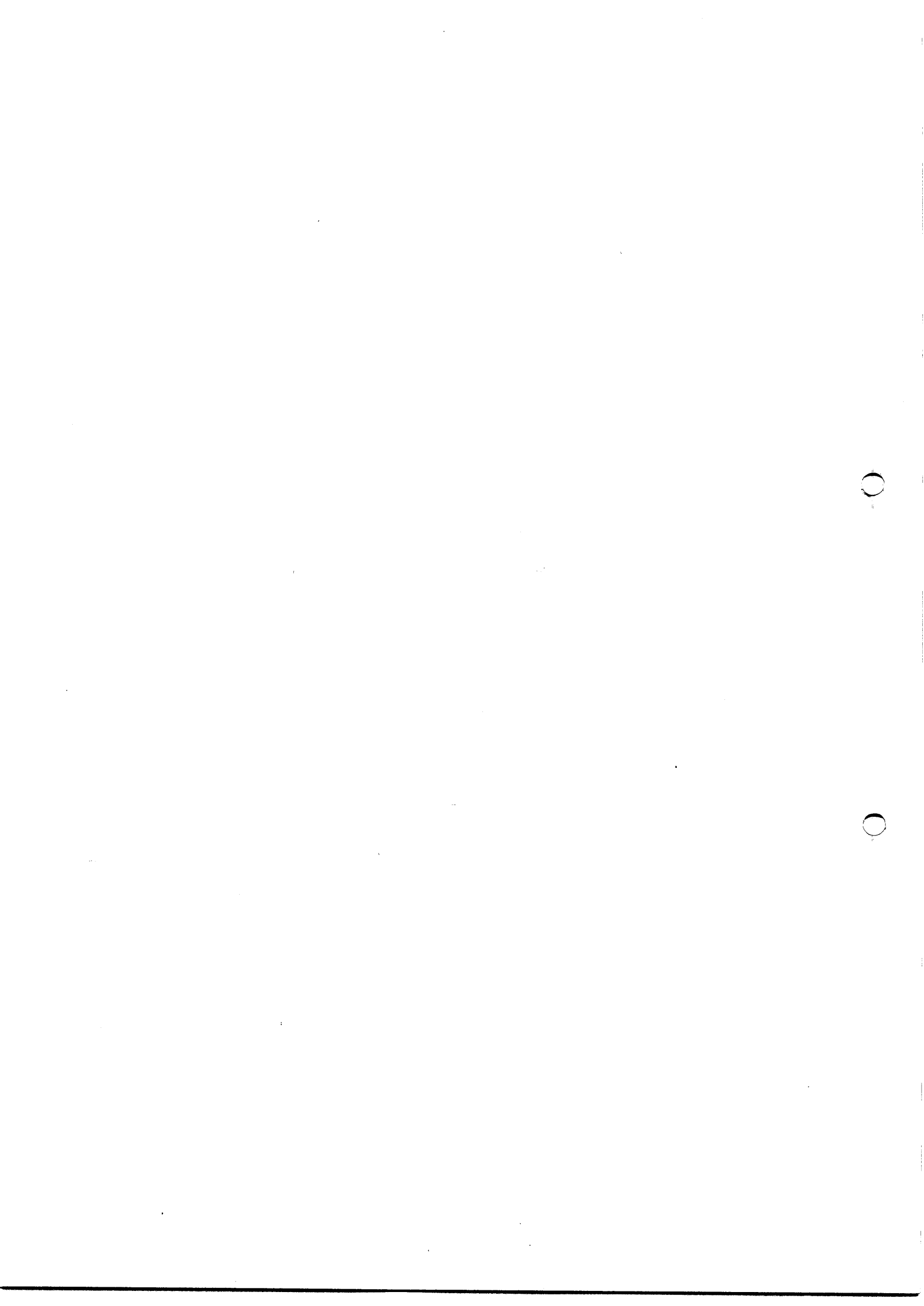


表 4.14: 学習文書数毎の 1 位認識率 (置換法)

社説	学習文書数										
	1	2	3	4	5	6	7	8	9	10	
21	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92	93.92
22	99.25	99.57	99.68	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51	91.51
23	99.75	99.75	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77	93.77
24	99.61	99.71	99.90	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60	90.60
25	99.80	99.90	99.90	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41	92.41
26	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00	94.00
27	99.30	99.30	99.82	100.00	100.00	100.00	100.00	99.72	99.91	100.00	100.00
	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58	91.58
28	99.80	99.80	99.80	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55	91.55
29	99.90	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64	93.64
30	99.57	99.89	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28	92.28

上段 ひらがな 65 種の 1 位認識率 (%)

下段 ひらがな 65 種以外の 1 位認識率 (%)





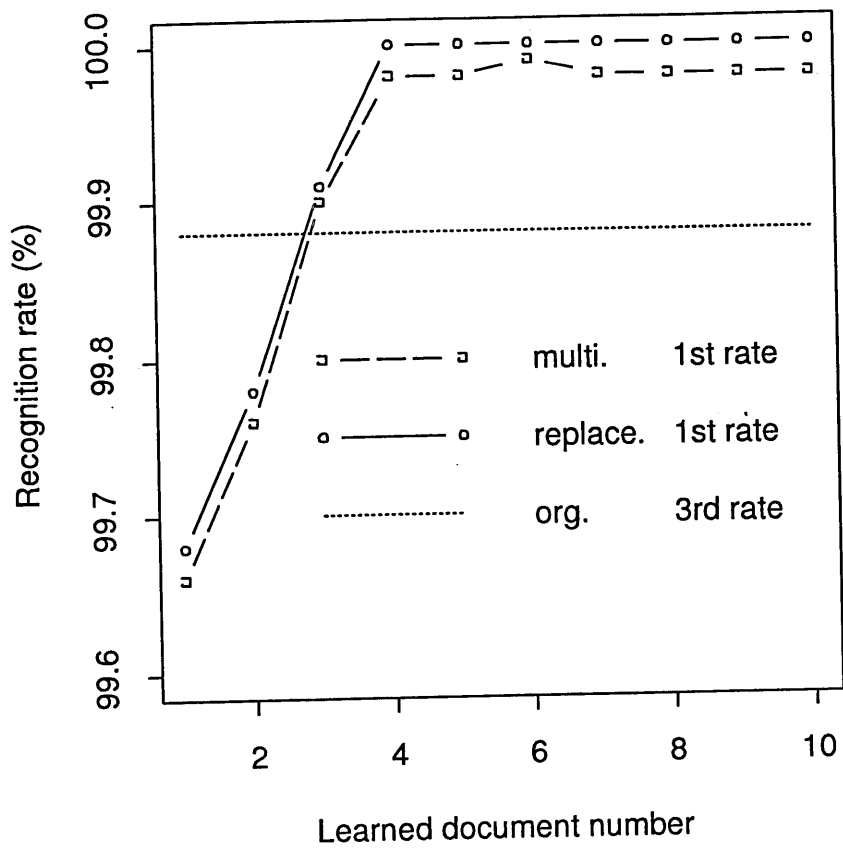


図 4.7: 学習文書数毎の 1 位認識率



## 第 5 章

### 結論

本研究では既存辞書を基にして、認識文書からフォントを抽出し、フォントの利用によって文書認識の高精度化を目的とした。その結果について記す。

#### 1. フォントの抽出

新聞社説 (1~10) の 10 部を用いて学習させた類似文字の利用により、候補数 6 までのデータの使用で新聞社説 (11~20) の 10 部から抽出対象となるひらがな 65 種を全て抽出することができ、今回提案したフォントの抽出のアルゴリズムが有効であることがわかった。

#### 2. 文書認識の高精度化

- 再認識による誤りの自動修正

抽出したフォントから作成した新辞書を用いて、抽出した辞書を再認識することで、フォントを抽出する際には 1 位で認識できなかった 370 文字 (10 部の合計) を全て正しく認識でき、誤りの自動修正に成功した。

- 新辞書による認識

社説 1 部から抽出した抽出辞書を用いて、他の新聞社説の認識を行ったが、辞書が抽出文書中に含まれる対象文字に依存するため認識結果は従来法よりは向上したが誤認識の全てを修正するまでには至らなかった。しかし、新聞社説 (11~20) を用いて新辞書の更新を行った実験では、新聞社説 (21~30) の 10 部に対して従来法での誤認識を全て修正できた (置換法)。

従来法と比較して上記の点で、今回提案したフォントの抽出を基にした文書認識システムの有効性が確認できたと考える。



しかし、本研究で提案した文書認識システムは未完成であり、今後の課題として下記のような項目が挙げられる。

#### 1. フォントの抽出法の高精度化

- パターンマッチング法だけでなく、構造解析法を用いる。
- 除去候補文字を検討する。

以上のような手段によりフォントの抽出法を高精度化する。特に、除去候補文字については、本研究中では事前に学習させたものを固定し、実験の結果をフィードバックしていない。これを抽出結果により、除去候補文字を学習・更新するシステムへ発展させていくことが多種フォント文書の自動認識のために有効であると考えられる。

#### 2. 文字の出現頻度を利用した文書認識システムの構築

今回のフォントの抽出では、事前の調査から対象をひらがな 65 文字にし、認識文書の約 50% 弱の文字をカバーし、認識精度の向上に成功したが、文書認識の自動化を目指すには、文字の出現頻度を考慮し、フォントを抽出する対象を拡大することが必要である。



## 謝辞

本研究を進めるにあたり、多大な御指導、御助言をいただいた東北大学工学部 阿曾弘具教授、北陸先端技術大学院大学 木村正行教授に心より感謝致します。また、本論文をまとめるにあたり、貴重な御意見をいただいた東北工学部 丸岡章教授、白鳥則郎教授に深く感謝致します。

また、東北大学大型計算機センター 孫寧助手、東北大学大学院博士課程 大町真一郎氏、同修士課程 越後和徳氏、セイコーエプソン株式会社 内山喜照氏には、文字認識の研究全般にわたって数多くの御指導、御助言をいただきました。ここに感謝の意を表します。

さらに、北陸先端技術大学院大学 下平博 助教授、東北大学情報処理教育センター 阿部亨助手、東北大学工学部 瀧本英二助手、東北大学大学院博士課程 成富敬氏、沼田一成氏、後藤英昭氏、同修士課程 福田大氏、中井満氏には、研究全般にわたる有益な御示唆をいただくとともに、計算機環境の整備を計っていただきました。厚く感謝いたします。

最後に御討論、御協力を戴き、また日頃の生活においてお世話になった旧木村研究室、丸岡研究室阿曾グループの皆様に感謝いたします。

0

0



## 参考文献

- [1] 飯島泰蔵：「パターン認識理論」  
森北出版 (1989 年)
- [2] 森健一：「パターン認識」  
電子情報通信学会編 (1988 年)
- [3] 長尾真：「パターン情報処理」  
コロナ社 (1983 年)
- [4] 孫寧、田原透、阿曾弘具、木村正行：「方向線素特徴量を用いた高精度文字認識」  
電子情報通信学会論文誌 (D-II), Vol.J74-D-II, No.3, pp. 330-339 (平成 3 年 3 月)
- [5] C. J. Hilditch : “Linear Skeleton from Square Cupboards”  
In Machine Intelligence 6, B. Meltzer & D. Michie, Eds., Univ. Press, Edinburgh, pp. 403-420 (1969)
- [6] 江島俊朗, 中村洋介, 木村正行：「構造情報を含む手書き漢字認識のための特徴量」  
電子情報通信学会論文誌 (D), Vol.J68-D, No.4, pp.789-796 (1985 年 4 月)
- [7] 八代博昭：「手書き漢字認識の高精度化に関する基礎研究」  
東北大学大学院工学研究科情報工学専攻 修士学位論文 (1988 年 2 月)
- [8] 池田啓明, 阿曾 弘具, 木村 正行：「文書認識における辞書の半自動作成システム」  
電気関係学会東北支部連合大会, 1H-3 (1991 年 8 月).
- [9] 池田啓明, 阿曾 弘具：「未知文書からの文字フォントの自動抽出」  
電気関係学会東北支部連合大会, 2C-1 (1992 年 8 月).

0

1

# 研究業績一覧

- 文書認識における辞書の半自動作成システム  
池田 啓明、阿曾 弘具、木村 正行  
平成 3 年度電気関係学会東北支部連合大会 1H3 (1991-08)
- 未知文書からの文字フォントの自動抽出  
池田 啓明、阿曾 弘具  
平成 4 年度電気関係学会東北支部連合大会 2C1 (1992-08)

3

0

# 第 A 章

## 社説の構成字種

今回用いた社説(1~20)の20部の構成字種について次に示す。なお、辞書に含まれていない字種(数字、英字、かな文字の一部)もある。

表 A.1: 社説 20 部中の各字種の出現頻度

	字種	文字数(累計)	割合(%)
記号	14/ 107	3650	8.02
数字	1/ 10	13	0.03
英字	18/ 52	206	0.45
かな	70/ 83	21138	46.45
カナ	72/ 86	1550	3.41
漢字	1276/ 2965	18937	41.61
合計	1451/ 3303	45504	100.00
その他	7/ -	10	0.03

字種の項はセット中の字種数 / 存在した字種数

0

1

表 A.2: 各字種毎の高出現頻度の文字

記号	、 1698	。 1075	－ 231	「 191	」 189
数字	0 13				
英字	N 46	A 32	T 26	S 22	E 18
かな	の 2010	い 1164	る 1145	に 1124	を 931
カナ	ン 105	ル 91	ト 88	ア 79	ス 77
漢字	国 336	政 233	日 201	一 186	大 155
全体	の 2010 。 1075	、 1698 を 931	い 1164 は 887	る 1145 な 862	に 1124 が 850

( ) 内は 20 部中での出現数

表 A.3: 高出現頻度の文字の全体に占める割合

文字数	5	10	20	30	40	50	100	150	200	250
割合 (%)	15.7	25.8	39.5	46.6	51.2	54.6	64.7	71.0	75.9	79.6

表 A.4: 上位 30 位までに含まれる文字

	ひらがな		記号		漢字	
	字種	文字数	字種	文字数	字種	文字数
10 位	8	8973	2	2773	—	—
20 位	18	15213	2	2773	—	—
30 位	27	20366	2	2773	1	336

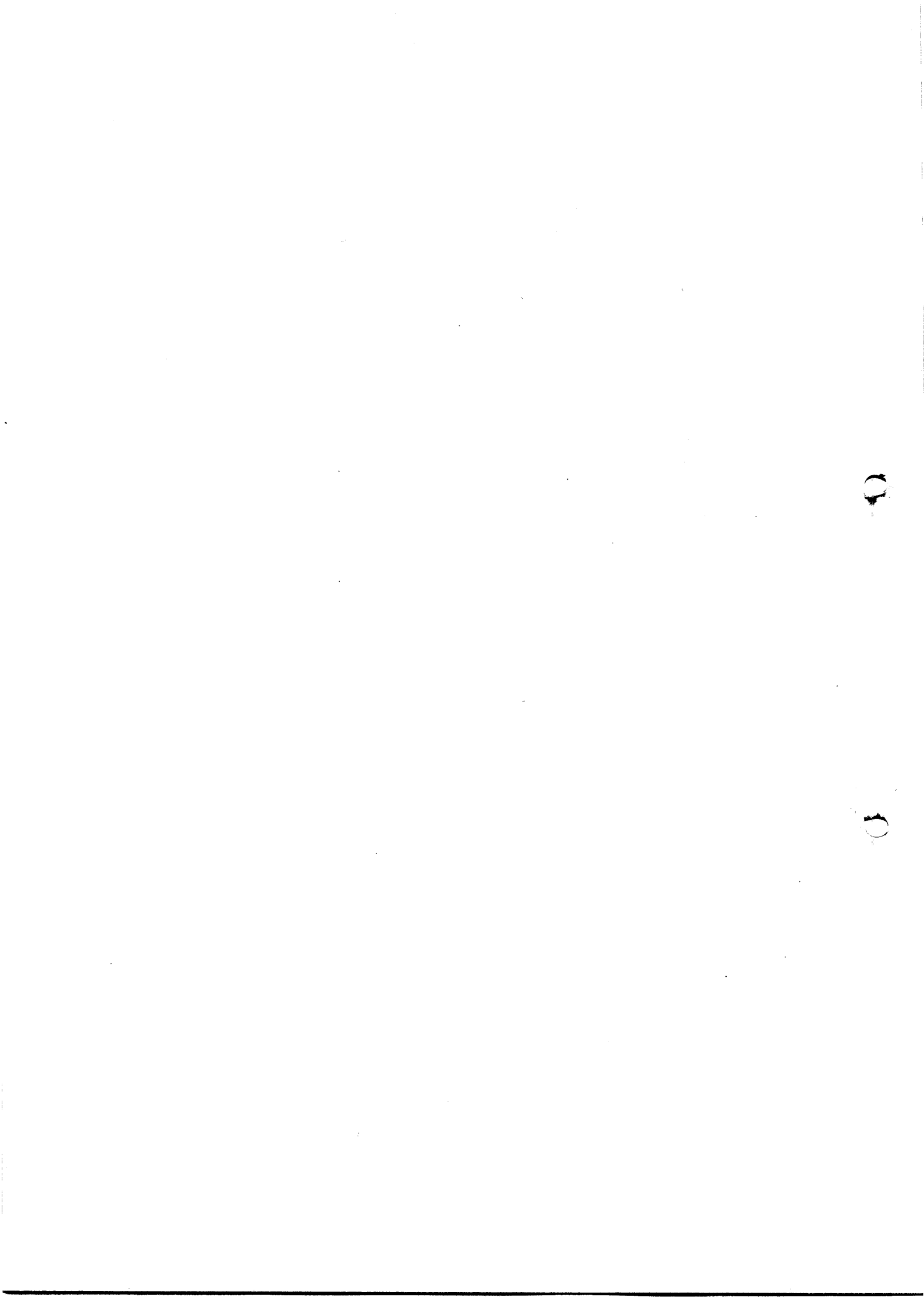




表 A.5: 各字種の出現頻度

社説	記号	数字	英字	ひらがな	カタカナ	漢字	その他	計
1	158	0	48	1134	58	877	0	2275
2	167	0	0	1057	59	974	1	2257
3	179	0	2	1126	66	876	1	2249
4	189	0	3	1073	90	948	1	2303
5	168	3	0	1062	190	885	0	2308
6	176	0	0	1050	66	1026	0	2318
7	186	0	9	1086	79	950	2	2310
8	168	4	0	1029	103	995	0	2299
9	189	0	4	1020	135	928	2	2276
10	211	3	0	1034	90	1005	0	2343
11	183	1	8	1077	40	973	1	2282
12	201	1	24	1078	40	992	3	2236
13	200	0	22	930	75	917	0	2144
14	178	0	8	1032	63	997	0	2278
15	161	0	0	1149	57	898	0	2265
16	216	0	3	1028	47	956	0	2250
17	173	0	68	1109	64	957	0	2371
18	176	0	0	1126	96	894	0	2292
19	207	0	0	1044	74	965	0	2290
20	164	1	7	899	53	924	0	2048

