

オフライン草書体文字認識システムの開発

西田 茂生 辻野 学

Development of an off-line cursive style body characters recognition system

Shigeki NISHIDA and Manabu TUJINO

There are two significant procedures to recognize off-line cursive style characters. One is recognition of the infinite form characters as separate of the character connected. Another is recognition of each separated character. In this paper, we propose an algorithm of off-line cursive style body characters using hidden Markov model (HMM). A characters databank was constructed using researcher's characters. To begin with the each character image is reduced to 50×50 pixels, it makes to four grey scale values. Next it has expanded as it is compared with the length and breadth of the character in the image. Finally the HMM parameter of each character was calculated. In this experiment, four sets of characters were used. The result of this, the recognition rate has become 43.5%. However, there was so little individual difference. It is thought that it becomes a system that can deal with everyone, and it is possible to correspond to various documents.

1. 緒 言

文字認識技術において、楷書体による定型文字やオンライン手書き文字などの認識技術はすでに実用化されているが、オフライン手書き文字については実験段階にあるものが多く、さらに日本文学の古典に多く用いられるオフライン草書体文字の認識技術は開発されていない。その理由として草書体文字は字形が不定形であり、文字間を1本の線で書くという特徴をもつこと、古い文献では文字の擦れや欠損が多くこれらが草書体文字の認識を困難にしている。オフライン草書体文字認識システム技術が開発されれば専門的な知識や訓練の必要なく草書体で書かれた文献を読解することができ、古典の研究に役立つと考えられる。本研究ではオフライン草書体文字認識システムの開発を最終目標とする。草書体文字を認識する場合は2つの重要なプロセスが必要となる。1つは文字間の分離であり、もう1つは不定形文字の認識である。そこで本論文ではその第1段階としてオフライン不定形文字認識手法の提案を行う。不定形文字のオフライン認識手法として、システムの発展性を考え学習機能を持ち、データベースを増やすほど認識率が上

がると考えられる隠れマルコフモデル（以下HMM）を用いた。また不定形文字としてデータベースには筆者の手書き文字、認識実験用には4名の被験者の手書き文字を用い実験を用いた。

2. 認識アルゴリズム

2.1 HMMの概要

HMMとは確率的な状態遷移と記号出力を備えたオートマトンである。以下にHMMの概要を述べる。

まずHMMは以下5項目により定義される。

(1) $Q=\{q_1, \dots, q_N\}$:状態の有限集合

(2) $\Sigma=\{o_1, \dots, o_M\}$:出力記号の有限集合

(3) $A=\{a_{ij}\}$:状態遷移確率分布

a_{ij} は状態 q_i から q_j への遷移確率であり $\sum a_{ij}=1$

(4) $B=\{b_i(o_t)\}$:記号出力確率分布

$b_i(o_t)$ は状態 q_i で記号 o_t を出力する確率であり、 $\sum b_i(o_t)=1$

(5) $\pi=\{\pi_i\}$:初期状態確率分布

π_i は状態 q_i が初期状態である確率 $P(X_1=q_i)$

4つの状態からなるHMMの例をFig.1に示す¹⁾。

Fig.1は状態 $q_1 \sim q_4$ があり、その状態から記号 $o_1 \sim o_4$ を出力することを表している。円の中にある $\pi_1 \sim \pi_4$ は添

番号の状態からはじまる初期状態確率である。

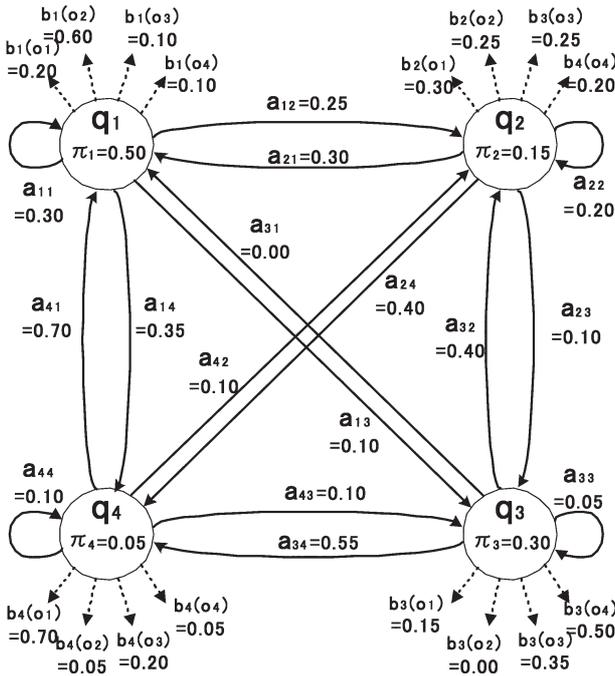


Fig.1 4つの状態にからなるHMM

矢印は状態の遷移を表しており $a_{11} \sim a_{44}$ はその際の状態遷移確率である。そして状態から記号 $o_1 \sim o_4$ が出力される記号出力確率が $b_i(o_1) \sim b_i(o_4)$ となる。例えば状態が $1 \rightarrow 1 \rightarrow 2 \rightarrow 4 \rightarrow 3$ (q の添え字を表す) と遷移するとき、記号が $3 \rightarrow 2 \rightarrow 4 \rightarrow 2 \rightarrow 1$ (o の添え字を表す) と出力されるとすると生成確率 P は次式で求めることができる。

$$P = \pi_1 \times b_1(o_3) \times a_{11} \times b_1(o_2) \times a_{12} \times b_2(o_4) \times a_{24} \times b_4(o_2) \times a_{43} \times b_3(o_1) \\ = 0.50 \times 0.10 \times 0.30 \times 0.60 \times 0.25 \times 0.20 \times 0.40 \times 0.05 \times 0.10 \times 0.15 \\ = 1.35 \times 10^{-7} \quad (1)$$

したがって、生成確率は 1.35×10^{-7} となる。また記号出力するとき状態遷移が数種類ある場合は、各生成確率の和がその記号出力の生成確率となる。

2.2 HMMの学習アルゴリズム

データベースを作成する際与えられた記号系列からモデルのパラメータを決定しなければならない。しかしHMMでは記号系列を生成した状態遷移系列は非観測であるため、直接最尤推定を行うことができない。このためEMアルゴリズムを用いた繰り返しアルゴリズムによりパラメータを推定する。

最初にモデル M が $O_1^T = O_1 \cdots O_T$ を生成して、時刻 t で状態 q_i に到達する確率 $\alpha_t(i)$ を導入する。

$$\alpha_t(i) = P(o_1^t, X_t = q_i | M) \quad (2)$$

$\alpha(\cdot)$ を前向き確率と呼び次のように再帰的に計算できる。

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}) \quad (3)$$

次に前向き確率とは双対な考えである後向き確率を導入する。後向き確率 $\beta_i(i)$ は状態 q_i から始まる状態遷移によって $o_{t+1}^T = o_{t+1} \cdots o_T$ が生成される確率である。

$$\beta_i(i) = P(o_{t+1}^T, X_t = q_i | M) \quad (4)$$

後向き確率も前向き確率と同様に再帰的に計算できる。

$$\beta_i(j) = \sum_{k=1}^N a_{jk} b_k(o_{t+1}) \beta_{t+1}(k) \quad (5)$$

そして与えられた記号系列 $O_1^T = O_1 \cdots O_T$ に対し状態 q_i から状態 q_j への遷移が時刻 t で生じた確率 $\gamma_t(i, j)$ を考える。

$$\gamma_t(i, j) = P(X_t = q_i, X_{t+1} = q_j | o_1^T, M) \\ = \frac{P(X_t = q_i, X_{t+1} = q_j, o_1^T | M)}{P(o_1^T | M)} \\ = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \alpha_t(i)} \quad (6)$$

また $\gamma_t(i)$ を以下のように定義する。

$$\gamma_t(i) = \sum_{j=1}^N \gamma_t(i, j) \quad (7)$$

$\gamma_t(i)$ は、時刻 t に状態 i に滞在した確率である。

$\gamma_t(i, j)$ および $\gamma_t(i)$ を用いてパラメータの再推定を以下のように行うことができる。

$$\bar{\pi}_i = \gamma_1(i) \quad (8)$$

$$\bar{a}_{ij} = \frac{\text{状態}i\text{から状態}j\text{へ遷移する回数の期待値}}{\text{状態}i\text{から遷移する回数の期待値}} \\ = \frac{\sum_{t, o_t=k} \gamma_t(i, j)}{T-1} \\ = \frac{\sum_{t=1}^{T-1} \gamma_t(i)}{T-1} \quad (9)$$

$$\bar{b}_i(k) = \frac{\text{状態}i\text{に滞在し記号}k\text{を出力する回数の期待値}}{\text{状態}i\text{に滞在する回数の期待値}} \\ = \frac{\sum_{t, o_t=k} \gamma_t(i)}{T} \\ = \frac{\sum_{t=1}^T \gamma_t(i)}{T} \quad (10)$$

上の再推定式を用いてパラメータを求める方法を本研究の学習アルゴリズムとして使用する。

2.3 HMMの文字認識への応用

今回の実験では1画素を1つの状態とし背景と同色(白)の場合、記号出力は無しとした。また入力画像を2値化ではなく、背景を除いて4値化し文字の濃淡にも

対応できるように考慮した．入力画像を M として 4 値化に用いた式(11)を以下に示す．

$$\left. \begin{aligned} (\text{median}M)/4 \times 3 - 0.1 \leq o_1 < (\text{median}M)/4 \times 4 - 0.1 \\ (\text{median}M)/4 \times 2 - 0.1 \leq o_2 < (\text{median}M)/4 \times 3 - 0.1 \\ (\text{median}M)/4 - 0.1 \leq o_3 < (\text{median}M)/4 \times 2 - 0.1 \\ 0 \leq o_4 < (\text{median}M)/4 - 0.1 \end{aligned} \right\} (11)$$

ここで median は中央値を意味する．画像の大部分が背景であることから，中央値と背景の値は近似していると考え，誤差を 0.1 に設定し背景を除去した．残った値で 4 値化を行う．その時の色調を Fig.2 に示す．

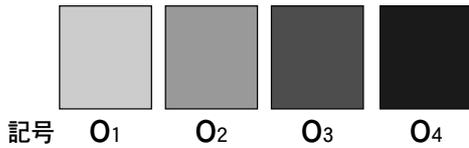


Fig.2 4 値化の色調

これらの記号 $o_1 \sim o_4$ を出力する場合の状態 q_i ，初期状態確率 π_i ，状態遷移確率 a_{ij} ，記号出力確率 $b_i(o_t)$ を Fig.3 に示すように設定した．

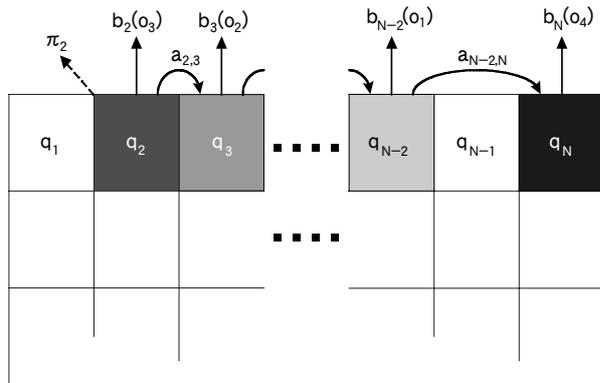


Fig.3 パラメータの設定方法

まず HMM パラメータを 1 行ごとに定義し，計算を行った．画素の左端を状態 q_1 ，右端を状態 q_N とし左から右にむかって状態が遷移すると定める．その時，最初の記号を状態 i で発見した場合に初期状態確率 π_i に 1 を代入し，その時の記号出力確率 $b_i(o_t)$ に 1 を代入する．次に記号が出力されている状態 j へ向かい状態遷移確率 a_{ij} に 1 を代入する．後は初期状態確率を除いて同様の手順で記号出力確率，状態遷移確率に 1 を代入していく．これで 1 枚の画像各行に HMM を定義することができる．この方法でデータベースを作成する場合，1 文字につき L 枚画像を取り込んだとする．そしてデータベースの $\pi_i, a_{ij}, b_i(o_t)$ を $(\pi_i)_D, (a_{ij})_D, \{b_i(o_t)\}_D$ ，取り込んだ K 枚目の画像の $\pi_i, a_{ij}, b_i(o_t)$ を $(\pi_i)_K, (a_{ij})_K, \{b_i(o_t)\}_K$ と表して書くと，式(12)により各パラメータを求めることができる．

$$\left. \begin{aligned} (\pi_i)_D &= \{(\pi_i)_1 + (\pi_i)_2 + \dots + (\pi_i)_K + \dots + (\pi_i)_{L-1} + (\pi_i)_L\} / L \\ (a_{ij})_D &= \{(a_{ij})_1 + (a_{ij})_2 + \dots + (a_{ij})_K + \dots + (a_{ij})_{L-1} + (a_{ij})_L\} / L \\ \{b_i(o_t)\}_D &= \{b_i(o_t)_1 + b_i(o_t)_2 + \dots + b_i(o_t)_K + \dots + b_i(o_t)_{L-1} + b_i(o_t)_L\} / L \end{aligned} \right\} (12)$$

本来であれば式(12)を用いるが，今回の HMM の状態・記号出力設定の場合，このまま後に示す式(15)を計算すると全ての値が少数であり，似ている文字ほどデータベースに含まれるパラメータ数が多いことから，式(15)の計算結果は小さくなる．しかし，似ている文字ほどパラメータ値は大きいので，パラメータ数が同じであると式(15)の計算値は大きくなる．このように計算値の大小と認識したい文字が一致しづらいため式(13)を用いる．

$$\left. \begin{aligned} (\pi_i)_D &= \{(\pi_i)_1 + (\pi_i)_2 + \dots + (\pi_i)_K + \dots + (\pi_i)_{L-1} + (\pi_i)_L\} \\ (a_{ij})_D &= \{(a_{ij})_1 + (a_{ij})_2 + \dots + (a_{ij})_K + \dots + (a_{ij})_{L-1} + (a_{ij})_L\} \\ \{b_i(o_t)\}_D &= \{b_i(o_t)_1 + b_i(o_t)_2 + \dots + b_i(o_t)_K + \dots + b_i(o_t)_{L-1} + b_i(o_t)_L\} \end{aligned} \right\} (13)$$

式(13)は，画像を取り込んだ枚数で割らず，HMM パラメータを全て整数で扱うようにした．パラメータ数が多い場合でも式(15)の計算結果が大きくなるので，計算値と認識したい文字の不一致を大幅に改善することができた．この式(13)を用い平仮名 46 文字のデータベースを作成する．ここまでの作業が終了したら認識実験を行う．まず実験画像を Fig.3 の方法でパラメータ化し作成した平仮名 46 字のデータベースとひとつひとつ重ね合わせる．その時用いた式で入力画像を M とし， $\pi_i, a_{ij}, b_i(o_t)$ を $(\pi_i)_M, (a_{ij})_M, \{b_i(o_t)\}_M$ と表してものを式(14)に示す．

$$\left. \begin{aligned} \pi_i &= (\pi_i)_M \times (\pi_i)_D \\ a_{ij} &= (a_{ij})_M \times (a_{ij})_D \\ b_i(o_t) &= \{b_i(o_t)\}_M \times \{b_i(o_t)\}_D \end{aligned} \right\} (14)$$

これにより，実験画像に対する平仮名 46 文字の各 HMM パラメータを求めることができる．

後は生成パラメータを式(15)によって計算する．

$$P = \sum_{i=1}^N \sum_{j=1}^N \sum_{t=1}^4 \pi_i a_{ij} b_i(o_t) \quad (15)$$

ただし $\pi_i, a_{ij}, b_i(o_t)$ の中に 0 が含まれている場合，乗算であるため他の値に関係なく結果が 0 になるので，0 を飛ばして計算している．また HMM パラメータは各行ごとに計算されるので実験画像の行数の数だけ計算される．

3 データベース

3.1 データベースの作成

データベースの作成手順を Fig.4 に示し，各手順の処理方法を述べる．

(i) データベース画像を取り込み平仮名各 46 文字の手書き文字のデータベースを作成する．600dpi のスキャ

ナで文字をパソコンに8ビットのグレースケールで取り込む。1文字につき10個とりこむ。文字は研究者のものを用いた。

(ii) 手作業により一文字ずつ切り出す。

(iii) 取り込んだ画像を50×50ピクセルになるようにサイズを調整する。これは計算機のメモリ容量に制限があるために行っている。「あ」を例にあげFig.5(a)に示す。これらの画像はデータベース作成のために取り込んだ研究者の文字画像である。この時点では画像は8ビットのグレースケールで表されている。

(iv) 式(11)を用いて背景、ノイズを除去し4値化した画像をFig.5(b)に示す。この時点で画像の各画素の光強度は2ビットに減少するが文字の濃淡の情報が残っていることを確認することができる。しかし文字ごとに大きさのばらつきがあるため次の作業を行う。

(v) Fig.5(b)で画像内の文字を縦横比は保持したまま50×50ピクセルに拡大しFig.5(c)に示す。文字の左右の余白は計算により均等になるように調節している。この作業は文字の大小による個人差をなくすためである。

(vi) Fig.5(c)の画像を式(13)によりパラメータ化し「あ」のデータベースが完成する。同様に平仮名46字のHMMパラメータを設定する。このとき1文字につき取り込む文字数を増やせばより緻密なHMMパラメータを設定できる。

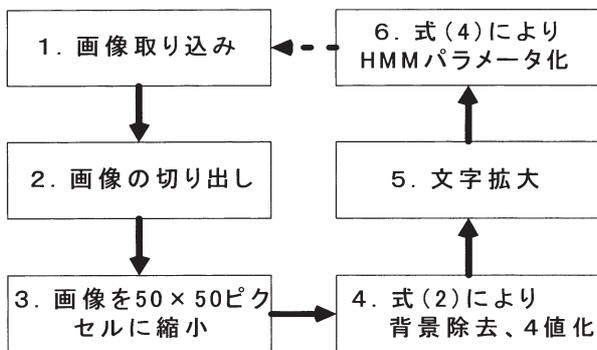


Fig.4 データベース作成手順



(a)取り込み画像 (b)4値化後画像 (c)拡大後画像

Fig.5 データベース文字画像「あ」の一例

3.2 データベース内文字の認識実験

認識実験手順をFig.6に示し各手順の処理方法を述べる。

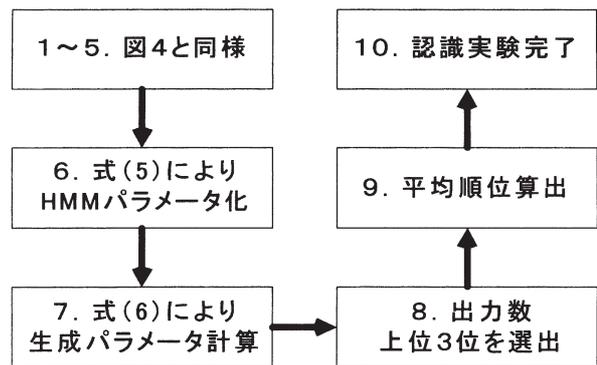


Fig.6 認識実験手順

(i) 認識対象画像を8ビットのグレースケールで取り込む。

(ii) 手作業により一文字ずつ切り出す。

(iii) 取り込んだ画像を50×50ピクセルになるようにサイズを調整する。調整後の文字画像をFig.7(a)に示す。

(iv) 式(11)を用いて背景、ノイズを除去し4値化する。処理後の文字画像をFig.7(b)に示す。

(v) Fig.10で画像内の文字を縦横比は保持したまま50×50ピクセルに拡大しFig.7(c)に示す。



(a)取り込み直後 (b)4値化画像 (c)拡大後画像

Fig.7 認識実験に用いた文字「あ」の一例

(vi) (c)をパラメータ化し式(14)によりデータベース46字と重ね合わせ(c)に対する平仮名46文字の各HMMパラメータを求める。

(vii) HMMパラメータより式(15)を用いて46文字の生成パラメータを計算する。

(viii) 各文字につき生成パラメータが0以外になる時の出力数をカウントする。その結果をFig.8に示す。この結果から出力数の1位「あ」、2位「え」、3位「ん」となった。

(ix) この3文字の生成パラメータを計算し、カウント数と式(16)で示す平均順位を算出したものをTable1に示す。

$$\text{平均順位} = \frac{\text{カウント数順位} + \text{生成パラメータ順位}}{2} \quad (16)$$

(x) 表からわかるようにこの認識実験結果から認識対象文字は「あ」を示しており、文字認識が行えたといえる。なお、データベース内の文字による認識実験結果は認識率100%であった。

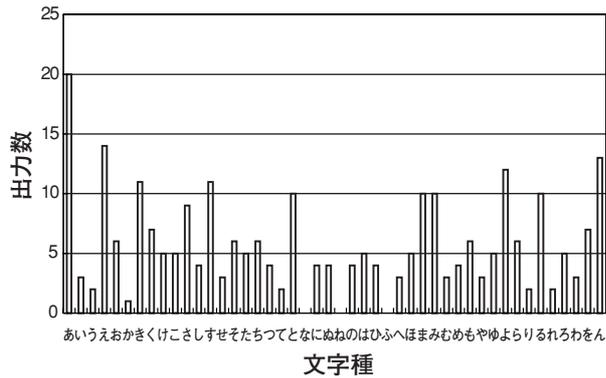


Fig.8 各文字のカウンタ数 (生成パラメータ 0 を除く)

Table 1 平均順位算出

	カウンタ数	生成パラメータ	平均順位
あ	20	2.07×10^{10}	1.5
え	14	9.36×10^9	2.5
ん	13	5.03×10^{10}	2

4. 認識実験

4.1 認識実験結果

前章の実験例で示したような実験を、Fig.9に示す被験者4人(データベース作成者と他の被験者3名)が手書きした平仮名46文字を用いて認識実験を行った。



(a)研究者 (b)被験者A (c)被験者B (d)被験者C

Fig.9 認識実験に使用した文字

この図以外にも濁点や半濁点を付け加えることができる文字があるが、本研究はデータベースにある文字種しか認識できないため、データベース内の文字を選んで認識させる必要がある。認識実験結果をTable2に示す。表の認識率は式(17)で定義したものをを用いる。

$$\text{認識率}(\%) = (\text{認識文字数} / \text{合計文字数}) \times 100 \quad (17)$$

認識率は全体で43.5%と実際に文字認識システムとして使用するには厳しい数字となってしまった。特にデー

Table2 認識実験結果 (認識率)

	認識 (文字種)	認識不可 (文字種)	認識率 (%)
研究者	20	26	43.5
被験者A	19	28	41.3
被験者B	21	25	45.7
被験者C	20	26	43.5
平均	20.0	26.3	43.5

タベース本人の認識率が予想に反し低く、逆に他の被験者と同じような認識率になった。この事実は本研究で提案するHMMを用いた認識手法が筆記者を選ばない不特定筆記者の文字認識の可能性を示していると考えられる。

4.2 考察

次に本認識手法による認識率が低い原因を考察する。まず4人の認識結果から各文字の認識数をFig.10に示す。

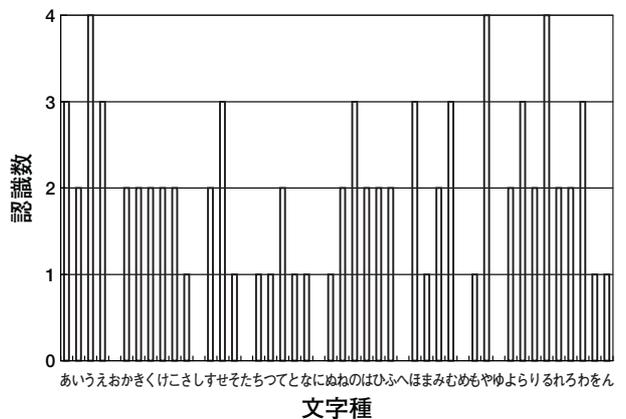
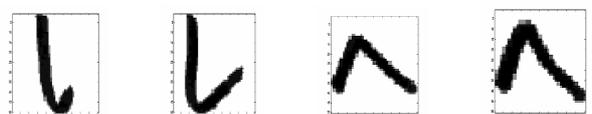


Fig.10 各文字の認識数 (被験者4名の総数)

図より「う」「や」「る」に関しては全員認識できているが、「お」「し」「た」「に」「へ」「め」「ゆ」は全員認識できなかった。認識できなかった文字について考察する。「し」「へ」に関しては文字の縦横比のばらつきにあると思われる。データベースに用いた「し」「へ」について幾つかをFig.11に示す。



(a)し;縦長 (b)し;横長 (c)へ;横長 (d)へ;縦長
Fig.11 縦横比のばらつき

このように文字の形にばらつきが出やすいため、データベースのHMMパラメータも低くなり式(5)により求められる生成パラメータも小さくなる。そのために認識されづらいと考えられる。対策としてはデータベースの拡大などがあると考えられる。

「お」「に」「た」の関しては、形の似ている文字の誤認識がほとんどであった。「お」は「む」と間違えやすく、「に」であれば「け」「ほ」「は」のどれかと間違えていた。また「た」では「む」「ん」と間違えている。

現在のHMM設定では背景を記号出力に含まずに設定していて、文字部だけを評価する手法をとっている。そのため、字形の似ている部分しか評価できずこのような誤認識が起こると考えられる。対策としては字形の違う部分も評価するために背景(白)も記号出力に設定し、画像全体を評価することで解決できると考える。

「め」「ゆ」に関しては上記に示したものの複合的な理由であると考えられる。対策方法も上記によると思われる。またTable3に上位選出数と非選出数、選出率を示す。

Table3 認識実験結果(選出率)

	上位選出 (文字種)	非選出 (文字種)	選出率 (認識率)(%)
研究者	28	18	60.9 (43.5)
被験者 A	27	19	58.7 (41.3)
被験者 B	28	18	60.9 (45.7)
被験者 C	26	20	56.5 (43.5)
平均	27.3	18.8	59.3 (43.5)

選出率の算出は式(18)による。

$$\text{選出率}(\%) = (\text{上位選出数} / \text{合計文字数}) \times 100 \quad (18)$$

この表より選出率と認識率の差が15.8%もあることがわかった。選出されても認識文字と一致しなかった文字が多かったことから、新たな判定基準を設け、この差をうめなければならない。

5. 結 言

本研究では、オフライン草書体文字認識システムの構築を行うため、その第1段階として、HMMを用いた認識アルゴリズムを提案した。認識率は全体で43.5%と決して高い数字ではないが、個人差がほとんどなかった。このことから改良を加え認識率を上げれば不特定筆記者に対応できる草書体文字認識システムの可能性が示された。

参考文献

- 1) T.Artieres, N.Gauthier, et al., A Hidden Markov Models combination framework for handwriting recognition, IJDAR, No.5 233-243(2002)
- 2) 北研二, 確率的言語モデル, 1999, 東京大学出版会