

The Optimum Classifier and the Performance Evaluation by Bayesian Approach

Xuexian HAN, Tetsushi WAKABAYASHI, and Fumitaka Kimura

Faculty of Engineering, Mie University
1515 Kamihama, Tsu 514-8507, JAPAN

Tel: 81-59-231-9457

Fax: 81-59-231-9456

e-mail: kimura@hi.info.mie-u.ac.jp

Abstract This paper deals with the optimum classifier and the performance evaluation by the Bayesian approach. Gaussian population with unknown parameters is assumed. The conditional density given a limited sample of the population has a relationship to the multivariate t -distribution. The mean error rate of the optimum classifier is theoretically evaluated by the quadrature of the conditional density. To verify the optimality of the classifier and the correctness of the mean error calculation, the results of Monte Carlo simulation employing a new sampling procedure are shown. It is also shown by the comparative study that the Bayesian formulas of the mean error rate have the following characteristics.

- 1) The unknown population parameters are not required in its calculation.
- 2) The expression is simple and clearly shows the limited sample effect on the mean error rate.
- 3) The relationship between the prior parameters and the mean error rate is explicitly expressed.

Keywords Statistical pattern recognition, optimum classifier, Monte Carlo simulation, Bayesian approach

1 INTRODUCTION

The Bayesian approach deals with unknown parameters as random variables and assumes their *a priori* distributions. The essential role of the *a priori* distribution has not been well known, and the validity of the Bayesian approach and its application has been long argued [1]. The fact that the Bayesian approach enables us to design the optimum classifier based on limited sample and to evaluate the mean error rate using known parameters alone clearly demonstrates the valuableness of this approach.

This paper deals with the optimum classifier and the performance evaluation by the Bayesian approach. Gaussian population with unknown parameters is assumed. The conditional density given a limited sample of the population has a relationship to the multivariate t -distribution. As a result, the obtained optimum classifier is different from the quadratic classifier known to be optimum for Gaussian distributions with known parameters. Especially when the sample size of classes are not equal, the optimum discriminant function is not quadratic, and the decision surface is not hyperquadratics.

The mean error rate of the optimum classifier is theoretically evaluated by the quadrature of the conditional density. For univariate case, the mean error rate of two-class problem with different sample size and different sample covariance matrixes is evaluated. For multivariate case, the one with common sample size, common sample covariance matrixes, and common *a priori* probabilities is evaluated. Since these mean error rates are obtained by taking the expectation of the error rate over unknown population parameters dealt as random variables, they only depend on known parameters such as sample parameters, sample size, and the dimensionality. In this point, the Bayesian mean error rate has its own interpretation and significance different from those of non-Bayesian mean error rate which requires the unknown population parameters in its calculation. To verify the optimality of the classifier and the correctness of the mean error calculation, the results of Monte Carlo simulation employing a new sampling procedure are shown.

The optimum classifier based on the Bayesian approach was first derived by Keehn [2]. He studied the asymptotic properties of the optimum classifier and calculated type I error, which is the rejection rate for a given threshold value

of the likelihood. However the mean error rate for two-class problem was not evaluated, and the properties of the optimum classifier except for the asymptotic properties were not studied.

In subsequent sections, a case with unknown covariance matrix (with known mean vector) is described in Section 2 to 4. A new sampling procedure and the result of Monte Carlo simulation are described in Section 5.

2 SAMPLE CONDITIONAL DENSITY OF GAUSSIAN POPULATION

Sample conditional density of d -dimensional feature vector X of Gaussian population with unknown covariance matrix given a sample $\chi = \{X_1, X_2, \dots, X_n\}$ is expressed by

$$p(X|\chi) = \int_S p(X|K)p(K|\chi)dK, \quad (1)$$

where K is the inverse of the population covariance matrix and S is $d(d+1)/2$ dimensional subspace on which K is positive definite. Since the mean vector is known, it can be assumed to be zero vector without loss of generality. Then the density $p(X|K)$ is the d -variate Gaussian distribution, and the density $p(K|\chi)$ is the Wishart distribution of n_n degrees of freedom [2], [5].

Performing the integration (1), we have

$$p(X|\chi) = (n_n\pi)^{-\frac{d}{2}} |\Sigma_n|^{-\frac{1}{2}} \frac{\Gamma\left(\frac{n_n+1}{2}\right)}{\Gamma\left(\frac{n_n-d+1}{2}\right)} \left\{ 1 + \frac{1}{n_n}(X-M)'\Sigma_n^{-1}(X-M) \right\}^{-\frac{n_n+1}{2}}$$

$$\Sigma_n = \frac{n_0\Sigma_0 + n\Sigma}{n_0 + n}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n X_i X_i'$$

$$n_n = n_0 + n, \quad (2)$$

where Σ_0 and n_0 are an initial estimate of the population covariance matrix, and the confidence constant, respectively. When n_0 is set to zero, n_n and Σ_n coincide to n and Σ respectively, and no knowledge about the prior distribution is utilized. $\Gamma(x)$ is the gamma function.

By variable transformation

$$X - M = \sqrt{\frac{n_n}{n_n - d + 1}} T \quad (3)$$

T leads to the multivariate elliptical t -distribution with $n_n - d + 1$ degrees of freedom [6].

3 OPTIMUM DISCRIMINANT FUNCTION

The optimum discriminant function for general case is derived from (2) as

$$g(X) = -2\log\{p(X|\chi)P(\omega)\}$$

$$= (n_n + 1)\log\left\{ 1 + \frac{1}{n_n}(X-M)'\Sigma_n^{-1}(X-M) \right\} + \log|\Sigma_n| - 2\log D - 2\log P(\omega)$$

$$D = (n_n\pi)^{-\frac{d}{2}} \frac{\Gamma\left(\frac{n_n+1}{2}\right)}{\Gamma\left(\frac{n_n-d+1}{2}\right)}. \quad (4)$$

4 EVALUATION OF MEAN ERROR RATE

4.1 Univariate Case

4.1.1 Case with Common Sample Size

Since the discriminant function for a case of equal sample size is univariate quadratic, the decision boundary and the mean error rate are easily calculated. For simplicities sake, the *a priori* probabilities are assumed to be common to classes. The extension to unequal *a priori* probability case is straight forward.

The discriminant function

$$g_i(x) = \sigma_{ii}^{2/(n_n+1)} \left\{ 1 + \frac{1}{n_n} \left(\frac{x-m_i}{\sigma_{ii}} \right)^2 \right\} \quad (i = 1, 2) \quad (5)$$

is derived from (4). Setting $h(x) = 0$ for

$$h(x) = g_1(x) - g_2(x)$$

$$= (a-b)x^2 - 2(am_1 - bm_2)x + am_1^2 - bm_2^2 + c$$

$$a = \frac{1}{n_n} \sigma_{n1}^{2/(n_n+1)-2}, \quad b = \frac{1}{n_n} \sigma_{n2}^{2/(n_n+1)-2}, \quad c = \sigma_{n1}^{2/(n_n+1)-2} - \sigma_{n2}^{2/(n_n+1)-2} \quad (6)$$

the decision boundaries are determined as

$$\begin{aligned} \alpha = \beta &= \frac{m_1 + m_2}{2} & (\sigma_{n1} = \sigma_{n2}) \\ \alpha, \beta &= \frac{am_1 - bm_2 \mp \sqrt{(m_1 - m_2)^2 ab - (a - b)c}}{a - b} & (\sigma_{n1} \neq \sigma_{n2}). \end{aligned} \quad (7)$$

The mean error rate for $\sigma_{n1} \neq \sigma_{n2}$ is given by

$$\begin{aligned} \varepsilon &= P(\omega_1)\varepsilon_1 + P(\omega_2)\varepsilon_2 = P(\omega_1)P(\text{error}|\chi, \omega_1) + P(\omega_2)P(\text{error}|\chi, \omega_2) \\ &= B + A + C \\ A &= \int_{-\infty}^{\alpha} p(x|\chi, \omega_2)P(\omega_2)dx = \frac{1}{2}\Phi_{n_n}\left(\frac{\alpha - m_2}{\sigma_{n2}}\right) \\ B &= \int_{\alpha}^{\beta} p(x|\chi, \omega_1)P(\omega_1)dx = \frac{1}{2}\Phi_{n_n}\left(\frac{\beta - m_1}{\sigma_{n1}}\right) - \frac{1}{2}\Phi_{n_n}\left(\frac{\alpha - m_1}{\sigma_{n1}}\right) \\ C &= \int_{\beta}^{\infty} p(x|\chi, \omega_2)P(\omega_2)dx = \frac{1}{2}\left\{1 - \Phi_{n_n}\left(\frac{\beta - m_2}{\sigma_{n2}}\right)\right\} \end{aligned} \quad (8)$$

(Fig.1), where $\Phi_n(x_0)$ is defined by

$$\Phi_n(x_0) = \int_{-\infty}^{x_0} t_n(x)dx. \quad (9)$$

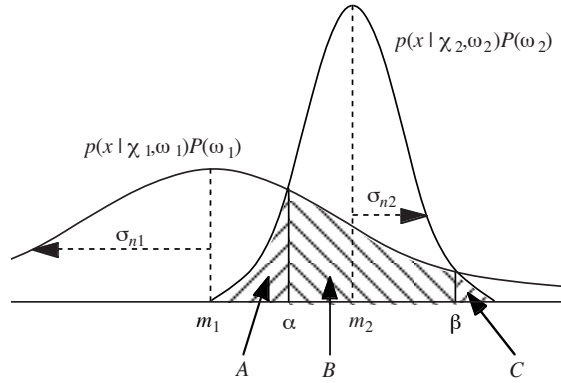


Fig.1. Calculation of mean error rate (univariate case).

The mean error rate obtained by the quadrature of the conditional density has an interpretation shown by

$$\begin{aligned} \int_{x_0}^{\infty} p(x|\chi)dx &= \int_{x_0}^{\infty} \int_0^{\infty} p(x|k)p(k|\chi)dkdx \\ &= \int_0^{\infty} \int_{x_0}^{\infty} p(x|k)dx p(k|\chi)dk = \int_0^{\infty} P(\text{error}|k)p(k|\chi)dk, \end{aligned} \quad (10)$$

where $p(x|\chi)$ is a sample conditional density of a class, x_0 is a decision boundary, and the region of the class is assumed to be $[-x_0, \infty)$. This equation implies that the mean error rate is the expectation of the error rate $p(\text{error}|k)$ given unknown population parameter k . Fig.2 shows the relationship between the error rate for individual population with variance k^{-1} and the expectation, when the sample variance is 1. When sample size is n and $n_0 = 0$, k is subject to the chi-squared distribution with n degrees of freedom, and the sample conditional density is the t -distribution with n degrees of freedom. The expectation of the hatched area of Gaussian distributions in upper row is equal to the hatched area of the t -distribution in lower row. This relation also holds for cases where there are multiple decision boundaries and the region of a class is separated into multiple segments.

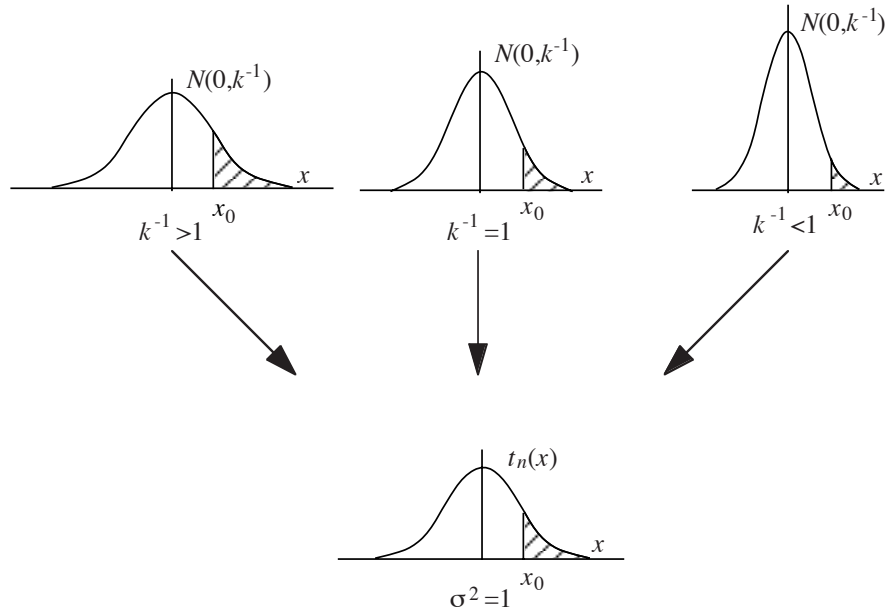


Fig.2. Relationship between the error rate for individual population with variance k^{-1} and the expectation.

4.1.2 Case with Different Sample Size

When the sample size of relevant classes are different, the optimum discriminant function is given by

$$g_i(x) = \left(\frac{\sigma_{ni}}{D_i} \right)^2 \left\{ 1 + \frac{1}{n_{ni}} \left(\frac{x - m_i}{\sigma_{ni}} \right)^2 \right\}^{n_{ni} + 1}$$

$$D_i = (n_{ni} \pi)^{-\frac{1}{2}} \frac{\Gamma\left(\frac{n_{ni} + 1}{2}\right)}{\Gamma\left(\frac{n_{ni}}{2}\right)} \quad (i = 1, 2). \quad (11)$$

The *a priori* probabilities are assumed to be common to all classes. The extension to unequal *a priori* probability case is straight forward.

Unfortunately the equation $h(x) = 0$ can not be solved analitically. In the experiment in Section 5, the decision boundary is calculated employing Newton's iterative formula. The mean error rate is calculated in the similar way as (8).

4.2 Multivariate Case

The sample size, the covariance matrixes and the *a priori* probabilities are assumed to be common to two classes. The logarithm of the likelihood ratio is given by

$$h(X) = (M_2 - M_1)' \Sigma_n^{-1} \sqrt{\frac{n_n - d + 1}{n_n}} X$$

$$+ \frac{1}{2} \sqrt{\frac{n_n - d + 1}{n_n}} (M_1' \Sigma_n^{-1} M_1 - M_2' \Sigma_n^{-1} M_2) \quad (12)$$

The distribution of $((n_n - d + 1)/n_n)^{1/2} X$ is d -variate elliptical t -distribution with $n_n - d + 1$ degrees of freedom, and the distribution of $h(X)$ is univariate t -distribution with the same degrees of freedom. The means of $h(X)$ are given by

$$\eta_1 = -\frac{1}{2} \sqrt{\frac{n_n - d + 1}{n_n}} \delta_n^2$$

$$\eta_2 = \frac{1}{2} \sqrt{\frac{n_n - d + 1}{n_n}} \delta_n^2$$

$$\delta_n^2 = (M_2 - M_1)' \Sigma_n^{-1} (M_2 - M_1). \quad (13)$$

The variances of $h(X)$ is given by

$$\sigma_1^2 = (M_2 - M_1)' \Sigma_n^{-1} E \left\{ \frac{n_n - d + 1}{n_n} (X - M_1)(X - M_1)' | \omega_1 \right\} \Sigma_n^{-1} (M_2 - M_1)$$

$$= \frac{n_n - d + 1}{n_n - d - 1} (M_2 - M_1)' \Sigma_n^{-1} (M_2 - M_1) = \frac{n_n - d + 1}{n_n - d - 1} \delta_n^2$$

$$\sigma_2^2 = \frac{n_n - d + 1}{n_n - d - 1} \delta_n^2. \quad (14)$$

Using these parameters the mean error rate is given by

$$\varepsilon = P(\omega_1) \varepsilon_1 + P(\omega_2) \varepsilon_2$$

$$\begin{aligned}
&= \frac{1}{2} \Phi_{n_n-d+1} \left(\frac{\eta_1}{\delta_n} \right) + \frac{1}{2} \left\{ 1 - \Phi_{n_n-d+1} \left(\frac{\eta_2}{\delta_n} \right) \right\} \\
&= 1 - \Phi_{n_n-d+1} \left(\frac{1}{2} \sqrt{\left(1 - \frac{d-1}{n} \right) \delta^2} \right).
\end{aligned} \tag{15}$$

When $n_0 = 0$,

$$\varepsilon = 1 - \Phi_{n-d+1} \left(\frac{1}{2} \sqrt{\left(1 - \frac{d-1}{n} \right) \delta^2} \right). \tag{16}$$

The Bayesian formulas of the mean error rate (15) and (16) have the following characteristics when compared with the non-Bayesian formulas.

- 1) The unknown population parameters are not required in its calculation.
- 2) The expression is simple and clearly shows the limited sample effect on the mean error rate.
- 3) The relationship between the prior parameters n_0 , Σ_0 and the mean error rate is explicitly expressed.

It should be noted that the Mahalanobis distance δ in (16) is an apparent one which is calculated using the known population mean vector and the sample covariance matrix. (16) reveals two causes which increase the mean error rate due to the limited sample effect. One is that the area of the tail of t -distribution increases due to the reduction of the degrees of freedom. The other is that the apparent squared Mahalanobis distance between two classes shrinks by $(d-1)/n$, and increases the mean error rate (Fig.3). The affection of the former is marginal and is negligible if $n-d+1$ is greater than 20 or so, because the t -distribution with this degrees of freedom can be approximated by the Gaussian distribution, which is the t -distribution with infinite degrees of freedom. On the other hand, the affection of the latter is so severe and is not negligible unless the sample size is much larger than the dimensionality. Such shrinkage of the apparent Mahalanobis distance has its origin in the variable transformation by (3), and causes a problem so called "peaking phenomenon" or "curse of dimensionality"[3], [4], [7]. This undesirable phenomenon is caused and aggravated by neglecting the prior distribution by setting $n_0 = 0$. The case for $n_0 = 0$ is discussed in Section 6.

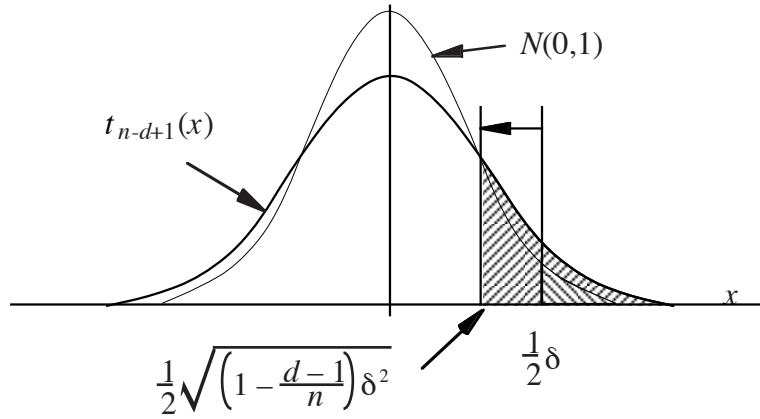


Fig.3. Increase of mean error rate due to limited sample effect.

5 COMPUTER SIMULATION

5.1 Bayesian Sampling

In the following computer simulation, a new sampling procedure called Bayesian sampling is employed together with the ordinary sampling procedure. Fig.4 illustrates the relationship between the ordinary sampling (a) and the Bayesian sampling (b). In the ordinary sampling, specified size of sample are drawn from a specified population and the sample parameters are calculated. Fig.4 (a) illustrates the case with a Gaussian population $N(0, I)$ and three samples of size five with the sample covariance matrixes Σ_a, Σ_b , and Σ_c . The classifiers are designed using these sample parameters and the mean error rate for the population is evaluated. Since the sample parameters are random variables, the expectation of the error rate is taken by repeating the sampling and the design and test of the classifier. On the contrary the Bayesian sampling generates populations from which a sample with specified parameter, e.g. $N(0, I)$, is extracted. When a sample of specified size is drawn from a temporal population $N(0, I)$, and the sample covariance matrix is Σ_a , the actual population is determined to be $N(0, \Sigma_a^{-1})$. Since the population parameters are random variables in this case, the expectation of the error rate is taken by repeating the Bayesian sampling and the test of the classifier. The design of the classifier need not be repeated because the design sample is fixed through the experiment. In this example, the sample mean vector and the sample covariance matrix are assumed to be zero vector and identity matrix, respectively.

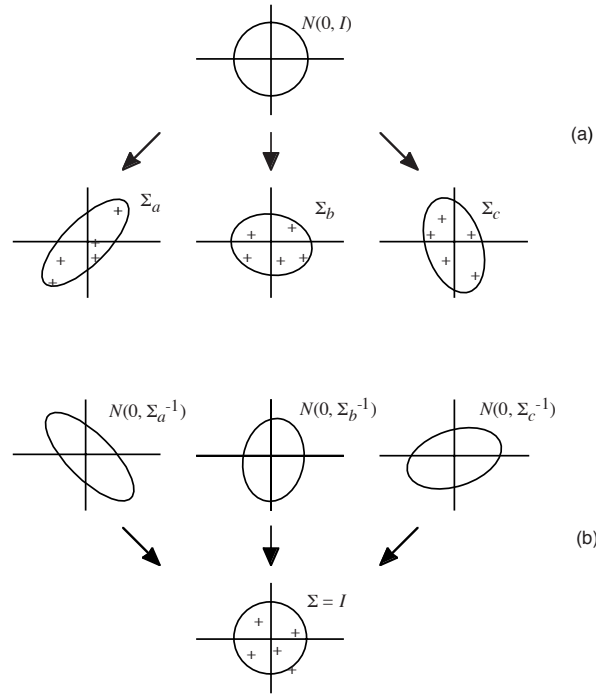


Fig.4. Relationship between ordinary sampling (a), and Bayesian sampling (b).

The general procedure is described below. The population parameters are determined so that the parameters of a sample drawn from the population is (μ_2, Σ_2) . The parameters of a sample of size n drawn from a temporal population $N(0, I)$ are denoted by (μ_1, Σ_1) , i.e.

$$\begin{aligned}\mu_1 &= \frac{1}{n} \sum_{i=1}^n X_i \\ \Sigma_1 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu_1)(X_i - \mu_1)'.\end{aligned}\quad (17)$$

By setting

$$Y = \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' (X - \mu_1) \quad (\Sigma_1 \Phi_1 = \Phi_1 \Lambda_1) \quad (18)$$

the sample parameters are transformed to $(0, I)$, i.e.

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n Y_i &= 0 \\ \frac{1}{n-1} \sum_{i=1}^n Y_i Y_i' &= \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' \Sigma_1 \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' = \Phi_1 \Phi_1' = I\end{aligned}\quad (19)$$

and the population parameters of Y are given by

$$\begin{aligned}E(Y) &= -\Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' \mu_1 \\ V(Y) &= E\{[Y - E(Y)][Y - E(Y)]'\} = \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' = \Sigma_1^{-1},\end{aligned}\quad (20)$$

where Λ_1 and Φ_1 are the eigenvalue matrix and eigenvector matrix of Σ_1 , respectively.

Further by setting

$$\begin{aligned}Z &= \Phi_2 \Lambda_2^{\frac{1}{2}} Y + \mu_2 \\ &= \Phi_2 \Lambda_2^{\frac{1}{2}} \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' X + \mu_2 - \Phi_2 \Lambda_2^{\frac{1}{2}} \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' \mu_1 \\ &\quad (\Sigma_2 \Phi_2 = \Phi_2 \Lambda_2)\end{aligned}\quad (21)$$

the sample parameters are transformed to (μ_2, Σ_2) , i.e.

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n Z_i &= \Phi_2 \Lambda_2^{\frac{1}{2}} \frac{1}{n} \sum_{i=1}^n Y_i + \mu_2 = \mu_2 \\ \frac{1}{n-1} \sum_{i=1}^n Y_i Y_i' &= \Phi_2 \Lambda_2^{\frac{1}{2}} I \Lambda_2^{\frac{1}{2}} \Phi_2' = \Phi_2 \Lambda_2 \Phi_2' = \Sigma_2\end{aligned}\quad (22)$$

and the population parameters of Z are given by

$$\begin{aligned}E(Z) &= \mu_2 - \Phi_2 \Lambda_2^{\frac{1}{2}} \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' \mu_1 \\ V(Z) &= \Phi_2 \Lambda_2^{\frac{1}{2}} \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' \Phi_1 \Lambda_1^{-\frac{1}{2}} \Phi_1' \Lambda_2^{\frac{1}{2}} \Phi_2' = \Phi_2 \Lambda_2^{\frac{1}{2}} \Sigma_1^{-1} \Lambda_2^{\frac{1}{2}} \Phi_2'.\end{aligned}\quad (23)$$

When the population mean vector M is known, (23) is replaced by

$$\begin{aligned}E(Z) &= M \\ V(Z) &= \Phi_2 \Lambda_2^{\frac{1}{2}} \Sigma_1^{-1} \Lambda_2^{\frac{1}{2}} \Phi_2' \\ \Sigma_1 &= \frac{1}{n} \sum_{i=1}^n X_i X_i'\end{aligned}\quad (24)$$

In the following experiments, n_0 is set to zero and the population is assumed to have known mean vector and unknown covariance matrix.

5.1.1 Univariate Case

TABLE I shows the mean error rate for a case of equal variances ($\sigma_1^2 = \sigma_2^2 = 1.0$). The *a priori* probabilities are common to two classes, and the two means are -1.0 and 1.0, respectively. In the table n_1 and n_2 denote the sample size of each class, and the columns *opt.* and *qdf.* show the mean error rate of the optimum classifier and the quadratic classifier, respectively. The rows *sim.* show the result by Monte Carlo simulation, and the rows *t* by quadrature of *t*-distribution. In the Monte Carlo simulation, the Bayesian sampling was employed and 2000 tests each of which used a test sample of size 1000 were repeated to calculate the mean error rate. When the sample variance and the sample size are common to classes, the optimum classifier and the quadratic classifier give the same results. For the rest of the case, the optimum classifier always outperforms the quadratic classifier. The mean error rates calculated by the simulation and the theoretical prediction are accurately coincident mutually. Fig.5 shows the relationship between the mean error rate and the sample size of class 1 when total sample size is fixed to 4 in TABLE I. When the sample variances are common to two classes, the mean error rate of the quadratic classifier is minimized for the case of equal sample size, while the one of the optimum classifier is maximized for the case. The similar relationship for a case of different variances ($\sigma_1^2 = 4.0, \sigma_2^2 = 0.25$) is shown in Fig.6. In this case, the mean error rates of the both classifiers are lower for the case where the class with less sample variance has greater sample size.

TABLE I

Mean error rate (%) v.s. sample size in univariate two-class problem with common sample variances ($\sigma_1^2 = \sigma_2^2 = 1.0$)

| $n_1 \backslash n_2$ | | 1 | | 2 | | 3 | |
|----------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | <i>opt.</i> | <i>qdf.</i> | <i>opt.</i> | <i>qdf.</i> | <i>opt.</i> | <i>qdf.</i> |
| 1 | <i>sim.</i> | 24.98 | 24.98 | 21.62 | 23.06 | 20.14 | 22.26 |
| | <i>t</i> | 25.00 | 25.00 | 21.67 | 23.07 | 20.19 | 22.27 |
| 2 | <i>sim.</i> | 21.77 | 23.08 | 21.16 | 21.16 | 20.13 | 20.37 |
| | <i>t</i> | 21.67 | 23.07 | 21.13 | 21.13 | 20.17 | 20.34 |
| 3 | <i>sim.</i> | 20.21 | 22.27 | 20.19 | 20.35 | 19.56 | 19.56 |
| | <i>t</i> | 20.19 | 22.27 | 20.17 | 20.34 | 19.55 | 19.55 |

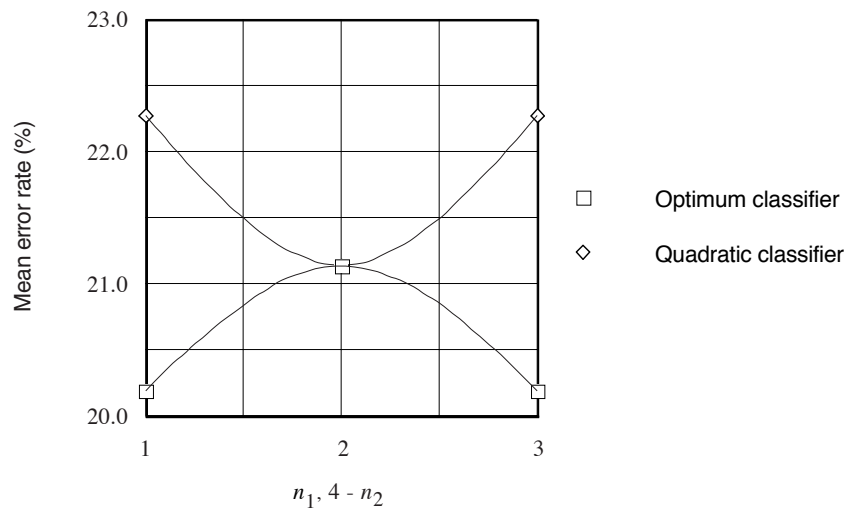


Fig.5. Theoretical mean error rate (%) vs. sample size with fixed total sample size ($n_1 + n_2 = 4, \sigma_{12} = \sigma_{22} = 1.0$).

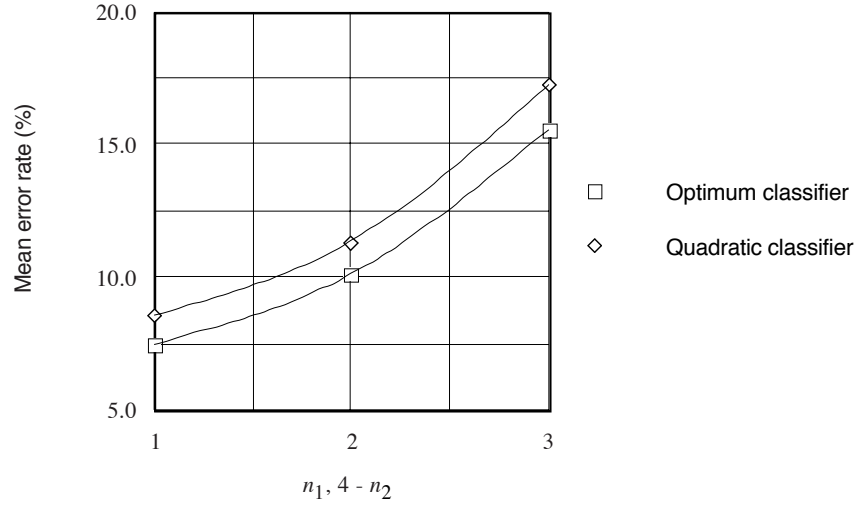


Fig.6. Theoretical mean error rate (%) vs. sample size with fixed total sample size ($n_1+n_2=4$, $\sigma_{12}=4.0$, $\sigma_{22}=0.25$).

5.1.2 Multivariate Case with Common Sample Covariance Matrix

TABLE II and Fig.7 show the results of experiments for multivariate case where the sample size, the sample covariance matrixes, and the *a priori* probabilities are all common to two classes. The rows *sim.* are the results by the Monte Carlo simulation employing the Bayesian sampling, where the size of test sample is 1000, and the number of iteration is 5000. The row *t* shows the mean error rate calculated as described in 4.2.

The optimum discriminant function employed in the simulation is derived from (4). The sample covariance matrix is $d \times d$ identity matrix, and the population mean vectors are

$$M_1 = (0, 0, 0, \dots, 0),$$

$$M_2 = (1, 1, 1, \dots, 1).$$

For these parameters, the Mahalanobis distance $\delta^2 = n$ and (16) is minimized when $d = (n + 1)/2$.

Because the sample size and the sample covariance matrixes are common to classes, the optimum classifier and the quadratic classifier give the same results. The mean error rates predicted by the *t*-quadrature is well coincident to those by Monte Carlo simulation.

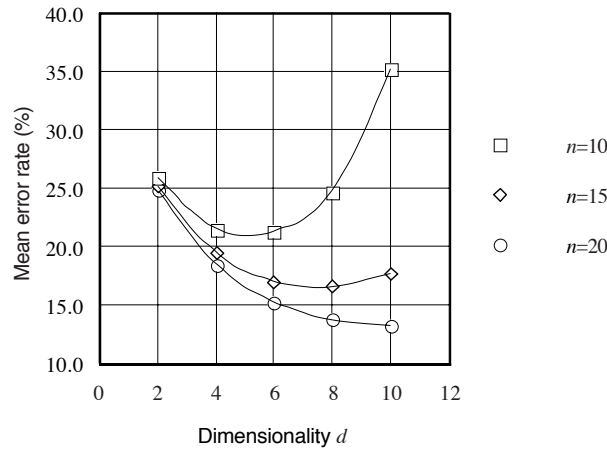


Fig.7. Theoretical mean error rate (%) vs. dimensionality.

TABLE II

Mean error rate (%) v.s. dimensionality in multivariate two-class problem with common sample covariance matrixes

| $d \backslash n$ | | 10 | 15 | 20 |
|------------------|-------------|--------------------|--------------------|--------------------|
| | | <i>opt. (qdf.)</i> | <i>opt. (qdf.)</i> | <i>opt. (qdf.)</i> |
| 2 | <i>sim.</i> | 25.97 | 25.29 | 24.95 |
| | <i>t</i> | 25.96 | 25.28 | 24.95 |
| 4 | <i>sim.</i> | 21.49 | 19.44 | 18.45 |
| | <i>t</i> | 21.52 | 19.43 | 18.47 |
| 6 | <i>sim.</i> | 21.26 | 16.94 | 15.17 |
| | <i>t</i> | 21.30 | 17.04 | 15.28 |
| 8 | <i>sim.</i> | 24.65 | 16.49 | 13.58 |
| | <i>t</i> | 24.75 | 16.59 | 13.74 |
| 10 | <i>sim.</i> | 35.32 | 17.65 | 13.22 |
| | <i>t</i> | 35.24 | 17.80 | 13.29 |

5.1.3 Multivariate Case with Different Sample Covariance

Fig.8 shows the mean error rates of the optimum classifier and the quadratic classifier for two classes with different sample covariance matrixes. The mean error rates were evaluated by Monte Carlo simulation employing the Bayesian sampling, where the size of test sample and the number of iteration are 5000. The size of design sample and the *a priori* probabilities are common to the classes. The sample covariance matrix of class1 is 8 x 8 identity matrix, and the one of class2 is 8 x 8 diagonal matrix with diagonal elements

$$\text{diag}\Sigma_2 = (8.41, 12.06, 0.12, 0.22, 1.49, 1.77, 0.35, 2.73)$$

The mean vectors are given by

$$M_1 = (-1, 0, 0, \dots, 0),$$

$$M_2 = -M_1$$

The mean error rates of the quadratic classifier approach to those of the optimum classifier as the sample size n increases, however the optimum classifier outperforms the quadratic classifier for all sample size.

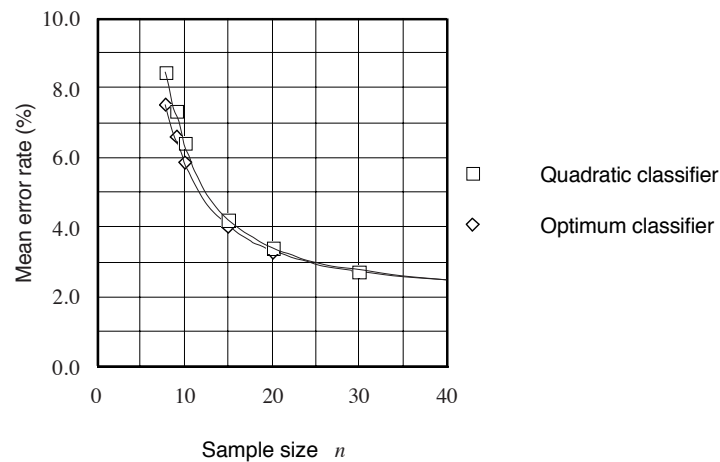


Fig.8. Mean error rate of quadratic classifier and optimum classifier vs. sample size in 8-variate two-class problem with individual sample covariance matrixes.

6 CONCLUSION and DISCUSSION

This paper dealt with the optimum classifier design and the performance evaluation by the Bayesian approach. To verify the optimality of the classifier and the correctness of the mean error calculation, the results of Monte Carlo simulation employing the Bayesian sampling were shown. It was also shown by the comparative study that the Bayesian formulas of the mean error rate have the following characteristics.

- 1) The unknown population parameters are not required in its calculation.
- 2) The expression is simple and clearly shows the limited sample effect on the mean error rate.
- 3) The relationship between the prior parameters and the mean error rate is explicitly expressed.

In the Monte Carlo simulation, the property of the optimum classifier was studied when n_0 was set to zero and the prior distribution was completely neglected. When n_0 is not zero, the mean error rate is expressed by (15) and is further minimized by selecting optimum n_0 which maximizes

$$f(n_0) = \left(1 - \frac{d-1}{n+n_0}\right) (M_2 - M_1)' \left\{ \frac{n}{n+n_0} \Sigma + \frac{n_0}{n+n_0} \Sigma_0 \right\}^{-1} (M_2 - M_1). \quad (25)$$

The increase of n_0 has similar effect as the increase of the sample size to add the degrees of freedom of the t -distribution, and to reduce the shrinkage of the apparent Mahalanobis distance. Therefore complete ignorance of the prior distribution by setting n_0 to zero does not lead the best possible classifier.

In most of the real world application, given sample parameters are fixed and the population parameters are unknown. The Bayesian sampling agrees better with these realities than non-Bayesian sampling, and provides us a new way of the Monte Carlo simulation such as the analysis of multi-category classification problems beginning with real world sample parameters at hand.

REFERENCES

- [1] R. O. Duda and P. E. Hart, *Pattern Classification and Scene Analysis*. New York: Wiley, 1973, p.76, p.68.
- [2] D. G. Keehn, "A Note on Learning for Gaussian Properties," *IEEE Trans. Inform. Theory*, vol. IT-11, no. 1, pp.126-132, Jan 1965.
- [3] S. Raudys and V. Pikelis, "On Dimensionality, Sample Size, Classification Error, and Complexity of Classification Algorithm in Pattern Recognition," *IEEE Trans. PAMI*, vol. PAMI-2, no.3, pp.242-252, May 1980.
- [4] K. Fukunaga and R. R. Hayes, "Effects of Sample Size in Classifier Design," *IEEE Trans. PAMI*, vol.11, no.8, pp.873-885, Aug 1989.
- [5] K. Fukunaga, *Introduction to Statistical Pattern Recognition, Second Edition*, New York: Academic Press, 1990, pp.392-393, pp.91-92.
- [6] R. J. Muirhead, *Aspects of Multivariate Statistical Theory*. New York: Wiley, 1982, p.32-49.
- [7] S. J. Raudys and A.K.Jain, "Small Sample Size Effects in Statistical Pattern Recognition: Recommendations for Practitioners," *IEEE Trans. PAMI*, vol. 13, no. 3, pp.252-264, Mar 1991.