

# AN APPLICATION OF INFORMATION THEORY TO MULTIVARIATE ANALYSIS

BY S. KULLBACK

*The George Washington University*

**0. Summary.** The problem considered is that of finding the “best” linear function for discriminating between two multivariate normal populations,  $\pi_1$  and  $\pi_2$ , without limitation to the case of equal covariance matrices. The “best” linear function is found by maximizing the divergence,  $J'(1, 2)$ , between the distributions of the linear function. Comparison with the divergence,  $J(1, 2)$ , between  $\pi_1$  and  $\pi_2$  offers a measure of the discriminating efficiency of the linear function, since  $J(1, 2) \geq J'(1, 2)$ . The divergence, a special case of which is Mahalanobis’s Generalized Distance, is defined in terms of a measure of information which is essentially that of Shannon and Wiener. Appropriate assumptions about  $\pi_1$  and  $\pi_2$  lead to discriminant analysis (Sections 4, 7), principal components (Section 5), and canonical correlations (Section 6).

**1. Introduction.** The following extract from Section (4), “Scientific Method,” of Cherry [15] is pertinent. “. . . the idea of information has existed in early times and has gradually entered into a great variety of sciences, to a certain extent integrating them together. Nowadays the concept of information would seem to be essential to all research workers, and as universal and fundamental as the concepts of energy or entropy. Speaking most generally, every time we make any observation, or perform any ‘experiment’, we are seeking for information; the question thus arises: How much can we know from a particular set of observations or experiments? The modern mathematical work, at which we have glanced, seeks to answer in precise terms this very question which, in its origin, is an epistemological one. But first a word of caution: the term ‘information’ has been used by different authors as having different meanings. . . . The information supplied by an experiment may perhaps be thought of as a ratio of a *posteriori* to the *a priori* probabilities (strictly, the logarithm of this ratio).”

R. A. Fisher’s measure of information (intrinsic accuracy) was introduced to compare the merits of different estimates. Shannon and Wiener’s measure of information was introduced to define and measure that which is being conveyed by a communication system, the latter considered as a stochastic process. (See references in [9], [15].)

The author and Leibler, in [9], generalized Shannon and Wiener’s definition to the abstract case and showed that it and Fisher’s definition are not unrelated. Properties of a measure of divergence between statistical populations, defined in terms of the measure of information, were also derived in [9].

Other approaches to a definition of the distance or divergence between two populations, and the applications of such a concept, have been made by Mahalanobis [10], Bhattacharyya [16], [17], and Rao [18].

The measure of divergence and its properties, as derived in [9], may be applied in particular to the problem of discrimination between certain multivariate normal populations by means of linear functions. We do not limit ourselves to the case of equal covariance matrices. Although no essentially new results are derived in the following, it is believed that the methods and underlying uniformity of approach may be of pedagogical interest.<sup>1</sup>

Matrix notation, methods and results are used and assumed known to the reader. For the purpose of this paper we limit ourselves to a discussion using population parameters and do not consider problems of estimation or distribution. Attention is invited to Bartlett [2], Brown [3], Cochran and Bliss [4], Kendall [8], Penrose [11], Smith [12], Tintner [13] and Wilks [14] for discussions of related problems and additional references to the literature.

**2. Divergence.**

(a) *Definition.* If two multivariate normal populations  $\pi_1$  and  $\pi_2$  have the respective probability densities  $f_i(x_1, x_2, \dots, x_k)$ ,  $i = 1, 2$ , then the divergence between  $\pi_1$  and  $\pi_2$  is defined by [9]

$$(2.1) \quad J(1, 2) = \int (f_1(x_1, \dots, x_k) - f_2(x_1, \dots, x_k)) \log \frac{f_1(x_1, \dots, x_k)}{f_2(x_1, \dots, x_k)} dx_1 \dots dx_k.$$

The mean information for discrimination between  $\pi_1$  and  $\pi_2$  per observation from  $\pi_1$  is defined by [9]

$$(2.2) \quad \begin{aligned} I(1:2) &= \int f_1(x_1, \dots, x_k) \log \frac{f_1(x_1, \dots, x_k)}{f_2(x_1, \dots, x_k)} dx_1 \dots dx_k \\ &= \int f_1(x_1, \dots, x_k) \log \frac{P(\pi_1 | x_1, \dots, x_k)}{P(\pi_2 | x_1, \dots, x_k)} dx_1 \dots dx_k - \log \frac{P(\pi_1)}{P(\pi_2)}, \end{aligned}$$

where  $P(\pi_i)$  and  $P(\pi_i | x_1, \dots, x_k)$  are respectively the a priori and a posteriori probabilities for  $\pi_i$ ,  $i = 1, 2$ , and a corresponding definition for  $I(2:1)$ . It is seen that  $J(1, 2) = I(1:2) + I(2:1)$ .

If

$$(2.3) \quad y_\alpha = y_\alpha(x_1, x_2, \dots, x_k), \quad \alpha = 1, 2, \dots, r, r \leq k,$$

are functions of the random variables  $x_1, x_2, \dots, x_k$ , such that the distribution of the  $y$ 's is given by the probability density function

$$(2.4) \quad g_i(y_1, y_2, \dots, y_r), \quad i = 1, 2,$$

---

<sup>1</sup>This approach has been found helpful in presenting certain aspects of multivariate analysis to a class at the George Washington University.

according as the  $x$ 's come from  $\pi_1$  or  $\pi_2$ , then the divergence between the population of  $y$ 's is defined by [9]

$$(2.5) \quad J'(1, 2) = \int (g_1(y_1, \dots, y_r) - g_2(y_1, \dots, y_r)) \log \frac{g_1(y_1, \dots, y_r)}{g_2(y_1, \dots, y_r)} dy_1 \cdots dy_r,$$

and the mean information for discrimination between  $\pi_1$  and  $\pi_2$  per observation from  $g_1(y_1, \dots, y_r)$  is defined by [9]

$$(2.6) \quad I'(1:2) = \int g_1(y_1, \dots, y_r) \log \frac{g_1(y_1, \dots, y_r)}{g_2(y_1, \dots, y_r)} dy_1 \cdots dy_r,$$

and a corresponding definition for  $I'(2:1)$ . It is seen that  $J'(1, 2) = I'(1:2) + I'(2:1)$ .

(b) *Properties.* The following properties of  $I$ ,  $I'$ ,  $J$  and  $J'$  will be utilized. (For proofs see [9].)

- i  $I(1:2) \geq 0$ ;  $J(1, 2) \geq 0$ , with equality if and only if  $f_1 = f_2$  a.e.;
- ii  $I, J$  are additive for independent random variables;
- iii  $I(1:2) \geq I'(1:2)$ ;  $J(1, 2) \geq J'(1, 2)$ , with equality if and only if
- iv  $\frac{f_1(x_1, \dots, x_k)}{f_2(x_1, \dots, x_k)} = \frac{g_1(y_1, \dots, y_r)}{g_2(y_1, \dots, y_r)}$  a.e.,

in which case we say the functions  $y_1, y_2, \dots, y_r$  are sufficient. The ratio  $J'(1, 2)/J(1, 2)$  is the discrimination efficiency of the  $y$ 's in the sense that  $N$  observations of the  $y$ 's will in the mean discriminate as well as  $n$  observations of the  $x$ 's where  $NJ' = nJ$ .

(c) *Particular cases.* If we denote the one-column matrix of the means of population  $\pi_i$  by  $\mu_{(i)}$ ,  $i = 1, 2$ , and the matrix of variances and covariances of population  $\pi_i$  by  $\sigma_{(i)}$ ,  $i = 1, 2$ , then evaluating (2.1) and (2.2) leads respectively to

$$(2.7) \quad J(1, 2) = \frac{1}{2} \text{tr}[(\sigma_{(1)} - \sigma_{(2)})(\sigma_{(2)}^{-1} - \sigma_{(1)}^{-1})] + \frac{1}{2}(\mu_{(1)} - \mu_{(2)})'(\sigma_{(1)}^{-1} + \sigma_{(2)}^{-1})(\mu_{(1)} - \mu_{(2)}),$$

$$(2.8) \quad I(1:2) = \frac{1}{2} \log \left| \frac{\sigma_{(2)}}{\sigma_{(1)}} \right| - \frac{k}{2} + \frac{1}{2} \text{tr} \sigma_{(1)} \sigma_{(2)}^{-1} + \frac{1}{2}(\mu_{(1)} - \mu_{(2)})' \sigma_{(2)}^{-1} (\mu_{(1)} - \mu_{(2)}),$$

where  $\text{tr } A$  is the trace (or spur) of the matrix  $A$  and the prime on a matrix denotes the transpose.

If  $\sigma_{(1)} = \sigma_{(2)} = \sigma$ , then (2.7) becomes

$$(2.9) \quad J(1, 2) = (\mu_{(1)} - \mu_{(2)})' \sigma^{-1} (\mu_{(1)} - \mu_{(2)}) = \delta' \sigma^{-1} \delta,$$

where  $\delta = \mu_{(1)} - \mu_{(2)}$ , and the last member in (2.9) is  $k$  times Mahalanobis's Generalized Distance ([4] p. 162, [10]), and (2.8) becomes

$$(2.10) \quad I(1:2) = \frac{1}{2} \delta' \sigma^{-1} \delta.$$

If  $\mu_{(1)} = \mu_{(2)}$ , then (2.7) becomes

$$(2.11) \quad \begin{aligned} J(1, 2) &= \frac{1}{2} \text{tr}[(\sigma_{(1)} - \sigma_{(2)})(\sigma_{(2)}^{-1} - \sigma_{(1)}^{-1})] \\ &= \frac{1}{2} \text{tr} \sigma_{(1)} \sigma_{(2)}^{-1} + \frac{1}{2} \text{tr} \sigma_{(2)} \sigma_{(1)}^{-1} - k, \end{aligned}$$

which for the single variate case is

$$(2.12) \quad J(1, 2) = \frac{1}{2}(\sigma_1^2 - \sigma_2^2) \left( \frac{1}{\sigma_2^2} - \frac{1}{\sigma_1^2} \right) = \frac{1}{2} \frac{\sigma_1^2}{\sigma_2^2} + \frac{1}{2} \frac{\sigma_2^2}{\sigma_1^2} - 1,$$

and (2.8) becomes

$$(2.13) \quad I(1:2) = \frac{1}{2} \log \left| \frac{\sigma_{(2)}}{\sigma_{(1)}} \right| - \frac{k}{2} + \frac{1}{2} \text{tr} \sigma_{(1)} \sigma_{(2)}^{-1},$$

which for the single variate case is

$$(2.14) \quad I(1:2) = \frac{1}{2} \log \frac{\sigma_2^2}{\sigma_1^2} - \frac{1}{2} + \frac{1}{2} \frac{\sigma_1^2}{\sigma_2^2}.$$

**3. Linear discriminant function.** Let us consider the following problem: Determine the values of the coefficients  $\alpha_1, \dots, \alpha_k$ , such that for

$$(3.1) \quad y = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k$$

the value of  $J'(1, 2)$  is a maximum, when  $x_1, x_2, \dots, x_k$  come from  $\pi_1$  and  $\pi_2$ . Depending on the assumptions regarding  $\pi_1$  and  $\pi_2$  we are led to a number of now classical results.

**4. Equal covariance matrices.** Assume that  $\sigma_{(1)} = \sigma_{(2)} = \sigma$ ; since  $y$  in (3.1) is normally distributed, we have as the single variate case of (2.9)

$$(4.1) \quad J'(1, 2) = (E(y_{(1)}) - E(y_{(2)}))^2 / \sigma_y^2 = (\alpha' \delta)^2 / \alpha' \sigma \alpha,$$

where  $\alpha$  is the one-column matrix of the  $\alpha_i, i = 1, 2, \dots, k$ , and  $\delta$  is defined as in (2.9). By selecting the  $\alpha$ 's such that

$$(4.2) \quad \alpha \sigma = \delta, \quad \alpha = \sigma^{-1} \delta,$$

$$(4.3) \quad J'(1, 2) = \frac{\delta' \sigma^{-1} \delta \delta' \sigma^{-1} \delta}{\delta' \sigma^{-1} \sigma \sigma^{-1} \delta} = \delta' \sigma^{-1} \delta = J(1, 2),$$

so that with the  $\alpha$ 's as given by (4.2), the linear function  $y$  of (3.1) is sufficient, and  $J'(1, 2)$  attains its maximum possible value (cf. [2], [3], [4], [5]).

**5. Principal components.** Let us assume that  $\mu_{(1)} = \mu_{(2)}$ , in which case, as we have seen,  $J(1, 2)$  is given by (2.11). For the linear function (3.1) we then derive, in view of (2.12),

$$(5.1) \quad J'(1, 2) = \frac{1}{2} \frac{\alpha' \sigma_{(1)} \alpha}{\alpha' \sigma_{(2)} \alpha} + \frac{1}{2} \frac{\alpha' \sigma_{(2)} \alpha}{\alpha' \sigma_{(1)} \alpha} - 1.$$

To find the values of the  $\alpha$ 's which will maximize (5.1), the usual calculus procedures yield the result that the  $\alpha$ 's must satisfy

$$(5.2) \quad \sigma_{(1)} \alpha = \lambda \sigma_{(2)} \alpha,$$

where  $\lambda$  is a root of the determinantal equation

$$(5.3) \quad | \sigma_{(1)} - \lambda \sigma_{(2)} | = 0,$$

all roots of which are real and positive. Let these roots be  $\lambda_1, \lambda_2, \dots, \lambda_k$ , arranged in ascending order. Corresponding to the root  $\lambda_i$ , and using (5.2), it is found that (5.1) may be written as

$$(5.4) \quad J'(1, 2; \lambda_i) = \frac{1}{2} \lambda_i + \frac{1}{2\lambda_i} - 1,$$

and that

$$(5.5) \quad \sum_{i=1}^k J'(1, 2; \lambda_i) = \frac{1}{2} \sum \lambda_i + \frac{1}{2} \sum \frac{1}{\lambda_i} - k = J(1, 2),$$

since

$$(5.6) \quad \sum \lambda_i = \text{tr } \sigma_{(1)} \sigma_{(2)}^{-1}, \quad \sum \frac{1}{\lambda_i} = \text{tr } \sigma_{(2)} \sigma_{(1)}^{-1};$$

also, using (2.13) and (2.14), we have

$$(5.7) \quad I'(1:2; \lambda_i) = -\frac{1}{2} \log \lambda_i - \frac{1}{2} + \frac{\lambda_i}{2},$$

$$(5.8) \quad I(1:2) = -\frac{1}{2} \log \lambda_1 \lambda_2 \cdots \lambda_k - \frac{k}{2} + \frac{1}{2} \sum \lambda_i = \sum_{i=1}^k I'(1:2; \lambda_i).$$

To determine the value for which (5.4) is a maximum, proceed as follows. Consider the function

$$(5.9) \quad f(\lambda) = \frac{1}{2} \lambda + \frac{1}{2\lambda} - 1, \quad \lambda > 0.$$

By examining the derivatives of  $f(\lambda)$ , it is readily determined that  $f(\lambda)$  is a minimum for  $\lambda = 1$ , is monotonically increasing for  $\lambda > 1$ , is monotonically decreasing for  $0 < \lambda < 1$ , and  $f(\lambda) = f(1/\lambda)$ . Thus, the maximum of (5.4) occurs for  $\lambda_1$  or  $\lambda_k$  according as

$$(5.10) \quad \lambda_1 \lambda_k < 1 \quad \text{or} \quad \lambda_1 \lambda_k > 1,$$

and the best linear discriminant function (3.1) is the one for which the  $\alpha$ 's respectively satisfy

$$(5.11) \quad \sigma_{(1)}\alpha = \lambda_1\sigma_{(2)}\alpha \quad \text{or} \quad \sigma_{(1)}\alpha = \lambda_2\sigma_{(2)}\alpha.$$

To illustrate, let us take

$$(5.12) \quad \sigma_{(1)} = \frac{1}{7} \begin{pmatrix} 12,496.8 & -6,786.6 \\ -6,786.6 & 32,985.0 \end{pmatrix}, \quad \sigma_{(2)} = \frac{1}{49} \begin{pmatrix} 136,972.6 & 58,549.0 \\ 58,549.0 & 71,496.1 \end{pmatrix},$$

which are respectively the treatments and residual values of Table I, p. 177 of Bartlett [2]. Using the values of (5.12), we find as the roots of (5.3)

$$(5.13) \quad \lambda_1 = 0.44158, \quad \lambda_2 = 6.38381,$$

so that (5.4) yields

$$(5.14) \quad J'(1, 2; \lambda_1) = .35309; \quad J'(1, 2; \lambda_2) = 2.27023;$$

and from (5.5) we have that

$$(5.15) \quad J(1, 2) = .35309 + 2.27023 = 2.62332.$$

Since  $\lambda_1\lambda_2 > 1$ , the best linear discriminant function is that associated with  $\lambda_2$  (as is also evident from (5.14)), and (5.11) becomes

$$(5.16) \quad \frac{1}{7} \begin{pmatrix} 12,496.8 & -6,786.6 \\ -6,786.6 & 32,985.0 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \frac{6.38381}{49} \begin{pmatrix} 136,972.6 & 58,549.0 \\ 58,549.0 & 71,496.1 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix},$$

or

$$(5.17) \quad \begin{cases} 112418.1\alpha_1 + 60181.5\alpha_2 = 0, \\ 60181.5\alpha_1 + 32217.3\alpha_2 = 0, \end{cases}$$

leading to

$$(5.18) \quad \alpha_1 = -.535\alpha_2,$$

that is, the linear function (see p. 179 of [2])

$$(5.19) \quad y = x_2 - 0.535x_1$$

is 86.5% efficient, since

$$(5.20) \quad J'(1, 2; \lambda_2)/J(1, 2) = \frac{2.27023}{2.62332} = .865.$$

Using the values in (5.13), we have from (5.8)

$$(5.21) \quad I(1:2) = 1.89449,$$

and since  $J(1, 2) = I(1:2) + I(2:1)$  [9],

$$(5.22) \quad I(2:1) = .72883.$$

For the linear function in (5.19) associated with  $\lambda_2$ , we have from (5.7)

$$(5.23) \quad I'(1:2; \lambda_2) = -\frac{1}{2} \log \lambda_2 - \frac{1}{2} + \frac{\lambda_2}{2} = 1.76502,$$

and since  $J'(1, 2; \lambda_2) = I'(1:2; \lambda_2) + I'(2:1; \lambda_2)$ ,

$$(5.24) \quad I'(2:1; \lambda_2) = .50521.$$

These values are summarized in Table 1.

TABLE 1

|           | $\lambda_1 = .44158$ | $\lambda_2 = 6.38381$ |                    |
|-----------|----------------------|-----------------------|--------------------|
| $I'(1:2)$ | .12947               | 1.76502               | 1.89449 = $I(1:2)$ |
| $I'(2:1)$ | .22362               | .50521                | .72883 = $I(2:1)$  |
| $J'(1,2)$ | .35309               | 2.27023               | 2.62332 = $J(1,2)$ |

From Table 1 it seems reasonable to infer that the linear function (5.19) is affected by the treatments.

If now we assume in particular that

$$(5.25) \quad \sigma_{(1)} = P, \quad \sigma_{(2)} = I_k,$$

where  $P$  is the matrix of population correlation coefficients, and  $I_k$  is the identity matrix of order  $k$ , then

$$(5.26) \quad \begin{aligned} J(1,2) &= \frac{1}{2} \text{tr} [(P - I_k)(I_k - P^{-1})] \\ &= \frac{1}{2} \text{tr} (P + P^{-1} - 2I_k) \\ &= \frac{1}{2} \sum_{i=1}^k (\rho^{ii} - 1) = \frac{1}{2} \sum_{i=1}^k \frac{\rho_{i,12 \dots (i-1)(i+1) \dots k}^2}{1 - \rho_{i,12 \dots (i-1)(i+1) \dots k}^2}, \end{aligned}$$

where  $\rho^{ii}$  are the diagonal elements of  $P^{-1}$ , and  $\rho_{i,12 \dots (i-1)(i+1) \dots k}$ ,  $i = 1, 2, \dots, k$  are the population multiple correlation coefficients, and

$$(5.27) \quad I(1:2) = -\frac{1}{2} \log |P|.$$

The best linear discriminant function (3.1) is that for which the  $\alpha$ 's satisfy

$$(5.28) \quad P\alpha = \lambda_1\alpha \quad \text{or} \quad P\alpha = \lambda_k\alpha,$$

according as  $\lambda_1\lambda_k < 1$  or  $\lambda_1\lambda_k > 1$ , where  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_k$  are the roots of

$$(5.29) \quad |P - \lambda I_k| = 0,$$

all of which are real and positive. It is easily verified ([6], [14]) that

$$(5.30) \quad \begin{cases} X'P^{-1}X = y_1^2/\lambda_1 + \dots + y_k^2/\lambda_k, \\ X'X = y_1^2 + \dots + y_k^2, \end{cases}$$

where  $y_i$  is the value of (3.1) corresponding to

$$(5.31) \quad P\alpha = \lambda_i\alpha.$$

For the bivariate case in particular, we have

$$(5.32) \quad P = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix},$$

$$(5.33) \quad I(1:2) = -\frac{1}{2} \log(1 - \rho^2), \quad J(1, 2) = \rho^2 / (1 - \rho^2),$$

$$(5.34) \quad |P - \lambda I_2| = \lambda^2 - 2\lambda + 1 - \rho^2 = 0,$$

$$(5.35) \quad \lambda_1 = 1 - \rho, \quad \lambda_2 = 1 + \rho, \quad \rho > 0,$$

$$(5.36) \quad \lambda_1 = 1 + \rho, \quad \lambda_2 = 1 - \rho, \quad \rho < 0,$$

$$(5.37) \quad J'(1, 2; \lambda_1) = \frac{1 - \rho}{2} + \frac{1}{2(1 - \rho)} - 1 = \frac{\rho^2}{2(1 - \rho)},$$

$$(5.38) \quad J'(1, 2; \lambda_2) = \frac{1 + \rho}{2} + \frac{1}{2(1 + \rho)} - 1 = \frac{\rho^2}{2(1 + \rho)},$$

$$(5.39) \quad J'(1, 2; \lambda_1) + J'(1, 2; \lambda_2) = \frac{\rho^2}{2(1 - \rho)} + \frac{\rho^2}{2(1 + \rho)} = \frac{\rho^2}{1 - \rho^2} = J(1, 2),$$

$$(5.40) \quad J'(1, 2; \lambda_1) / J(1, 2) = (1 + \rho) / 2,$$

$$(5.41) \quad J'(1, 2; \lambda_2) / J(1, 2) = (1 - \rho) / 2,$$

$$(5.42) \quad I'(1:2; \lambda_1) = -\frac{1}{2} \log(1 - \rho) - \frac{1}{2} + \frac{1}{2}(1 - \rho) = -\frac{1}{2} \log(1 - \rho) - \frac{1}{2}\rho,$$

$$(5.43) \quad I'(1:2; \lambda_2) = -\frac{1}{2} \log(1 + \rho) - \frac{1}{2} + \frac{1}{2}(1 + \rho) = -\frac{1}{2} \log(1 + \rho) + \frac{1}{2}\rho,$$

$$(5.44) \quad y_1 = (x_1 - x_2) / \sqrt{2}; \quad y_2 = (x_1 + x_2) / \sqrt{2},$$

$$(5.45) \quad \frac{1}{1 - \rho^2} (x_1^2 - 2\rho x_1 x_2 + x_2^2) = \frac{(x_1 - x_2)^2}{2(1 - \rho)} + \frac{(x_1 + x_2)^2}{2(1 + \rho)}.$$

Note that if  $\rho > 0$ , the best linear discriminant function corresponds to  $\lambda_1$ , as is evident from (5.40), and also from the fact that  $\lambda_1 \lambda_2 = 1 - \rho^2 < 1$ .

**6. Canonical correlation.** Let us assume that  $\mu_{(1)} = \mu_{(2)}$ , and that

$$(6.1) \quad \sigma_{(1)} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}, \quad \sigma_{(2)} = \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix},$$

where

$$\Sigma_{11} = (\sigma_{ij}), \quad i, j = 1, 2, \dots, k_1,$$

$$\Sigma_{22} = (\sigma_{rs}), \quad r, s = k_1 + 1, \dots, k_1 + k_2 = k,$$

$$\Sigma_{12} = (\sigma_{is}), \quad \Sigma_{21} = \Sigma'_{12}.$$



Since, as may be readily verified,

$$(6.2) \quad \begin{pmatrix} I_{k_1} & 0 \\ -\Sigma_{21} \Sigma_{11}^{-1} & I_{k_2} \end{pmatrix} \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} I_{k_1} - \Sigma_{11}^{-1} \Sigma_{12} \\ 0 & I_{k_2} \end{pmatrix} = \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \end{pmatrix},$$

we have

$$(6.3) \quad \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} = \begin{pmatrix} I_{k_1} - \Sigma_{11}^{-1} \Sigma_{12} & \Sigma_{11}^{-1} & 0 \\ 0 & I_{k_2} & \Sigma_{22.1}^{-1} \end{pmatrix} \begin{pmatrix} I_{k_1} & 0 \\ -\Sigma_{21} \Sigma_{11}^{-1} & I_{k_2} \end{pmatrix} \\ = \begin{pmatrix} \Sigma_{11}^{-1} + \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} & -\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \\ -\Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} & \Sigma_{22.1}^{-1} \end{pmatrix},$$

where  $\Sigma_{22.1} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ . Thus (cf. [1] p. 182)

$$(6.4) \quad J(1,2) = \frac{1}{2} \text{tr} \left[ \left\{ \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} - \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix} \right\} \left\{ \begin{pmatrix} \Sigma_{11}^{-1} & 0 \\ 0 & \Sigma_{22}^{-1} \end{pmatrix} - \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} \right\} \right] \\ = \frac{1}{2} \text{tr} \left[ \begin{pmatrix} 0 & \Sigma_{12} \\ \Sigma_{21} & 0 \end{pmatrix} \begin{pmatrix} -\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} & \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \\ \Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} & \Sigma_{22}^{-1} - \Sigma_{22.1}^{-1} \end{pmatrix} \right] \\ = \frac{1}{2} \text{tr} \begin{pmatrix} \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} & \cdot \\ \cdot & \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \end{pmatrix} \\ = \text{tr} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} = \text{tr} \Sigma_{22} \Sigma_{22.1}^{-1} - k_2,$$

where the dots indicate matrices which are not needed, and

$$(6.5) \quad I(1:2) = \frac{1}{2} \log \frac{\begin{vmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{vmatrix}}{\begin{vmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{vmatrix}} - \frac{k}{2} + \frac{1}{2} \text{tr} \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \Sigma_{11}^{-1} & 0 \\ 0 & \Sigma_{22}^{-1} \end{pmatrix} \\ = \frac{1}{2} \log \frac{\begin{vmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{vmatrix}}{\begin{vmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{vmatrix}} = \frac{1}{2} \log \frac{|\Sigma_{11}| |\Sigma_{22}|}{\begin{vmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{vmatrix}}.$$

If we write the linear function of (3.1) as

$$(6.6) \quad y = \beta_1 x_1 + \cdots + \beta_{k_1} x_{k_1} + \gamma_1 x_{k_1+1} + \cdots + \gamma_{k_2} x_{k_1+k_2},$$

then (5.2) may be written as

$$(6.7) \quad \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \beta \\ \gamma \end{pmatrix} = \lambda \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \beta \\ \gamma \end{pmatrix},$$

where  $\beta$  and  $\gamma$  are respectively the one-column matrices of  $\beta_1, \dots, \beta_{k_1}$  and  $\gamma_1, \dots, \gamma_{k_2}$ , and (5.3) may be written as

$$(6.8) \quad \begin{vmatrix} (1-\lambda)\Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & (1-\lambda)\Sigma_{22} \end{vmatrix} = 0.$$

Since (6.7) is equivalent to

$$(6.9) \quad \begin{aligned} \Sigma_{11}\beta + \Sigma_{12}\gamma &= \lambda\Sigma_{11}\beta, \\ \Sigma_{21}\beta + \Sigma_{22}\gamma &= \lambda\Sigma_{22}\gamma, \end{aligned}$$

or

$$(6.10) \quad \begin{aligned} \beta &= -\frac{1}{1-\lambda}\Sigma_{11}^{-1}\Sigma_{12}\gamma, \\ -\frac{1}{1-\lambda}\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\gamma + (1-\lambda)\Sigma_{22}\gamma &= 0, \end{aligned}$$

we conclude that (6.8) is equivalent to

$$(6.11) \quad |\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} - \rho^2\Sigma_{22}| = 0,$$

where  $\rho^2 = (1-\lambda)^2$ . According to (5.3), the roots of (6.8) are all real and positive. If we take  $k_2 \leq k_1$ , then since  $k = k_1 + k_2$ , and the determinant of (6.11) is of order  $k_2$ ,

$$(6.12) \quad \begin{aligned} \lambda_i &= 1 - \rho_i, \quad \lambda_{k_1+i} = 1 + \rho_{k_2+1-i}, \quad i = 1, 2, \dots, k_2, \\ \lambda_{k_2+1} &= \dots = \lambda_{k_2+(k_1-k_2)} = 1, \end{aligned}$$

where  $\rho_1 \geq \rho_2 \geq \dots \geq \rho_{k_2}$ . We may also conclude that  $-1 \leq \rho_i \leq 1$  since the  $\lambda$ 's cannot be negative. The  $\rho_i$  are Hotelling's canonical correlations [7]. The results of (5.4) now become

$$(6.13) \quad \begin{aligned} J'(1, 2; \lambda_i) &= \frac{1}{2}(1 + \rho_i) + \frac{1}{2(1 + \rho_i)} - 1 = \frac{1}{2}\rho_i^2/(1 + \rho_i), \\ J'(1, 2; \lambda_{k_1+i}) &= \frac{1}{2}(1 - \rho_{k_2+1-i}) + \frac{1}{2(1 - \rho_{k_2+1-i})} - 1 = \frac{\rho_{k_2+1-i}^2}{2(1 - \rho_{k_2+1-i})}, \\ J'(1, 2; \lambda_j) &= \frac{1}{2} + \frac{1}{2} - 1 = 0, \\ & \quad i = 1, \dots, k_2, j = k_2 + 1, \dots, k_2 + (k_1 - k_2), \end{aligned}$$

or

$$(6.14) \quad J'(1, 2; \lambda_i) + J'(1, 2; \lambda_{k_1+1-i}) = \rho_i^2/(1 - \rho_i^2), \quad i = 1, 2, \dots, k_2,$$

and from (5.5) and (5.8) respectively

$$(6.15) \quad J(1, 2) = \sum_{i=1}^{k_2} \rho_i^2/(1 - \rho_i^2),$$

$$(6.16) \quad I(1:2) = -\frac{1}{2} \log(1 - \rho_1^2)(1 - \rho_2^2) \dots (1 - \rho_{k_2}^2).$$

Since

$$\lambda_1\lambda_k = (1 - \rho_1)(1 + \rho_1) = 1 - \rho_1^2 < 1,$$

the best linear discriminant function (6.6) corresponds to the value  $\lambda_1$ , or to the largest canonical correlation.

If we pose the problem of finding the best pair of linear discriminant functions

$$(6.17) \quad \begin{cases} \mu = \beta_1 x_1 + \cdots + \beta_{k_1} x_{k_1}, \\ \nu = \gamma_1 x_{k_1+1} + \cdots + \gamma_{k_2} x_{k_1+k_2}, \end{cases}$$

then we want to maximize (see (5.33))

$$(6.18) \quad J'(1, 2; \mu, \nu) = \frac{\rho_{\mu\nu}^2}{1 - \rho_{\mu\nu}^2} = \frac{(\beta' \Sigma_{12} \gamma)^2}{(\beta' \Sigma_{11} \beta)(\gamma' \Sigma_{22} \gamma) - (\beta' \Sigma_{12} \gamma)^2}.$$

The usual methods lead us again to the condition (6.9), where  $(1 - \lambda)^2 = \rho_{\mu\nu}^2$ . We thus see that the canonical correlation coefficients are the values of  $\rho_{\mu\nu}$ , and

$$(6.19) \quad J'(1, 2; \mu_i, \nu_i) = J'(1, 2; \lambda_i) + J'(1, 2; \lambda_{k_1+1-i}), \quad i = 1, 2, \dots, k_2.$$

The best pair of linear discriminant functions thus corresponds to  $\rho_1^2$ , that is, the largest of the canonical correlations.

To illustrate, let us take as  $\sigma_{(1)}$  the matrix

$$(6.20) \quad \begin{pmatrix} 1.0000 & .6328 & | & .2412 & .0586 \\ .6328 & 1.0000 & | & -.0553 & .0655 \\ \hline .2412 & -.0553 & | & 1.0000 & .4248 \\ .0586 & .0655 & | & .4248 & 1.0000 \end{pmatrix},$$

which is Kelley's data discussed on p. 342 of [7]. As the roots of (6.8) we find (cf. Ex. 28.4, p. 351, of [8])

$$(6.21) \quad \lambda_1 = 0.6055, \quad \lambda_2 = 0.9312, \quad \lambda_3 = 1.0688, \quad \lambda_4 = 1.3945,$$

$$(6.22) \quad \begin{cases} \rho_1^2 = (1 - .6055)(1.3945 - 1) = .1556, \\ \rho_2^2 = (1 - .9312)(1.0688 - 1) = .0047, \end{cases}$$

$$(6.23) \quad \begin{aligned} J'(1, 2; \lambda_1) &= .1285, & J'(1, 2; \lambda_2) &= .0025, \\ J'(1, 2; \lambda_3) &= .0022, & J'(1, 2; \lambda_4) &= .0558, \end{aligned}$$

$$(6.24) \quad \frac{\rho_1^2}{1 - \rho_1^2} = \frac{.1556}{.8444} = .1843, \quad \frac{\rho_2^2}{1 - \rho_2^2} = \frac{.0047}{.9953} = .0047,$$

$$(6.25) \quad J(1, 2) = .1843 + .0047 = .1890.$$

The linear function associated with  $\lambda_1$  is 68.0% efficient. The pair of linear functions (6.17) related to the correlation  $\rho_1^2 = .1556$  (see [7], [8] loc.cit.),

$$(6.26) \quad \begin{cases} \mu_1 = -2.7772x_1 + 2.2655x_2, \\ \nu_1 = -2.4404x_3 + x_4, \end{cases}$$

are 97.5% efficient (and thus practically sufficient) since

$$(6.27) \quad \frac{J'(1, 2; \lambda_1)}{J(1, 2)} = \frac{.1285}{.1890} = .680, \quad \frac{J'(1, 2; \mu_1 \nu_1)}{J(1, 2)} = \frac{.1843}{.1890} = .975.$$

Using the values in (6.22), we have from (6.16)

$$(6.28) \quad I(1:2) = -\frac{1}{2} \log (.8444)(.9953) = .0869,$$

and therefore

$$(6.29) \quad I(2:1) = .1890 - .0869 = .1021.$$

Similarly

$$(6.30) \quad I'(1:2; \mu_1, \nu_1) = -\frac{1}{2} \log(.8444) = .0846,$$

$$(6.31) \quad I'(2:1; \mu_1, \nu_1) = .1843 - .0846 = .0997.$$

These values are summarized in Table 2.

TABLE 2

|            | $\mu_1, \nu_1$ | $\mu_2, \nu_2$ |                   |
|------------|----------------|----------------|-------------------|
| $I'(1:2)$  | .0846          | .0023          | .0869 = $I(1:2)$  |
| $I'(2:1)$  | .0997          | .0024          | .1021 = $I(2:1)$  |
| $J'(1, 2)$ | .1843          | .0047          | .1890 = $J(1, 2)$ |

From Table 2 it seems reasonable to infer that the linear functions (6.26) are the only such components. (See p. 342 [7].)

**7. Discriminant functions with covariance.** Assume that

$$(7.1) \quad \sigma_{(1)} = \sigma_{(2)} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix},$$

with  $\Sigma_{11}, \Sigma_{12}, \Sigma_{21}, \Sigma_{22}$  as defined in (6.1). Let the one-row matrix of means be

$$(7.2) \quad (\mu'_{(i)}, \nu'_{(i)}), \quad i = 1, 2,$$

where the  $\mu$ 's are the means of the first  $k_1$  variables, and the  $\nu$ 's the means of the last  $k_2$  variables. Assume that  $\nu_{(1)} = \nu_{(2)}$ . Then

$$(7.3) \quad J(1, 2) = (\delta, ' 0) \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} \begin{pmatrix} \delta \\ 0 \end{pmatrix},$$

where  $\delta$  is defined as in (2.9). Using the value of the inverse matrix as given in (6.3) it follows that

$$(7.4) \quad J(1, 2) = \delta' \Sigma_{11}^{-1} \delta + \delta' \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \delta.$$

Let

$$(7.5) \quad y_i = x_i, \quad i = 1, 2, \dots, k_1,$$

then

$$(7.6) \quad J'(1, 2) = \delta' \Sigma_{11}^{-1} \delta.$$

The gain due to the use of the covariance variates  $x_{k_1+1}, \dots, x_{k_1+k_2}$  in the linear discriminant function is thus given by (see [4])

$$(7.7) \quad J(1, 2)/J'(1, 2) = 1 + \lambda,$$

where

$$(7.8) \quad \lambda = \frac{\delta' \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \delta}{\delta' \Sigma_{11}^{-1} \delta}$$

and  $\lambda$  will take on a value between the smallest and largest root of the determinantal equation

$$(7.9) \quad | \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} \Sigma_{11}^{-1} - \lambda \Sigma_{11}^{-1} | = 0.$$

Indeed, since the quadratic form in the denominator of  $\lambda$  is positive definite, there exists a real nonsingular transformation

$$(7.10) \quad \delta = A\gamma,$$

such that

$$(7.11) \quad \lambda = \frac{\lambda_1 \gamma_1^2 + \lambda_2 \gamma_2^2 + \dots + \lambda_{k_1} \gamma_{k_1}^2}{\gamma_1^2 + \dots + \gamma_{k_1}^2},$$

where  $\lambda_1, \lambda_2, \dots, \lambda_{k_1}$ , are the roots of (7.9), or

$$(7.12) \quad | \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} - \lambda \Sigma_{11} | = 0.$$

If  $k_1 \geq k_2$ , then since

$$(7.13) \quad \begin{aligned} \begin{vmatrix} \lambda \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22.1} \end{vmatrix} &= \begin{vmatrix} \lambda \Sigma_{11} - \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} & \Sigma_{12} \\ & 0 & \Sigma_{22.1} \end{vmatrix} \\ &= \begin{vmatrix} \lambda \Sigma_{11} & 0 \\ \Sigma_{21} & \Sigma_{22.1} - \frac{1}{\lambda} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \end{vmatrix}, \end{aligned}$$

it follows that

$$(7.14) \quad \left| \Sigma_{12} \Sigma_{22.1}^{-1} \Sigma_{21} - \lambda \Sigma_{11} \right| = 0 = \left| \Sigma_{22.1} - \frac{1}{\lambda} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \right|.$$

Since  $\Sigma_{22,1} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$ , the right member of (7.14) reduces to

$$(7.15) \quad | \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} - \rho^2\Sigma_{22} | = 0,$$

where  $\rho^2 = \lambda/(1 + \lambda)$ . Comparison with (6.11) shows that the roots of (7.15) are the canonical correlations.

In the case  $k_2 = 1$ , there is only one canonical correlation and it is the multiple correlation of  $x_k$  on  $x_1, x_2, \dots, x_{k-1}$ , so that

$$(7.16) \quad \frac{J(1, 2)}{J'(1, 2)} = 1 + \lambda \leq 1 + \frac{\rho_{k \cdot 12 \dots (k-1)}^2}{1 - \rho_{k \cdot 12 \dots (k-1)}^2} = \frac{1}{1 - \rho_{k \cdot 12 \dots (k-1)}^2}.$$

Thus, using the values for the Rabbits  $\times$  doses components in Table 2 p. 157 of Cochran and Bliss [4], we have for the matrix (7.1),

$$(7.17) \quad \frac{1}{33} \left( \begin{array}{cc|c} 3223 & 1200 & 1259 \\ 1200 & 3137 & 1340 \\ \hline 1259 & 1340 & 2351 \end{array} \right),$$

and

$$(7.18) \quad \frac{1259 \left| \begin{array}{cc} 1259 & 1200 \\ 1340 & 1373 \end{array} \right| + 1340 \left| \begin{array}{cc} 3223 & 1259 \\ 1200 & 1340 \end{array} \right|}{\left| \begin{array}{cc} 3223 & 1200 \\ 1200 & 3137 \end{array} \right|} = 774,$$

$$(7.19) \quad \rho_{3 \cdot 12}^2 = \frac{774}{2351} = .33,$$

$$(7.20) \quad J/J' \leq \frac{1}{1 - .33} = 1.50.$$

On p. 162 of [4] it was concluded that the use of covariance gives 50% more information.

Solving for the coefficients of the discriminant function in the equations (see [4], p. 157)

$$(7.21) \quad \begin{aligned} 3223\alpha_1 + 1200\alpha_2 + 1259\alpha_3 &= -1197.2 \times 33, \\ 1200\alpha_1 + 3137\alpha_2 + 1340\alpha_3 &= -844.3 \times 33, \\ 1259\alpha_1 + 1340\alpha_2 + 2351\alpha_3 &= 0, \end{aligned}$$

it is found that (see (4.2), (4.3))

$$(7.22) \quad 33J = 1197.2 \times .41848 + 844.3 \times .27070 = 729.556.$$

If we solve (omitting the covariance variable)

$$(7.23) \quad \begin{aligned} 3223\beta_1 + 1200\beta_2 &= -1197.2 \times 33 \\ 1200\beta_1 + 3137\beta_2 &= -844.3 \times 33 \end{aligned}$$

for the coefficients of the linear discriminant function, it is found that (see (4.2), (4.3))

$$(7.24) \quad 33J' = 1197.2 \times .31629 + 844.3 \times .14815 = 503.845,$$

and

$$(7.25) \quad J/J' = \frac{729.556}{503.845} = 1.45 < 1.50.$$

**8. Conclusion.** It is seen that the multivariate analysis techniques of discriminant analysis, principal components and canonical correlations are indeed closely related concepts associated with a linear discriminant function, and differing primarily in the assumption about the underlying populations.

#### REFERENCES

- [1] M. S. BARTLETT, "A note on tests of significance in multivariate analysis," *Proc. Cambridge Philos. Soc.*, Vol. 35 (1939), pp. 180-185.
- [2] M. S. BARTLETT, "Multivariate analysis," *Jour. Roy. Stat. Soc., Suppl.*, Vol. 9 (1947), pp. 176-197.
- [3] G. W. BROWN, "Discriminant functions," *Annals of Math. Stat.*, Vol. 18 (1947), pp. 514-528.
- [4] W. G. COCHRAN AND C. I. BLISS, "Discriminant functions with covariance," *Annals of Math. Stat.*, Vol. 19 (1948), pp. 151-176.
- [5] R. A. FISHER, "The statistical utilization of multiple measurements," *Annals of Eugenics*, Vol. 8 (1938), pp. 376-386.
- [6] H. HOTELLING, "Analysis of a complex of statistical variables into principal components," *Jour. Educ. Psych.*, Vol. 24 (1933), pp. 417-441; 498-520.
- [7] H. HOTELLING, "Relations between two sets of variates," *Biometrika*, Vol. 28 (1936), pp. 321-377.
- [8] M. G. KENDALL, *The Advanced Theory of Statistics*, Vol. 2, C. Griffin & Co. Ltd, London, 1946.
- [9] S. KULLBACK AND R. A. LEIBLER, "On information and sufficiency," *Annals of Math. Stat.*, Vol. 22 (1951), pp. 79-86.
- [10] P. C. MAHALANOBIS, "On the generalized distance in statistics," *Proc. Nat. Inst. Sci. India*, Vol. 12 (1936), pp. 49-55.
- [11] L. S. PENROSE, "Some notes on discrimination," *Annals of Eugenics*, Vol. 13 (1947), pp. 228-237.
- [12] C. A. B. SMITH, "Some examples of discrimination," *Annals of Eugenics*, Vol. 13 (1947), pp. 272-282.
- [13] GERHARD TINTNER, "Some formal relations in multivariate analysis," *Jour. Roy. Stat. Soc., Ser. B*, Vol. 12 (1950), pp. 95-101.
- [14] S. S. WILKS, *Mathematical Statistics*, Princeton University Press, 1943.
- [15] E. C. CHERRY, "A history of the theory of information," *Proc. Inst. Elec. Engrs. Part III*, Vol. 98 (1951), pp. 383-393.
- [16] A. BHATTACHARYYA, "On a measure of divergence between two statistical populations defined by their probability distributions," *Bull. Calcutta Math. Soc.*, Vol. 35 (1943), pp. 99-109.
- [17] A. BHATTACHARYYA, "On a measure of divergence between two multinomial populations," *Sankhyā*, Vol. 7 (1946), pp. 58-61.
- [18] C. R. RAO, "On the distance between two populations," *Sankhyā*, Vol. 9 (1949), pp. 246-248.