# ESTIMATION ASSOCIATED WITH LINEAR DISCRIMINANTS[1]

By Seymour Geisser

State University of New York at Buffalo

**1. Introduction.** Suppose we have two $p$-variate normal populations where $\Pi_1$ is $N(\mu_1, \Sigma)$ and $\Pi_2$ is $N(\mu_2, \Sigma)$ and sample estimates $\bar{x}_1$ of $\mu_1$, $\bar{x}_2$ of $\mu_2$ based on $n_1$ and $n_2$ observations respectively. Further assume we have an independent estimate $S$ of the common covariance matrix $\Sigma$ based on $\nu$ degrees of freedom i.e., $\nu S$ is Wishart, $W(\Sigma, \nu)$ and $\nu = n_1 + n_2 - 2$ if all the information on $\Sigma$ is supplied by the two samples. Now the samples provide us with a linear discriminant, Anderson [1], page 138,

$$(1.1) \qquad V = [z - \tfrac{1}{2}(\bar{x}_1 + \bar{x}_2)]'S^{-1}(\bar{x}_1 - \bar{x}_2)$$

as an estimate of the population discriminant

$$(1.2) \qquad U = [z - \tfrac{1}{2}(\mu_1 + \mu_2)]'\Sigma^{-1}(\mu_1 - \mu_2)$$

where the new observation $z$ has prior probability $q_1$ of being from $\Pi_1$ and $q_2$ from $\Pi_2$, $q_1 + q_2 = 1$. For $r = q_2/q_1$,

$$(1.3) \qquad U > \log r \quad \text{assigns} \quad z \quad \text{to} \quad \Pi_1,$$
$$U < \log r \quad \text{assigns} \quad z \quad \text{to} \quad \Pi_2;$$

or if $\mu_1$, $\mu_2$ and $\Sigma$ are unknown

$$(1.4) \qquad V > \log r \quad \text{assigns} \quad z \quad \text{to} \quad \Pi_1,$$
$$V < \log r \quad \text{assigns} \quad z \quad \text{to} \quad \Pi_2.$$

There are then three questions that naturally arise here. Firstly, the estimation of the population discriminant $U$. Secondly the estimation of the true errors of misclassification, $\epsilon_1$ and $\epsilon_2$; and lastly, though most importantly, the estimation of the "index" errors of misclassification, $\beta_1$ and $\beta_2$ i.e., the errors incurred in using the sample discriminant $V$ on future observations. These problems will be investigated by the Bayes approach previously outlined by the author [4], [5], [6], [7]. In essence it is asserted that in the absence of any prior objective knowledge on $\mu_1$, $\mu_2$ and $\Sigma$ it is convenient to assume that the prior density for these parameters may be represented by

$$(1.5) \qquad g(\mu_1, \mu_2, \Sigma^{-1}) \propto |\Sigma|^{(p+1)/2}.$$

This leads to

$$(1.6) \quad P(\mu_1, \mu_2, \Sigma^{-1}) \propto |\Sigma^{-1}|^{(\nu-p+1)/2}$$
$$\cdot \exp\{-\tfrac{1}{2}\operatorname{tr}\Sigma^{-1}[\nu S + n_1(\bar{x}_1 - \mu_1)(\bar{x}_1 - \mu_1)' + n_2(\bar{x}_2 - \mu_2)(\bar{x}_2 - \mu_2)']\}$$

for the joint posterior densities of $\mu_1$, $\mu_2$ and $\Sigma^{-1}$ as the starting point for our investigation to the aforementioned estimation problems. What then will be offered is a method for assessing the errors involved in discriminatory problems for those investigators and statisticians who find linear discriminants useful and appealing but prefer a Bayesian orientation in their interpretation. As a secondary feature this method can then be compared with the complete Bayesian approach of [4], [5] (see Section 5) and with orthodox frequency techniques where available or as outlined in a paper by John [8]. At this point an exact comparison between the sampling theory methods and the approach advocated here is difficult because of the rather complicated nature of the distributions involved in the frequentist procedures. Suffice it to say that they will differ for finite sample sizes but tend to be close as the sample sizes increase.

**2. The distribution of $U$.** We now seek the posterior density of $U$ first under $z \varepsilon \Pi_1$ and then under $z \varepsilon \Pi_2$. From Anderson [1], page 135, it is clear that, conditional on $\mu_1$, $\mu_2$ and $\Sigma$, $U$ is $N(\frac{1}{2}\alpha, \alpha)$ under $\Pi_1$ and $N(-\frac{1}{2}\alpha, \alpha)$ under $\Pi_2$ with the "distance" between $\Pi_1$ and $\Pi_2$ given as

$$(2.1) \qquad \alpha = (\mu_1 - \mu_2)'\Sigma^{-1}(\mu_1 - \mu_2).$$

For fixed $\Sigma$, $c\alpha$ is non-central chi-square $\chi^2(p, 2\lambda)$ from (1.6), where $c = n_1 n_2/(n_1 + n_2)$ and

$$(2.2) \qquad \lambda = \tfrac{1}{2}c(\bar{x}_1 - \bar{x}_2)'\Sigma^{-1}(\bar{x}_1 - \bar{x}_2).$$

Thus for $\alpha > 0$

$$(2.3) \qquad P(\alpha \mid \lambda) = e^{-c\alpha/2 - \lambda} \sum_{j=0}^{\infty} \frac{\lambda^j c(c\alpha)^{\frac{1}{2}(p+2j)-1}}{j!2^{(p+2j)/2}\Gamma((p+2j)/2)}.$$

Further from (1.6) it is clear that $b\lambda$ is central $\chi^2$ with $\nu$ degrees of freedom where

$$(2.4) \qquad b = 2\nu/c(\bar{x}_1 - \bar{x}_2)'S^{-1}(\bar{x}_1 - \bar{x}_2)$$

so that

$$P(\lambda) = b(b\lambda)^{\nu/2-1}e^{-b\lambda/2}/2^{\nu/2}\Gamma(\nu/2).$$

Now

$$(2.5) \qquad P(U, \alpha, \lambda) = P(U \mid \alpha, \lambda)P(\alpha \mid \lambda)P(\lambda)$$

and

$$(2.6) \qquad P(U, \alpha) = P(U \mid \alpha)P(\alpha)$$

and from (2.3) and (2.4) we find

$$(2.7) \qquad P(\alpha) = \sum_{j=0}^{\infty} \frac{c(c\alpha)^{\frac{1}{2}(p+2j)-1}e^{-c\alpha/2}\Gamma(\nu/2+j)b^{\nu/2}}{j!2^{p/2}\Gamma((p+2j)/2)\Gamma(\nu/2)(2+b)^{\nu/2+j}},$$

which is incidentally the posterior density of the non-centrality parameter associated with Hotelling's $T^2$ statistic. It may also be expressed as

$$(2.8) \qquad P(c\alpha) = \sum_{j=0}^{\infty} w_j f(\chi^2 \mid p + 2j)$$

where the weights

$$(2.9) \qquad w_j = (b/(2 + b))^{\nu/2} (2/(2 + b))^j \binom{\nu/2+j-1}{j}$$

are the individual terms of a negative binomial density and $f(\chi^2 \mid p + 2j)$ represents the density of $\chi^2 = c\alpha$ with $p + 2j$ degrees of freedom. This is similar to the non-central $\chi^2$, which can be put into the same form except that the weights $w_j$ are the individual terms of a Poisson density. From (2.7) we find

$$(2.10) \qquad E(\alpha) = c^{-1}[p + 2\nu b^{-1}],$$

$$(2.11) \qquad \mathrm{Var}\,(\alpha) = 2c^{-2}[p + 4\nu(b + 1)b^{-2}]$$

and

$$(2.12) \qquad E[e^{it\alpha}] = (1 - 2c^{-1}it)^{(\nu-p)/2}(1 - 2c^{-1}b^{-1}(b + 2)it)^{-\nu/2}.$$

Since $U$, given $\alpha$, is $N(\tfrac{1}{2}\alpha, \alpha)$ under $\Pi_1$ and $N(-\tfrac{1}{2}\alpha, \alpha)$ under $\Pi_2$, we may find the density of $U$ in the following way:

$$(2.13) \qquad P(U \mid \Pi_i) = \int_0^{\infty} P(U \mid \alpha, \Pi_i)P(\alpha)\,d\alpha, \qquad i = 1, 2.$$

This density can be evaluated by reference to Table 98, page 143 of Bierens De Haan [2] but since it is somewhat complicated and not too important from the Bayesian viewpoint, it will not be given here. However we can easily find, without explicit evaluation of (2.13), that $E(U \mid \Pi_1) = \tfrac{1}{2}E(\alpha) = -E(U \mid \Pi_2)$ so that

$$(2.14) \qquad E(U \mid \Pi_1) = \tfrac{1}{2}p(n_1^{-1} + n_2^{-1}) + \tfrac{1}{2}(\bar{x}_1 - \bar{x}_2)'S^{-1}(\bar{x}_1 - \bar{x}_2) = -E(U \mid \Pi_2).$$

In either case, $V(U) = E(U^2) - E^2(U) = E(\alpha) + \tfrac{1}{4}\mathrm{Var}\,(\alpha)$, which results in

$$(2.15) \qquad \mathrm{Var}\,(U) = pc^{-1}[1 + (2c)^{-1}] + [1 + c^{-1} + (cb)^{-1}](\bar{x}_1 - \bar{x}_2)'S^{-1}(\bar{x}_1 - \bar{x}_2).$$

This yields $p(n_1^{-1} + n_2^{-1}) + (\bar{x}_1 - \bar{x}_2)'S^{-1}(\bar{x}_1 - \bar{x}_2)$ as an estimate of the "distance" between the two populations. More generally the characteristic functions of $U$ under $\Pi_1$ or $\Pi_2$ is also obtained without any trouble by integrating first with respect to $U$ and then $\alpha$ in the joint density of $U$ and $\alpha$. This yields

$$(2.16) \qquad E[e^{itU} \mid \Pi_i] = g_i^{-p/2}[((2 + b)g_i - 2)/bg_i]^{-\nu/2}, \qquad i = 1, 2,$$

where

$$(2.17) \qquad g_1 = 1 - c^{-1}it(1 + it); \qquad g_2 = 1 - c^{-1}it(it - 1).$$

It may be shown that from (2.16) and (2.17) that $U$ in either case is asymptotically normal for increasing $\nu$, when $\nu = n_1 + n_2 - 2$. Hence as a rough approximation in finding limits on $U$ we may use $X = [U - E(U)]/[\mathrm{Var}\,(U)]^{\frac{1}{2}}$ as $N(0, 1)$.

Thus far we have considered the posterior distribution of $U$ for random $z$. This of course does not provide us with an estimate of the "discriminant" $U$.

To find such an estimate requires that we fix $z$ and then find the posterior distribution of $U$. This can be accomplished by writing the algebraic equivalent of (1.2) which is

$$(2.18) \qquad U = \tfrac{1}{2}(z - \mu_2)'\Sigma^{-1}(z - \mu_2) - \tfrac{1}{2}(z - \mu_1)'\Sigma^{-1}(z - \mu_1)$$

and noticing from (1.6) that for fixed $z$ the two terms on the right are each distributed as a constant times non-central $\chi^2$ conditional on $\Sigma^{-1}$. In other words $U$, for fixed $z$, is a linear combination of non-central $\chi^2$ variables conditional on $\Sigma^{-1}$. One then could combine this with the marginal distribution of $\Sigma^{-1}$ and then integrate out $\Sigma^{-1}$ to find the density of $U$ for fixed $z$. This is somewhat complicated and since we are mainly interested in this density only as it provides us with a sample estimate of $U$, we shall not attempt its explicit evaluation. However we shall return to this point at the end of Section 4.

**3. Estimation of the "true" misclassification errors.** Now for fixed $\mu_1$, $\mu_2$ and $\Sigma^{-1}$

$$(3.1) \quad \Pr\left[U < \log r \mid \mu_1, \mu_2, \Sigma^{-1}; z \,\varepsilon\, \Pi_1\right] = \epsilon_1 = \int_{-\infty}^{\tau_1} \phi(v)\, dv = \Phi(\tau_1),$$

$$(3.2) \quad \Pr\left[U > \log r \mid \mu_1, \mu_2, \Sigma^{-1}; z \,\varepsilon\, \Pi_2\right] = \epsilon_2 = \int_{\tau_2}^{\infty} \phi(v)\, dv = 1 - \Phi(\tau_2)$$

where $\phi(v) = (2\pi)^{-\frac{1}{2}}e^{-\frac{1}{2}v^2}$ is the standard normal density and

$$(3.3) \qquad \tau_1 = (\log r - \tfrac{1}{2}\alpha)/\alpha^{\frac{1}{2}}, \qquad \tau_2 = (\log r + \tfrac{1}{2}\alpha)/\alpha^{\frac{1}{2}}.$$

In other words we define $\epsilon_1$ and $\epsilon_2$ to be functions of the random variable $\alpha$ whose distribution is conditioned on the sample values $\bar{x}_1$, $\bar{x}_2$ and $S$.

Now for $\log r \geqq 0$ (we may choose $q_2 \geqq q_1$ and label $\Pi_1$ and $\Pi_2$ accordingly) we note that $\tau_1$ is a monotone decreasing function of $\alpha$ and is a monotone increasing function $\epsilon_1$. Hence for $\log r > 0$

$$(3.4) \qquad \Pr\left[\alpha_1 < \alpha < \alpha_2\right] = 1 - 2p$$

is equivalent

$$(3.5) \qquad \Pr\left[\tau_1(\alpha_2) < \tau_1 < \tau_1(\alpha_1)\right] = 1 - 2p$$

and consequently

$$(3.6) \qquad \Pr\left\{\Phi[\tau_1(\alpha_2)] < \epsilon_1 < \Phi[\tau_1(\alpha_1)]\right\} = 1 - 2p.$$

Hence limits on $\epsilon_1$ can be obtained from the posterior density of $\alpha$. For $\log r > 0$, $\tau_2$ decreases monotonically with $\alpha$ until $\alpha = 2 \log r$ where $\tau_2$ attains its minimum value, $(2 \log r)^{\frac{1}{2}}$, then $\tau_2$ increases monotonically as $\alpha$ increases. Therefore limits on $\tau_2$ can be computed from $\alpha$ in the following way: Let

$$(3.7) \qquad \Pr\left[\tau_{21} < \tau_2 < \tau_{22}\right] = 1 - 2p$$

for $(2 \log r)^{\frac{1}{2}} < \tau_{21} < \tau_{22}$, then

$$(3.8) \qquad \Pr\left[\alpha_0 < \alpha < \alpha_1\right] + \Pr\left[\alpha_2 < \alpha < \alpha_3\right] = 1 - 2p$$

where $\alpha_0 < \alpha_1 < 2 \log r < \alpha_2 < \alpha_3$, i.e.,

$$(3.9) \quad \begin{aligned} \alpha_0 &= 2(\tau_{22}^2 - \log r) - 2\tau_{22}(\tau_{22}^2 - 2 \log r)^{\frac{1}{2}}, \\ \alpha_1 &= 2(\tau_{21}^2 - \log r) - 2\tau_{21}(\tau_{21}^2 - 2 \log r)^{\frac{1}{2}}, \\ \alpha_2 &= 2(\tau_{21}^2 - \log r) + 2\tau_{21}(\tau_{21}^2 - 2 \log r)^{\frac{1}{2}}, \\ \alpha_3 &= 2(\tau_{22}^2 - \log r) + 2\tau_{22}(\tau_{22}^2 - 2 \log r)^{\frac{1}{2}}. \end{aligned}$$

Further since $\epsilon_2$ is a monotone decreasing function of $\tau_2$ we have

$$(3.10) \qquad \Pr\,[1 - \Phi(\tau_{22}) < \epsilon_2 < 1 - \Phi(\tau_{21})] = 1 - 2p.$$

As a rough approximation to the distribution of $\alpha$ given by (2.7) we may equate the first two moments of $\alpha$ to a constant times a chi-square distribution with $d$ degree of freedom resulting in

$$(3.11) \qquad (n_1 n_2/(n_1 + n_2))(p + 2\nu/b)(p + 4\nu(b + 1)/b^2)^{-1}\alpha \sim \chi_d^2$$

where

$$(3.12) \qquad d = (p + 2\nu/b)^2/[p + 4\nu(b + 1)/b^2].$$

For $q_2 = q_1$, so that $\log r = 0$, we note that $-\tau_1 = \tau_2 = \frac{1}{2}\alpha^{\frac{1}{2}}$ with $0 < \tau_2 < \infty$. Hence we obtain identical limits for $\epsilon_1$ and $\epsilon_2$ in this special case.

Confidence limits on $\epsilon_i$ may be obtained depending on confidence limits on $\alpha$ through the sampling distribution of $c(\bar{x}_1 - \bar{x}_2)'S^{-1}(\bar{x}_1 - \bar{x}_2) = cQ$ which is $\nu p(\nu - p - 1)^{-1}F(p, \nu - p - 1)$ where $F$ is non-central $F$ with non-centrality parameter $c\alpha$. As $\nu \to \infty$, the sampling distribution of $cQ$ tends to non-central $\chi^2(p)$ with non-centrality parameter $c\alpha$, thus limits on $c\alpha$ are obtained from the confidence inversion of the non-central $\chi^2(p)$. However as $\nu \to \infty$ the posterior distribution of $c\alpha$ tends to non-central $\chi^2(p)$ with non-centrality parameter $cQ$, thus limits on $c\alpha$ are obtained directly from the non-central $\chi^2(p)$ distribution. Hence it is clear that the confidence limits and the posterior limits will differ somewhat.

**4. The "index" errors of misclassification.** We now turn to the problem which is paramount from the practical point of view and that is the estimation of the index misclassification errors when using the sample or "index" discriminant, $V$, on future observations. In actual practice, an investigator is often interested in how well his particular index discriminant will do on future observations since this sample is what he must work with. He is probably more interested in these index errors of misclassification, $\beta_1$ and $\beta_2$, than in the "true" errors, $\epsilon_1$ and $\epsilon_2$. Now for the fixed values $\bar{x}_1$, $\bar{x}_2$ and $S$

$$(4.1) \quad \Pr\,(V < \log r \mid \mu_1, \mu_2, \Sigma; z \,\varepsilon\, \Pi_1) = \beta_1 = \int_{-\infty}^{\theta_1} \phi(v)\,dv = \Phi(\theta_1),$$

$$(4.2) \quad \Pr\,(V > \log r \mid u_1, \mu_2, \Sigma; z \,\varepsilon\, \Pi_2) = \beta_2 = \int_{\theta_2}^{\infty} \phi(v)\,dv = 1 - \Phi(\theta_2),$$

where

(4.3)   $\theta_1 = \{[\frac{1}{2}(\bar{x}_1 + \bar{x}_2) - \mu_1]'S^{-1}(\bar{x}_1 - \bar{x}_2) + \log r\}$

$$\cdot [(\bar{x}_1 - \bar{x}_2)'S^{-1}\Sigma S^{-1}(\bar{x}_1 - \bar{x}_2)]^{-\frac{1}{2}},$$

(4.4)   $\theta_2 = \{[\frac{1}{2}(\bar{x}_1 + \bar{x}_2) - \mu_2]'S^{-1}(\bar{x}_1 - \bar{x}_2) + \log r\}$

$$\cdot [(\bar{x}_1 - \bar{x}_2)'S^{-1}\Sigma S^{-1}(\bar{x}_1 - \bar{x}_2)]^{-\frac{1}{2}}$$

and $\theta_1$ and $\theta_2$ are random variables that are functions of $\mu_1$, $\mu_2$ and $\Sigma$. Hence we have defined $\beta_1$ and $\beta_2$ as functions of the random variables $\mu_1$, $\mu_2$, $\Sigma$ for fixed values of $\bar{x}_1$, $\bar{x}_2$ and $S$ which differs from the sampling interpretation where $\beta_1$ and $\beta_2$ are considered either as functions of the fixed parameters $\mu_1$, $\mu_2$, $\Sigma$ obtained from the unconditional sampling distribution of $V$ in terms of the random variables $\bar{x}_1$, $\bar{x}_2$, and $S$, or defined as functions of the random variables $\bar{x}_1$, $\bar{x}_2$, $S$, in particular see John [8].

First we note that (4.1) and (4.2) follow from the fact that $V$ is normal, conditional on $\mu_1$, $\mu_2$ and $\Sigma$ for the fixed values of $\bar{x}_1$, $\bar{x}_2$ and $S$. Further from (1.6) $\theta_1$ and $\theta_2$ are normally distributed conditional on $\Sigma$ with

(4.5)          $E(\theta_1 \mid \Sigma) = [\log r - \frac{1}{2}(\bar{x}_1 - \bar{x}_2)'a]/[a'\Sigma a]^{\frac{1}{2}},$

$$E(\theta_2 \mid \Sigma) = [\log r + \frac{1}{2}(\bar{x}_1 - \bar{x}_2)'a]/[a'\Sigma a]^{\frac{1}{2}};$$

(4.6)          $\text{Var}(\theta_1 \mid \Sigma) = n_1^{-1}, \qquad \text{Var}(\theta_2 \mid \Sigma) = n_2^{-1}$

where $a = S^{-1}(\bar{x}_1 - \bar{x}_2)$. Moreover it is easy to show that $y = \nu a'Sa/a'\Sigma a$ is distributed as $\chi^2$ with $\nu - p + 1$ degrees of freedom. Hence the joint density of $\theta_1$, $\theta_2$ and $y$ is

(4.7)   $P(\theta_1, \theta_2, y) = (n_1 n_2)^{\frac{1}{2}}(2\pi)^{-1} \exp\{-\frac{1}{2}[n_1(\theta_1 - t_1 y^{\frac{1}{2}})^2 + n_2(\theta_2 - t_2 y^{\frac{1}{2}})^2]\}$

$$\cdot e^{-y/2} y^{(\nu-p-1)/2} / 2^{(\nu-p+1)/2} \Gamma((\nu - p + 1)/2)$$

where

(4.8)          $t_1 = [\log r - \frac{1}{2}(\bar{x}_1 - \bar{x}_2)'a]/(\nu a'Sa)^{\frac{1}{2}} = \nu^{-\frac{1}{2}} k_1,$

$$t_2 = [\log r + \frac{1}{2}(\bar{x}_1 - \bar{x}_2)'a]/(\nu a'Sa)^{\frac{1}{2}} = \nu^{-\frac{1}{2}} k_2.$$

We then integrate out $y$ in (4.7) and obtain

(4.9)   $P(\theta_1, \theta_2) = (n_1 n_2)^{\frac{1}{2}} \exp[-\frac{1}{2}(n_1\theta_1^2 + n_2\theta_2^2)]/$

$$2\pi\Gamma((\nu - p + 1)/2)[1 + t_1^2 n_1 + t_2^2 n_2]^{(\nu-p+1)/2}$$

$$\times \sum_{j=0}^{\infty} [(t_1 n_1 \theta_1 + t_2 n_2 \theta_2)2^{\frac{1}{2}}(1 + t_1^2 n_1 + t_2^2 n_2)^{-\frac{1}{2}}]^j$$

$$\cdot \Gamma((\nu + j - p + 1)/2)/j!.$$

The marginal density for $\theta_i$, $i = 1, 2$, is

(4.10)   $P(\theta_i) = (n_i/2\pi)^{\frac{1}{2}} \exp[\frac{1}{2}n_i\theta_i^2][\Gamma((\nu - p + 1)/2)(1 + n_i t_i^2)^{(\nu-p+1)/2}]^{-1}$

$$\cdot \sum_{j=0}^{\infty} [t_i n_i \theta_i 2^{\frac{1}{2}}(1 + n_i t_i^2)^{-\frac{1}{2}}]^j \Gamma((\nu + j - p + 1)/2)/j!.$$

We may also obtain most conveniently by conditional and unconditional expectations the following:

$$(4.11) \qquad E(\theta_1) = k_1 m, \qquad E(\theta_2) = k_2 m$$

where

$$(4.12) \quad m = \Gamma((\nu - p + 2)/2)(2/\nu)^{\frac{1}{2}}/\Gamma((\nu - p + 1)/2) \to ((\nu - p + \tfrac{1}{2})/\nu)^{\frac{1}{2}}$$
$$\to 1 \quad \text{as} \quad \nu \to \infty$$

and

$$(4.13) \quad \text{Var}(\theta_1) = n_1^{-1} + k_1^2 f, \quad \text{Var}(\theta_2) = n_2^{-1} + k_2^2 f, \quad \text{cov}(\theta_1, \theta_2) = k_1 k_2 f$$

where

$$(4.14) \quad f = \nu^{-1}[\nu - p + 1 - 2\Gamma^2((\nu - p + 2)/2)/\Gamma^2((\nu - p + 1)/2)]$$
$$\to 1/2\nu \to 0 \quad \text{as} \quad \nu \to \infty.$$

Marginal limits on $\theta_1$ and $\theta_2$ can readily be converted from

$$(4.15) \qquad \text{Pr}\,[\theta_{11} < \theta_1 < \theta_{12}] = 1 - 2p$$

to

$$(4.16) \qquad \text{Pr}\,[\Phi(\theta_{11}) < \beta_1 < \Phi(\theta_{12})] = 1 - 2p$$

and from

$$(4.17) \qquad \text{Pr}\,[\theta_{21} < \theta_2 < \theta_{22}] = 1 - 2p$$

to

$$(4.18) \qquad \text{Pr}\,[1 - \Phi(\theta_{22}) < \beta_2 < 1 - \Phi(\theta_{21})] = 1 - 2p$$

or

$$(4.19) \qquad \text{Pr}\,[\Phi(-\theta_{22}) < \beta_2 < \Phi(-\theta_{21})] = 1 - 2p$$

due to the monotone character of the transformations from $\theta_1$ and $\theta_2$ to $\beta_1$ and $\beta_2$.

Further the simultaneous probability

$$(4.20) \quad \text{Pr}\,[\Phi(\theta_{11}) < \beta_1 < \Phi(\theta_{12}); 1 - \Phi(\theta_{22}) < \beta_2 < 1 - \Phi(\theta_{21})] = P$$

is equivalent to

$$(4.21) \qquad \text{Pr}\,[\theta_{11} < \theta_1 < \theta_{12}; \theta_{21} < \theta_2 < \theta_{22}] = P.$$

The joint characteristic function of $\theta_1$ and $\theta_2$ is calculated to be

$$(4.22) \quad E[\exp(iu_1\theta_1 + iu_2\theta_2)] = \exp[-\tfrac{1}{2}(n_1^{-1}u_1^2 + n_2^{-1}u_2^2)]$$
$$\cdot \sum_{j=0}^{\infty} [i(t_1 u_1 + t_2 u_2)2^{\frac{1}{2}}]^j j!^{-1}$$
$$\cdot \Gamma((\nu + j - p + 1)/2)/\Gamma((\nu - p + 1)/2)$$

with the marginal characteristic function of $\theta_i$

(4.23) $\quad E[\exp(i\theta_i u_i)] = \exp(-u_i^2/2n_i) \sum_{j=0}^{\infty} [it_i u_i 2^{\frac{1}{2}}]^j j!^{-1}$

$$\cdot \Gamma((\nu - p + 1 + j)/2)/\Gamma((\nu - p + 1)/2).$$

From (4.22) and (4.23) it can be shown that the joint distribution of

(4.24) $\qquad\qquad v_i = (\theta_i - k_i m)(n_i^{-1} + k_i^2 f)^{-\frac{1}{2}}, \qquad\qquad i = 1, 2,$

tends to the bivariate normal distribution $N(0, 0, 1, 1, \rho)$ where

(4.25) $\qquad\qquad \rho = k_1 k_2 f/(n_1^{-1} + k_1^2 f)^{\frac{1}{2}}(n_2^{-1} + k_2^2 f)^{\frac{1}{2}}$

for increasing $\nu$. In particular the marginal distribution of $v_i$ tends to the $N(0, 1)$ distribution. Hence approximate limits on $\beta_1$ and $\beta_2$ are given by

(4.26) $\quad \Pr\{\Phi[k_1 m - (n_1^{-1} + k_1^2 f)^{\frac{1}{2}} v_p] < \beta_1 < \Phi[k_1 m + (n_1^{-1} + k_1^2 f)^{\frac{1}{2}} v_p]\}$

$$\cong 1 - 2p,$$

(4.27) $\quad \Pr\{\Phi[-k_2 m - (n_2^{-1} k_2^2 f)^{\frac{1}{2}} v_p] < \beta_2 < \Phi[-k_2 m + (n_2^{-1} + k_2^2 f)^{\frac{1}{2}} v_p]\}$

$$\cong 1 - 2p$$

where

(4.28) $\qquad\qquad\qquad\qquad p = \int_{v_p}^{\infty} \phi(v)\, dv.$

It is to be noted that the normal approximation given here for the posterior density of $\beta_i$ differs slightly from the normal approximation to the sampling distribution of this quantity given by John [8] in formula (97). The basic difference being in the variance. If we are interested in the probability for the joint rectangular region $R$ specified in (4.26) and (4.27) then

(4.29) $\quad \Pr([\beta_1, \beta_2]\ \varepsilon\ R] \cong 2[L(v_p, v_p; \rho) + L(v_p, v_p; -\rho) + 1 - 2p] - 1$

where

(4.30) $\qquad\qquad L(v_p, v_p; \rho) = \int_{v_p}^{\infty} \int_{v_p}^{\infty} \phi(x, y \mid 0, 0, 1, 1, \rho)\, dx\, dy$

and here $\phi(x, y \mid 0, 0, 1, 1, \rho)$ is the bivariate normal density with means $(0, 0)$, variances $(1, 1)$ and correlation $\rho$. The evaluation of $L(v_p, v_p; \rho)$ can be obtained from tables, or very conveniently from the charts of Zelen and Severo [10].

An inspection of the marginal characteristic function of $\theta_i$ (4.23) reveals that $\theta_i$ is distributed as the sum of two independent variables $X$ and $k_i Y$ where $X$ is $N(0, n_i^{-1})$ and $Y$ is distributed as $[\nu^{-1} \chi_{\nu-p+1}^2]^{\frac{1}{2}}$. To have an idea of the usefulness of the normal approximations (4.26), (4.27) and (4.29) we recall that $(2\chi_d^2)^{\frac{1}{2}} - (2d - 1)^{\frac{1}{2}}$ is well approximated by the standard normal distribution for a moderate number of degrees of freedom $d$. Therefore $\theta_i$, which is the sum of a normal component $X$, and $Y$ which tends to $N([k_i^2(\nu - p + \frac{1}{2})/\nu]^{\frac{1}{2}}, k_i^2/2\nu)$ reasonably quickly as $\nu$ increases, should tend even more rapidly to $N(k_i[(\nu - p + \frac{1}{2})/\nu]^{\frac{1}{2}}, 1/n_i + k_i^2/2\nu)$ then does the approximation $(2\chi_d^2)^{\frac{1}{2}} -$

$(2d - 1)^{\frac{1}{2}}$ to $N(0, 1)$. Further the approximation suggested (4.26) and (4.27) uses the exact means and variances rather than the asymptotic ones which should perhaps even enhance the value of the approximation. Generally $\nu = n_1 + n_2 - 2$ if the estimate of the covariance matrix is computed only from the sample (otherwise $\nu$ will be even greater) so that the approximation should be fairly accurate for $\nu - p + 1 \geqq 20$. The same reasoning will hold in the bivariate case because from an inspection of the joint characteristic function (4.22) we determine that we can represent $\theta_1 = X_1 + k_1 Y$, $\theta_2 = X_2 + k_2 Y$ where $X_1$, $X_2$ and $Y$ are mutually independent and are $N(0, n_1^{-1})$, $N(0, n_2^{-1})$ and $[\nu^{-1}\chi^2_{\nu-p+1}]^{\frac{1}{2}}$ variables, respectively.

We note for $p = 1$, the univariate case, that $n_i^{\frac{1}{2}}\theta_i = n_i^{\frac{1}{2}}X_i + n_i^{\frac{1}{2}}k_i Y$ has a posterior density which is the fiducial density with a suitable change in notation obtained by Fisher [3], page 123, formula 91. Although he cast his problem in a different framework it can readily be transformed into the discriminatory setting thus implying that the limits on $\beta_i$ are fiducial limits in the sense of Fisher, at least in the univariate case.

It is to be noted that $\beta_1$ and $\beta_2$, the conditional probabilities of misclassification, are functions of the random variables $\mu_1$, $\mu_2$ and $\Sigma^{-1}$ within the previous framework. What has been presented are probability limits on these random variables. The unconditional or posterior predictive probabilities of misclassification

$$(4.31) \qquad \Pr [V < \log r \mid z \,\varepsilon\, \Pi_1]$$

$$(4.32) \qquad \Pr [V > \log r \mid z \,\varepsilon\, \Pi_2]$$

may also be obtained. For example

$$(4.33) \quad E(\beta_1) = \int \Pr [V < \log r \mid \mu_1, \mu_2; \Sigma, z \,\varepsilon\, \Pi_1] P(\mu_1, \mu_2, \Sigma^{-1}) \, d\mu_1 \, d\mu_2 \, d\Sigma^{-1}$$

$$= \int_{-\infty}^{\log r} \int f(V \mid \mu_1, \mu_2, \Sigma; z \,\varepsilon\, \Pi_1) P(\mu_1, \mu_2, \Sigma^{-1}) \, d\mu_1 \, d\mu_2 \, d\Sigma^{-1} dV$$

where $f(V \mid \mu_1, \mu_2, \Sigma; z \,\varepsilon\, \Pi_1)$ represents the conditional density of $V$. Hence

$$(4.34) \qquad E(\beta_1) = \int_{-\infty}^{\log r} f(V \mid z \,\varepsilon\, \Pi_1) \, dV = \Pr [V < \log r \mid z \,\varepsilon\, \Pi_1]$$

where $f(V \mid z \,\varepsilon\, \Pi_1)$ represents the unconditional or predictive density of $V$. The evaluation of (4.34) is accomplished by noting that for $Q = (\bar{x}_1 - \bar{x}_2)' \cdot S^{-1}(\bar{x}_1 - \bar{x}_2)$,

$$(4.35) \quad [\nu(n_1 + 1)Q/n_1]^{-\frac{1}{2}}V = [(\nu(n_1 + 1)/n_1)Q]^{-\frac{1}{2}}(z - \bar{x}_1)'S^{-1}(\bar{x}_1 - \bar{x}_2)$$
$$+ \tfrac{1}{2}(n_1 Q/\nu(n_1 + 1))^{\frac{1}{2}}.$$

It has previously been shown, [6] or [7], that the predictive distribution of the first term on the right hand side of (4.35) involving $z$ is $(\nu + 1 - p)^{-\frac{1}{2}}t_{\nu+1-p}$ where $t_{\nu+1-p}$ is the $t$ distribution with $\nu + 1 - p$ degrees of freedom. Therefore the predictive density of $V$ is

$$(4.36) \qquad [\nu(n_1 + 1)Q/n_1(\nu + 1 - p)]^{\frac{1}{2}}t_{\nu+1-p} + \tfrac{1}{2}Q.$$

Thus we obtain

$$(4.37) \quad E(\beta_1) = \Pr\left[V < \log r \mid z \; \varepsilon \; \Pi_1\right]$$

$$= \Pr\left[t_{\nu+1-p} < (\log r - \tfrac{1}{2}Q)[\nu(n_1 + 1)Q/(\nu + 1 - p)n_1]^{-\frac{1}{2}}\right]$$

which may be evaluated directly from tables of the $t$-distribution. Similarly

$$(4.38) \quad E(\beta_2) = \Pr\left[V > \log r \mid z \; \varepsilon \; \Pi_2\right]$$

$$= \Pr\left[t_{\nu+1-p} > (\log r + \tfrac{1}{2}Q)[\nu(n_2 + 1)Q/(\nu + 1 - p)n_2]^{-\frac{1}{2}}\right].$$

These expectations, the unconditional or predictive misclassification errors, are in a sense the optimum point estimates of the conditional misclassification errors. From one point of view the limits on $\beta_1$ and $\beta_2$ provide a more comprehensive guide to errors incurred in using the linear discriminant $V$ than their expectations, although the latter are also extremely valuable as "optimum" point estimates.

It is to be noted that we so far have made no attempt to justify the use of $V$ as the appropriate discriminant other than as Anderson states, "It seems intuitively reasonable." In support of this we may further add that for any fixed $z$ it is easily shown that the posterior mean of $U$, which minimizes the squared-error loss function, is

$$(4.39) \qquad\qquad E[U \mid z] = \tfrac{1}{2}p(n_2^{-1} - n_1^{-1}) + V.$$

This is a consequence of the fact that the expectation of $U$ is the sum of the expectations of the right hand side of (2.18). These can be obtained conditional on $\Sigma^{-1}$ and then unconditionally by integrating over the posterior distribution of $\Sigma^{-1}$, thus avoiding the exact evaluation of the density of $U$. Now if we wish to incorporate this bias $\tfrac{1}{2}p(n_2^{-1} - n_1^{-1})$ into our analysis this is easily accomplished by substituting throughout this section $\log r - \tfrac{1}{2}p(n_2^{-1} - n_1^{-1})$ for $\log r$.

**5. Some remarks.** It is suggested that use be made of the results in this paper in the following way. First the investigator, for some values $n_1$ and $n_2$ already chosen, should compute limits on the true errors $\epsilon_1$ and $\epsilon_2$ to assess the discriminatory power of the variables under consideration. If he finds the discriminatory power unsatisfactory, he may wish to include more variables until he is satisfied that the true errors are small enough for his purposes. Then he calculates his estimate of the index errors—if his estimated index errors are much larger than his estimated true errors, he might, if possible, collect more data on the two populations until he can drive his estimated index errors down to what he thinks is close enough to his estimated true errors for his purposes. If new classifiable data is unavailable, he at least knows how well his sample discriminant will do and what its possibilities for improvement are, given the variables he has to work with.

We note here that the Bayesian approach taken in this paper is somewhat different from the one previously taken in [5] and [6]. There we presented what we conceive of as a theory for a complete Bayesian approach to the problem of

classifying new observations when the type and number of populations is known with the number being exhaustive and samples are available on each of them. Classification of the observation depended on the product of the prior probability and the predictive density of $z$ under the various and exhaustive possibilities. The predictive discriminant therein obtained from the ratio of the predictive densities was linear only when $n_1 = n_2$, in particular see [5]. In addition, misclassification errors were defined there in terms of the predictive densities and here in terms of the linear discriminants. It turns out that $E(\beta_1)$ and $E(\beta_2)$ are exactly the predictive errors of misclassification if $n_1 = n_2$ for the complete Bayesian approach. For $n_1 \neq n_2$, the predictive discriminant is non-linear and the regions over which the predictive misclassification errors are computed are consequently more complicated. However the differential error will diminish as the difference between the sample means increases and as the sample sizes increase.

The Bayesian (or perhaps better termed semi-Bayesian) approach presented in this paper adheres more closely to the orthodox discriminatory procedure due to its emphasis on linearity and parametric estimation. This semi-Bayesian approach is somewhat akin in spirit to the previous presentation of a pseudo-posterior density for the correlation coefficient in [7] and to Pratt's incomplete Bayesian approach [9].

Finally, there is little doubt that for large sample sizes, the two Bayesian classification methods will yield substantially equivalent results. Moreover, linearity has much to recommend it due to its intrinsic simplicity and important interpretive uses. Work is also currently in progress on the problem of quadratic discriminants via this semi-Bayesian approach, i.e., for $\Sigma_1 \neq \Sigma_2$.

## REFERENCES

[1] ANDERSON, T. W. (1958). *An Introduction to Multivariant Statistical Analysis*. Wiley, New York.

[2] BIERENS, DE HAAN, D. (1957). *Nouvelles Tables D'Intégrales Définies*. Hafner, New York.

[3] FISHER, R. A. (1959). *Statisical Methods and Scientific Inference*, 2nd edition. Oliver and Boyd, Edinburgh.

[4] GEISSER, S. (1965). Bayesian estimation in multivariate analysis. *Ann. Math. Statist.* **36** 150–159.

[5] GEISSER, S. (1964). Posterior odds for multivariate normal classification. *J. Roy. Statist. Soc. Ser. B* **26** 69–76.

[6] GEISSER, S. (1966). Predictive discrimination. *Proceedings of the International Symposium on Multivariate Analysis*, 149–163. Academic Press, New York.

[7] GEISSER, S. and CORNFIELD, J. (1963). Posterior distributions for multivariate normal parameters. *J. Roy. Statist. Soc. Ser. B* **25** 368–376.

[8] JOHN, S. (1961). Errors in discrimination. *Ann. Math. Statist.* **32** 1125–1144.

[9] PRATT, J. W. (1965). Bayesian interpretation of standard inference statements. *J. Roy. Statist. Soc. Ser B* **27** 169–203.

[10] ZELEN, M. and SEVERO, N. (1960). Graphs for bivariate normal probabilities. *Ann. Math. Statist.* **31** 619–624.