

# 固有ベクトルのずれの影響を軽減する固有値補正法

岩村雅一      大町真一郎      阿曾弘具

東北大学大学院工学研究科

〒 980-8579 仙台市青葉区荒巻字青葉 05

Tel:022-217-7088      Fax:022-263-9419

E-mail:masa@aso.ecei.tohoku.ac.jp

あらまし マハラノビス距離などの統計的手法を用いたパターン認識ではパターンの分布が既知であることが求められるが、一般には未知であり、学習サンプルから推定された標本共分散行列を用いる。このとき推定に用いる学習サンプルが少ないと推定誤差が生じ、認識性能が低下することが知られている。固有値の推定誤差についてはよく調査され、固有値を修正する方法が提案されているが、固有ベクトルの推定誤差についてはほとんど考慮されてこなかった。本論文では、固有ベクトルの推定誤差について調査し、固有ベクトルの推定誤差についても認識性能に悪影響を与える可能性があることを示す。そして、固有ベクトルの推定誤差の影響を固有値の修正により軽減する手法を提案する。

キーワード マハラノビス距離, 固有ベクトル, 推定誤差, パターン認識, NIST Special Database 19

## A Modification Method of Eigenvalues to Avoid Recognition Accuracy Reduction Caused by Estimation Errors of Eigenvectors.

Masakazu IWAMURA, Shin'ichiro OMACHI, and Hiroto ASO

Department of Electrical and Communication Engineering,

Graduate School of Engineering, Tohoku University

Aoba 05, Aramaki, Aoba-ku, Sendai-shi, 980-8579 Japan

Tel:022-217-7088      Fax:022-263-9419

E-mail:masa@aso.ecei.tohoku.ac.jp

Abstract Statistical pattern recognition methods such as Mahalanobis distance need to know true distributions of the patterns, but it is usually impossible. As a result, the covariance matrices estimated from training samples are used. If training samples are not enough, estimation errors in both eigenvalues and eigenvectors will occur and they will cause recognition accuracy reduction. Many researchers have studied on estimation errors of eigenvalues and have proposed various methods to modify eigenvalues. However, estimation errors of eigenvectors have not been considered enough. In this paper, we investigate estimation errors of eigenvectors and show the possibility that they make recognition performance down. We propose a new method estimating eigenvalues so as to reduce bad influence on estimation errors of eigenvectors.

key words Mahalanobis distance, eigenvector, estimation error, pattern recognition, NIST Special Database 19

# 1 はじめに

多次元正規分布の確率密度関数から導かれる2次識別関数は、サンプルの分布が  $n$  次元正規分布に従うときに最適である。マハラノビス距離は、サンプルの分布が  $n$  次元正規分布に従い、各カテゴリの共分散行列が同じであるときに最適である。これらの統計的手法は Bayes 決定と呼ばれる、損失の期待値を最小限にとどめる分類法を導く。

2次識別関数やマハラノビス距離のような統計的手法を用いて認識システムを構築するには真の分布が既知であることが必要である。しかし、一般には未知であるため、分布形状を正規分布と仮定し、そのパラメータを学習サンプルを用いて推定する。このとき、パラメータの推定に用いられる学習サンプル数が十分でないと、推定された分布は誤差を含み、認識性能が著しく低下することが知られている。

分布の推定誤差に関する問題は統計的パターン認識において避けては通れない問題であり、古くから議論されてきた。この問題への対処法は、信頼性の低い推定値を用いないで認識を行なう手法と、推定誤差の傾向から推定値を補正する手法に大別できる。後者は前者に比べて難しい反面、分布が正しく推定できれば認識性能がより良くなることが期待できる。

前者には修正2次識別関数 (MQDF) [1, 2], 改良型マハラノビス距離 (MMD)[3], 簡素化マハラノビス距離 (SMD)[4] などがある。MQDF は、木村らにより提案された手法で、式の形の上ではマハラノビス距離とユークリッド距離の加重和からなる。MMD は、加藤らにより提案された手法で、マハラノビス距離を計算する際、固有値にバイアスを加える。SMD は、孫らにより提案された手法で、MQDF のユークリッド距離に相当する項の重みを比較的信頼性の高い固有値から求める。また、竹下らは MQDF で使用するための信頼できる固有値・固有ベクトル数を調査した [5]。

一方、後者には酒井らによる2次識別関数 (RQDF)[6] がある。酒井らは、Fukunaga が行なった、摂動法を用いて固有値・固有ベクトルの推定誤差の分散を求める解析 [7] を2次にまで拡張して標本共分散行列の固有値の偏りを調査し、標本固有値の偏りを補正する手法を提案している。しかし、この手法は固有値の偏りの補正のみで固有ベクトルの推定誤差の影響については考慮していない。

そこで本論文では、まず固有ベクトルの推定誤差について調査する。ここでは実際の文字画像から求めた特徴ベクトルとパラメータが既知の正規分布に

従う人工サンプルからそれぞれ推定した固有ベクトルのずれを2次元平面上にプロットすることによって行なう。その結果、固有ベクトルの推定誤差(ずれ)が2次識別関数やマハラノビス距離などの統計的手法に影響し、認識性能の低下が引き起こされる可能性があることを示す。さらに、固有ベクトルの推定誤差の影響を固有値に吸収させることにより軽減する手法を提案する。

なお、本論文では統計的手法のうち、マハラノビス距離を用いるが、本質的な部分は他の手法にも同様に適用できると思われる。

## 2 マハラノビス距離を用いたパターン認識

### 2.1 主成分分析

$N$  個の学習サンプルを  $t_1, t_2, \dots, t_N$  とすると、標本平均ベクトル  $\hat{\mu}$  と標本共分散行列  $\hat{\Sigma}$  は

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^N t_k \quad (1)$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{k=1}^N (t_k - \hat{\mu})(t_k - \hat{\mu})^T \quad (2)$$

のように推定することができる。ここで標本共分散行列  $\hat{\Sigma}$  は第  $i$  固有値  $\hat{\lambda}_i$  とそれに対応する固有ベクトル  $\hat{\phi}_i$  を用いて

$$\hat{\Sigma} = \hat{\Phi} \hat{\Lambda} \hat{\Phi}^T \quad (3)$$

$$= \sum_{k=1}^d \hat{\lambda}_k \hat{\phi}_k \hat{\phi}_k^T \quad (4)$$

と表わすことができる。ここで、 $\hat{\Lambda}$  と  $\hat{\Phi}$  は

$$\hat{\Lambda} = \text{diag}(\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_d) \quad (5)$$

$$\hat{\Phi} = (\hat{\phi}_1 \quad \hat{\phi}_2 \quad \dots \quad \hat{\phi}_d) \quad (6)$$

であり、 $d$  は特徴量の次元数である。なお、固有値の並び順に本質的な意味はないが、本論文では全て降順であるとする。

### 2.2 マハラノビス距離

マハラノビス距離  $d^2(x)$  は  $x$  を未知入力ベクトルとすると、真の共分散行列  $\Sigma$  を用いて

$$d^2(x) = (x - \mu)^T \Sigma^{-1} (x - \mu) \quad (7)$$

と表わされる．しかし，一般に真の平均ベクトル  $\mu$  と真の共分散行列  $\Sigma$  は未知であるため，標本平均ベクトル  $\hat{\mu}$  と標本共分散行列  $\hat{\Sigma}$  を用いて

$$d^2(x) = (x - \hat{\mu})^T \hat{\Sigma}^{-1} (x - \hat{\mu}) \quad (8)$$

のように計算する．また，標本共分散行列  $\hat{\Sigma}$  を固有値と固有ベクトルに置き換えることにより

$$d^2(x) = \sum_{k=1}^d \frac{1}{\hat{\lambda}_k} \left( \hat{\phi}_k \cdot (x - \hat{\mu}) \right)^2 \quad (9)$$

と表わすこともできる．実際の認識は全てのカテゴリについて未知入力  $x$  とのマハラノビス距離  $d^2(x)$  を求めて，未知入力  $x$  を距離が最も小さいカテゴリに決定することで行なわれる．

### 3 学習サンプル数と固有ベクトルの推定誤差に関する調査

#### 3.1 固有ベクトルの推定誤差の評価方法

ここでは固有ベクトルのずれを

1. ずれの大きさ
2. ずれの方向

という観点で調査する．しかし，固有ベクトルの推定誤差は固有値の推定誤差と同様，確率的な変動でしか得ることができない．このため，以下の定義を行ない，視覚的に調査する．

固有ベクトルの推定誤差は 2 本の固有ベクトルの内積を用いて表わすことができる．学習サンプル  $N$  個から推定した標本共分散行列の第  $k$  固有ベクトルを  $\hat{\phi}_k^N$  とおく． $N_1, N_2$  を分布の推定に用いる学習サンプル数とすれば（ただし， $N_1 > N_2$  であるとする）， $\hat{\phi}_i^{N_1}$  と  $\hat{\phi}_j^{N_2}$  の内積を  $\hat{\phi}_i^{N_1} \cdot \hat{\phi}_j^{N_2}$  と表わす． $i = j$  のときは  $|\hat{\phi}_i^{N_1} \cdot \hat{\phi}_j^{N_2}|$  が 1 に近いほど推定誤差が小さいことを意味し， $i \neq j$  のときは  $|\hat{\phi}_i^{N_1} \cdot \hat{\phi}_j^{N_2}|$  がずれの方向と大きさを表わす．

例えば学習サンプル  $N_1$  個から求めた共分散行列の固有ベクトルと  $N_2$  個から求めた共分散行列の固有ベクトルが同じであれば，対応する固有ベクトルは全て同じであるので，2 組のベクトルのずれは無く，

$$|\hat{\phi}_i^{N_1} \cdot \hat{\phi}_j^{N_2}| = \begin{cases} 1, & i = j, \\ 0, & i \neq j. \end{cases} \quad (10)$$

が成り立つ．

0	1	2	3	4
40363	44704	40072	41112	39154
5	6	7	8	9
36606	39937	41893	39579	39533

表 1: NIST Special Database 19 に含まれる数字サンプルの数

さらに  $F(N_1, N_2)$  を

$$F(N_1, N_2) = \hat{\Phi}_1^T \hat{\Phi}_2 \quad (11)$$

$$= \begin{pmatrix} \hat{\phi}_1^{N_1} \cdot \hat{\phi}_1^{N_2} & \dots & \hat{\phi}_1^{N_1} \cdot \hat{\phi}_d^{N_2} \\ \vdots & \ddots & \vdots \\ \hat{\phi}_d^{N_1} \cdot \hat{\phi}_1^{N_2} & \dots & \hat{\phi}_d^{N_1} \cdot \hat{\phi}_d^{N_2} \end{pmatrix} \quad (12)$$

とおく．ここで， $\hat{\Phi}_1$  と  $\hat{\Phi}_2$  の各列は  $N_1$  個， $N_2$  個の学習サンプルから求めた共分散行列の固有ベクトルある．

$F(N_1, N_2)$  は対角成分に注目すれば各固有ベクトルのずれの大きさがわかり，非対角成分に注目すればずれの方向がわかる．この  $F(N_1, N_2)$  を 2 次元平面上にプロットすることで，確率的にしか得られない固有ベクトルの推定誤差の傾向をずれの大きさ，ずれの方向という面から視覚的に調査することができる．ただし， $F(N_1, N_2)$  の各要素の絶対値をプロットすることにする．今後の議論においても  $F(N_1, N_2)$  の各要素は大きさのみ用いるので影響はない．

#### 3.2 固有ベクトルの推定誤差の調査

実際の文字サンプルと人工サンプルを用いて固有ベクトルのずれの傾向について調査を行なう．

##### 3.2.1 実際の文字サンプルによる実験

実際の文字サンプルから固有ベクトルを推定し，学習サンプル数の違いによって推定誤差がどのように変化するかを調べた．

実験には NIST(National Institute of Standards and Technology) の NIST Special Database 19(以降，NIST19) を用いた．このデータベースには数字とアルファベットが含まれているが，ここでは，よりサンプル数の多い数字サンプル 10 字種を用いることにした．このデータベースに含まれる数字のサンプル数は表 1 の通りである．このデータベースの文字サンプルの大きさは  $128 \times 128$  画素である．これ

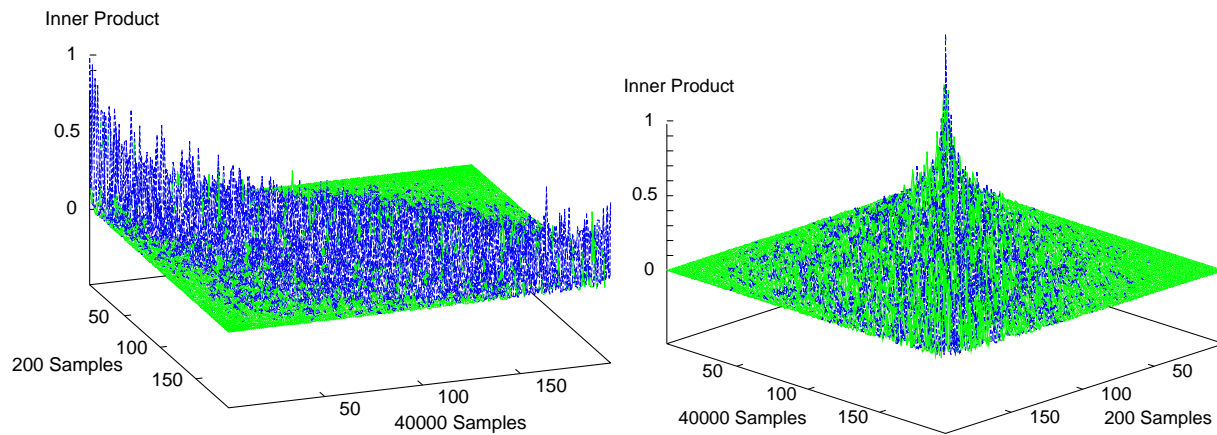


図 1: 実際の文字サンプルの  $F(200, 40000)$

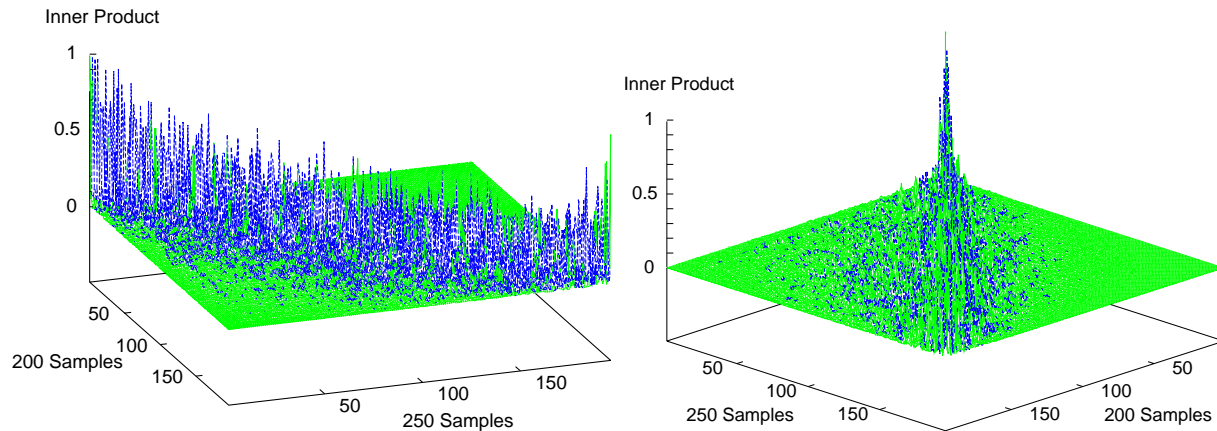


図 2: 実際の文字サンプルの  $F(200, 250)$

を  $64 \times 64$  画素に非線形正規化 [8] した後, 196 次元の方向線素特徴量 [9] を抽出した.

実験は学習サンプル数  $N_1, N_2$  の組み合わせを変えて行なった. 文字 “0” の結果の一部を図 1 から図 4 に示す. 実験結果は字種によらずほぼ同様であった. この実験から

1. 推定に用いる学習サンプル数が違うと, 推定された固有ベクトルにはずれが生じる. これは推定誤差の大きさがサンプル数に依存していることを示している.
2. 用いた学習サンプル数の差が大きいほど推定誤差が大きい.
3. 用いた学習サンプル数の差が同程度の場合, 絶対値が小さいほど推定誤差が大きい.

(例:  $F(1000, 2000)$  と  $F(39000, 40000)$  では前者のほうがずれが大きい)

ということがわかる.

### 3.2.2 人工サンプルによる実験

実際の文字サンプルによる実験では, 有限個の学習サンプルを用いた調査しかできない. そこで人工サンプルを作成し, 実験を行なった. これにより, 学習サンプル数が無限大に相当する真の値との関係も調査することができる.

実験に用いた人工サンプルは NIST19 の数字サンプル 10 字種から求めた. この人工サンプルは 3.2.1 の実際の文字サンプルと同様に 196 次元の方向線素特徴量を抽出し, 一定数の文字サンプルから推定した平均, 分散を真の値とみなし, 乱数を用いて作成した. 文字 “0” の場合は母集団として 40,000 個の文字サンプルから推定した分布を用いた.

文字 “0” の結果の一部を図 5 から図 8 に示す. 実験結果は他の字種でもほぼ同様であった. 人工サンプルはその作成法ゆえ, 正規分布に従う. しかし, 実際の文字サンプルは必ずしも正規分布にしたがうとは限らない. 人工サンプルを用いた場合より実際の

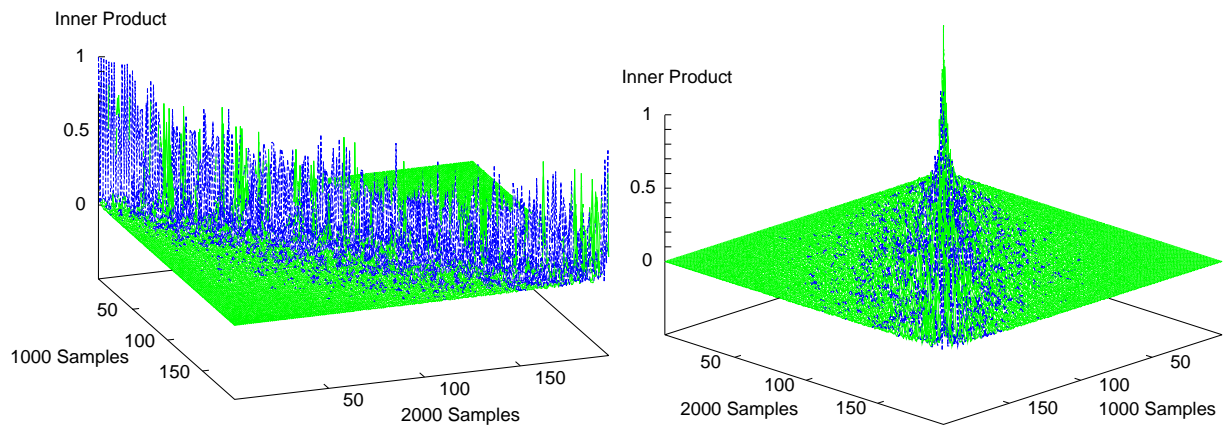


図 3: 実際の文字サンプルの  $F(1000, 2000)$

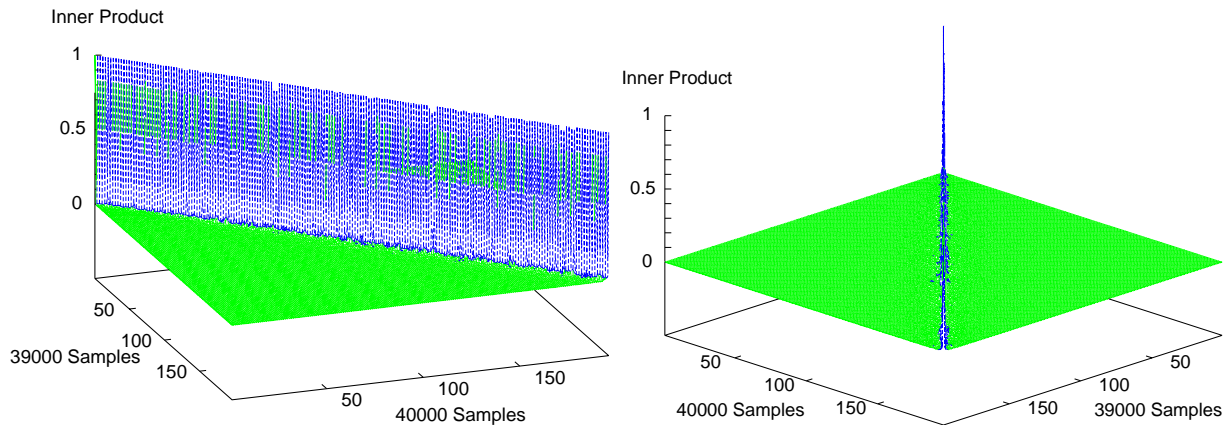


図 4: 実際の文字サンプルの  $F(39000, 40000)$

文字サンプルを用いた場合のほうがずれが若干大きいですが、人工サンプルの結果は 3.2.1 の実際の文字サンプルとほとんど同じ傾向を示している。また、真の値と比較した場合、10,000,000 個の学習サンプルから推定した固有ベクトルもずれを含むことがわかる。

#### 4 固有ベクトルの推定誤差を補正する方法の検討

3 章の実験は固有ベクトルがずれる方向についての興味深い事実を示している。まず、固有ベクトルは Fukunaga の理論式 [7] にあるように、固有値の大きさが同程度である固有ベクトル軸にずれ易いことがわかる。しかし、そのみでなく、固有ベクトルは固有値の大きさが極端に違う軸に対してもわずかながらずれていることがわかる。この傾向は、正規分布に従う人工サンプルのみでなく、必ずしも正規分布に従う保証がない実際の文字サンプルでも見

られた。

マハラノビス距離は (9) 式のように表わすことができるが、この式は“各固有ベクトル軸  $\hat{\phi}_k$  方向の重みを  $\frac{1}{\hat{\lambda}_k}$  にして距離を計算する”と解釈することができる。

これまでの研究で固有値の推定誤差を考慮した固有値補正法が提案され、認識性能の向上に効果を上げているが、固有値は固有ベクトルに依存している。固有ベクトルのずれを考慮しないで固有値のみ補正するという方法では固有値と固有ベクトルのバランスが崩れ、補正の効果が半減してしまう恐れがある。これらのことから、固有値のみでなく固有ベクトルの推定誤差についても補正が必要であると言える。

##### 4.1 固有ベクトルを補正する方法の検討

今、学習サンプルとして利用可能なサンプルが  $N_2$  個あるとする。その  $N_2$  個の学習サンプルを用いて

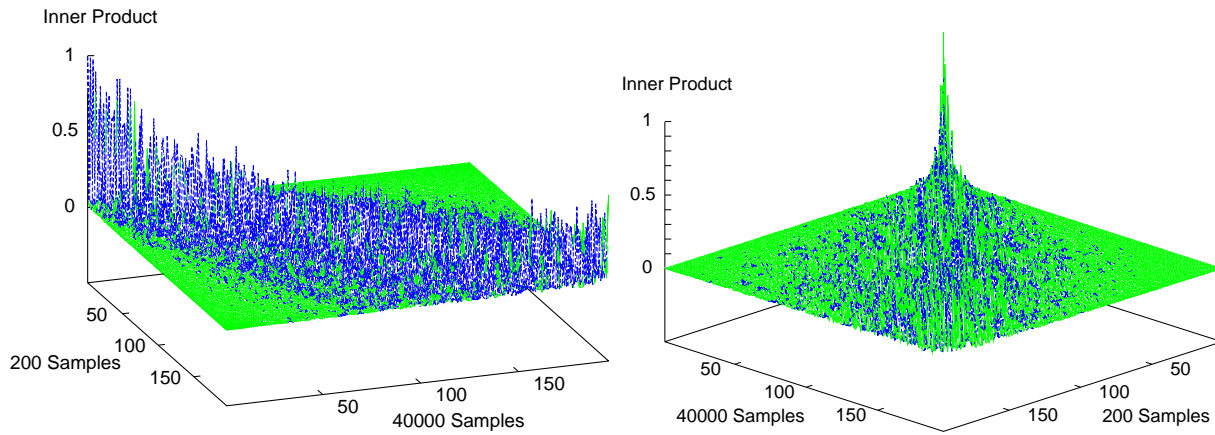


図 5: 人工サンプルの  $F(200, 40000)$

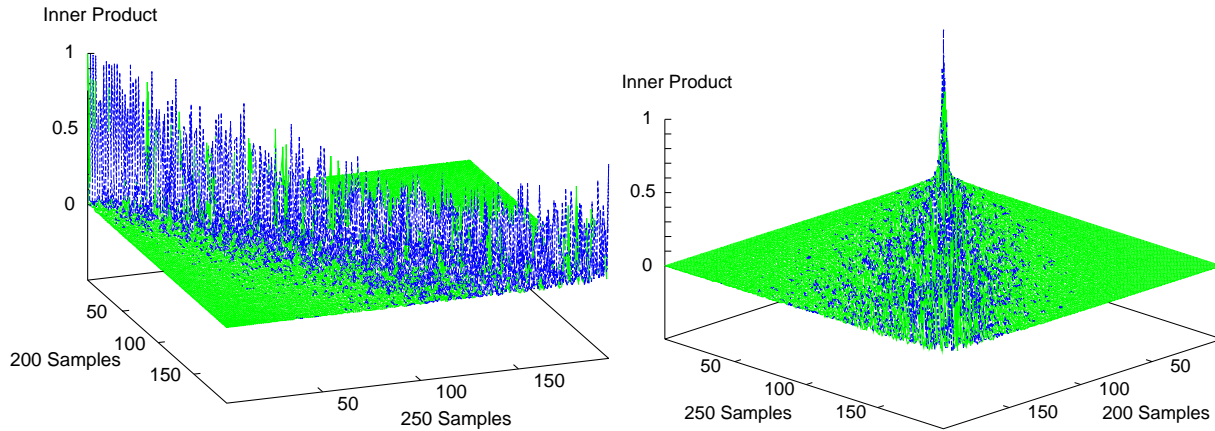


図 6: 人工サンプルの  $F(200, 250)$

辞書用に分布 (平均・分散) を推定する。ここで、固有値の推定誤差を考慮して固有値を補正した結果、 $N_1$  個の学習サンプルから推定した固有値と同様の値が得られたとする。ところが、固有値の推定誤差のみ減少したとしても認識の際には固有ベクトルも使用するため、先に示した理由により必ずしも十分な結果が得られないことが予想できる。

そこで、固有ベクトルの推定誤差も固有値の推定誤差同様、学習サンプル  $N_1$  個から推定したときに含まれる程度のずれの量に修正すればよいが、固有ベクトルは固有値に比べて自由度が高く、確率的な固有ベクトルの変動からでは補正するのが困難である。しかし、固有値の自由度は固有ベクトルに比べて著しく小さいため、確率的な固有ベクトルの変動からでも固有値を補正することは可能である。そこで、固有ベクトル自身を補正するのではなく、固有ベクトルのずれを考慮して固有値をさらに補正する手法を提案する。

## 4.2 固有値を補正する方法

前提条件として、学習サンプル数  $N_1$  個から推定された固有ベクトルと学習サンプル数  $N_2$  個から推定された固有ベクトルとのずれを表わす  $F(N_1, N_2)$  が既知であるとする (ここでも  $N_1 > N_2$  であるとする)。また、学習サンプル数  $N_2$  個から推定された固有値を任意の固有値補正法で補正して得られた値を  $\lambda'_k$  とし、学習サンプル数  $N_1$  個から推定された固有値と同程度の推定誤差を含むとする。このとき、次式を用いて、 $F(N_1, N_2)$  と  $\lambda'_k$  から固有ベクトルの推定誤差を吸収する固有値  $\hat{\lambda}_k$  を推定する。

$$\hat{\lambda}_k = \sum_{i=1}^d (\phi_i^{N_1} \cdot \phi_k^{N_2})^2 \lambda'_i \quad (13)$$

この式の導出は以下のように行なう。学習サンプル数  $N_1$  個から推定された標本共分散行列を  $\hat{\Sigma}_1$  とし、その固有値、固有ベクトルを  $\hat{\Lambda}_1, \hat{\Phi}_1$  とする。学習サンプル数が  $N_2$  個のときも同様に  $\hat{\Sigma}_2, \hat{\Lambda}_2, \hat{\Phi}_2$

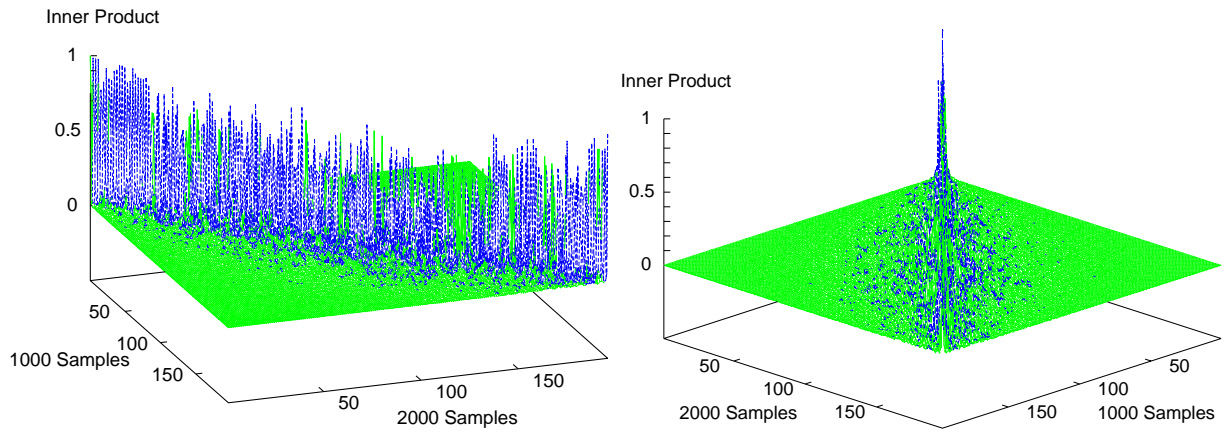


図 7: 人工サンプルの  $F(1000, 2000)$

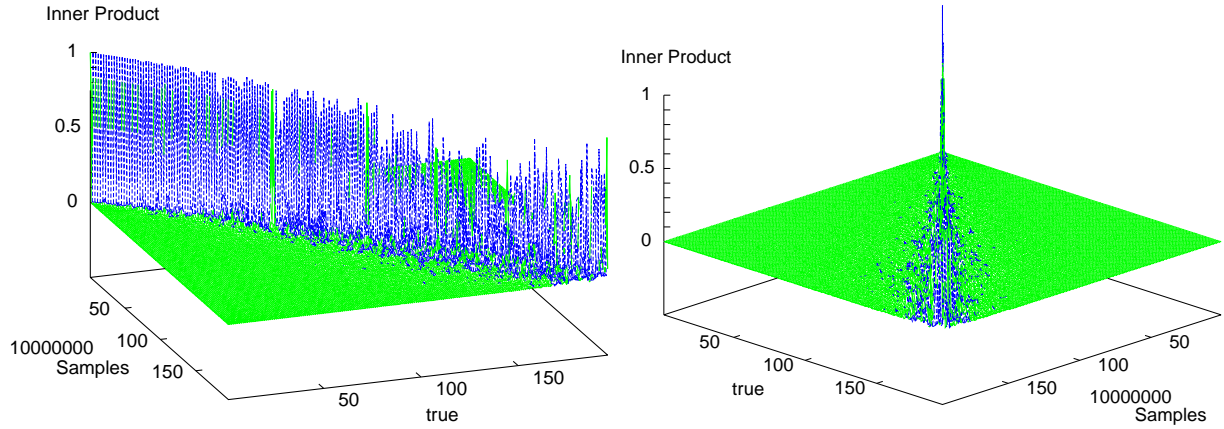


図 8: 人工サンプルの  $F(10000000, \infty)$

とすると、次の 2 式が成り立つ。

$$\hat{\Sigma}_1 = \hat{\Phi}_1 \hat{\Lambda}_1 \hat{\Phi}_1^T \quad (14)$$

$$\hat{\Sigma}_2 = \hat{\Phi}_2 \hat{\Lambda}_2 \hat{\Phi}_2^T \quad (15)$$

固有ベクトルをそのままに  $\hat{\Sigma}_2$  を  $\hat{\Sigma}_1$  に近づけるように固有値を修正することを考え、 $\hat{\Sigma}_2 = \hat{\Sigma}_1$  とおくと、

$$\hat{\Phi}_1 \hat{\Lambda}_1 \hat{\Phi}_1^T = \hat{\Phi}_2 \hat{\Lambda}_2 \hat{\Phi}_2^T \quad (16)$$

$$\hat{\Lambda}_2 = \hat{\Phi}_2^T \hat{\Phi}_1 \hat{\Lambda}_1 \hat{\Phi}_1^T \hat{\Phi}_2 \quad (17)$$

$$= F(N_1, N_2)^T \hat{\Lambda}_1 F(N_1, N_2) \quad (18)$$

となり、ここから (13) 式が導出される。

ここでは、 $N_1$  を有限個のサンプルとしたが、 $F$  を求めるときに人工サンプルを用いるなどの方法で  $N_1$  個から求めた固有値・固有ベクトルを擬似的に真の分布にすることで学習サンプルが無限個の場合をシミュレートすることができる。

## 5 人工データによる認識実験

提案手法の有効性を示すために、人工サンプルを用いた認識実験を行なった。人工サンプルを使用するので、真の分布がわかる。固有値のみ真の値にしたものと、提案手法により固有値を補正したものの二つを比較する。提案手法で固有値を補正するとき用いた固有ベクトルのずれを表わす  $F$  は、 $N_1 \rightarrow \infty$  としたものである。

実験に用いたサンプルは NIST19 の数字 10 字種である。作成方法は、3.2.2 のときと同様に NIST19 の文字サンプルから 196 次元の方向線素特徴量を抽出し、各字種 36,000 個の学習サンプルから分布 (平均・分散) を推定し、これを母集団として字種毎に各 1,000 個のサンプルを認識用に、各 1,000,000 個のサンプルを学習用に作成した。

辞書作成用に用いる学習サンプル数を 1,000,000 個の範囲内で変え、学習サンプル数と認識率の関係を調査した。認識用サンプルは常に 1,000 個使用した。

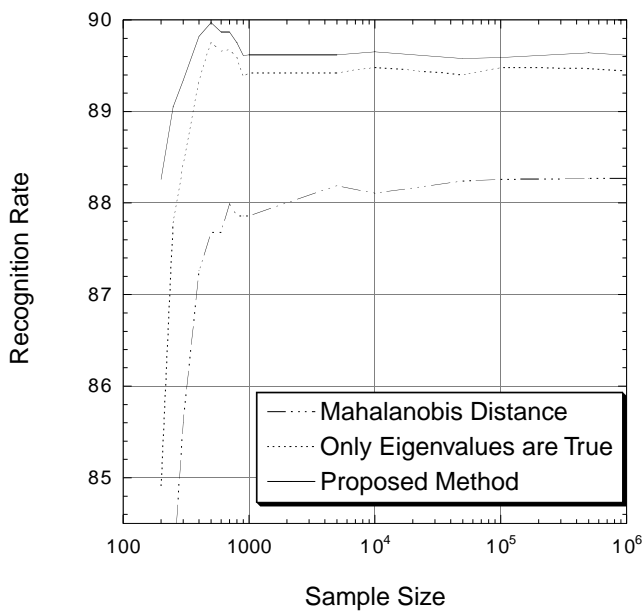


図 9: 提案手法の有効性を示す実験

結果を図 9 に示す．固有値のみ真の値にしたものを “Only Eigenvalues are True”，提案手法により固有値を補正したものを “Proposed Method” とした．また，参考のため標本共分散行列をそのまま用いたマハラノビス距離による認識結果を “Mahalanobis Distance” として示した．

この結果，再推定に用いる学習サンプル数によらず，提案手法のほうが真の固有値を用いた場合よりも認識率が高かった．真の固有値を用いた場合は，学習サンプル数によらず標本共分散行列をそのまま用いたマハラノビス距離よりも認識率が高かった．これにより，提案手法の有効性が確認された．

## 6 結び

本論文では，まず実際の文字サンプルと人工サンプルを用いて固有ベクトルのずれについて傾向を調べた．その結果，固有ベクトルのずれがマハラノビス距離に影響し，認識性能の低下が引き起こされる可能性があることを示した．そして，固有ベクトルの推定誤差の影響を固有値に吸収させることにより軽減する手法を提案した．

また，今回は，固有値を推定するときに用いた固有ベクトルのずれ  $F(N_1, N_2)$  を実験的に求めたが，これを理論的に導くことは今後の課題である．

## 参考文献

- [1] 木村文隆，高階健治，鶴岡信治，三宅康二，“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)，vol.J69-D, no.9, pp.1328–1334, Sep. 1986.
- [2] F. Kimura and M. Shridhar, “Handwritten numerical recognition based on multiple algorithms,” Pattern Recognition, vol.24, no.10, pp.969–983, 1991.
- [3] 加藤 寧，安倍正人，根元義章，“改良型マハラノビス距離を用いた高精度な手書き文字認識”，信学論 (D-II)，vol.J79-D-II, no.1, pp.45–52, Jan. 1996.
- [4] F. Sun, S. Omachi, and H. Aso, “Precise selection of candidates for handwritten character recognition using feature regions,” IEICE Trans. Inf. & Syst., vol.E79-D, no.5, pp.510–515, May, 1996.
- [5] 竹下鉄夫，木村文隆，三宅康二，“マハラノビス距離の推定誤差に関する考察”，信学論 (D)，vol.J70-D, no.3, pp.567–573, Mar. 1987.
- [6] 酒井 充，米田政明，長谷博行，丸山 博，直江美知子，“固有値の偏り補正に基づく 2 次識別関数”，信学論 (D-II)，vol.J82-D-II, no.4, pp.631–640, Apr. 1999.
- [7] K. Fukunaga, “Introduction to statistical pattern recognition,” 2nd edition, Academic Press, pp.425–435, 1990.
- [8] 山田博三，斉藤泰一，山本和彦，“線密度イコライゼーション—相関法のための非線形正規化法”，信学論 (D)，vol.J67-D, no.11, pp.1379–1383, Nov. 1984.
- [9] 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”，信学論 (D-II)，vol.J74-D-II, no.3, pp.330–339, Mar. 1991.