

## *Discrimination for Variance Matrices*

By Masashi OKAMOTO

### CONTENTS

	<i>Page</i>
0. Introduction ... ..	1
PART I: COMPLETELY SPECIFIED POPULATIONS ... ..	2
1. Problem of discrimination ... ..	2
2. Bayes discrimination for variance matrices ...	4
3. Minimax discrimination for variance matrices ...	7
4. Reduction in dimensions ... ..	9
5. Weighted sum of $\chi^2$ variates ... ..	10
PART II: INCOMPLETELY SPECIFIED POPULATIONS ... ..	18
6. Discrimination when parameters are estimated	18
7. Reduction in dimensions ... ..	22
8. Asymptotic distribution of the eigenvalues and the eigenvectors ... ..	22
9. Asymptotic distribution of the quadratic discrimi- nant function ... ..	28
PART III: AN APPLICATION ... ..	31
10. Discrimination of zygoty of twins ... ..	31
11. Analysis of the data of twins ... ..	34

**0. Introduction.** The problem of statistical discrimination has been hitherto investigated with respect to mean vectors of several multi-dimensional normal populations with a common variance matrix by Fisher [7], [8], [9], Wald [31], Rao [23], [24], Anderson [3], [4] and others. We shall consider in this paper the problem with respect to variance matrices of two multi-dimensional normal populations with a common mean vector. The reason why this problem has not been taken up before seems to be the complexity of its theory on the one hand and the scantiness of its application to practical sciences on the other, compared with that for mean vectors. But the theory can be developed to some extent and there has been found at least one interesting application in the field of biometry.

The paper is divided into three parts. Part I is concerned with the case when the populations are completely specified, or the common mean

vector together with two variance matrices are all known. Part II deals with the case when the populations are incompletely specified, while the information for unknown parameters is provided by a random sample taken from each population. Part III finally illustrates the theory by a practical example.

## PART I. COMPLETELY SPECIFIED POPULATIONS

**1. Problem of discrimination.** Let  $\Pi_i$  ( $i=1, 2$ ) be two populations in a space  $\mathcal{X}$  and let  $P_i$  be their probability functions. We do not in this section set up any assumption for  $P_i$  such as normality. If an observation is performed on either  $\Pi_1$  or  $\Pi_2$  and the result  $x$  belonging to  $\mathcal{X}$  is informed to us, we are confronted with the problem of discrimination. We may take either the decision  $d_1$ , judging that  $x$  came from the population  $\Pi_1$ , or the alternative decision  $d_2$  in favor of the other possibility.

Let  $w_i$  ( $i=1, 2$ ) be the loss incurred in adopting  $d_{3-i}$  when  $x$  comes in fact from  $\Pi_i$ . Let  $\varphi$ ,  $0 \leq \varphi(x) \leq 1$  for any  $x$ , be a randomized decision function to the effect that when  $x$  is observed we adopt the decision  $d_2$  with the assigned probability  $\varphi(x)$ . We call  $\varphi$  a *discrimination function*, distinguishing it from the term discriminant function originated by R. A. Fisher. The error that we adopt  $d_2$  when  $\Pi_1$  is true or the error that we adopt  $d_1$  when  $\Pi_2$  is true has the probability

$$(1.1) \quad P(2|1, \varphi) = \int \varphi(x)P_1(dx) \quad \text{or} \quad P(1|2, \varphi) = \int (1-\varphi(x))P_2(dx),$$

respectively. And the expected loss or risk when  $\Pi_i$  is true is given by

$$(1.2) \quad R_\varphi(P_i) = w_i P(3-i|i, \varphi).$$

If by some random device  $\Pi_i$  ( $i=1, 2$ ) is chosen with probability  $\pi_i$  ( $\pi_1 + \pi_2 = 1$ ) to give rise to an observed value  $x$ , then the average risk is

$$(1.3) \quad R_\varphi = \sum_{i=1}^2 \pi_i R_\varphi(P_i).$$

A discrimination function  $\varphi$  minimizing this expression is called a *Bayes discrimination (function)* with respect to  $(\pi_1, \pi_2)$  and is often denoted later by  $\varphi_B$ . Similarly,  $\varphi$  minimizing the maximum risk

$$(1.4) \quad R_\varphi^* = \max (R_\varphi(P_1), R_\varphi(P_2))$$

is called a *minimax discrimination (function)* and is denoted by  $\varphi_M$ .

We assume now that the probability functions  $P_i$  ( $i=1, 2$ ) have

density functions  $p_i(x)$  with respect to a certain measure  $\nu$  in  $\mathcal{X}$ . Welch [33], Rao [24] and Anderson [4] states a theorem on non-randomized Bayes discrimination which is easily adapted to the randomized case such as

**Theorem 1.** *A necessary and sufficient condition that a discrimination function  $\varphi$  is Bayes with respect to  $(\pi_1, \pi_2)$  is that  $\varphi$  satisfies*

$$(1.5) \quad \varphi(x) = \begin{cases} 1 & \text{if } \pi_1 w_1 p_1(x) < \pi_2 w_2 p_2(x) \\ 0 & \text{if } \pi_1 w_1 p_1(x) > \pi_2 w_2 p_2(x). \end{cases}$$

*There exists at least one Bayes discrimination.*

The theorem is obtained from the equations

$$\begin{aligned} R_\varphi &= \pi_1 w_1 \int \varphi(x) P_1(dx) + \pi_2 w_2 \int (1 - \varphi(x)) P_2(dx) \\ &= \pi_2 w_2 + \int \varphi(x) [\pi_1 w_1 p_1(x) - \pi_2 w_2 p_2(x)] \nu(dx). \end{aligned}$$

There exist many Bayes discrimination functions so long as the set of  $x$  determined by  $\pi_1 w_1 p_1(x) = \pi_2 w_2 p_2(x)$  has positive  $\nu$  measure but of course the average risk  $R_\varphi$  for each of them coincides.

Similarly, corresponding to the theorem of Mises [18] and Anderson [4] on the non-randomized minimax discrimination, we get under less assumptions

**Theorem 2.** *A necessary and sufficient condition that a discrimination function  $\varphi$  is minimax is that  $\varphi$  satisfies*

$$(1.6) \quad R_\varphi(P_1) = R_\varphi(P_2)$$

*as well as (1.5) for some  $(\pi_1, \pi_2)$ . There exists at least one minimax discrimination function.*

The theorem is implied in the following Lemmas 1 and 2.

**Lemma 1.** *Two properties below are equivalent:*

- (i)  $\varphi$  is minimax.
- (ii)  $\varphi$  minimizes  $R_\varphi(P_1)$  under the restriction (1.6).

*Proof.* It suffices to show that any minimax discrimination function satisfies the equation (1.6). Now suppose that for some  $\varphi^*$  we have  $R_{\varphi^*}(P_1) < R_{\varphi^*}(P_2)$ . For  $\varepsilon$ ,  $0 < \varepsilon < 1$ , put  $\varphi(x) = \varepsilon + (1 - \varepsilon)\varphi^*(x)$ , then

$$R_\varphi(P_1) = \varepsilon w_1 + (1 - \varepsilon)R_{\varphi^*}(P_1), \quad R_\varphi(P_2) = (1 - \varepsilon)R_{\varphi^*}(P_2).$$

Choosing  $\varepsilon$  such that  $R_\varphi(P_1) = R_\varphi(P_2)$ , we know that

$$\max_{i=1,2} R_\varphi(P_i) < \max_{i=1,2} R_{\varphi^*}(P_i),$$

which implies that  $\varphi^*$  is not minimax.

**Lemma 2.** *The property (ii) in Lemma 1 is equivalent to (iii)  $\varphi$  is Bayes for some  $(\pi_1, \pi_2)$  and satisfies (1.6).*

Proof. The equation (1.6) is written in the form

$$\int (1 - \varphi(x)) [w_1 p_1(x) + w_2 p_2(x)] \nu(dx) = w_1$$

and the condition  $R_\varphi(P_1) = \min.$  is equivalent to

$$\int (1 - \varphi(x)) w_1 p_1(x) \nu(dx) = \max.$$

Thus the problem is quite analogous to that of testing hypothesis; hence Neyman-Pearson's fundamental lemma gives the solution

$$1 - \varphi(x) = \begin{cases} 1, & \text{if } c_1 w_1 p_1(x) > c_2 [w_1 p_1(x) + w_2 p_2(x)], \\ 0, & \text{if } c_1 w_1 p_1(x) < c_2 [w_1 p_1(x) + w_2 p_2(x)] \end{cases}$$

for some  $c_1$  and  $c_2$ . This is equivalent to

$$(1.7) \quad \varphi(x) = \begin{cases} 1, & \text{if } c_1^* w_1 p_1(x) < c_2^* w_2 p_2(x), \\ 0, & \text{if } c_1^* w_1 p_1(x) > c_2^* w_2 p_2(x) \end{cases}$$

for some  $c_1^*$  and  $c_2^*$ . Since  $c_1^*$  and  $c_2^*$  are readily seen to be non-negative, they determine a prior probability  $(\pi_1, \pi_2)$  rendering (1.7) identical with (1.5).

Existence of a minimax discrimination function results from the condition (iii).

**2. Bayes discrimination for variance matrices.** From this section through the last let  $\Pi_i$  ( $i=1, 2$ ) denote a  $p$ -dimensional normal population  $N(\mu, \Sigma_i)$  with a common mean vector  $\mu$  and a variance matrix  $\Sigma_i$ . Then the space  $\mathcal{X}$  of observation is an Euclidean  $p$ -space. Furthermore throughout Part I we assume that  $\mu$  and  $\Sigma_i$ 's are known completely. The  $i$ th density function (with respect to Lebesgue measure) is

$$p_i(\mathbf{x}) = (2\pi)^{-p/2} |\Sigma_i|^{-1/2} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu)' \Sigma_i^{-1} (\mathbf{x} - \mu) \right],$$

a prime superfix to any vector or any matrix denoting always the transpose. From Theorem 1 in the preceding section any Bayes discrimination function is given by

$$(2.1) \quad \varphi_B(\mathbf{x}) = \begin{cases} 1 & \text{if } Q > k_B \\ 0 & \text{if } Q < k_B, \end{cases}$$

where

$$(2.2) \quad Q = (\mathbf{x} - \boldsymbol{\mu})'(\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})(\mathbf{x} - \boldsymbol{\mu})$$

and

$$(2.3) \quad k_B = 2 \log \frac{\pi_1 w_1}{\pi_2 w_2} + \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|}.$$

The function  $Q$  corresponds to the linear discriminant function appearing in the discrimination for means and hence will be called the *quadratic discriminant function*. Since the set of  $x$  satisfying  $Q = k_B$  has Lebesgue measure zero, the Bayes discrimination is determined uniquely with probability one.

There exists a non-singular matrix  $\mathbf{F}$  as well as a diagonal  $\mathbf{A}$  with the diagonal elements in descending order in magnitude such that

$$(2.4) \quad \mathbf{F}'\boldsymbol{\Sigma}_1\mathbf{F} = \mathbf{I}, \quad \mathbf{F}'\boldsymbol{\Sigma}_2\mathbf{F} = \mathbf{A},$$

where  $\mathbf{I}$  denotes the identity matrix (cf. for example Roy [28]). The diagonal elements  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  of  $\mathbf{A}$  are the roots of the determinantal equation

$$(2.5) \quad |\boldsymbol{\Sigma}_2 - \lambda \boldsymbol{\Sigma}_1| = 0,$$

or in other words the eigenvalues with respect to the pair  $(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2)$ . They are all positive since both  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$  are positive definite. Two equations in (2.4) together imply that  $\boldsymbol{\Sigma}_2\mathbf{F} = \boldsymbol{\Sigma}_1\mathbf{F}\mathbf{A}$ , and therefore  $\boldsymbol{\Sigma}_2\mathbf{f}_i = \lambda_i\boldsymbol{\Sigma}_1\mathbf{f}_i$  for the  $i$ th column  $\mathbf{f}_i$  of  $\mathbf{F}$  ( $i=1, \dots, p$ ). This means that  $\mathbf{f}_i$  is an eigenvector with respect to the eigenvalue  $\lambda_i$  and the pair  $(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2)$ . If all the  $\lambda_i$  are distinct, then  $\mathbf{F}$  is determined uniquely except for the sign of every column. Uniqueness will be required later but not for the present.

Define

$$(2.6) \quad \mathbf{y} = \mathbf{F}'(\mathbf{x} - \boldsymbol{\mu}),$$

then we have

$$(2.7) \quad Q = \mathbf{y}'(\mathbf{I} - \mathbf{A}^{-1})\mathbf{y} = Q(\mathbf{y}) \quad (\text{say}).$$

If the observation  $\mathbf{x}$  comes from the population  $\Pi_1$  or from  $\Pi_2$ , then  $\mathbf{y}$  may be regarded as coming from the normal population  $N(\mathbf{0}, \mathbf{I})$  or  $N(\mathbf{0}, \mathbf{A})$ , respectively, to which we refer as  $P_1^Y$  or  $P_2^Y$ . The latter is also referred to as  $P_\lambda^Y$  indicating explicitly the dependence on  $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_p)$ . Thus the distribution of  $Q$  depends only on  $\boldsymbol{\lambda}$ , whether  $\mathbf{x}$  comes from  $\Pi_1$  or from  $\Pi_2$ . This is the canonical reduction used frequently in the multivariate analysis and we call  $\mathbf{y}$  a *canonical variate*. We could but

did not indeed start from the reduced populations, intending to make a correspondence with the discussion in Part II where the  $\Sigma_i$  are unknown and so the reduction is not permitted. By the way the reason why the discrimination for variance matrices is more complicated than that for mean vectors is that the canonical parameter is  $p$ -dimensional  $\boldsymbol{\lambda}$  for the former, while it is one-dimensional Mahalanobis' distance  $D^2$  for the latter. Hence results also the difficulty of dealing with more than two normal populations in this paper.

In terms of the canonical variate  $\mathbf{y}$  the Bayes discrimination is performed as follows: we adopt the decision  $d_1$  or  $d_2$  according as  $\mathbf{y}$  belongs to the set

$$(2.8) \quad D_1(\boldsymbol{\lambda}) = \{\mathbf{y}; Q(\mathbf{y}) < k_B(\boldsymbol{\lambda})\} \quad \text{or} \quad D_2(\boldsymbol{\lambda}) = \{\mathbf{y}; Q(\mathbf{y}) > k_B(\boldsymbol{\lambda})\},$$

where

$$(2.9) \quad k_B(\boldsymbol{\lambda}) = 2 \log \frac{\pi_1 w_1 + \sum_{i=1}^p \log \lambda_i}{\pi_2 w_2}.$$

We may take either one of  $d_1$  and  $d_2$  whenever  $Q(\mathbf{y}) = k_B(\boldsymbol{\lambda})$ . Thus the discrimination is essentially determined by the pair  $(D_1(\boldsymbol{\lambda}), D_2(\boldsymbol{\lambda}))$  and the probabilities of error of two kinds are given by

$$(2.10) \quad \begin{aligned} P(2|1, \varphi_B) &= P_1^Y(D_2(\boldsymbol{\lambda})) = Pr \left( \sum_{i=1}^p \left(1 - \frac{1}{\lambda_i}\right) Z_i^2 > k_B(\boldsymbol{\lambda}) \right) \\ P(1|2, \varphi_B) &= P_2^Y(D_1(\boldsymbol{\lambda})) = Pr \left( \sum_{i=1}^p (\lambda_i - 1) Z_i^2 < k_B(\boldsymbol{\lambda}) \right), \end{aligned}$$

where  $Z_i$  ( $i=1, 2, \dots, p$ ) are independent  $N(0, 1)$  variates.

We shall now investigate the behavior of the average risk  $R_{\varphi_B}$  given by (1.3) for the Bayes discrimination  $\varphi_B$  when  $\lambda_i$  varies with the  $\pi_i$  and the  $w_i$  being fixed. We write  $R_B(\boldsymbol{\lambda})$  for  $R_{\varphi_B}$ .

**Theorem 3.** For each  $i$  ( $i=1, 2, \dots, p$ )  $R_B(\boldsymbol{\lambda})$  is strictly monotonically

(i) increasing in  $\lambda_i$  if  $0 < \lambda_i \leq 1$ ,

and

(ii) decreasing in  $\lambda_i$  if  $1 \leq \lambda_i$ .

Both this and Theorem 4 in Section 3 are based on the

**Lemma 3.** Let  $D_1(\boldsymbol{\lambda})$  be the set of  $\mathbf{y}$  satisfying  $Q(\mathbf{y}) < k$  (const). Then for each  $i$  ( $i=1, 2, \dots, p$ ) the function  $P_{\lambda_i^*}^Y(D_1(\boldsymbol{\lambda}))$  of  $\boldsymbol{\lambda}^* = (\lambda_1^*, \dots, \lambda_p^*)$  is monotonically

(i) *non-decreasing in  $\lambda_i^*$  if  $0 < \lambda_i \leq 1$ ,*

*and*

(ii) *non-increasing in  $\lambda_i^*$  if  $1 \leq \lambda_i$ .*

Proof of (ii). Suppose  $\lambda_i \geq 1$  and let  $\boldsymbol{\lambda}^{**}$  be a vector which differs by one component  $\lambda_i^{**} > \lambda_i^*$  from  $\boldsymbol{\lambda}^*$ . Since  $(1 - 1/\lambda_j)\lambda_j^{**} \geq (1 - 1/\lambda_j)\lambda_j^*$  for every  $j$ , we have as in (2.10)

$$\begin{aligned} P_{\boldsymbol{\lambda}^{**}}^Y(D_1(\boldsymbol{\lambda})) &= \Pr\left(\sum_{j=1}^p\left(1 - \frac{1}{\lambda_j}\right)\lambda_j^{**}Z_j^2 < k\right) \\ &\leq \Pr\left(\sum_{j=1}^p\left(1 - \frac{1}{\lambda_j}\right)\lambda_j^*Z_j^2 < k\right) = P_{\boldsymbol{\lambda}^*}^Y(D_1(\boldsymbol{\lambda})). \end{aligned}$$

And (i) is proved similarly.

Proof of Theorem 3. To prove (ii) suppose  $1 \leq \lambda_i$  and let  $\boldsymbol{\lambda}^*$  be a vector which differs by one component  $\lambda_i^* > \lambda_i$  from  $\boldsymbol{\lambda}$ . We must show

$$(2.11) \quad R_B(\boldsymbol{\lambda}) > R_B(\boldsymbol{\lambda}^*).$$

It is seen from (2.10) that

$$R_B(\boldsymbol{\lambda}) = \pi_1 w_1 P_1^Y(D_2(\boldsymbol{\lambda})) + \pi_2 w_2 P_{\boldsymbol{\lambda}^*}^Y(D_1(\boldsymbol{\lambda})).$$

Lemma 3 then yields that

$$R_B(\boldsymbol{\lambda}) \geq \pi_1 w_1 P_1^Y(D_2(\boldsymbol{\lambda})) + \pi_2 w_2 P_{\boldsymbol{\lambda}^*}^Y(D_1(\boldsymbol{\lambda})).$$

Obviously the Bayes discrimination  $(D_1(\boldsymbol{\lambda}^*), D_2(\boldsymbol{\lambda}^*))$  corresponding to the pair of populations  $(P_1^Y, P_{\boldsymbol{\lambda}^*}^Y)$  does not coincide (a.e.) with  $(D_1(\boldsymbol{\lambda}), D_2(\boldsymbol{\lambda}))$ ; hence the right-hand side of the last relation is larger than

$$\pi_1 w_1 P_1^Y(D_2(\boldsymbol{\lambda}^*)) + \pi_2 w_2 P_{\boldsymbol{\lambda}^*}^Y(D_1(\boldsymbol{\lambda}^*))$$

which is equal to  $R_B(\boldsymbol{\lambda}^*)$ . This implies (2.11). The proof of (i) is quite analogous.

Theorem 3 means that as eigenvalues  $\lambda_i$  are the more distant from 1, larger or smaller, the smaller grows the average risk  $R_B(\boldsymbol{\lambda})$  or the more efficient the Bayes discrimination becomes.

**3. Minimax discrimination for variance matrices.** We shall now consider the minimax discrimination for two normal populations  $\Pi_1$  and  $\Pi_2$ . From Theorem 2 in Section 1 we obtain the minimax discrimination function

$$(3.1) \quad \varphi_M(\boldsymbol{x}) = \begin{cases} 1 & \text{if } Q > k_M, \\ 0 & \text{if } Q < k_M, \end{cases}$$

where

$$(3.2) \quad Q = (\mathbf{x} - \boldsymbol{\mu})'(\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})(\mathbf{x} - \boldsymbol{\mu})$$

and  $k_M = k_M(\boldsymbol{\lambda})$  is the constant depending only on  $\boldsymbol{\lambda}$ , determined by the equation  $R_{\varphi_M}(P_1) = R_{\varphi_M}(P_2)$  or by

$$(3.3) \quad w_1 Pr \left( \sum_{i=1}^p \left( 1 - \frac{1}{\lambda_i} \right) Z_i^2 > k_M(\boldsymbol{\lambda}) \right) = w_2 Pr \left( \sum_{i=1}^p (\lambda_i - 1) Z_i^2 < k_M(\boldsymbol{\lambda}) \right)$$

on account of (2.10). On each side of (3.3) there appears a weighted sum of  $\chi^2$  variates and so the equation cannot be solved to represent  $k_M(\boldsymbol{\lambda})$  in such a simple formula as (2.9) for the Bayes case. In particular the minimax discrimination function when  $w_1 = w_2$  does not coincide with the Bayes one when  $w_1 = w_2$  and  $\pi_1 = \pi_2$ , while for the discrimination for means two functions coincide with each other. The value of  $k_M(\boldsymbol{\lambda})$  will be obtained by numerical computation as will be explained in Section 5 but we note here that it is continuous in  $\boldsymbol{\lambda}$ .

Let us study the behavior of the risk  $R_{\varphi_M}^* = R_{\varphi_M}(P_i)$  of the minimax discrimination  $\varphi_M$  when  $\lambda_i$  varies with the  $w_i$  being fixed. We write  $R_M(\boldsymbol{\lambda})$  for  $R_{\varphi_M}^*$  to indicate its dependence on  $\boldsymbol{\lambda}$ .

**Theorem 4.** *For each  $i$  ( $i = 1, 2, \dots, p$ ) the function  $R_M(\boldsymbol{\lambda})$  is strictly monotonically*

(i) *increasing in  $\lambda_i$  if  $0 < \lambda_i \leq 1$*

*and*

(ii) *decreasing in  $\lambda_i$  if  $1 \leq \lambda_i$ .*

*Proof.* To prove (ii) suppose  $1 \leq \lambda_i$  and let  $\boldsymbol{\lambda}^*$  be a vector which differs by one component  $\lambda_i^* > \lambda_i$  from  $\boldsymbol{\lambda}$ . we show that

$$(3.4) \quad R_M(\boldsymbol{\lambda}) > R_M(\boldsymbol{\lambda}^*).$$

Indeed we get from (2.10)

$$R_M(\boldsymbol{\lambda}) = \max(w_1 P_1^Y(D_2(\boldsymbol{\lambda})), w_2 P_{\lambda}^Y(D_1(\boldsymbol{\lambda}))),$$

where  $D_1(\boldsymbol{\lambda})$  and  $D_2(\boldsymbol{\lambda})$  are given by (2.8) with  $k_M(\boldsymbol{\lambda})$  in place of  $k_B(\boldsymbol{\lambda})$  there. Then from Lemma 3

$$R_M(\boldsymbol{\lambda}) \geq \max(w_1 P_1^Y(D_2(\boldsymbol{\lambda})), w_2 P_{\lambda^*}^Y(D_1(\boldsymbol{\lambda}))).$$

Obviously again the minimax discrimination  $(D_1(\boldsymbol{\lambda}^*), D_2(\boldsymbol{\lambda}^*))$  corresponding to the pair of populations  $(P_1^Y, P_{\lambda^*}^Y)$  does not coincide (a.e.) with  $(D_1(\boldsymbol{\lambda}), D_2(\boldsymbol{\lambda}))$ , and hence the right-hand side of the last expression is larger than

$$\max (w_1 P_Y^Y(D_2(\boldsymbol{\lambda}^*)), w_2 P_{\boldsymbol{\lambda}^*}^Y(D_1(\boldsymbol{\lambda}^*))),$$

which is equal to  $R_M(\boldsymbol{\lambda}^*)$ . Thus (3.4) holds as asserted and the proof of (ii) is complete. (i) is proved similarly.

It should be remarked that this theorem is not included in Theorem 3 in despite of the fact that the minimax discrimination is Bayes for a particular choice of the prior probability.

We find that the more distant are the  $\lambda$ 's from 1, the more efficient are the minimax as well as the Bayes discrimination.

**4. Reduction in dimensions.** In the preceding two sections we have discussed the Bayes and the minimax discrimination utilizing the whole information of a  $p$ -dimensional observation  $\boldsymbol{x}$  in the space  $\mathcal{X}$ . What problem will arise if any discrimination is to be performed utilizing only a projection of  $\boldsymbol{x}$  on a certain  $q$ -dimensional ( $q \leq p$ ) subspace  $\mathcal{X}^*$  of  $\mathcal{X}$ ? There occur two cases:  $\mathcal{X}^*$  is given to us a priori or it can be so chosen by us as to enjoy in some sense optimal property, with only the number  $q$  of dimensions being fixed. While the former case involves no new problem, the latter does. Assume for simplicity that the eigenvalues defined by (2.5) are distinct and let  $\boldsymbol{y} = (y_1, y_2, \dots, y_p)'$  be the canonical variate defined uniquely by (2.6). The problem is then how to choose a new variate

$$(4.1) \quad \boldsymbol{x}^* = \boldsymbol{A}\boldsymbol{y},$$

where  $\boldsymbol{A}$  denotes a  $q \times p$  constant matrix, as a basis of the subspace  $\mathcal{X}^*$  in order to get the most efficient discrimination.

Now there is a well-known (cf. for example Hamburger & Grimshaw [10], p. 75)

**THEOREM** (Cauchy's inequality). *Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  be the eigenvalues of a real symmetric  $p \times p$  matrix  $\boldsymbol{A}$  and let  $\lambda_1^* \geq \lambda_2^* \geq \dots \geq \lambda_q^*$  be those of  $\boldsymbol{A}\boldsymbol{A}'$ ,  $\boldsymbol{A}$  denoting any  $q \times p$  matrix such that  $\boldsymbol{A}\boldsymbol{A}' = \boldsymbol{I}_q$  (identity). Then it holds that*

$$(4.2) \quad \lambda_{i+p-q} \leq \lambda_i^* \leq \lambda_i \quad (i=1, 2, \dots, q).$$

Now we state

**Theorem 5.** *Given the number  $q$  of dimensions, a basis  $\boldsymbol{x}^*$  of the subspace which minimizes the average (or maximum) risk of the Bayes (or minimax) discrimination is given by one of  $(q+1)$  variates  $(y_1, \dots, y_s, y_{p-q+s+1}, \dots, y_p)$  where  $s=0, 1, \dots, q$ . The corresponding discrimination function is*

$$(4.3) \quad \varphi^*(\mathbf{x}^*) = \begin{cases} 1 & \text{if } Q^* > k \\ 0 & \text{if } Q^* < k, \end{cases}$$

where

$$(4.4) \quad Q^* = \left( \sum_{i=1}^s + \sum_{i=p-q+s+1}^n \right) \left( 1 - \frac{1}{\lambda_i} \right) y_i^2$$

and

$$(4.5) \quad k = k_B(\boldsymbol{\lambda}^*) \quad \text{or} \quad k = k_M(\boldsymbol{\lambda}^*),$$

which is obtained from (2.9) or (3.3) by replacing  $\boldsymbol{\lambda}$  there by  $\boldsymbol{\lambda}^* = (\lambda_1, \dots, \lambda_s, \lambda_{p-q+s+1}, \dots, \lambda_p)$ , according as  $\varphi^*$  is Bayes or minimax.

The solution is determined uniquely:  $s=q$  when all  $\lambda_i$  are  $\geq 1$ , or  $s=0$  when all  $\lambda_i$  are  $\leq 1$ .

Proof. We shall consider only the Bayes case since the proof is quite similar for the minimax one. We may assume that  $\mathbf{A}\mathbf{A}' = \mathbf{I}_q$ . If  $\mathbf{x}$  comes from the population  $\Pi_1$  or from  $\Pi_2$ , then  $\mathbf{x}^*$  is regarded as coming from the population  $\Pi_1^* : N(\mathbf{0}, \mathbf{I}_q)$  or  $\Pi_2^* : N(\mathbf{0}, \mathbf{A}\mathbf{A}')$ , respectively. Denote by  $\lambda_i^*$  ( $\lambda_1^* \geq \lambda_2^* \geq \dots \geq \lambda_q^*$ ) the roots of the determinantal equation  $|\mathbf{A}\mathbf{A}' - \lambda^* \mathbf{I}_q| = 0$ .

Applying the discussion in Section 2 to the pair of populations ( $\Pi_1^*, \Pi_2^*$ ), we know that the average risk of the Bayes discrimination is a function  $R_B(\boldsymbol{\lambda}^*)$  of  $\boldsymbol{\lambda}^* = (\lambda_1^*, \dots, \lambda_q^*)$ , which grows smaller as the  $\lambda_i^*$  are more distant from 1. For the eigenvalues  $\lambda_i$  and  $\lambda_i^*$  the Cauchy inequality (4.2) holds, whence it is readily seen that  $R_B(\boldsymbol{\lambda}^*)$  attains its minimum at  $\boldsymbol{\lambda}^* = (\lambda_1, \dots, \lambda_q)$  if all the  $\lambda_i$  are  $\geq 1$  or at  $\boldsymbol{\lambda}^* = (\lambda_{p-q+1}, \dots, \lambda_p)$  if  $\lambda_i \leq 1$  for all  $i$ . If among the  $\lambda$ 's there are some  $> 1$  and some  $< 1$ , then the minimum point  $\boldsymbol{\lambda}^*$  picks up some largest  $\lambda_i$  among those which are greater than 1 and some smallest among those which are less than 1 to fill the dimension  $q$ . Summing up,  $R_B(\boldsymbol{\lambda}^*)$  attains its minimum at one of  $(q+1)$  points  $\boldsymbol{\lambda}^* = (\lambda_1, \dots, \lambda_s, \lambda_{p-q+s+1}, \dots, \lambda_p)$ , where  $s=0, 1, \dots, q$ . To this  $\boldsymbol{\lambda}^*$  corresponds the variate  $\mathbf{x}^* = (y_1, \dots, y_s, y_{p-q+s+1}, \dots, y_p)$  and this proves the theorem.

**5. Weighted sum of  $\chi^2$  variates.** In applying the results in the preceding sections to any actual data it is necessary to calculate the probabilities of the form

$$Pr \left( \sum_{i=1}^n (\lambda_i - 1) Z_i^2 < k \right) \quad \text{or} \quad Pr \left( \sum_{i=1}^n \left( 1 - \frac{1}{\lambda_i} \right) Z_i^2 > k \right),$$

where  $Z_1, Z_2, \dots, Z_p$  are independent  $N(0, 1)$  variates. We give in this section a practical procedure convenient for evaluating such probabilities.

We shall consider the distribution of a weighted sum

$$(5.1) \quad W = \sum_{i=1}^p a_i Z_i^2$$

of  $\chi^2$  variates  $Z_i^2$ , where we assume that all the coefficients  $a_i$  are of the same sign, positive in fact. This requires that in the preceding sections either all the  $\lambda_i \geq 1$  or all  $\leq 1$  but this, it seems to the author, is not a serious restriction for any practical application.

The distribution of  $W$  has been studied occasionally by several authors. For instance, Satterthwaite [29] approximates it by the distribution of  $k\chi_n^2$ , where  $k = \Sigma a_i^2 / \Sigma a_i$  and  $n = (\Sigma a_i)^2 / \Sigma a_i^2$ . The notation  $\chi_n^2$  will be used henceforth consistently for a  $\chi^2$  variate with  $n$  degrees of freedom. The paper [20] of the present author gives an inequality

$$Pr(\sum_{i=1}^p a_i Z_i^2 < c) \leq Pr(a\chi_p^2 < c),$$

where  $a = (\Pi a_i)^{1/p}$  and  $c$  is any constant. It may be available as an approximation provided that the  $a_i$ 's differ relatively little. Robbins & Pitman [27] and Box [5] also dealt with this problem. The former obtained an interesting result by means of the method of mixture due to Robbins [26], which will be improved as follows from the point of view of accelerating the convergence.

The characteristic function of  $W$  is

$$(5.2) \quad \varphi(t) = \prod_{j=1}^p (1 - 2ia_j t)^{-1/2}.$$

Using an arbitrary positive constant  $x$ , we get

$$1 - 2ia_j t = \frac{a_j}{x} \left[ (1 - 2ixt) - \left( 1 - \frac{x}{a_j} \right) \right] = \frac{a_j}{x} w^{-2} (1 - c_j w^2),$$

where

$$(5.3) \quad c_j = 1 - \frac{x}{a_j}, \quad w = (1 - 2ixt)^{-1/2}.$$

We suppose here that the complex function  $w$  denotes the branch which takes the value 1 when  $t=0$ . Putting

$$(5.4) \quad f(w^2) = \prod_{j=1}^p (1 - c_j w^2),$$

we rewrite (5.2) in the form

$$(5.5) \quad \varphi(t) = x^{p/2} \left( \prod_{j=1}^p a_j \right)^{-1/2} w^p [f(w^2)]^{-1/2}.$$

For  $x$ , an arbitrary constant, Robbins & Pitman [27] chose

$$(5.6) \quad x = \min_j a_j = x_0 \quad (\text{say}),$$

which implies  $0 \leq c_j < 1$  in view of (5.3) for every  $j$ . This and the fact that  $|w^2| = |1 - 2ix_0 t|^{-1} \leq 1$  for any real  $t$  together guarantee the absolute convergence of the expansion of  $(1 - c_j w^2)^{-1/2}$  in  $w^2$  and hence we find

$$[f(w^2)]^{-1/2} = \sum_{n=0}^{\infty} f_n w^{2n},$$

$f_n$  being constant coefficients. Then we have from (5.5)

$$(5.7) \quad \varphi(t) = x_0^{p/2} \left( \prod_{j=1}^n a_j \right)^{-1/2} \sum_{n=0}^{\infty} f_n w^{2n+p}.$$

Since  $w^m = (1 - 2ix_0 t)^{-m/2}$  is the characteristic function of a variate  $x_0 \chi_m^2$ , we obtain Robbins-Pitman's expansion

$$(5.8) \quad Pr\left(\sum_{j=1}^n a_j Z_j^2 < c\right) = x_0^{p/2} \left( \prod_{j=1}^n a_j \right)^{-1/2} \sum_{n=0}^{\infty} f_n Pr(x_0 \chi_{2n+p}^2 < c)$$

for any  $c$ .

Though the convergence of (5.8) results necessarily from the choice (5.6) of  $x$ , its speed proves to be very slow in many cases. Thus it is required to find an alternative choice of  $x$  capable of accelerating the convergence. We recommend here one which makes the coefficient of  $w^2$  in the expansion of  $f(w^2)$  vanish, although unfortunately we have not succeeded yet in proving the convergence in general. The condition stated implies

$$(5.9) \quad \frac{p}{x} = \sum_{j=1}^n \frac{1}{a_j},$$

or  $x$  is the harmonic mean of the  $a_i$ 's. From the definition of  $x$ , however, it will be expected that the expansion of  $[f(w^2)]^{-1/2}$  in  $w^2$  converges about twice as fast as that of the Robbins-Pitman expansion provided the convergence is assured anyhow for the former. This expectation is really met in the following example.

Consider the special case where the number  $p$  of dimensions is even ( $=2r$ ) and the coefficients  $a_j$  are divided into two groups such that

$$a_1 = a_2 = \cdots = a_r \quad \text{and} \quad a_{r+1} = \cdots = a_p.$$

Then from (5.9), (5.3) and (5.4) we have

$$(5.10) \quad x = \frac{2a_1 a_p}{a_1 + a_p}, \quad c_1 = \cdots = c_r = -c_{r+1} = \cdots = -c_p = \frac{a_1 - a_p}{a_1 + a_p}$$

and

$$(5.11) \quad f(w^2) = (1 - c_1 w^2)^r (1 - c_p w^2)^r = (1 - c_1^2 w^4)^r.$$

Since  $0 \leq c_1^2 < 1$  and  $|w^4| = |1 - 2ixt|^{-2} \leq 1$  for any real  $t$ , we obtain the expansion converging absolutely

$$(5.12) \quad [f(w^2)]^{-1/2} = (1 - c_1^2 w^4)^{-r/2} = \sum_{n=0}^{\infty} b_n c_1^{2n} w^{4n},$$

where

$$(5.13) \quad b_0 = 1, \quad b_n = r(r+2) \cdots (r+2(n-1))/2^n n! \quad (n = 1, 2, \dots).$$

Substitution of (5.12) into (5.5) yields

$$\varphi(t) = x^r (a_1 a_p)^{-r/2} \sum_{n=0}^{\infty} b_n c_1^{2n} w^{4n+2r},$$

which implies

$$(5.14) \quad Pr(a_1 \chi_r^2 + a_p \chi_r'^2 < c) = x^r (a_1 a_p)^{-r/2} \sum_{n=0}^{\infty} b_n c_1^{2n} Pr(\chi_{4n+2r}^2 < c)$$

for any  $c$ , where  $\chi_r'^2$  denotes a  $\chi^2$  variate with  $r$  degrees of freedom distributed independently of  $\chi_r^2$ . We may suppose  $a_1 \geq a_p$  without loss of generality and put  $a_1 = a a_p$  ( $a \geq 1$ ). Replacing  $c$  in (5.14) by  $a_p c$  and substituting (5.10), we have the result

$$(5.15) \quad Pr(a \chi_r^2 + \chi_r'^2 < c) = \left(\frac{2\sqrt{a}}{a+1}\right)^r \sum_{n=0}^{\infty} b_n \left(\frac{a-1}{a+1}\right)^{2n} Pr\left(\chi_{4n+2r}^2 < \frac{a+1}{2a} c\right)$$

for any  $c$ .

On the other hand if we set  $x = \min a_j = a_p$  after Robbins & Pitman, then we have

$$[f(w^2)]^{-1/2} = \left[1 - \left(1 - \frac{1}{a}\right) w^2\right]^{-r/2} = \sum_{n=0}^{\infty} b_n \left(1 - \frac{1}{a}\right)^n w^{2n}$$

with the coefficients  $b_n$  defined in (5.13). Hence

$$(5.16) \quad Pr(a \chi_r^2 + \chi_r'^2 < c) = \left(\frac{1}{\sqrt{a}}\right)^r \sum_{n=0}^{\infty} b_n \left(\frac{a-1}{a}\right)^n Pr(\chi_{2n+2r}^2 < c)$$

for any  $c$ . The speed of convergence of either (5.15) or (5.16) depends on  $a$ ,  $c$  and  $r$ . Especially it becomes slower as  $a$  varies from 1 to infinity because of the factor  $\left(\frac{a-1}{a+1}\right)^{2n}$  or  $\left(\frac{a-1}{a}\right)^n$ .

We shall now compare two expansions (5.15) and (5.16) with each other. Specifically, let us compare the  $n$ th term of (5.15) with the  $(2n)$ th term of (5.16) in the following four items:

- (i) the coefficients  $b_n$  and  $b_{2n}$  are of comparable magnitude;
- (ii) the ratio  $\left(\frac{a-1}{a+1}\right)^2$  of the geometric series of (5.15) is smaller than  $\left(\frac{a-1}{a}\right)^2$  of (5.16);
- (iii) each number of degrees of freedom of  $\chi^2$  variates is equal ( $=4n+2r$ );
- (iv) the abscissa  $\frac{a+1}{2a}c$  of the  $\chi^2$  distribution of (5.15) is smaller than  $c$  of (5.16).

From these comparisons it might be said that (5.15) converges twice or more as speedily as (5.16) does. And when  $a$  is very large the advantage is more than twice as seen from the item (iv) above, which property is very profitable because the convergence then slows down for either expansion.

Note that the formula (5.15) can be easily evaluated numerically by using

$$(5.17) \quad \begin{aligned} Pr(\chi_{2m}^2 < A) &= \frac{1}{(m-1)!} \int_0^{A/2} x^{m-1} e^{-x} dx \\ &= 1 - e^{-A/2} \sum_{i=0}^{m-1} \frac{1}{i!} \left(\frac{A}{2}\right)^i. \end{aligned}$$

In particular, if  $p$  equals two, then as a special case ( $r=1$ ) of (5.15) we get

$$(5.18) \quad Pr(aZ_1^2 + Z_2^2 < c) = \frac{2\sqrt{a}}{a+1} \sum_{n=0}^{\infty} b_n \left(\frac{a-1}{a+1}\right)^{2n} Pr\left(\chi_{4n+2}^2 < \frac{a+1}{2a}c\right)$$

for any  $c$ , where from (5.13)

$$(5.19) \quad b_0 = 1, \quad b_n = \frac{1 \cdot 3 \cdot \dots \cdot (2n-1)}{2 \cdot 4 \cdot \dots \cdot (2n)} \quad \text{for } n \geq 1.$$

In Section 3 as well as in Section 4 we have reduced the problem of the minimax discrimination to the equation (3.3), which is written in the form

$$(5.20) \quad Pr((\lambda_1 - 1)Z_1^2 < k_M(\lambda_1)) = Pr\left(\left(1 - \frac{1}{\lambda_1}\right) Z_1^2 > k_M(\lambda_1)\right)$$

or

$$(5.21) \quad \begin{aligned} Pr((\lambda_1 - 1)Z_1^2 + (\lambda_2 - 1)Z_2^2 < k_M(\boldsymbol{\lambda})) \\ = Pr\left(\left(1 - \frac{1}{\lambda_1}\right) Z_1^2 + \left(1 - \frac{1}{\lambda_2}\right) Z_2^2 > k_M(\boldsymbol{\lambda})\right), \quad \boldsymbol{\lambda} = (\lambda_1, \lambda_2) \end{aligned}$$

when  $w_1=w_2$  and  $p=1$  or  $2$ , respectively. The value of  $k_M(\lambda_1)$  satisfying (5.20) is easily calculated from any table of normal probability functions. As for (5.21) we may apply the expansion (5.18) to both sides, replace the probabilities appearing there by the formula (5.17) and utilize any numerical method for solving an equation.

Table 1 gives the probability (5.20) or (5.21) of the error committed by the minimax discrimination procedure and Table 2 the corresponding critical value  $k_M(\lambda)$  for some typical values of  $\lambda_1$  or of the pair  $(\lambda_1, \lambda_2)$ . The first column in each table designated as  $\lambda_2=1$  corresponds to (5.20) and others to (5.21). Actual computation was carried out by the automatic computer NEAC 2203 at Electronic Equipment Industry Division, Nippon Electric Company, Tokyo.

Table 1. Probabilities of error of the minimax discrimination when  $q=1$  or 2.

$\lambda_2 \backslash \lambda_1$	1	2	4	6	8	10	12	14	16	18	20	25	30	35	40	45	50	60	70	80	90	100	
1	.5000																						
2	.4263	.3820																					
4	.3557	.3238	.2755																				
6	.3171	.2904	.2473	.2219																			
8	.2911	.2676	.2280	.2045	.1883																		
10	.2718	.2506	.2136	.1915	.1763	.1649																	
12	.2567	.2372	.2022	.1812	.1668	.1560	.1474																
14	.2443	.2262	.1929	.1728	.1590	.1486	.1405	.1338															
16	.2340	.2169	.1850	.1657	.1524	.1425	.1346	.1282	.1228														
18	.2251	.2090	.1783	.1596	.1468	.1372	.1296	.1234	.1182	.1138													
20	.2173	.2020	.1724	.1543	.1419	.1326	.1252	.1192	.1142	.1099	.1061												
25	.2016	.1878	.1604	.1435	.1319	.1232	.1163	.1107	.1060	.1020	.0984	.0913											
30	.1894	.1768	.1510	.1351	.1241	.1159	.1094	.1041	.0996	.0958	.0925	.0858	.0805										
35	.1795	.1678	.1434	.1283	.1178	.1100	.1038	.0986	.0945	.0909	.0877	.0813	.0763	.0723									
40	.1713	.1604	.1371	.1226	.1126	.1051	.0991	.0943	.0902	.0868	.0837	.0776	.0728	.0690	.0658								
45	.1643	.1540	.1317	.1178	.1081	.1009	.0952	.0905	.0866	.0833	.0803	.0744	.0698	.0661	.0631	.0604							
50	.1582	.1484	.1270	.1136	.1042	.0972	.0917	.0872	.0835	.0802	.0774	.0717	.0672	.0637	.0607	.0582	.0560						
60	.1481	.1392	.1192	.1066	.0978	.0912	.0860	.0818	.0782	.0752	.0725	.0671	.0630	.0596	.0568	.0545	.0524	.0490					
70	.1400	.1318	.1130	.1010	.0926	.0864	.0814	.0774	.0740	.0711	.0686	.0635	.0595	.0564	.0537	.0515	.0495	.0463	.0437				
80	.1333	.1256	.1077	.0963	.0883	.0823	.0776	.0738	.0705	.0678	.0654	.0605	.0567	.0537	.0511	.0490	.0471	.0441	.0416	.0396			
90	.1276	.1204	.1033	.0923	.0847	.0789	.0744	.0707	.0676	.0649	.0626	.0579	.0543	.0514	.0489	.0469	.0451	.0422	.0398	.0379	.0362		
100	.1227	.1158	.0995	.0889	.0815	.0760	.0716	.0680	.0650	.0625	.0602	.0557	.0522	.0494	.0471	.0451	.0434	.0405	.0382	.0364	.0348	.0334	

Table 2. Critical values of the minimax discrimination procedure when  $q=1$  or 2.

$\lambda_2$ $\lambda_1$	1	2	4	6	8	10	12	14	16	18	20	25	30	35	40	45	50	60	70	80	90	100
1	0.000																					
2	0.316	0.962																				
4	0.640	1.397	1.934																			
6	0.834	1.629	2.211	2.509																		
8	0.976	1.790	2.400	2.711	2.922																	
10	1.087	1.913	2.545	2.865	3.081	3.244																
12	1.179	2.014	2.661	2.989	3.209	3.376	3.509															
14	1.259	2.100	2.760	3.093	3.317	3.486	3.622	3.735														
16	1.328	2.174	2.844	3.183	3.410	3.581	3.718	3.833	3.932													
18	1.390	2.240	2.919	3.262	3.491	3.644	3.803	3.919	4.018	4.106												
20	1.446	2.298	2.986	3.332	3.564	3.738	3.878	3.995	4.096	4.184	4.262											
25	1.565	2.424	3.128	3.481	3.718	3.895	4.038	4.157	4.259	4.348	4.428	4.596										
30	1.665	2.527	3.244	3.603	3.843	4.023	4.168	4.288	4.391	4.482	4.562	4.732	4.870									
35	1.750	2.615	3.342	3.706	3.949	4.131	4.277	4.399	4.503	4.595	4.676	4.847	4.986	5.104								
40	1.825	2.691	3.427	3.796	4.041	4.225	4.372	4.495	4.600	4.692	4.774	4.947	5.087	5.205	5.307							
45	1.892	2.759	3.503	3.875	4.122	4.308	4.456	4.580	4.686	4.779	4.861	5.035	5.176	5.294	5.397	5.487						
50	1.952	2.820	3.571	3.946	4.195	4.382	4.531	4.656	4.762	4.856	4.939	5.113	5.255	5.374	5.477	5.568	5.649					
60	2.057	2.927	3.688	4.069	4.322	4.511	4.662	4.787	4.895	4.989	5.073	5.249	5.392	5.512	5.616	5.708	5.789	5.931				
70	2.147	3.018	3.789	4.173	4.429	4.620	4.772	4.899	5.007	5.102	5.187	5.364	5.508	5.629	5.734	5.826	5.908	6.050	6.170			
80	2.226	3.097	3.876	4.264	4.522	4.714	4.868	4.995	5.105	5.200	5.285	5.464	5.609	5.730	5.835	5.928	6.011	6.154	6.274	6.379		
90	2.296	3.167	3.953	4.344	4.604	4.798	4.952	5.080	5.191	5.287	5.373	5.552	5.697	5.820	5.925	6.018	6.101	6.245	6.366	6.471	6.564	
100	2.359	3.231	4.022	4.417	4.678	4.873	5.028	5.157	5.268	5.365	5.451	5.631	5.777	5.900	6.006	6.099	6.183	6.327	6.448	6.553	6.646	6.729

## PART II. INCOMPLETELY SPECIFIED POPULATIONS

**6. Discrimination when parameters are estimated.** We turn to the case when for two  $p$ -dimensional populations  $\Pi_i: N(\mu, \Sigma_i)$  ( $i=1, 2$ ) the variance matrices  $\Sigma_i$  are unknown whereas the common mean vector  $\mu$  is either (1°) known or (2°) unknown, while a random sample taken from each population provides the information about unknown parameters.

Let  $\mathbf{X}_i = (\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}, \dots, \mathbf{x}_{n_i}^{(i)})$  be a sample of size  $n_i$  from each  $\Pi_i$ . Suppose that given an observation  $\mathbf{x}$  we wish to decide whether  $\mathbf{x}$  comes from  $\Pi_1$  or  $\Pi_2$ , utilizing also the knowledge of  $\mathbf{X}_1$  and  $\mathbf{X}_2$ . Then any discrimination function should be a function of  $\mathbf{x}$ ,  $\mathbf{X}_1$  and  $\mathbf{X}_2$ . In this general formulation, however, there is unfortunately no such simple results as in Section 1 of Part I. It might be possible to introduce the notion such as the invariance and the power for the purpose of obtaining any optimal discrimination function, as was done by Kudô [16], [17] in order to justify the classical discriminant function for the problem of mean vectors. But we shall look to another occasion for this approach and confine ourselves in the present paper to a study of the sampling distribution and of estimation, adopting the discrimination function obtained by applying a conventional modification to that used in Part I.

Let us assume for simplicity throughout this Part that the eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_p$  given by

$$(6.1) \quad |\Sigma_2 - \lambda \Sigma_1| = 0$$

are all distinct. Then the matrix  $\mathbf{F}$  of eigenvectors defined by

$$(6.2) \quad \mathbf{F}'\Sigma_1\mathbf{F} = \mathbf{I}, \quad \mathbf{F}'\Sigma_2\mathbf{F} = \Lambda$$

is determined uniquely up to the sign of every column. To determine  $\mathbf{F}$  completely we require that the first row of  $\mathbf{F}$  consists of positive elements. Now two cases should be distinguished.

CASE (1°):  $\mu$  is known. Put for each  $i$

$$(6.3) \quad \hat{\Sigma}_i = \frac{1}{n_i} \sum_{\alpha=1}^{n_i} (\mathbf{x}_\alpha^{(i)} - \mu)(\mathbf{x}_\alpha^{(i)} - \mu)',$$

the maximum likelihood estimate of  $\Sigma_i$ . Since the statistic  $(\hat{\Sigma}_1, \hat{\Sigma}_2)$  is sufficient for the unknown parameter  $(\Sigma_1, \Sigma_2)$  we have only to consider a discrimination function which depends only on  $\mathbf{x}$ ,  $\hat{\Sigma}_1$  and  $\hat{\Sigma}_2$ , provided that our concern lies in minimizing the risk. In fact we adopt

$$(6.4) \quad \varphi(\mathbf{x}, \hat{\Sigma}_1, \hat{\Sigma}_2) = \begin{cases} 1 & \text{if } \hat{Q} > \hat{k} \\ 0 & \text{if } \hat{Q} < \hat{k}, \end{cases}$$

where

$$(6.5) \quad \hat{Q} = (\mathbf{x} - \boldsymbol{\mu})' (\hat{\Sigma}_1^{-1} - \hat{\Sigma}_2^{-1}) (\mathbf{x} - \boldsymbol{\mu}),$$

which will be called the (*estimated*) *quadratic discriminant function*.  $\hat{Q}$  is defined by replacing the  $\Sigma_i$ 's in the formula (2.2) by their estimates  $\hat{\Sigma}_i$ 's. The critical value  $\hat{k}$  will be discussed later on.

First we consider the sampling distribution of  $\hat{Q}$ . We show that it depends only on  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$ . Indeed, putting

$$(6.6) \quad \Sigma_1^* = \mathbf{F}' \hat{\Sigma}_1 \mathbf{F}, \quad \Sigma_2^* = \mathbf{F}' \hat{\Sigma}_2 \mathbf{F}$$

and

$$(6.7) \quad \mathbf{y} = \mathbf{F}'(\mathbf{x} - \boldsymbol{\mu})$$

we have

$$(6.8) \quad Q = \mathbf{y}'(\Sigma_1^{*-1} - \Sigma_2^{*-1})\mathbf{y}.$$

As is well-known the random matrix  $n_i \hat{\Sigma}_i$  follows the Wishart distribution  $W(\Sigma_i, n_i)$  for each  $i$  (for the notation see Anderson [4]); hence  $n_1 \Sigma_1^*$  and  $n_2 \Sigma_2^*$  are distributed according to  $W(\mathbf{I}, n_1)$  and  $W(\boldsymbol{\Lambda}, n_2)$ , respectively. On the other hand  $\mathbf{y}$  follows  $N(\mathbf{0}, \mathbf{I})$  or  $N(\mathbf{0}, \boldsymbol{\Lambda})$  according as  $\mathbf{x}$  comes from  $\Pi_1$  or from  $\Pi_2$ . Thus from (6.8) it is seen that the distribution of  $\hat{Q}$  depends only on  $\boldsymbol{\lambda}$  as contended.

Next, consider the estimation of  $\boldsymbol{\Lambda}$ ,  $\mathbf{F}$  and  $\mathbf{y}$ . As in Section 2 we know that there exist a matrix  $\hat{\mathbf{F}}$  and a diagonal matrix  $\hat{\boldsymbol{\Lambda}}$  with the diagonal elements in descending order such that

$$(6.9) \quad \hat{\mathbf{F}}' \hat{\Sigma}_1 \hat{\mathbf{F}} = \mathbf{I}, \quad \hat{\mathbf{F}}' \hat{\Sigma}_2 \hat{\mathbf{F}} = \hat{\boldsymbol{\Lambda}}.$$

The diagonal elements  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$  of  $\hat{\boldsymbol{\Lambda}}$  are the roots of

$$(6.10) \quad |\hat{\Sigma}_2 - \hat{\lambda} \hat{\Sigma}_1| = 0.$$

Then  $\hat{\boldsymbol{\Lambda}}$  consists of eigenvalues and  $\hat{\mathbf{F}}$  of eigenvectors with respect to the pair  $(\hat{\Sigma}_1, \hat{\Sigma}_2)$ . We adopt  $(\hat{\boldsymbol{\Lambda}}, \hat{\mathbf{F}})$  as an estimate of  $(\boldsymbol{\Lambda}, \mathbf{F})$ . It is the maximum likelihood estimate because of the corresponding property of  $(\hat{\Sigma}_1, \hat{\Sigma}_2)$ . The probability that the equation (6.10) has equal roots is zero, and hence  $\hat{\mathbf{F}}$  is determined uniquely with probability one if we adopt again the convention used for  $\mathbf{F}$ .

Quite similarly there exists a canonical reduction  $(\boldsymbol{\Lambda}^*, \mathbf{F}^*)$  such that

$$(6.11) \quad \mathbf{F}^{*'} \Sigma_1^* \mathbf{F}^* = \mathbf{I}, \quad \mathbf{F}^{*'} \Sigma_2^* \mathbf{F}^* = \boldsymbol{\Lambda}^*.$$

The diagonal elements of  $\Lambda^*$  are the roots of

$$(6.12) \quad |\Sigma_2^* - \lambda^* \Sigma_1^*| = 0$$

and the uniqueness of  $F^*$  can be ensured if we require that every diagonal element of  $F^*$  is non-negative. Substitution of (6.6) into (6.11) yields

$$(\mathbf{F}\mathbf{F}^*)' \hat{\Sigma}_1(\mathbf{F}\mathbf{F}^*) = \mathbf{I}, \quad (\mathbf{F}\mathbf{F}^*)' \hat{\Sigma}_2(\mathbf{F}\mathbf{F}^*) = \Lambda^*.$$

Comparing this with (6.9) we know from the uniqueness of the solution of (6.9) that with probability one

$$(6.13) \quad \hat{\Lambda} = \Lambda^*, \quad \hat{F} = \mathbf{F}\mathbf{F}^* \mathbf{D}_\pm,$$

where  $\mathbf{D}_\pm$  denotes a certain diagonal matrix with elements  $+1$  or  $-1$ . The problem is thus reduced to  $(\Lambda^*, F^*)$  given by (6.11) and will be taken up in the following section. But we note here that the distribution of  $(\Lambda^*, F^*)$  depends only on  $\lambda$ .

We turn to the "estimation" of the canonical variate  $\mathbf{y}$ . Note that the quantity to be estimated is not a constant but a variate. We may take as an estimate of  $\mathbf{y}$

$$(6.14) \quad \hat{\mathbf{y}} = \hat{F}'(\mathbf{x} - \mu)$$

with  $\hat{F}$  being defined by (6.9). In view of (6.13) and (6.7) this is written in the form

$$(6.15) \quad \hat{\mathbf{y}} = \mathbf{D}_\pm \mathbf{F}'^* \mathbf{y},$$

which reduces the problem again to  $F^*$ . Since the distribution of both  $F^*$  and  $\mathbf{y}$  depend only on  $\lambda$ , so is that of  $\hat{\mathbf{y}}$  up to the sign of every component.

As for  $\hat{k}$  in the formula (6.4) we set

$$(6.16) \quad \hat{k} = k_B(\hat{\lambda}) \quad \text{or} \quad \hat{k} = k_M(\hat{\lambda})$$

according as we wish to get an asymptotically Bayes or a minimax discrimination, where  $k_B(\lambda)$  is given by (2.9),  $k_M(\lambda)$  by (3.3) and  $\hat{\lambda}$  by (6.10).

CASE (2°):  $\mu$  is unknown. For each  $i$  ( $i=1, 2$ ) the sample mean vector and the sample variance matrix are defined as follows:

$$(6.17) \quad \bar{\mathbf{x}}^{(i)} = \frac{1}{n_i} \sum_{\alpha=1}^{n_i} \mathbf{x}_\alpha^{(i)},$$

$$(6.18) \quad \hat{\Sigma}_i = \frac{1}{n_i - 1} \sum_{\alpha=1}^{n_i} (\mathbf{x}_\alpha^{(i)} - \bar{\mathbf{x}}^{(i)}) (\mathbf{x}_\alpha^{(i)} - \bar{\mathbf{x}}^{(i)})'.$$

Since the statistic  $(\bar{\mathbf{x}}^{(1)}, \bar{\mathbf{x}}^{(2)}, \hat{\Sigma}_1, \hat{\Sigma}_2)$  is sufficient for the unknown parameter  $(\boldsymbol{\mu}, \Sigma_1, \Sigma_2)$ , it suffices to consider a discrimination function which is dependent only on this sufficient statistic and  $\mathbf{x}$ . We adopt in fact

$$(6.19) \quad \varphi(\mathbf{x}, \bar{\mathbf{x}}^{(1)}, \bar{\mathbf{x}}^{(2)}, \hat{\Sigma}_1, \hat{\Sigma}_2) = \begin{cases} 1 & \text{if } \hat{Q} > \hat{k}, \\ 0 & \text{if } \hat{Q} < \hat{k}, \end{cases}$$

where

$$(6.20) \quad \hat{Q} = (\mathbf{x} - \bar{\mathbf{x}})' (\hat{\Sigma}_1^{-1} - \hat{\Sigma}_2^{-1}) (\mathbf{x} - \bar{\mathbf{x}})$$

$$(6.21) \quad \bar{\mathbf{x}} = \frac{n_1 \bar{\mathbf{x}}^{(1)} + n_2 \bar{\mathbf{x}}^{(2)}}{n_1 + n_2}.$$

In doing this we have replaced  $\boldsymbol{\mu}, \Sigma_1$  and  $\Sigma_2$  in (2.2) by their estimates  $\bar{\mathbf{x}}, \hat{\Sigma}_1$  and  $\hat{\Sigma}_2$ , which are *not* the maximum likelihood estimates contrary to Case (1°).

The problem again concerns with the sampling distribution of  $\hat{Q}$  and the estimation of  $\mathbf{A}, \mathbf{F}$  and  $\mathbf{y}$ . Put

$$(6.22) \quad \Sigma_1^* = \mathbf{F}' \hat{\Sigma}_1 \mathbf{F}, \quad \Sigma_2^* = \mathbf{F}' \hat{\Sigma}_2 \mathbf{F}$$

and

$$(6.23) \quad \bar{\mathbf{y}} = \mathbf{F}'(\bar{\mathbf{x}} - \boldsymbol{\mu}),$$

then it follows that

$$(6.24) \quad \hat{Q} = (\mathbf{y} - \bar{\mathbf{y}})' (\Sigma_1^{*-1} - \Sigma_2^{*-1}) (\mathbf{y} - \bar{\mathbf{y}}).$$

We know that  $(n_i - 1) \hat{\Sigma}_i$  obeys the Wishart distribution  $W(\Sigma_i, n_i - 1)$  for each  $i$  and hence  $(n_1 - 1) \Sigma_1^*$  and  $(n_2 - 1) \Sigma_2^*$  obey  $W(\mathbf{I}, n_1 - 1)$  and  $W(\mathbf{A}, n_2 - 1)$  respectively. Furthermore the distribution of  $\bar{\mathbf{y}}$  is  $N(\mathbf{0}, (n_1 \mathbf{I} + n_2 \mathbf{A}) / (n_1 + n_2)^2)$  and that of  $\mathbf{y}$  is  $N(\mathbf{0}, \mathbf{I})$  or  $N(\mathbf{0}, \mathbf{A})$  according as  $x$  comes from  $\Pi_1$  or from  $\Pi_2$ . This time again the distribution of  $\hat{Q}$  depends only on  $\boldsymbol{\lambda}$ .

As regards the estimation of  $(\mathbf{A}, \mathbf{F})$  and the choice of  $\hat{k}$  in (6.19) the situation is quite analogous with Case (1°), so that the equation from (6.9) through (6.13) as well as (6.16) hold with only the alteration that the number of degrees of freedom of  $\Sigma_i$  is reduced from  $n_i$  to  $n_i - 1$  for each  $i$ .

As for the estimation of  $\mathbf{y}$  we take the estimate

$$(6.25) \quad \hat{\mathbf{y}} = \mathbf{F}'(\mathbf{x} - \bar{\mathbf{x}})$$

in place of (6.14). This can be written in the form

$$(6.26) \quad \hat{\mathbf{y}} = D_{\pm} \mathbf{F}^{*'} (\mathbf{y} - \bar{\mathbf{y}}),$$

whose distribution depends only on  $\boldsymbol{\lambda}$ .

**7. Reduction in dimensions.** We have utilized in the preceding section the whole information provided by  $p$ -dimensional variates  $\mathbf{x}$  and  $\mathbf{x}_a^{(i)}$ . In parallel with Section 4 of Part I it is conceivable to reduce the number of dimensions from  $p$  to  $q(\leq p)$ , holding then the efficiency of discrimination as high as possible. Along the line of Theorem 5 we consider a discrimination function written in the form

$$(7.1) \quad \varphi(\mathbf{x}, \hat{\Sigma}_1, \hat{\Sigma}_2) = \begin{cases} 1 & \text{if } \hat{Q}^* > \hat{k}, \\ 0 & \text{if } \hat{Q}^* < \hat{k}. \end{cases}$$

Then we have only to determine  $\hat{Q}^*$  and  $\hat{k}$  as respective estimate of  $Q^*$  in (4.4) and of  $k$  in (4.5). We set in fact

$$(7.2) \quad \hat{Q}^* = \left( \sum_{i=1}^s + \sum_{i=p-q+s+1}^n \right) \left( 1 - \frac{1}{\hat{\lambda}_i} \right) \hat{y}_i^2,$$

where the  $\hat{\lambda}_i$  are given by (6.9), while the  $\hat{y}_i$  are defined either by (6.14) for Case (1°) or by (6.25) for Case (2°). We set further

$$(7.3) \quad \hat{k} = k_B(\hat{\boldsymbol{\lambda}}^*) \quad \text{or} \quad \hat{k} = k_M(\hat{\boldsymbol{\lambda}}^*),$$

according as we aim at an asymptotically Bayes or a minimax discrimination, where the functions  $k_B(\boldsymbol{\lambda}^*)$  and  $k_M(\boldsymbol{\lambda}^*)$  are the same as in (4.5) and

$$\hat{\boldsymbol{\lambda}}^* = (\hat{\lambda}_1, \dots, \hat{\lambda}_s, \hat{\lambda}_{p-q+s+1}, \dots, \hat{\lambda}_p).$$

Clearly, from the discussion of the preceding section the distribution of  $Q^*$  depends only on  $\boldsymbol{\lambda}$  but not necessarily on  $\boldsymbol{\lambda}^* = (\lambda_1, \dots, \lambda_s, \lambda_{p-q+s+1}, \dots, \lambda_p)$ . Now from (6.8), (6.11) and (6.15) it holds that

$$\hat{Q} = \hat{\mathbf{y}}'(\mathbf{I} - \hat{\mathbf{A}}^{-1})\hat{\mathbf{y}} = \sum_{i=1}^p \left( 1 - \frac{1}{\hat{\lambda}_i} \right) \hat{y}_i^2.$$

A comparison of this and (7.2) shows that  $\hat{Q}^*$  picks up dominant terms,  $q$  in number, of the canonical form of  $\hat{Q}$ .

**8. Asymptotic distribution of the eigenvalues and the eigenvectors.** Though we are interested in the asymptotic distribution of the random variable  $(\mathbf{A}^*, \mathbf{F}^*)$  defined by (6.11), let us change the notation for convenience. Let  $\mathbf{A}_{n_1}$  and  $\mathbf{B}_{n_2}$  be random matrices such that  $n_1\mathbf{A}_{n_1}$  and  $n_2\mathbf{B}_{n_2}$  follow independently the Wishart distributions  $W(\mathbf{I}, n_1)$  and  $W(\mathbf{A}, n_2)$ , respectively, where  $\mathbf{A}$  is diagonal with diagonal elements satisfying

$$(8.1) \quad \lambda_1 > \lambda_2 > \dots > \lambda_p > 0.$$

Then the set of equations

$$(8.2) \quad \mathbf{F}'_n \mathbf{A}_{n_1} \mathbf{F}_n = \mathbf{I}, \quad \mathbf{F}'_n \mathbf{B}_{n_2} \mathbf{F}_n = \mathbf{A}_n$$

determines uniquely with probability one the diagonal matrix  $\mathbf{A}_n = (\lambda_i(n) \delta_{ij})$  of eigenvalues and the matrix  $\mathbf{F}_n = (f_{ij}(n))$  of eigenvectors subject to the additional conditions

$$(8.3) \quad \lambda_1(n) \geq \lambda_2(n) \geq \dots \geq \lambda_p(n)$$

and

$$(8.4) \quad f_{ii}(n) \geq 0 \quad (i=1, 2, \dots, p).$$

For simplicity we suppose  $n_1 = n$ , regarding  $n_2$  as a function of  $n_1$ . The problem is to find the asymptotic distribution of  $(\mathbf{A}_n, \mathbf{F}_n)$  when  $n_1, n_2 \rightarrow \infty$  in such a way that

$$(8.5) \quad \sqrt{\frac{n_1}{n_2}} \rightarrow c \quad (\text{a positive constant}).$$

Analogous problems are dealt with in Hsu [12], [13], [14] and Anderson [1], [2]. The following derivation is based on Rubin's theorem as is the case with Anderson [2]:

**THEOREM (Rubin).** *Let  $\mathbf{X}_n$  ( $n=1, 2, \dots$ ) and  $\mathbf{X}$  be  $p$ -dimensional random vectors and let  $\mathbf{f}_n$  ( $n=1, 2, \dots$ ) and  $\mathbf{f}$  be mappings from a Euclidean  $p$ -space to a  $q$ -space. Suppose that*

- (i)  $\mathbf{X}_n$  converges in law to  $\mathbf{X}$  as  $n \rightarrow \infty$ ,
- (ii) for every continuity point  $\mathbf{x}$  of  $\mathbf{f}$ , it holds that  $\mathbf{f}_n(\mathbf{x}_n) \rightarrow \mathbf{f}(\mathbf{x})$  whenever  $\mathbf{x}_n \rightarrow \mathbf{x}$ , and
- (iii) the probability that  $\mathbf{X}$  falls in the set of discontinuities of  $\mathbf{f}$  is zero.

Then  $\mathbf{f}_n(\mathbf{x}_n)$  converges in law to  $\mathbf{f}(\mathbf{X})$  as  $n \rightarrow \infty$ .

We now state the

**Theorem 6.** (1) For each  $i$  ( $i=1, 2, \dots, p$ ) the random vector  $(\lambda_i(n), f_{ii}(n))$  follows asymptotically a normal distribution with

$$\begin{pmatrix} \lambda_i \\ 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 2\lambda_i^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right) & \frac{\lambda_i}{n_1} \\ \frac{\lambda_i}{n_1} & \frac{1}{2n_1} \end{pmatrix}$$

as the mean vector and the variance matrix, respectively.

(2) For each pair  $(i, j)$  ( $i, j=1, 2, \dots, p; i < j$ ) the random vector  $(f_{ij}(n), f_{ji}(n))$  is distributed asymptotically normally with

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \text{ and } \frac{1}{(\lambda_i - \lambda_j)^2} \begin{pmatrix} \lambda_j \left( \frac{\lambda_j + \lambda_i}{n_1} + \frac{\lambda_i}{n_2} \right) & -\lambda_i \lambda_j \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \\ -\lambda_i \lambda_j \left( \frac{1}{n_1} + \frac{1}{n_2} \right) & \lambda_i \left( \frac{\lambda_i + \lambda_j}{n_1} + \frac{\lambda_j}{n_2} \right) \end{pmatrix}$$

as the mean vector and the variance matrix, respectively.

(3) *These random vectors are distributed asymptotically independently.*

As a consequence of this theorem we know the following properties concerning the quantities defined by (6.9), (6.11) and (6.13). First

$$(8.6) \quad \text{plim } \mathbf{A}^* = \mathbf{A}, \quad \text{plim } \mathbf{F}^* = \mathbf{I}$$

as  $n_1, n_2 \rightarrow \infty$  under (8.5). Accordingly,  $\text{plim } \mathbf{D}_\pm = \mathbf{I}$  and hence

$$(8.7) \quad \text{plim } \hat{\mathbf{A}} = \mathbf{A}, \quad \text{plim } \hat{\mathbf{F}} = \mathbf{F}.$$

Thus  $(\hat{\mathbf{A}}, \hat{\mathbf{F}})$  is a consistent estimate of  $(\mathbf{A}, \mathbf{F})$ .

Furthermore, from (1) of Theorem 6 it is seen that  $\log \hat{\lambda}_i$  follows asymptotically

$$(8.8) \quad N\left(\log \lambda_i, 2\left(\frac{1}{n_1} + \frac{1}{n_2}\right)\right),$$

whence a confidence interval of  $\log \lambda_i$  or of  $\lambda_i$  itself can be obtained.

As for  $\hat{\mathbf{y}}$ , the estimated canonical variate, we know from (6.15) or (6.26) that

$$(8.9) \quad \hat{\mathbf{y}} \rightarrow \mathbf{y} \quad (\text{in law})$$

for either of Cases (1°) and (2°).

Proof of Theorem 6. Put

$$(8.10) \quad \mathbf{U}_n = \sqrt{n_1}(\mathbf{A}_{n_1} - \mathbf{I}), \quad \mathbf{V}_n = \sqrt{n_2}(\mathbf{B}_{n_2} - \mathbf{A}),$$

then from the central limit theorem their elements  $u_{ij}(n)$  and  $v_{ij}(n)$  ( $i, j=1, 2, \dots, p; i \leq j$ ) follow in the limit independent normal distributions with the mean zero in common and the variances

$$(8.11) \quad \begin{cases} E(u_i^2) = 2, & E(v_{ii}^2) = 2\lambda_i^2, \\ E(u_j^2) = 1, & E(v_{ij}^2) = \lambda_i \lambda_j \quad (i < j). \end{cases}$$

Put

$$(8.12) \quad \boldsymbol{\theta}_n = \sqrt{n}(\mathbf{A}_n - \mathbf{A}), \quad \mathbf{Z}_n = \sqrt{n}(\mathbf{F}_n - \mathbf{I}),$$

or, in terms of elements,

$$(8.13) \quad \theta_i(n) = \sqrt{n}(\lambda_i(n) - \lambda_i), \quad z_{ij}(n) = \sqrt{n}(f_{ij}(n) - \delta_{ij}),$$

then we have to prove that  $(\theta_i(n), z_{ii}(n))$  and  $(z_{ij}(n), z_{ji}(n))$  follows in the limit independent normal distributions with the mean  $(0, 0)$  in common and the variance matrices

$$(8.14) \quad \begin{pmatrix} 2\lambda_i^2(1+c^2) & \lambda_i \\ \lambda_i & \frac{1}{2} \end{pmatrix} \text{ and } \frac{1}{(\lambda_i - \lambda_j)^2} \begin{pmatrix} \lambda_j(\lambda_j + c^2\lambda_i) & -\lambda_i\lambda_j(1+c^2) \\ -\lambda_i\lambda_j(1+c^2) & \lambda_i(\lambda_i + c^2\lambda_j) \end{pmatrix},$$

respectively.

Substituting  $A_{n_1}$ ,  $B_{n_2}$ ,  $A_n$  and  $F_n$  given by (8.10) and (8.12) into (8.2), (8.3) and (8.4), we get respectively

$$(8.15) \quad \begin{cases} \left(I + \frac{1}{\sqrt{n}} Z'_n\right) \left(I + \frac{1}{\sqrt{n_1}} U_n\right) \left(I + \frac{1}{\sqrt{n}} Z_n\right) = I, \\ \left(I + \frac{1}{\sqrt{n}} Z'_n\right) \left(A + \frac{1}{\sqrt{n_2}} V_n\right) \left(I + \frac{1}{\sqrt{n}} Z_n\right) = A + \frac{1}{\sqrt{n}} \theta_n, \end{cases}$$

$$(8.16) \quad \lambda_1 + \frac{1}{\sqrt{n}} \theta_1(n) \geq \lambda_2 + \frac{1}{\sqrt{n}} \theta_2(n) \geq \dots \geq \lambda_p + \frac{1}{\sqrt{n}} \theta_p(n)$$

and

$$(8.17) \quad 1 + \frac{1}{\sqrt{n}} z_{ii}(n) \geq 0 \quad (i = 1, 2, \dots, p).$$

These equations together determine  $(\theta_n, Z_n)$  uniquely with probability one. If we neglect in them the terms of order  $n^{-1}$  or higher, then, omitting the subscript  $n$  of  $\theta_n$ ,  $Z_n$ ,  $U_n$  and  $V_n$ , we get from (8.15)

$$(8.18) \quad \begin{cases} Z + Z' + U = \mathbf{0}, \\ AZ + Z'A + cV = \theta, \end{cases}$$

while both (8.16) and (8.17) reduce to triviality.

It is easily seen that (8.18) determines  $(\theta, Z)$  uniquely in terms of  $(A, U, V)$ . Indeed, putting  $U = (u_{ij})$ ,  $V = (v_{ij})$ ,  $\theta = (\theta_i \delta_{ij})$  and  $Z = (z_{ij})$ , we get the solution

$$(8.19) \quad \begin{cases} \theta_i = cv_{ii} - \lambda_i u_{ii}, & z_{ii} = -\frac{1}{2} u_{ii}, \\ z_{ij} = \frac{\lambda_j u_{ij} - cv_{ij}}{\lambda_i - \lambda_j}, & z_{ji} = \frac{cv_{ij} - \lambda_i u_{ij}}{\lambda_i - \lambda_j} \quad (i < j). \end{cases}$$

Suppose that  $u_{ij}$  and  $v_{ij}$  ( $i, j = 1, 2, \dots, p; i \leq j$ ) follow independent normal distributions with the mean zero and the variances given by (8.11), so that  $U_n \rightarrow U$  and  $V_n \rightarrow V$  in law as  $n \rightarrow \infty$ . Then (8.19) implies that  $(\theta_i, z_{ii})$  and  $(z_{ij}, z_{ji})$  are distributed independently normally with the

mean  $(0, 0)$  and the variance matrices given by (8.14). Thus we have only to prove that  $(\theta_n, \mathbf{Z}_n)$  defined by the equations (8.15), (8.16) and (8.17) converges in law to  $(\theta, \mathbf{Z})$  defined by (8.18). Substituting  $\mathbf{X}_n = (\mathbf{U}_n, \mathbf{V}_n)$ ,  $\mathbf{X} = (\mathbf{U}, \mathbf{V})$ ,  $\mathbf{f}_n(\mathbf{X}_n) = (\theta_n, \mathbf{Z}_n)$  and  $\mathbf{f}(\mathbf{X}) = (\theta, \mathbf{Z})$  in Rubin's theorem, we see that the condition (i) is already ascertained and (iii) is trivial since the mapping  $\mathbf{f}$  is continuous as is seen from (8.19). It remains thus to verify the condition (ii).

From now on we regard all variates in the equations from (8.15) to (8.18) as non-stochastic and prove that for any pair of matrices  $(\mathbf{U}, \mathbf{V})$  the convergence

$$(8.20) \quad (\mathbf{U}_n, \mathbf{V}_n) \rightarrow (\mathbf{U}, \mathbf{V}) \quad \text{as } n \rightarrow \infty$$

implies that

$$(8.21) \quad (\theta_n, \mathbf{Z}_n) \rightarrow (\theta, \mathbf{Z}) \quad \text{as } n \rightarrow \infty.$$

We rewrite (8.15) in the form

$$(8.22) \quad \begin{cases} \mathbf{F}'_n \left( \mathbf{I} + \frac{1}{\sqrt{n}} \mathbf{U}_n \right) \mathbf{F}_n = \mathbf{I}, \\ \mathbf{F}'_n \left( \mathbf{A} + \frac{1}{\sqrt{n_2}} \mathbf{V}_n \right) \mathbf{F}_n = \mathbf{A}_n \end{cases}$$

and we begin by showing that

$$(8.23) \quad \mathbf{F}_n \rightarrow \mathbf{I}, \quad \mathbf{A}_n \rightarrow \mathbf{A} \quad \text{as } n \rightarrow \infty.$$

To verify this we have to prove that for any subsequence  $\{n'\}$  of  $\{n\}$  there exists a subsequence  $\{n''\}$  of  $\{n'\}$  such that

$$(8.24) \quad \mathbf{F}_{n''} \rightarrow \mathbf{I}, \quad \mathbf{A}_{n''} \rightarrow \mathbf{A} \quad \text{as } n'' \rightarrow \infty.$$

Now from the first equation of (8.22) we see that  $\mathbf{F}_n$  ( $n=1, 2, \dots$ ) are bounded and hence for any  $\{n'\}$  there exist a subsequence  $\{n''\}$  and a matrix  $\mathbf{F}^* = (f_{ij}^*)$  such that  $\mathbf{F}_{n''} \rightarrow \mathbf{F}^*$  as  $n'' \rightarrow \infty$ . This together with (8.20) and (8.22) implies that  $\mathbf{A}_{n''}$  converges to a certain diagonal matrix  $\mathbf{A}^* = (\lambda_i^* \delta_{ij})$  as  $n'' \rightarrow \infty$  such that

$$(8.25) \quad \mathbf{F}^* \mathbf{F}^* = \mathbf{I}, \quad \mathbf{F}^* \mathbf{A} \mathbf{F}^* = \mathbf{A}^*.$$

Passing to the limit in (8.3) and (8.4), we know that  $\lambda_1^* \geq \lambda_2^* \geq \dots \geq \lambda_p^*$  and  $f_{ii}^* \geq 0$  for every  $i$ . Accordingly, the uniqueness of the canonical reduction (8.25) implies  $\mathbf{F}^* = \mathbf{I}$  and  $\mathbf{A}^* = \mathbf{A}$ , which proves (8.24) as asserted.

Thus we have obtained the first terms of  $\mathbf{A}_n$  and  $\mathbf{F}_n$ . We now turn to the second terms  $\boldsymbol{\theta}_n$  and  $\mathbf{Z}_n$  defined by (8.12). The convergence  $\mathbf{F}_n \rightarrow \mathbf{I}$  implies that for sufficiently large  $n$   $\mathbf{F}_n$  is nonsingular; hence two equations in (8.22) together with the first equation of (8.12) yield

$$(8.26) \quad \left( \mathbf{A} + \frac{1}{\sqrt{n_2}} \mathbf{V}_n \right) \mathbf{F}_n = \left( \mathbf{I} + \frac{1}{\sqrt{n}} \mathbf{U}_n \right) \mathbf{F}_n \left( \mathbf{A} + \frac{1}{\sqrt{n}} \boldsymbol{\theta}_n \right).$$

The equation for the  $(i, i)$  elements is

$$\begin{aligned} \lambda_i f_{ii}(n) + \frac{1}{\sqrt{n_2}} \sum_{k=1}^p v_{ik}(n) f_{ki}(n) \\ = \left( \lambda_i + \frac{1}{\sqrt{n}} \theta_i(n) \right) \left[ f_{ii}(n) + \frac{1}{\sqrt{n}} \sum_{k=1}^p u_{ik}(n) f_{ki}(n) \right] \end{aligned}$$

or equivalently

$$(8.27) \quad \begin{aligned} \sum_{k=1}^p f_{ki}(n) \left[ \sqrt{\frac{n_1}{n_2}} v_{ik}(n) - \lambda_i u_{ik}(n) \right] \\ = \theta_i(n) \left[ f_{ii}(n) + \frac{1}{\sqrt{n}} \sum_{k=1}^p u_{ik}(n) f_{ki}(n) \right]. \end{aligned}$$

Using (8.20) and (8.23) or, in terms of elements,  $u_{ij}(n) \rightarrow u_{ij}$ , and  $v_{ij}(n) \rightarrow v_{ij}$  and  $f_{ij}(n) \rightarrow \delta_{ij}$ , we find from (8.27) that

$$(8.28) \quad \theta_i(n) \rightarrow cv_{ii} - \lambda_i u_{ii} \quad \text{for every } i.$$

The equation for the  $(i, j)$  elements ( $i \neq j$ ) of (8.26) is

$$\begin{aligned} \lambda_i f_{ij}(n) + \frac{1}{\sqrt{n_2}} \sum_{k=1}^p v_{ik}(n) f_{kj}(n) \\ = \left( \lambda_j + \frac{1}{\sqrt{n}} \theta_j(n) \right) \left[ f_{ij}(n) + \frac{1}{\sqrt{n}} \sum_{k=1}^p u_{ik}(n) f_{kj}(n) \right]. \end{aligned}$$

Substituting  $z_{ij}(n) = \sqrt{n} f_{ij}(n)$ , we get

$$\begin{aligned} (\lambda_i - \lambda_j) z_{ij}(n) + \sum_{k=1}^p f_{kj}(n) \left[ \sqrt{\frac{n_1}{n_2}} v_{ik}(n) - \lambda_j u_{ik}(n) \right] \\ = \theta_j(n) \left[ f_{ij}(n) + \frac{1}{\sqrt{n}} \sum_{k=1}^p u_{ik}(n) f_{kj}(n) \right] \end{aligned}$$

In view of the convergences (8.5), (8.20), (8.23) and (8.28) it holds that

$$(8.29) \quad z_{ij}(n) \rightarrow \frac{\lambda_j u_{ij} - cv_{ij}}{\lambda_i - \lambda_j} \quad \text{for } i \neq j.$$

Finally from the first equation of (8.22) we get for every  $i$

$$\sum_{k=1}^n f_{ki}^2(n) + \frac{1}{\sqrt{n}} \sum_{k=1}^n \sum_{l=1}^n f_{ki}(n) u_{kl}(n) f_{li}(n) = 1.$$

Substituting  $f_{ii}(n) = 1 + z_{ii}(n)/\sqrt{n}$ , we obtain after a straight-forward calculation

$$(8.30) \quad z_{ii}(n) \rightarrow -\frac{1}{2} u_{ii} \quad \text{for every } i.$$

The relations (8.28), (8.29), (8.30) together with (8.19) imply (8.21) and the proof of the theorem is complete.

### 9. Asymptotic distribution of the quadratic discriminant function.

We consider the asymptotic property of the discrimination procedure defined by (7.1) including then (6.14) and (6.19) as special cases. As is seen in the preceding section  $\hat{\lambda} = (\hat{\lambda}_1, \dots, \hat{\lambda}_p)$  converges in probability to  $\lambda = (\lambda_1, \dots, \lambda_p)$  and  $\hat{y} = (\hat{y}_1, \dots, \hat{y}_p)'$  converges in law to  $y = (y_1, \dots, y_p)'$  as  $n_1, n_2 \rightarrow \infty$  under the restriction (8.5). Therefore from the well-known theorem for stochastic convergence (cf. Cramér [6], p. 254) we find that

$$(9.1) \quad \hat{Q}^* \rightarrow Q^* \quad (\text{in law}),$$

where  $Q^*$  is given by (4.4). On the other hand, since the functions  $k_B(\lambda^*)$  and  $k_M(\lambda^*)$  are continuous in  $\lambda^* = (\lambda_1, \dots, \lambda_s, \lambda_{p-q+s+1}, \dots, \lambda_p)$ , we get

$$(9.2) \quad \text{plim } k_B(\hat{\lambda}^*) = k_B(\lambda^*), \quad \text{plim } k_M(\hat{\lambda}^*) = k_M(\lambda^*).$$

Two convergences (9.1) and (9.2) together imply that the discrimination function (7.1) is asymptotically equivalent to (4.3) and the probabilities of error of two kinds are approximated either by (2.10) for the Bayes case or by the probabilities appearing in (3.3) for the minimax case, using in either case  $\lambda^*$  in place of  $\lambda$ .

This is the first approximation, as it were, of the asymptotic distribution of  $\hat{Q}^*$ . What will be then the second approximation? For the problem of discrimination for means of two normal populations  $\Pi_i: N(\mu^{(i)}, \Sigma)$  with the common variance matrix, the forthcoming paper [21] by the author deals with the asymptotic distribution of the linear discriminant function  $V = (\mathbf{x}^{(1)} - \mathbf{x}^{(2)})' \hat{\Sigma}^{-1} \left[ \mathbf{x} - \frac{1}{2} (\mathbf{x}^{(1)} + \mathbf{x}^{(2)}) \right]$ , where  $\mathbf{x}^{(i)}$  denotes for each  $i$  ( $i=1, 2$ ) a sample mean of size  $n_i$  from  $\Pi_i$  and  $\hat{\Sigma}$  an unbiased estimate of  $\Sigma$  with  $f$  degrees of freedom. As  $n_1, n_2, f \rightarrow \infty$   $V$  follows asymptotically the same distribution as  $U = (\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} \left[ \mathbf{x} - \frac{1}{2} (\mu^{(1)} + \mu^{(2)}) \right]$  does, which is  $N\left(\frac{1}{2}D^2, D^2\right)$  or  $N\left(-\frac{1}{2}D^2, D^2\right)$  according as  $\mathbf{x}$  comes from

$\Pi_1$  or from  $\Pi_2$ , where  $D^2 = (\mu^{(1)} - \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)})$  is the Mahalanobis distance. Let  $\Phi(t)$  and  $\phi(t)$  be the cdf and pdf, respectively, of  $N(0, 1)$ , then in the expansion

$$(9.3) \quad Pr\left(V < \frac{1}{2} D^2 + tD \mid \Pi_1\right) = \Phi(t) + \phi(t) \left(\frac{A_1}{n_1} + \frac{A_2}{n_2} + \frac{A_3}{f}\right) + \dots$$

in the inverse power of  $n_1$ ,  $n_2$  and  $f$  the coefficients  $A_1$ ,  $A_2$  and  $A_3$  are given in terms of  $t$  and  $D^2$ . The probability  $P(2|1)$  of error is obtained by setting  $t = -D/2$ . It may be said that the formula (9.3) gives the second approximation of the distribution of  $V$ .

The similar formula for  $\hat{Q}^*$  may be obtainable but the derivation, it seems, is rather difficult for the general case and hence we give only a sketch of the derivation after the studentization method of Moriguti [19] and Wallace [32] for the special case where the mean vector is known and the reduced number  $q$  of dimensions equals one.

We evaluate the probability that the quadratic discriminant function  $\hat{Q}^* = (1 - 1/\hat{\lambda}_1) \hat{y}_1^2$  with  $\hat{y}_1$  given in (6.14) is larger than  $K(\hat{\lambda}_1)$ ,  $K(\lambda)$  being any function twice differentiable in  $\lambda$ , when  $\mathbf{x}$  comes from  $\Pi_1$ , that is,

$$(9.4) \quad Pr(\hat{Q}^* > K(\hat{\lambda}_1) \mid \Pi_1) = P(2|1) \quad (\text{say}).$$

Clearly

$$(9.5) \quad P(2|1) = E' [Pr(\hat{Q}^* > K(\hat{\lambda}_1) \mid \hat{\lambda}; \Pi_1)],$$

where  $E'$  denotes the expectation with regard to  $\hat{\lambda}$ , while  $Pr$  the conditional probability given  $\hat{\lambda}$ . From (6.15) it follows that  $\hat{y}_1 = \pm \sum_{i=1}^p f_{i1}^* y_i$  and hence the conditional distribution of  $\hat{y}_1$  given  $\hat{\lambda}$  is  $N(0, \sum_{i=1}^p f_{i1}^{*2})$ .

If we put

$$\Psi(t) = Pr(X_1^2 < t) = \int_0^t \frac{1}{\sqrt{2\pi t}} e^{-t/2} dt,$$

then (9.5) is written in the form

$$(9.6) \quad P(2|1) = 1 - E' \left[ \Psi \left( \frac{\hat{\lambda}_1 K(\hat{\lambda}_1)}{(\hat{\lambda}_1 - 1) \sum_{i=1}^p f_{i1}^{*2}} \right) \right].$$

If we have the expansion

$$\frac{\hat{\lambda}_1 K(\hat{\lambda}_1)}{(\hat{\lambda}_1 - 1) \sum_{i=1}^p f_{i1}^{*2}} = \frac{\lambda_1 K(\lambda_1)}{\lambda_1 - 1} + \frac{1}{\sqrt{n}} A + \frac{1}{n} B$$

in the power of  $n^{-1/2}$ , where  $n = n_1$  denotes the size of the first sample  $\mathbf{X}_1$ , neglecting the terms of higher order, and if  $E'(A) = 0$ , then it holds that

$$(9.7) \quad P(2|1) = 1 - \Psi\left(\frac{\lambda_1 K(\lambda_1)}{\lambda_1 - 1}\right) - \frac{1}{n} \left[ \psi_2 E'(B) + \frac{1}{2} \psi_2' E'(A^2) \right],$$

where

$$\psi(t) = \frac{d}{dt} \Psi(t) = \frac{1}{\sqrt{2\pi t}} e^{-t/2}, \quad \psi'(t) = \frac{d}{dt} \psi(t)$$

and

$$\psi_2 = \psi\left(\frac{\lambda_1 K(\lambda_1)}{\lambda_1 - 1}\right), \quad \psi_2' = \psi'\left(\frac{\lambda_1 K(\lambda_1)}{\lambda_1 - 1}\right).$$

The formula (9.7) gives the second approximation required. And it remains only to calculate  $E'(B)$  and  $E'(A^2)$ .

From the arguments in Section 7 it follows that

$$\hat{\lambda}_1 = \lambda_1 + \frac{1}{\sqrt{n}} \theta_1 + \frac{1}{n} \phi_1$$

$$f_{11}^* = 1 + \frac{1}{\sqrt{n}} z_{11} + \frac{1}{n} w_1, \quad f_{i1}^* = \frac{1}{\sqrt{n}} z_{i1} \quad (i \neq 1)$$

with the terms of higher order being neglected, where  $\theta_1$ ,  $z_{11}$  and  $z_{i1}$  are given by (8.19) and

$$\phi_1 = \sum_{i=2}^p \frac{(c v_{1i} - \lambda_1 u_{1i})^2}{\lambda_1 - \lambda_i} + u_{11} (\lambda_1 u_{11} - c v_{11}),$$

$$w_1 = -\frac{1}{2} \sum_{i=1}^p z_{i1}^2,$$

$u_{1i}$  and  $v_{1i}$  ( $i=1, \dots, p$ ) denoting independent normal variates with the mean zero in common and the variances (8.11). After a straightforward calculation we get

$$E'(A) = 0,$$

$$(9.8) \quad E'(A^2) = \frac{2\lambda^2}{(\lambda-1)^2} \left[ \lambda^2 \left( \frac{K}{\lambda-1} - K' \right)^2 + c^2 \left( \frac{K}{\lambda-1} - \lambda K' \right)^2 \right],$$

$$(9.9) \quad E'(B) = \frac{\lambda}{\lambda-1} \left\{ \lambda^2 K'' + \lambda K' \left( \lambda S - \frac{2}{\lambda-1} \right) - K \left( 2\lambda S + \frac{\lambda S}{\lambda-1} - \frac{2\lambda}{(\lambda-1)^2} \right) \right\}$$

$$+ c^2 \frac{\lambda}{\lambda-1} \left\{ \lambda^2 K'' + \lambda K' \left( \lambda S - \frac{2}{\lambda-1} - p + 1 \right) - K \left( \frac{\lambda S}{\lambda-1} - \frac{2\lambda}{(\lambda-1)^2} - \frac{p-1}{\lambda-1} \right) \right\},$$

where  $c$  is given by (8.5) and we have put  $\lambda = \lambda_1$ ,  $K = K(\lambda)$ ,  $K' = \frac{d}{d\lambda} K(\lambda)$ ,  $K'' = \frac{d}{d\lambda} K'(\lambda)$  and  $S = \sum_{i=2}^p (\lambda_1 - \lambda_i)^{-1}$ . Substitution of (9.8) and (9.9) into (9.7) yields the second approximation of the probability (9.4).

Similarly, for the probability

$$(9.10) \quad P(1|2) = Pr(\hat{Q}^* < K(\hat{\lambda}_1) | \Pi_2)$$

of the reversed error we have

$$(9.11) \quad P(1|2) = \Psi\left(\frac{K(\lambda_1)}{\lambda_1 - 1}\right) + \frac{1}{n} \left[ \psi_1 E'(D) + \frac{1}{2} \psi'_1 E'(C^2) \right]$$

where

$$\psi_1 = \psi\left(\frac{K}{\lambda_1 - 1}\right), \quad \psi'_1 = \psi'\left(\frac{K}{\lambda_1 - 1}\right)$$

and

$$(9.12) \quad E'(C^2) = \lambda^{-2} E'(A^2)$$

$$(9.13) \quad E'(D) = \left\{ \lambda^2 K'' + \lambda K' \left( \lambda S - \frac{2}{\lambda - 1} \right) - K \left( \frac{\lambda^2 S}{\lambda - 1} - \frac{2\lambda}{(\lambda - 1)^2} \right) \right\} \\ + c^2 \left\{ \lambda^2 K'' + \lambda K' \left( \lambda S - \frac{2}{\lambda - 1} - p + 1 \right) + K \left( (\lambda S - p + 1) \frac{\lambda - 2}{\lambda - 1} + \frac{2\lambda}{(\lambda - 1)^2} \right) \right\}.$$

If we wish to get the probabilities of the two kinds of error committed by the Bayes discrimination procedure, we may set  $K(\lambda) = k_B(\lambda) = \log \lambda$  together with  $K' = \lambda^{-1}$  and  $K'' = -\lambda^{-2}$  in the formulas (9.9) and (9.13), assuming  $\pi_1 = \pi_2$  and  $w_1 = w_2$ . For the minimax case we have the function  $K = K(\lambda) = k_M(\lambda)$  implicitly defined by

$$(9.14) \quad \Psi\left(\frac{K}{\lambda - 1}\right) + \Psi\left(\frac{\lambda K}{\lambda - 1}\right) = 1.$$

By differentiating (9.14) once or twice in  $\lambda$  we know

$$K' = \frac{\psi_1 + \psi_2}{(\lambda - 1)(\psi_1 + \lambda\psi_2)} K, \\ K'' = \frac{\psi_2}{2(\lambda - 1)(\psi_1 + \lambda\psi_2)^2} \left[ \frac{\psi_1(\psi_1 + \psi_2)}{\psi_1 + \lambda\psi_2} K^2 - \left( \frac{3\lambda + 1}{\lambda} \psi_1 + 4\psi_2 \right) K \right],$$

which we may substitute in (9.9) and (9.13).

### PART III. AN APPLICATION

**10. Discrimination of zygoty of twins.** As an illustration of the theory so far developed let us consider the problem of discrimination of zygoty of twins. It is recognized by biologists that there are two kinds of twins, monozygotic or dizygotic. Generally speaking, two members of a pair of monozygotic twins resemble each other in many respects, physical or mental, qualitative or quantitative, more closely than those

of dizygotic twins do. For instance, a pair of monozygotic twins is necessarily like-sexed but this is not true for a dizygotic pair. Conversely, if a pair of twins is opposite-sexed, then it must be dizygotic but if otherwise, we cannot assert anything. We encounter thus the problem of discrimination.

Consider some characteristics of a person,  $p$  in number, distributed continuously among the population of persons and suppose that discrimination is to be based on measurements of these characteristics performed for each member of a pair of twins. Denote by  $\xi$  and  $\eta$  the observed two  $p$ -vectors. It is noted here that we cannot specify which member of the pair  $\xi$  or  $\eta$  should be referred to;  $\xi$  may be referred to either member and  $\eta$  to the other. This situation accompanied with a kind of arbitrariness will be called an *intraclass property*. Now for the population of  $(2p)$ -dimensional variate  $\begin{pmatrix} \xi \\ \eta \end{pmatrix}$  we assume normality, which will not be so unrealistic. Then from the intraclass property above it will be reasonable to assume that the population has the form

$$(10.1) \quad \Pi'_i: N\left(\begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} \Sigma & \Gamma_i \\ \Gamma_i & \Sigma \end{bmatrix}\right) \quad (i = 1, 2),$$

where  $\Gamma_i$ , the covariance matrix of  $\xi$  and  $\eta$ , is symmetric and  $\Pi'_1$  corresponds to monozygotic population and  $\Pi'_2$  to dizygotic. In this formulation the problem is to discriminate between two normal populations having distinct variance matrices with the mean vector in common. The parameters are supposed to be unknown in general.

Of course it is possible at this stage to apply the theory of Part I or II but, conforming to the principle of economy, it is desirable to reduce the number of dimensions of variates considered as far as possible. For this purpose introduce a one-to-one linear transformation

$$(10.2) \quad \begin{pmatrix} \xi \\ \eta \end{pmatrix} \rightarrow \begin{pmatrix} x \\ \tilde{x} \end{pmatrix} = \begin{pmatrix} \xi - \eta \\ \xi + \eta \end{pmatrix}.$$

For each  $i$  ( $i=1, 2$ ) if  $\begin{pmatrix} \xi \\ \eta \end{pmatrix}$  comes from  $\Pi'_i$ , then the new variates  $x$  and  $\tilde{x}$  are distributed independently in

$$(10.3) \quad N(0, 2(\Sigma - \Gamma_i)) \quad \text{and} \quad N(2\mu, 2(\Sigma + \Gamma_i))$$

respectively. As is already seen in Section 2 or 3 the efficiency of the Bayes or minimax discrimination depends only on the eigenvalues with respect to the pair of the variance matrices. Now from the subject

matter consideration above it will be expected that both covariance matrices  $\Gamma_1$  and  $\Gamma_2$  are close to the variance matrix  $\Sigma$  and furthermore that  $\Gamma_1$  is much closer to  $\Sigma$  than  $\Gamma_2$  is. This situation may be represented symbolically in a figure below arranging the matrices in a linear order ;

$$\overline{\begin{array}{ccccccc} \mathbf{O} & \Sigma - \Gamma_1 & \Sigma - \Gamma_2 & & \Sigma & & \Sigma + \Gamma_2 & \Sigma + \Gamma_1 & 2\Sigma \end{array}}$$

hence we have a set of approximate inequalities

$$(10.4) \quad \Sigma < \Sigma + \Gamma_2 < \Sigma + \Gamma_1 < 2\Sigma,$$

where the sign of inequality means that the difference matrix is positive definite. If (10.4) holds exactly, then as is easily verified every eigenvalue with respect to the pair  $(2(\Sigma + \Gamma_2), 2(\Sigma + \Gamma_1))$  is smaller than 2. On the other hand it is quite probable that there exist large eigenvalues among those with respect to  $(2(\Sigma - \Gamma_1), 2(\Sigma - \Gamma_2))$ . This means that the contribution of the variate  $\mathbf{x}$  to the efficiency of discrimination occupies the major part of the efficiency provided by the whole information  $(\mathbf{x}, \tilde{\mathbf{x}})$ , which is equivalent to  $(\xi, \eta)$ . It will not therefore bring a heavy loss to utilize only the measurement  $\mathbf{x}$  instead of  $(\xi, \eta)$ . Then a discrimination should be made between two normal populations

$$(10.5) \quad \Pi_i : N(\mathbf{0}, \Sigma_i), \quad \text{where } \Sigma_i = 2(\Sigma - \Gamma_i) \quad (i = 1, 2).$$

This reduction indeed has been used by biologists from the intuitive ground. We mention further that it has other favorable properties as follows :

(i) The fact that the mean vector of  $\Pi_i$  is known, in fact constant zero, simplifies considerably the theory required and also reduces drastically the amount of numerical computation for obtaining the estimates  $\hat{\Sigma}_i$ , as is seen from a comparison of (6.3) with (6.18).

(ii) It might happen that in the real field of application there exists a certain variation of the population mean  $\mu$ , which we assumed constant. This possibility endangers any application of the theory to the populations  $\Pi'_i$  but not to  $\Pi_i$ . There may also exist a variation of  $(\Sigma, \Gamma_i)$  but most of its possible effect will be absorbed by taking the difference  $\Sigma_i = 2(\Sigma - \Gamma_i)$ .

(iii) Though less probable than (ii) the assumption that  $\mu$  is common in  $\Pi'_1$  and  $\Pi'_2$  may not represent well the situation. As for  $\Pi_i$  we are free also from such a difficulty.

**11. Analysis of the data of twins.** We shall analyze the data\* consisting of measurements of ten anthropological characteristics on a series of 143 pairs of like-sexed twins, aged from 10 to 13, from primary and junior high schools in Osaka City. The constitution is represented in the table below, where the discrimination of the data themselves between

	Male	Female	Total
Monozygotic	48	43	91
Dizygotic	18	34	52

mono- and dizygotic groups is performed by "polysymptomatische Aehnlichkeitsdiagnose" due to H. W. Siemens and O. v. Verschuer and hence there have been excluded several pairs of twins which were difficult to discriminate. Ten characteristics are (1) stature, (2) right iliospinal height, (3) biacromial breadth, (4) upper limb length, (5) maximum head length, (6) maximum head breadth, (7) maximum bizygomatic breadth, (8) bien-tokanthial breadth, (9) total facial length and finally (10) auricular height. All these variates are of continuous type and may be regarded as jointly normally distributed.

Though the data provide us  $p(=10)$ -dimensional variate  $\xi$  as well as  $\eta$  for every pair of twins, we utilize only the difference  $\alpha = \xi - \eta$ , neglecting the information of  $\tilde{\alpha} = \xi + \eta$  as is described in the preceding section. Recall now the property (ii) there. In the present context there exists certainly a variation of the mean  $\mu$  and perhaps also of  $(\Sigma, \Gamma_i)$  resulting from the variation of ages ranging from 10 to 13. But it is expected that such a variation does not exert any serious influence on the application of the theory if we consider only the variate  $\alpha$ . There may also exist most probably a definite difference between male and female populations within either the monozygotic group or the dizygotic, and hence to aim at rigor it is necessary to deal with the problem for each sex separately. But in this example we amalgamate two samples from male and female groups in order to enhance the precision of inference by enlarging the sample sizes and so we try the discrimination of zygosity with the sexes mixed. This procedure though somewhat rough may be tolerable because of the same reason as above.

---

\* These data have been provided to the author cordially by Professor Mototsugu Kohama, Department of Anatomy, School of Medicine, Osaka University and have been analyzed also from other points of view; the case of discrete characteristics is dealt with in Okamoto & Ishii [22] as intraclass contingency table, and Tanaka [30] investigates the inference concerning intraclass canonical correlation.

Tables 3.1 and 3.2 give the moment matrices  $\hat{\Sigma}_1$  and  $\hat{\Sigma}_2$ , respectively, computed by the formula (6.3) with regard to the variate  $x$ . Units of

Table 3.1. Moment matrix of ten characteristics for monozygotic twins

	<i>St</i>	<i>IH</i>	<i>B<sub>a</sub>B</i>	<i>LL</i>	<i>HL</i>	<i>HB</i>	<i>B<sub>2</sub>B</i>	<i>B<sub>c</sub>B</i>	<i>FL</i>	<i>AH</i>
<i>St</i>	3.425	1.948	0.568	1.154	1.965	0.910	1.331	0.880	1.154	1.513
<i>IH</i>		1.810	0.369	0.773	1.342	0.764	0.910	0.465	0.901	0.803
<i>B<sub>a</sub>B</i>			0.774	0.312	-0.158	0.154	0.458	0.186	0.147	-0.226
<i>LL</i>				1.009	0.499	0.480	0.467	0.389	0.353	0.240
<i>HL</i>					12.111	-1.714	-0.879	0.912	0.846	4.671
<i>HB</i>						8.429	3.517	0.033	0.517	1.418
<i>B<sub>2</sub>B</i>							5.616	0.330	1.012	0.231
<i>B<sub>c</sub>B</i>								1.747	0.143	-0.550
<i>FL</i>									4.847	-0.220
<i>AH</i>										19.857

Table 3.2. Moment matrix of ten characteristics for dizygotic twins

	<i>St</i>	<i>IH</i>	<i>B<sub>a</sub>B</i>	<i>LL</i>	<i>HL</i>	<i>HB</i>	<i>B<sub>2</sub>B</i>	<i>B<sub>c</sub>B</i>	<i>FL</i>	<i>AH</i>
<i>St</i>	48.248	28.510	9.286	20.294	1.696	9.135	15.244	1.129	16.550	14.027
<i>IH</i>		19.413	5.423	10.386	0.744	4.656	8.848	0.508	8.460	6.464
<i>B<sub>a</sub>B</i>			3.191	4.360	-0.192	1.875	3.517	-0.150	3.633	1.571
<i>LL</i>				9.941	-1.623	3.404	6.544	1.056	7.352	3.873
<i>HL</i>					50.443	-1.885	6.192	4.039	6.615	11.654
<i>HB</i>						27.635	14.250	3.789	10.654	3.962
<i>B<sub>2</sub>B</i>							19.693	5.173	11.250	4.558
<i>B<sub>c</sub>B</i>								6.404	3.904	0.827
<i>FL</i>									29.808	3.000
<i>AH</i>										35.885

measurement are 1 cm for the first four items and 1 mm for the others. As was expected, every diagonal element of  $\hat{\Sigma}_1$  is much smaller than the corresponding element of  $\hat{\Sigma}_2$ . The eigenvalues given by

$$|\hat{\Sigma}_2 - \lambda \hat{\Sigma}_1| = 0$$

were computed after the fashion of Rao [24]. Actual computation performed by the relay computer FACOM 128A under supervision of Mr. Tutomu Komazawa, the Institute of Statistical Mathematics, yields the largest eigenvalue

$$\lambda_1 = 19.631$$

and the corresponding eigenvector

$$y_1 = 0.5107x_1 + 0.0719x_2 - 0.0244x_3 + 0.2447x_4 - 0.1138x_5 - 0.0563x_6 - 0.0613x_7 - 0.3238x_8 - 0.0378x_9 - 0.0074x_{10},$$

where the coefficients are so standardized that  $y_1$  follows  $N(0, 1)$  for the

population  $\Pi_1$ , that is,  $y_1$  is the first canonical component. The second largest eigenvalue is

$$\lambda_2 = 8.990$$

and the corresponding eigenvector or the second canonical component is

$$v_2 = 0.1384x_1 - 0.7441x_2 + 0.0967x_3 + 0.5291x_4 + 0.0619x_5 \\ + 0.0673x_6 + 0.0579x_7 + 0.1337x_8 + 0.2836x_9 + 0.0122x_{10}.$$

It is seen that in  $y_1$  the term of  $x_1$  is dominant and is followed by  $x_5$  and  $x_8$ , while in  $y_2$  the term of  $x_2$  is dominant, followed by  $x_4$  and  $x_9$ . In either of  $y_1$  and  $y_2$  both  $x_3$  and  $x_{10}$  contribute almost negligible and  $x_6$  and  $x_7$  follow them. Most of these results agree well with those of Katō [15] who investigated the resemblance of twins by the method of "relative deviation", using the same data as ours.

The procedure of discrimination of an observation  $\mathbf{x}$  is as follows. We assume equal weight  $w_1 = w_2$  and equal prior probability  $\pi_1 = \pi_2$ . If the number  $q$  of dimensions utilized is one, we assign  $\mathbf{x}$  to  $\Pi_1$  or to  $\Pi_2$  according as

$$\left(1 - \frac{1}{\lambda_1}\right) y_1^2 < k \quad \text{or} \quad > k,$$

where  $k = k_M = 1.436$  for the minimax case and  $k = k_B = 2.977$  for the Bayes. For the former either probability of error of two kinds is 0.2187 while for the latter  $P(2|1) = 0.0766$  and  $P(1|2) = 0.3107$ . Note that two probabilities for the Bayes case differ markedly. If  $q = 2$ , then we assign  $\mathbf{x}$  to  $\Pi_1$  or to  $\Pi_2$  according as

$$\left(1 - \frac{1}{\lambda_1}\right) y_1^2 + \left(1 - \frac{1}{\lambda_2}\right) y_2^2 < k \quad \text{or} \quad > k,$$

where  $k = k_M = 3.642$  or  $k = k_B = 5.173$ . For the minimax case  $P(2|1) = P(1|2) = 0.1378$  and for the Bayes case  $P(2|1) = 0.0560$  and  $P(1|2) = 0.1894$ . Thus the efficiency improves considerably when  $q = 2$ , compared with the case  $q = 1$ . Such a sharp increase of efficiency will not be expected when  $q$  changes from 2 to 3, since the third eigenvalue  $\lambda_3 = 5.435$  is not so large.

These values of the probability of error are exact only if 19.631 and 8.990 are the true values of the population eigenvalues  $\lambda_1$  and  $\lambda_2$ . Since in fact they are nothing but estimates based on the random samples we can fathom the reliability of the data by calculating  $P(2|1)$  and  $P(1|2)$  along the line described in Section 9. From the theoretical limitation mentioned there we consider only the case  $q = 1$ , obtaining the results below.

$P(2 1)$			$P(1 2)$		
	Minimax	Bayes		Minimax	Bayes
1st approx.	0.2187	0.0765	1st approx.	0.2187	0.3107
Term due to $B$	0.0465	0.0305	Term due to $D$	-0.0039	-0.0062
Term due to $A^2$	0.0020	0.0022	Term due to $C^2$	-0.0004	-0.0006
Total (2nd approx.)	0.2672	0.1092	Total (2nd approx.)	0.2144	0.3039

Since each of the eight correction terms shows small value, compared with the first approximation, except perhaps for the term due to  $B$  in the Bayes  $P(2|1)$ , it might be said that discrimination procedure advocated here, either Bayes or minimax, enjoys the respective optimum property in a rather good approximation.

The author would like to thank Prof. J. Ogawa, Nihon University, for his keen interest and encouraging guidance during the development of this work, Prof. M. Kohama for permission to use his data and Dr. S. Katō as well as Mr. M. Tanaka for helpful suggestions. My thanks are also due to many members of the Institute of Statistical Mathematics and of Nippon Electric Company for computational assistance.

(Received November 28, 1960)

**References**

- [ 1 ] Anderson, T. W.: The asymptotic distributions of the roots of certain determinantal equations, *J. Roy. Statist. Soc. Ser. B.* **10** (1948), 132-139.
- [ 2 ] Anderson, T. W.: The asymptotic distribution of certain characteristic roots and vectors, *Proc. 2nd Berkeley Symposium* (1950), 103-130.
- [ 3 ] Anderson, T. W.: Classification by multivariate analysis, *Psychometrika* **16** (1951), 31-50.
- [ 4 ] Anderson, T. W.: *An Introduction to Multivariate Statistical Analysis*, New York, John Wiley and Sons, 1958.
- [ 5 ] Box, G. E. P.: Some theorems on quadratic forms applied in the study of analysis of variance problems, I. Effect of inequality of variance in the one-way classification, *Ann. Math. Statist.* **25** (1954), 290-302.
- [ 6 ] Cramér, H.: *Mathematical Methods of Statistics*, Princeton, Princeton University Press, 1946.
- [ 7 ] Fisher, R. A.: The use of multiple measurements in taxonomic problems, *Ann. Eugen.* **7** (1936), 179-188.
- [ 8 ] Fisher, R. A.: The statistical utilization of multiple measurements, *Ann. Eugen.* **8** (1938), 376-386.

- [9] Fisher, R. A.: The precision of discriminant functions, *Ann. Eugen.* **10** (1940), 422-429.
- [10] Hamburger, H. L. and M. E. Grimshaw: *Linear Transformations in  $n$ -dimensional Vector Space*, Cambridge University Press, 1951.
- [11] Hoel, P. G. and R. P. Peterson: A solution to the problem of optimum classification, *Ann. Math. Statist.* **20** (1949), 433-438.
- [12] Hsu, P. L.: On the distribution of roots of certain determinantal equations, *Ann. Eugen.* **9** (1939), 250-258.
- [13] Hsu, P. L.: On the limiting distribution of roots of determinantal equation, *J. London Math. Soc.* **16** (1941), 183-194.
- [14] Hsu, P. L.: On the limiting distribution of the canonical correlations, *Biometrika* **32** (1941), 38-45.
- [15] Katō, S.: Somatometrical studies on the resemblance of twins, *Anthrop. Reports*, No. **22** (1958), 1-18.
- [16] Kudō, A.: The classificatory problem viewed as a two-decision problem, *Mem. Fac. Sci. Kyushu Univ.* **13** (1959), 96-125.
- [17] Kudō, A.: The classificatory problem viewed as a two-decision problem, II. *Mem. Fac. Sci. Kyushu Univ.* **14** (1960), 63-83.
- [18] Mises, R. von: On the classification of observation data into distinct groups, *Ann. Math. Statist.* **16** (1945), 68-73.
- [19] Moriguti, S.: A note on Hartley's formula of studentization, *Rep. Statist. Appl. Res.* **2** (1953), 99-103.
- [20] Okamoto, M.: An inequality for the weighted sum of  $\chi^2$  variates, *Bull. Math. Statist.* **9** (1960), 69-70.
- [21] Okamoto, M.: The asymptotic distribution of the linear discriminant function when parameters are estimated, to be published.
- [22] Okamoto, M. and G. Ishii: Test of independence in intraclass  $2 \times 2$  tables, to be published in *Biometrika* **48** (1961).
- [23] Rao, C. R.: The utilization of multiple measurements in problems of biological classification, *J. Roy. Statist. Soc. Ser. B.* **10** (1948), 159-203.
- [24] Rao, C. R.: *Advanced Statistical Methods in Biometric Research*, New York, John Wiley and Sons, 1952.
- [25] Rao, C. R.: A general theory of discrimination when the information on alternative hypotheses is based on samples, *Ann. Math. Statist.* **25** (1954), 651-670.
- [26] Robbins, H.: Mixture of distributions, *Ann. Math. Statist.* **19** (1948), 360-369.
- [27] Robbins, H. and E. J. G. Pitman: Application of the method of mixtures to quadratic forms in normal variates, *Ann. Math. Statist.* **20** (1949), 552-560.
- [28] Roy, S. N.: *Some Aspects of Multivariate Analysis*, New York, John Wiley and Sons, 1957.
- [29] Satterthwaite, F. E.: An approximate distribution of estimates of variance components, *Biometrics* **2** (1946), 110-114.
- [30] Tanaka, M.: On the intraclass canonical correlation, *Osaka Tokei-danwakai Hokoku* **3** (1958), 93-113.

- [31] Wald, A.: On a statistical problem arising in the classification of an individual into one of two groups, *Ann. Math. Statist.* **15** (1944), 145-162.
- [32] Wallace, D. L.: Asymptotic approximations to distributions, *Ann. Math. Statist.* **29** (1958), 635-654.
- [33] Welch, B. L.: Note on discriminant functions, *Biometrika* **31** (1939), 218-220.

