

## 標本数が少ない状況下における識別器の評価

浜本義彦<sup>†</sup> 内村俊二<sup>‡</sup> 金岡泰保<sup>†</sup> 富田真吾<sup>†</sup>

<sup>†</sup> 山口大学工学部

<sup>‡</sup> 大島商船高等専門学校

755 山口県宇部市常盤台 2557

742-21 山口県大島郡大島町小松 1091-1

あらし 統計的パターン認識における主要な問題は、識別器を設計することである。識別器は、有限個の訓練サンプルを用いて学習される。特に、識別性能の高い識別器を得るためには、大量の訓練サンプルが必要である。ところが、実際のパターン認識問題では、特徴数に対する訓練サンプル数の比が小さい場合が多い。そこで、そのような状況下で、誤識別率の意味で、いずれの識別器が最良かという疑問が生じる。本論文では、人工データを用いて、よく知られている識別器を識別性能の観点から比較する。

和文キーワード パターン認識, 識別器, 誤識別率, サンプルサイズ, 特徴数

## Evaluation of classifiers in a small sample size case

Yoshihiko Hamamoto<sup>†</sup> Shunji Uchimura<sup>‡</sup> Taiho Kanaoka<sup>†</sup> Shingo Tomita<sup>†</sup>

<sup>†</sup> Yamaguchi University

<sup>‡</sup> Oshima National College of Maritime Technology

2557 Tokiwadai, Ube, 755, Japan

1091-1 Komatsu, Oshima-cho, 742-21, Japan

Abstract The main problem in statistical pattern recognition is to design a classifier. The classifier must be learned from the training samples. A large training sample is essential to design a classifier with a very low error probability. In many practical situations, the ratio of sample size to dimensionality is small. One may ask which of classifiers is best in terms of the error probability, when the ratio of sample size to dimensionality is small. We compare the classification performances of the well-known classifiers on three artificial data sets.

英文 key words pattern recognition, classifier, error probability, sample size, dimensionality

# 1 まえがき

パターン認識における識別器の目的は、任意に与えられたパターンを、しかるべきクラスに対応づけることである。その際、誤って識別されるパターンの数が少ない程、良い識別器と考えられる。統計的パターン認識論によると、誤識別率が最小となる識別器は、Bayes 識別器であることが知られている。この Bayes 識別器は、クラスの条件付き確率密度関数が既知であることを前提としている。しかし、現実のパターン認識問題では、確率密度関数が未知である。そこで、無作為に抽出された標本を用いて、識別器の学習というアプローチがとられた。学習には、例えば正規分布などの分布を仮定するパラメトリックな手法と、分布を仮定しないノンパラメトリックな手法がある。いずれにしても、有限個の標本を用いて、識別器の学習および評価を行わなければならない。特に、識別器の学習に用いられる標本は訓練サンプルと、また評価に用いられる標本はテストサンプルと呼ばれている。

さて、実際の認識においては、パターンから種々の特徴が抽出され、得られた特徴の組で構成される特徴空間上において識別器が学習される。このとき、特徴数に比べ大量の訓練サンプルを用いた識別器については、学習が十分なされていると考えられる。識別器の学習において、特徴数に対する訓練サンプル数の比が重要である、ということは多くの研究者により指摘されてきた [1][2][3][4]。Jain は、特徴数に対する訓練サンプル数の比が 10 以上あることが望ましいと述べている。上述の Bayes 識別器は、訓練サンプル数が無限にあるとき得られるものである。Bayes 識別器は特殊な例としても、従来のパターン認識論では訓練サンプル数が十分多いということを暗に仮定してきた。しかしながら、現実のパターン認識問題では、特徴数に対する訓練サンプル数の比が小さい場合が多い。例えば、手書き漢字認識問題では、訓練サンプル数が高々 160 であるのに対し特徴数は 64 で、その比が 2.5 という状況で認識実験がなされている [5]。一般に、特徴数に対する訓練サンプル数の比が小さい状況下では、十分な学習が行われなく、結果として誤識別率の増加という問題が生じている。

このようなパターン認識論と現実のパターン認識問題との間のギャップを縮小しようとする研究が、近年盛んに行われるようになった [6][1]。本論文もこれらの研究と立場を同じくするもので、より現実の状況下で真に有効な識別器の解明をその目的としている。

# 2 準備

本論文では、2 クラス問題のみを論じ、2 クラスの事前確率は等しいとする。まず、諸記号の説明を行う。

$n$	: 特徴空間の次元数 (特徴数)
$x_{ij}$	: クラス $i$ の $j$ 番目の訓練サンプル
$N_i$	: クラス $i$ の訓練サンプル数
$N = N_1 + N_2$	: 全訓練サンプル数
$\mu_i$	: クラス $i$ の母平均ベクトル
$\Sigma_i$	: クラス $i$ の母共分散行列
$\hat{\mu}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} x_{ij}$	: クラス $i$ の標本平均ベクトル
$\hat{\Sigma}_i = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (x_{ij} - \hat{\mu}_i)(x_{ij} - \hat{\mu}_i)^T$	: クラス $i$ の標本共分散行列
$\hat{W} = \frac{1}{2}(\hat{\Sigma}_1 + \hat{\Sigma}_2)$	: 標本クラス内共分散行列
$\lambda_i^s$	: $\hat{\Sigma}_i$ の $s$ 番目の固有値
$\phi_i^s$	: $\hat{\Sigma}_i$ の $s$ 番目の固有ベクトル

また、 $x^T$  は  $x$  の転置を、 $|A|$  は行列  $A$  の行列式を、 $A^{-1}$  は  $A$  の逆行列を、 $\|x - y\|$  は  $x$  と  $y$  との間のユークリッド距離を表わす。なお、固有値は  $\lambda_1^i \geq \dots \geq \lambda_n^i$  と順序づけられている。

## 2.1 誤識別率の推定法

誤識別率の推定法としては、Resubstitution 法、Leave-one-out 法、Holdout 法が提案されている [7]。

まず、Resubstitution 法は、識別器を設計するために用いる訓練サンプルと誤識別率を推定するために用いるテストサンプルとを同一とするもので、低い方へ偏った誤識別率が得られる、という問題を含んでいる。これは、訓練サンプルとテストサンプルとの独立性が保たれていないことによる。この独立性が重要であることを Fukunaga らは指摘している [8]。訓練サンプルとテストサンプルとの独立性を保つ手法が、Leave-one-out 法と Holdout 法である。

Leave-one-out 法は、利用できる標本数が限られるとき、それを最大限に活用する手法として提案されたものである [9]。これは、いま  $N$  個の標本が利用できるとすると、 $N$  個の標本の中からテストサンプルとして 1 個の標本を取り出し、残りの  $N - 1$  個の標本を訓練サンプルとして識別器の学習を行い、得られた識別器でテストサンプルの識別を行うものである。この処理を  $N$  回繰り返して誤識別率を求める。しかし、Leave-one-out 法により推定された誤識別率の分散は大きい、ということが指摘されている [10]。

一方、Holdout 法は、利用できる標本を訓練サンプルとテストサンプルとに分割し、誤識別率を求めるものである。この Holdout 法では、実在データを対象とした場合、利用できる標本数に限りがあるため、訓練サンプル数およびテストサンプル数の不足により偏りと分散の大きい誤識別率しか得られない。また、訓練サンプルとテストサンプルの分割法も問題となる。しかし、人工データを対象とした場合、利用できる標本数に制約がないため、上述の訓練サンプル数、テストサンプル数の不足および分割法の問題が解消される。また、実在データに対しては一般に Bayes 誤識別率が未知であるのに対し、人工データについては Bayes 誤識別率が既知である、という長所もある。このため、真の Bayes 誤識別率が既知という条件のもとで、識別器の評価を行うことができる。

以上のことから、本論文では人工データ集合を用いた Holdout 法により誤識別率の推定を行う。

## 2.2 人工データ集合

本論文では、母平均ベクトル  $\mu_i$ 、母共分散行列  $\Sigma_i$  ( $i = 1, 2$ ) の正規分布に従う 3 種類の人工データ集合 [7] を用いる。いずれの人工データ集合も 8 次元特徴空間上で記述されている。いま、 $\mu_i, \Sigma_i$  ( $i = 1, 2$ ) をそれぞれ

$$\mu_1 = [0 \cdots 0]^T, \quad \mu_2 = [m_1 \cdots m_8]^T$$

$$\Sigma_1 = \begin{pmatrix} 1 & & & & & & & \\ & \ddots & & & & & & \\ & & 1 & & & & & \\ 0 & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ 0 & & & & & & & 1 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} r_1 & & & & & & & \\ & \ddots & & & & & & \\ & & r_2 & & & & & \\ & & & \ddots & & & & \\ & & & & r_3 & & & \\ & & & & & \ddots & & \\ & & & & & & r_4 & \\ & & & & & & & \ddots & \\ & & & & & & & & r_5 & \\ & & & & & & & & & \ddots & \\ & & & & & & & & & & r_8 \end{pmatrix}$$

とおくと、 $m_i, r_i, (1 \leq i \leq 8)$  の値により、次の 3 種類のデータ集合が得られる。

### 1. データ集合 1

$$m_1 = 2.56, m_2 = \cdots = m_8 = 0,$$

$$r_1 = \cdots = r_8 = 1$$

このデータ集合は、2 クラスの母平均ベクトルのみが異なり、母共分散行列は共に単位行列である。データ集合 1 については、真の Bayes 誤識別率が 10% である。

### 2. データ集合 2

$$m_1 = \cdots = m_8 = 0,$$

$$r_1 = \cdots = r_8 = 4$$

このデータ集合は、2 クラスの母平均ベクトルが等しく、母共分散行列のみが異なる。データ集合 2 については、真の Bayes 誤識別率が 9% である。

### 3. データ集合 3

表 1: データ集合 3 のパラメータ

$i$	1	2	3	4	5	6	7	8
$m_i$	3.86	3.10	0.84	0.84	1.64	1.08	0.26	0.01
$r_i$	8.41	12.06	0.12	0.22	1.49	1.77	0.35	2.73

このデータ集合は、 $m_i, r_i$  ( $1 \leq i \leq 8$ ) が表 1 で示される値をとるもので、2 クラスの母平均ベクトル、母共分散行列が共に異なる。データ集合 3 は、Marill と Green によりアルファベット手書き文字 A, B のデータから得られたものもある [11]。このデータ集合 3 については、真の Bayes 誤識別率が 1.9% であることが報告されている [12]。

## 2.3 評価された識別器

(a) Fisher の線形識別関数  $G_j(x)$  [13]

これは、標本平均ベクトルの差異に基づいて識別がなされるもので、

$$G_j(x) = (x - \hat{\mu}_j)^T \hat{W}^{-1} (x - \hat{\mu}_j) \quad (1)$$

で与えられる。Fisher の線形識別関数では、

$$G_i(x) = \min\{G_1(x), G_2(x)\} \quad (2)$$

ならば、入力パターン  $x$  はクラス  $i$  に属するものと決定する。以下、この識別器を FLDF と呼ぶ。

(b) 二次識別関数  $H_j(x)$  [13]

これは、標本平均ベクトルおよび標本共分散行列の差異に基づいて識別がなされるもので、

$$H_j(x) = (x - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (x - \hat{\mu}_j) + \log |\hat{\Sigma}_j| \quad (3)$$

で与えられる。二次識別関数では、

$$H_i(x) = \min\{H_1(x), H_2(x)\} \quad (4)$$

ならば、入力パターン  $x$  はクラス  $i$  に属するものと決定する。以下、この識別器を QDF と呼ぶ。

(c) 修正二次識別関数  $M_j(x)$  [5]

これは、標本共分散行列の推定誤差を低減する工夫を取り入れたもので、

$$M_j(x) = \frac{1}{\lambda_j^{\xi+1}} \left[ \|x - \hat{\mu}_j\|^2 - \sum_{s=1}^{\xi} \left(1 - \frac{\lambda_j^{\xi+1}}{\lambda_j^s}\right) \{\phi_j^{sT} \cdot (x - \hat{\mu}_j)\}^2 \right] + \log \left( \prod_{s=1}^{\xi} \lambda_j^s \cdot \prod_{s=\xi+1}^n \lambda_j^{\xi+1} \right) \quad (5)$$

で与えられる。ここで、 $\xi (0 \leq \xi \leq n-1)$  は実験的に定められる定数である。修正二次識別関数では、

$$M_i(x) = \min\{M_1(x), M_2(x)\} \quad (6)$$

ならば、入力パターン  $x$  はクラス  $i$  に属するものと決定する。以下、この識別器を MQDF と呼ぶ。

(d) k 最近傍識別器 [13]

これは、パターン分布の形を仮定しないノンパラメトリック識別器の一つである。

$N$  個の全訓練サンプルの中で、入力パターン  $x$  とのユークリッド距離が  $k$  番目に小さい訓練サンプルを  $x^k$  とする。全訓練サンプルに対し、

$$\|x - x^*\| \leq \|x - x^k\| \quad (7)$$

なる訓練サンプル  $x^*$  のうち、クラス  $j$  に属する訓練サンプルの個数を  $v_k^j(x)$  とする。k 最近傍識別器では、

$$v_k^i(x) = \max\{v_k^i(x), v_k^2(x)\} \quad (8)$$

ならば、入力パターン $x$ はクラス $i$ に属するものと決定する。以下、この識別器をk-NNと呼ぶ。

(e) ユークリッド距離識別器 ED

これは、最も簡単な識別器で、これでは、

$$\|x - \hat{\mu}_i\| = \min\{\|x - \hat{\mu}_1\|, \|x - \hat{\mu}_2\|\} \quad (9)$$

ならば、入力パターン $x$ はクラス $i$ に属するものと決定する。以下、この識別器をEDと呼ぶ。

(f) Parzen 識別器 [14][15]

Parzen 識別器は、ノンパラメトリック識別器の一つである。

まず、訓練サンプル $x_{ij}$ を用いて、クラス $i$ の条件付き確率密度関数 $\hat{p}(x|i)$ を推定する:

$$\hat{p}(x|i) = \frac{1}{N_i} \sum_{j=1}^{N_i} K(x - x_{ij}) \quad (10)$$

ここで、 $K(\cdot)$ は核関数と呼ばれるものである。本論文では、核関数 $K$ に

$$K(x - x_{ij}) = \frac{1}{(2\pi)^{n/2} h^n |\hat{\Sigma}_i|^{1/2}} \exp\left\{\frac{-1}{2h^2} (x - x_{ij})^T \hat{\Sigma}_i^{-1} (x - x_{ij})\right\} \quad (11)$$

を用いる [15]。ここで、 $h$ は実験的に定められるウィンドウ幅である。Parzen 識別器では、

$$\hat{p}(x|i) = \max\{\hat{p}(x|1), \hat{p}(x|2)\} \quad (12)$$

ならば、入力パターン $x$ はクラス $i$ に属するものと決定する。以下、この識別器をPARZENと呼ぶ。

### 3 計算機シミュレーション

#### 3.1 手順

本論文で重要な役割を果たすパラメータ $\eta$ を次のように定義する。

$$\eta = N_i/n \quad (13)$$

すなわち、 $\eta$ は特徴数に対する訓練サンプル数の比である。本実験では、 $n = 8$ 、 $N_1 = N_2$ とし、 $\eta$ を1, 2, 4, 6, 8, 10と変え、それぞれについて識別器を設計し、その誤識別率を推定した。精度の高い誤識別率を推定するためには、大量のテストサンプルが必要である。そこで、訓練サンプルとは独立に発生させたテストサンプルを各クラスとも1000個用意した。以上の処理を独立に100回繰り返し、誤識別率の平均値、標準偏差および95%信頼区間を求めた。なお、 $\xi$ の最適値は次のように求めた。 $\xi$ の値を、0, 1, 2, 3, 4, 5, 6, 7と変え、それぞれについて誤識別率を求めた。それらの中で誤識別率を最小とする値を、 $\xi$ の最適値とした。 $h$ の最適値も、 $h$ の値を0.2から2.2まで0.2ずつ変えて求めた。

#### 3.2 実験結果と検討

データ集合1に対する実験結果を表2に、データ集合2に対する実験結果を表3に、データ集合3に対する実験結果を表4にそれぞれ示す。 $\xi$ の最適値と $\eta$ との関係、 $h$ の最適値と $\eta$ との関係をそれぞれ表5、表6に示す。

一般に、 $\eta$ の値が小さくなるにつれて、パターン分布の正規分布からのずれが大きくなっていく。特に、 $\eta$ の値が1.2の場合、正規性の仮定は適当でないと考えられる。

実験結果から以下のことが指摘される。

1.  $\eta$ の値が小さくなるにつれて、いずれの識別器の誤識別率もBayes誤識別率からの偏りおよび標準偏差が大きくなっていく。標準偏差が大きいくということは、識別器の信頼性が低い、ということ意味する。いずれのデータ集合についても、比較された中でもEDが $\eta$ の影響を受けにくく、QDFが一番 $\eta$ の影響を受ける。
2. 従来から指摘されているように、QDFよりFLDFの方がロバスト性が高い。QDFにおいては、共分散行列の推定誤差が問題となる。
3.  $\eta$ の値が小さいときは、MQDFが有効である。 $\eta = 2$ の場合における識別器の識別性能順位を表7に示す。
4.  $\eta$ の値が大きくなるにつれて、QDFの誤識別率はMQDFのそれに接近して行く。
5. 表5より、データ集合2と3については、 $\eta$ の値が小さくなるにつれて、 $\xi$ の最適値も小さくなっていく。木村ら[16]によれば、MQDFは、 $\xi = n - 1$ の場合、QDFに等しく、 $\xi = 0$ の場合ユークリッド距離と等価になり、 $0 < \xi < n - 1$ の場合、低次の $\xi$ 次元部分空間におけるQDFと、高次の $(n - \xi)$ 次元部分空間におけるユークリッド距離の荷重和となる。実験結果は、 $\eta$ の値が小さいときは、QDFよりむしろユークリッド距離に比重をおいた識別の方が良い、ということを示している。修正二次識別関数では、共分散行列の推定誤差を低減しているため、QDFに比べ、良い識別結果が得られている。
6. 表6より、データ集合2と3については、 $\eta$ の値が小さくなるにつれて、 $h$ の最適値は大きくなっていく。すなわち、訓練サンプル数が減少していくと、ウィンドウ幅を大きくとった方が、良い識別結果が得られる。
7. Van Ness [17]により、高次元のとき、すなわち $\eta$ の値が小さいとき、たとえパラメトリックな仮定が正しくても、ノンパラメトリック識別器であるPARZENの方が、パラメトリック識別器であるQDFよりも、良い識別結果を与える、ということが指摘された。実験結果は、 $\eta = 2$ の場合すべてのデータ集合について、PARZENの方がQDFより優れていることを示し、Van Nessの主張を支持している。

表 2: データ集合 1 に対する識別器の誤識別率 [%]

(上段:平均値, 中段:標準偏差, 下段:95%信頼区間, Bayes 誤識別率:10%)

識別器	$n$					
	1	2	4	6	8	10
1-NN	21.22	20.67	19.34	18.58	18.44	18.46
	4.44	4.16	2.61	1.95	1.80	1.76
	(20.35,22.09)	(19.86,21.49)	(18.82,19.85)	(18.20,18.96)	(18.09,18.79)	(18.11,18.80)
3-NN	18.31	17.24	15.83	15.12	15.08	14.95
	3.57	3.29	2.00	1.60	1.42	1.27
	(17.61,19.01)	(16.59,17.88)	(15.44,16.22)	(14.81,15.43)	(14.80,15.36)	(14.71,15.20)
5-NN	17.40	15.93	14.54	13.89	13.79	13.70
	3.60	2.71	1.65	1.44	1.19	1.18
	(16.70,18.11)	(15.40,16.46)	(14.22,14.86)	(13.61,14.17)	(13.55,14.02)	(13.47,13.93)
PARZEN	****	24.84	17.73	14.95	13.84	13.31
		4.84	3.00	2.40	1.49	1.54
		(23.89,25.79)	(17.14,18.32)	(14.48,15.42)	(13.55,14.13)	(13.01,13.61)
ED	13.74	12.14	11.04	10.76	10.58	10.69
	2.21	1.31	0.79	0.74	0.75	0.82
	(13.30,14.17)	(11.88,12.40)	(10.88,11.19)	(10.61,10.91)	(10.44,10.73)	(10.53,10.85)
FLDF	21.98	15.67	12.38	11.63	11.12	11.14
	5.68	3.08	1.29	1.00	0.91	0.88
	(20.86,23.09)	(15.07,16.28)	(12.13,12.64)	(11.44,11.83)	(10.94,11.30)	(10.97,11.31)
QDF	****	25.28	16.38	14.20	12.97	12.54
		5.93	2.18	1.43	1.19	1.05
		(24.12,26.45)	(15.95,16.81)	(13.92,14.49)	(12.74,13.20)	(12.33,12.74)
MQDF	15.92	12.98	11.50	10.96	10.75	10.79
	3.30	1.84	0.94	0.82	0.75	0.83
	(15.28,16.57)	(12.62,13.34)	(11.32,11.68)	(10.80,11.12)	(10.60,10.90)	(10.63,10.96)

表 3: データ集合 2 に対する識別器の誤識別率 [%]

(上段:平均値, 中段:標準偏差, 下段:95%信頼区間, Bayes 誤識別率:9%)

識別器	$n$					
	1	2	4	6	8	10
1-NN	38.74	36.44	33.34	31.67	30.66	29.70
	1.90	2.24	1.51	1.55	1.49	1.54
	(38.37,39.11)	(36.01,36.88)	(33.05,33.64)	(31.36,31.97)	(30.37,30.95)	(29.40,30.00)
3-NN	43.12	40.57	37.22	34.75	33.54	31.97
	1.86	1.53	1.50	1.36	1.23	1.15
	(42.75,43.48)	(40.27,40.87)	(36.93,37.51)	(34.48,35.01)	(33.30,33.78)	(31.74,32.19)
5-NN	45.99	43.95	40.90	38.37	36.82	34.99
	2.02	1.53	1.44	1.25	1.26	1.15
	(45.60,46.39)	(43.65,44.25)	(40.62,41.18)	(38.13,38.62)	(36.58,37.07)	(34.76,35.21)
PARZEN	****	23.30	16.18	14.35	13.86	13.83
		3.53	1.82	1.11	1.06	1.00
		(22.61,24.00)	(15.83,16.54)	(14.13,14.57)	(13.65,14.07)	(13.63,14.02)
ED	44.61	45.49	47.40	47.54	47.88	48.09
	2.62	2.38	2.09	1.81	1.50	1.48
	(44.10,45.13)	(45.02,45.96)	(46.99,47.81)	(47.18,47.89)	(47.58,48.17)	(47.80,48.38)
FLDF	46.53	46.02	47.53	47.72	47.95	48.14
	2.97	2.39	2.03	1.83	1.47	1.44
	(45.94,47.11)	(45.55,46.49)	(47.13,47.93)	(47.36,48.08)	(47.66,48.24)	(47.86,48.42)
QDF	****	29.93	17.50	13.97	12.32	11.56
		4.10	2.09	1.51	1.15	1.04
		(29.13,30.74)	(17.09,17.91)	(13.68,14.27)	(12.09,12.54)	(11.35,11.76)
MQDF	17.09	12.79	11.22	10.46	10.10	9.89
	4.57	1.75	0.97	0.80	0.74	0.70
	(16.19,17.98)	(12.45,13.13)	(11.03,11.41)	(10.30,10.61)	(9.95,10.24)	(9.76,10.03)

表 4: データ集合 3 に対する識別器の誤識別率 [%]

(上段:平均値, 中段:標準偏差, 下段:95%信頼区間, Bayes 誤識別率:1.0%)

識別器	$\eta$					
	1	2	4	6	8	10
1-NN	14.97 3.52 (14.28,15.66)	12.34 2.49 (11.85,12.83)	10.56 1.71 (10.23,10.90)	9.02 1.52 (8.73,9.32)	8.81 1.13 (8.55,9.06)	8.14 0.91 (7.96,8.32)
3-NN	15.48 4.34 (14.63,16.33)	11.40 2.23 (10.96,11.84)	9.34 1.38 (9.07,9.61)	8.36 1.24 (8.12,8.60)	7.76 1.09 (7.55,7.97)	7.37 0.92 (7.19,7.55)
5-NN	16.88 5.01 (15.90,17.86)	11.53 2.34 (11.07,11.99)	9.33 1.43 (9.05,9.61)	8.14 1.19 (7.90,8.37)	7.48 1.00 (7.28,7.68)	7.07 0.92 (6.88,7.25)
PARZEN	*****	6.58 2.58 (6.07,7.08)	3.69 0.66 (3.56,3.81)	3.00 0.53 (2.89,3.10)	2.83 0.53 (2.73,2.93)	2.64 0.42 (2.56,2.72)
ED	16.69 4.78 (15.75,17.63)	15.17 4.08 (14.37,15.97)	14.70 3.36 (14.04,15.36)	14.44 2.80 (13.89,14.98)	14.17 2.54 (13.67,14.66)	14.41 2.13 (13.99,14.82)
FLDF	14.28 5.41 (13.22,15.34)	9.33 2.14 (8.91,9.75)	7.43 1.27 (7.18,7.68)	6.67 0.88 (6.50,6.84)	6.28 0.65 (6.15,6.41)	6.13 0.63 (6.00,6.25)
QDF	*****	6.87 2.14 (6.45,7.29)	3.54 0.78 (3.39,3.69)	2.72 0.51 (2.62,2.82)	2.44 0.38 (2.36,2.51)	2.25 0.37 (2.18,2.32)
MQDF	11.29 3.13 (10.68,11.90)	6.10 1.94 (5.72,6.48)	3.39 0.74 (3.24,3.53)	2.68 0.53 (2.57,2.78)	2.44 0.38 (2.36,2.51)	2.25 0.37 (2.18,2.32)

表 5: MQDF の  $\xi$  の最適値

$\eta$	1	2	4	6	8	10
データ集合 1	2	1	2	1	1	1
データ集合 2	3	3	4	4	4	4
データ集合 3	3	6	6	6	7	7

表 6: PARZEN の  $h$  の最適値

$\eta$	1	2	4	6	8	10
データ集合 1	**	1.4	1.4	1.4	1.4	1.4
データ集合 2	**	1.2	0.8	0.8	0.8	0.8
データ集合 3	**	1.8	1.2	1.2	1.2	1.0

表 7: 識別性能による識別器の順位

	1位	2位	3位	4位	5位	6位	7位	8位
データ集合 1	ED	MQDF	FLDF	5-NN	3-NN	1-NN	PARZEN	QDF
データ集合 2	MQDF	PARZEN	QDF	1-NN	3-NN	5-NN	ED	FLDF
データ集合 3	MQDF	PARZEN	QDF	FLDF	3-NN	5-NN	1-NN	ED

8. 表中の\*は、共分散行列の非正則性のため、識別器が得られなかったことを示す。このように、訓練サンプル数が少ないと、共分散行列が正則でなくなるため、QDF および PARZEN は実現できない、ということになる。
9. 最近傍識別器 (1-NN) については、訓練サンプル数が無限にあるとき、その誤識率が Bayes 誤識率の 2 倍を越えない、ということが理論的に示されている [18]。実験結果は、訓練サンプル数の有限性のため、偏りのある誤識率が得られたことを示している。この偏りについては、Fukunaga らにより検討がなされている [19]。文献 [19] で、64 次元特徴空間上で、訓練サンプル数を 1000 個から 10000 個に増加させても、誤識率の偏りは 6.9% しか減少しない、ということが示された。このことから、Fukunaga は、最近傍識別器の誤識率を改善するためには、かなり大量の訓練サンプルが必要である、と述べている [7]。現実には、訓練サンプル数が限られているため、最近傍識別器の改良には、用いる距離を検討することが、重要と思われる。距離の改良については、文献 [20][21] がある。

#### 4 むすび

本論文では、特徴数に対する訓練サンプル数の比が小さい状況下で、識別器の性能比較を誤識率の観点から行った。比較された識別器は、Fisher の線形識別関数、二次識別関数、修正二次識別関数、k 最近傍識別器、ユークリッド距離識別器、Parzen 識別器である。これらの中では修正二次識別関数が有効であった。

実在データを用いて更に比較検討を行い、真に有効な識別器を設計することが、今後の課題である。

#### 参考文献

- [1] S. J. Raudys and A. K. Jain, "Small Sample Size Effects in Statistical Pattern Recognition", IEEE Trans. PAMI-13, 3, pp.252-264 (1991).
- [2] D. H. Foley, "Consideration of Sample and Feature Size", IEEE Trans. IT-18, 5, pp.618-626 (1972).
- [3] L. Kanal and B. Chandrasekaran, "On Dimensionality and Sample Size in Statistical Pattern Classification", in Proc. 1968 National Electronics Conf., pp.2-7 (1968).
- [4] A. K. Jain and B. Chandrasekaran, "Dimensionality and Sample Size Considerations in Pattern Recognition Practice", in Handbook of Statistics, Vol.2, P. R. Krishnaiah and L. N. Kanal, Eds., North-Holland, pp.835-855 (1982).
- [5] 鶴岡, 栗田, 原田, 木村, 三宅, "加重方向指数ヒストグラム法による手書き漢字・ひらがな認識", 信学論 (D), J70-D, 7, pp.1390-1397 (1987).
- [6] K. Fukunaga and R. R. Hayes, "Effects of Sample Size in Classifier Design", IEEE Trans. PAMI-11, 8, pp.873-885 (1989).
- [7] K. Fukunaga, "Introduction to Statistical Pattern Recognition", Second Edition, Academic Press (1990).
- [8] K. Fukunaga and R. R. Hayes, "Estimation of Classifier Performance", IEEE Trans. PAMI-11, 10, pp.1087-1101 (1989).
- [9] P. A. Lachenbruch and R. M. Mickey, "Estimation of Error Rates in Discriminant Analysis", Technometrics, 10, 1, pp.1-11 (1968).
- [10] A. K. Jain, "Advances in Statistical Pattern Recognition", in Pattern Recognition Theory and Applications, eds., P. A. Devijver and J. Kittler, pp.1-19, Springer-Verlag Berlin Heidelberg (1987).
- [11] T. Marill and D. M. Green, "On the Effectiveness of Receptors in Recognition Systems", IEEE Trans. IT-9, pp.11-27 (1963).
- [12] K. Fukunaga and T. F. Krile, "Calculation of Bayes Recognition Error for Two Multivariate Gaussian Distributions", IEEE Trans. C-18, pp.220-229 (1969).
- [13] A. K. Jain, R. C. Dubes and C. C. Chen, "Bootstrap Techniques for Error Estimation", IEEE Trans. PAMI-9, 5, pp.628-633 (1987).

- [14] K. Fukunaga and D. M. Hummels, "Bayes Error Estimation Using Parzen and k-NN Procedures", IEEE Trans. PAMI-9, 5, pp.634-643 (1987).
- [15] A. K. Jain and M. D. Ramaswami, "Classifier Design with Parzen Windows", Pattern Recognition and Artificial Intelligence, pp.211-228, Elsevier Science Publishers B. V. (1988).
- [16] 木村, 高階, 鶴岡, 三宅, "2次識別関数のピーキング現象とその防止に関する考察", 信学論 (D), J69-D, 9, pp.1328-1334 (1986).
- [17] J. Van Ness, "On the Dominance of Non-parametric Bayes Rule Discriminant Algorithms in High Dimensions", Pattern Recognition, 12, pp.355-368 (1980).
- [18] R. O. Duda and P. E. Hart, "Pattern Classification and Scene Analysis", John Wiley & Sons (1973).
- [19] K. Fukunaga and D. M. Hummels, "Bias of Nearest Neighbor Estimates", IEEE Trans. PAMI-9, pp.103-112 (1987).
- [20] R. D. Short and K. Fukunaga, "The Optimal Distance Measure for Nearest Neighbor Classification". IEEE Trans. IT-27, pp.622-627 (1981).
- [21] K. Fukunaga and T. E. Flick, "An Optimal Global Nearest Neighbor Metric", IEEE Trans. PAMI-6, pp.314-318 (1984).