

Inference on Discriminant Function Coefficients

C. RADHAKRISHNA RAO, *Indian Statistical Institute*

1. INTRODUCTION

Some years ago, the author developed some tests for examining hypotheses such as "the coefficients of some specified characters in the linear discriminant function are zero" (Rao, 1946; 1948; 1949) and "the coefficients of two given characters are in a specified ratio" (Rao, 1952). These tests were meant to be generalizations of an earlier test by Fisher (1940) for a proposed (assigned) discriminant function (i.e., when the proportions of all the coefficients are specified). During the last few years, the author has received a number of queries regarding the theory and application of these tests. The present paper is an attempt to answer these queries and to propose tests for other hypotheses of interest in the use of discriminant functions.

It must be mentioned that the individual coefficients in the linear discriminant function are not definite population parameters, as only the ratios of the coefficients are unique. For this reason estimates of the individual coefficients and their

standard errors are not meaningful. We can only draw inferences on the ratios of the coefficients.

All the tests considered in the paper are special cases of a test for examining the sufficiency of a given subset out of a larger set of variables, for purposes of discrimination between two populations. *The hypothesis of sufficiency* is explicitly defined as follows. Let (X_1, \dots, X_p) be a p dimensional random variable. Then the subset (X_1, \dots, X_q) is said to be sufficient for discrimination between two populations if the conditional distributions of (X_1, \dots, X_p) given (X_1, \dots, X_q) are the same for both the populations. We may also describe the hypothesis of sufficiency of (X_1, \dots, X_q) as the *absence of additional information* contained in (X_{q+1}, \dots, X_p) when the variables (X_1, \dots, X_q) are already available. We consider this general problem in Section 2.

The problem of inference on the coefficients of a genetic selection index as developed by Smith (1936) needs an entirely different approach. A test is described for examining the adequacy of a *straight selection function* or any proposed selection function. Other questions such as the adequacy of a subset of the phenotypic observations for assessing some well defined genetic worth of an individual need further study.

2. NOTATIONS AND PRELIMINARY RESULTS

Let $\mathbf{X}' = (X_1, \dots, X_p)$ be a p dimensional random variable and consider the partition $\mathbf{X}' = (\mathbf{X}'_1 : \mathbf{X}'_2)$ where \mathbf{X}'_1 consists of the first q components of \mathbf{X} and \mathbf{X}'_2 , the rest of the components of \mathbf{X} . Assuming the first two moments of \mathbf{X} to exist, we can write the corresponding partitions of $E(\mathbf{X})$, the expectation of \mathbf{X} , and $D(\mathbf{X})$, the dispersion matrix of \mathbf{X} , as

$$(1) \quad E(\mathbf{X}') = \boldsymbol{\mu}' = (E(\mathbf{X}'_1) : E(\mathbf{X}'_2)),$$

$$(2) \quad D(\mathbf{X}) = \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

Let H_1 and H_2 be two simple hypotheses specifying the means and dispersion matrices as

$$(3) \quad E(\mathbf{X}' | H_1) = \boldsymbol{\mu}'_1 = (\boldsymbol{\mu}'_{11} : \boldsymbol{\mu}'_{12}), \quad D(\mathbf{X}) = \boldsymbol{\Sigma},$$

$$(4) \quad E(\mathbf{X}' | H_2) = \boldsymbol{\mu}'_2 = (\boldsymbol{\mu}'_{21} : \boldsymbol{\mu}'_{22}), \quad D(\mathbf{X}) = \boldsymbol{\Sigma}.$$

If $\boldsymbol{\delta}' = \boldsymbol{\mu}'_1 - \boldsymbol{\mu}'_2 = (\boldsymbol{\delta}'_1 : \boldsymbol{\delta}'_2)$, then the linear discriminant function between H_1 and H_2 based on \mathbf{X} is defined as $\boldsymbol{\delta}' \boldsymbol{\Sigma}^{-1} \mathbf{X}$, while that based on \mathbf{X}_1 alone is $\boldsymbol{\delta}'_1 \boldsymbol{\Sigma}_{11}^{-1} \mathbf{X}_1$. We shall represent the linear discriminant function $\boldsymbol{\delta}' \boldsymbol{\Sigma}^{-1} \mathbf{X}$ by $\boldsymbol{\lambda}' \mathbf{X} = \lambda_1 X_1 + \dots + \lambda_p X_p$ and consider some hypotheses on $\boldsymbol{\lambda}$.

Note that the Mahalanobis distance between H_1 and H_2 based on \mathbf{X} is $\Delta_p^2 = \boldsymbol{\delta}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}$ while that based on \mathbf{X}_1 is $\Delta_q^2 = \boldsymbol{\delta}'_1 \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\delta}_1$.

With the above notations, the following statements are equivalent.

(a) $\boldsymbol{\delta}_2 - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\delta}_1 = \mathbf{0}$, i.e., the random variable $\mathbf{X}_2 - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \mathbf{X}_1$ obtained by subtracting from \mathbf{X}_2 its regression on \mathbf{X}_1 has the same expected value for both the populations.

(b) The coefficients $\lambda_{q+1}, \dots, \lambda_p$ of the components $\mathbf{X}_{q+1}, \dots, \mathbf{X}_p$, in the linear discriminant function based on \mathbf{X} are all zero.

(c) $\Delta_p^2 = \Delta_q^2$, i.e., there is no additional distance contributed by the variables $\mathbf{X}_{q+1}, \dots, \mathbf{X}_p$.

(d) Every linear function of \mathbf{X} uncorrelated with \mathbf{X}_1 has the same expected value for both the populations.

(e) If $\mathbf{X} \sim \mathcal{N}_p$ (i.e., distributed as p -variate normal), then the conditional distribution of \mathbf{X} given \mathbf{X}_1 is the same for both the populations, which is the same as saying that \mathbf{X}_1 is sufficient for discrimination between the populations.

The compounding vector $\boldsymbol{\lambda}$ of \mathbf{X} in the linear discrimination function $\boldsymbol{\delta}' \boldsymbol{\Sigma}^{-1} \mathbf{X}$ is

$$(5) \quad (\boldsymbol{\lambda}'_1 : \boldsymbol{\lambda}'_2) = (\boldsymbol{\delta}'_1 : \boldsymbol{\delta}'_2) \cdot \begin{pmatrix} \boldsymbol{\Sigma}^{11} & \boldsymbol{\Sigma}^{12} \\ \boldsymbol{\Sigma}^{21} & \boldsymbol{\Sigma}^{22} \end{pmatrix}$$

where $\lambda'_2 = (\lambda_{q+1}, \dots, \lambda_p)$ is the vector of coefficients of the components of X_2 . From (5) $\lambda'_2 = \delta'_1 \Sigma^{12} + \delta'_2 \Sigma^{22}$ and by the statement (b), $\lambda_2 = 0$. By virtue of the algebraic equivalence

$$(6) \quad \delta'_1 \Sigma^{12} + \delta'_2 \Sigma^{22} = 0 \iff \delta_2 - \Sigma_{21} \Sigma_{11}^{-1} \delta_1 = 0,$$

this proves the equivalence of statements (a) and (b).

It is easy to prove the identity

$$\begin{aligned} \Delta_p^2 &= \delta \Sigma^{-1} \delta \\ &= \delta'_1 \Sigma_{11}^{-1} \delta_1 + (\delta_2 - \Sigma_{21} \Sigma_{11}^{-1} \delta_1)' (\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})^{-1} (\delta_2 - \Sigma_{21} \Sigma_{11}^{-1} \delta_1) \end{aligned} \quad (7)$$

If $\Delta_p^2 = \Delta_q^2 = \delta'_1 \Sigma_{11}^{-1} \delta_1$, the second term in (7), which is a positive definite quadratic form is zero. Hence $(\delta_2 - \Sigma_{21} \Sigma_{11}^{-1} \delta_1) = 0$ and vice-versa, i.e., (a) \iff (c).

Let $L'_1 X_1 + L'_2 X_2$ be a linear function of X . Then the statement (d) implies $L'_1 \Sigma_{11} + L'_2 \Sigma_{21} = 0$, i.e., $L_1 = -\Sigma_{11}^{-1} \Sigma_{12} L_2$. Consider $L'_1 \delta_1 + L'_2 \delta_2 = L'_2 (-\Sigma_{21} \Sigma_{11}^{-1} \delta_1 + \delta_2) = 0$ for all $L_2 \iff \delta_2 - \Sigma_{21} \Sigma_{11}^{-1} \delta_1 = 0$, i.e., (d) \iff (a).

To prove (c) \iff (a) we need only consider the conditional distribution of X_2 given X_1 , which is $(p-q)$ variate normal with

$$(8) \quad E(X_2 | X_1, H_1) - E(X_2 | X_1, H_2) = \delta_2 - \Sigma_{21} \Sigma_{11}^{-1} \delta_1$$

$$(9) \quad D(X_2 | X_1, H_1) = D(X_2 | X_1, H_2) = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}.$$

Sufficiency of X_1 is equivalent to $\delta_2 - \Sigma_{21} \Sigma_{11}^{-1} \delta_1 = 0$, which proves the desired result.

3. TEST FOR SUFFICIENCY OF A SUBSET

We shall develop a test for the hypothesis that " X_1 is sufficient" or " X_2 has no additional information in the presence of X_1 " on the basis of samples of sizes n_1 and n_2 from the two populations. Let

$$(10) \quad d' = (d'_1 : d'_2)$$

represent the difference in sample means of \mathbf{X}_1 and \mathbf{X}_2 and

$$(11) \quad S = \begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix}$$

the pooled sum of products (S. P. matrix) within samples on (n_1+n_2-2) degrees of freedom (d.f.).

Let us observe that under the assumption of normality

$$(12) \quad E(\mathbf{X}_2 | \mathbf{X}_1, H_1) = \alpha_1 + \Gamma \mathbf{X}_1,$$

$$E(\mathbf{X}_2 | \mathbf{X}_1, H_2) = \alpha_2 + \Gamma \mathbf{X}_1,$$

where Γ represents the matrix of regression coefficients. If the dispersion matrix of \mathbf{X} is the same under H_1 and H_2 , then $D(\mathbf{X}_2 | \mathbf{X}_1, H_1) = D(\mathbf{X}_2 | \mathbf{X}_1, H_2)$. Hence the hypothesis under test is $\alpha_1 = \alpha_2$.

The formulation (12) may be recognized as the multivariate extension of the Gauss-Markoff set up, with $(p-k)$ variables (the components of \mathbf{X}_2). So, no new problem arises in the consideration of test criteria. One obtains, to begin with, the dispersion matrices due to deviation from hypothesis and due to error. We also notice that the set up (12) involves two sets of parameters, with the null hypothesis concerning only one of the sets. In such case the computations are simple, involving what is known as covariance adjustment (see the discussion on page 119 in Rao (1952) in the univariate case and on pages 468-69 in Rao (1965) in the multivariate case).

The S.P. matrix within populations jointly for $\mathbf{X}_1, \mathbf{X}_2$ is

$$(13) \quad W = \begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix} = S$$

on (n_1+n_2-2) d.f. and that between populations is

$$(14) \quad B = \frac{n_1 n_2}{n_1 + n_2} \begin{pmatrix} d_1 d_1' & d_1 d_2' \\ d_2 d_1' & d_2 d_2' \end{pmatrix}$$

on 1 d.f., giving the total S.P. matrix

$$(15) \quad T = W + B = \begin{pmatrix} T_{11} & T_{12} \\ T_{21} & B_{22} \end{pmatrix}.$$

We now compute the residual S.P. matrices for X_2 , making covariance adjustment for X_1 , obtaining

$$(16) \quad S_{22} - S_{21}S_{11}^{-1}S_{12} \quad \text{with d.f.} = n_1 + n_2 - 2 - q$$

for within, and

$$(17) \quad T_{22} - T_{21}T_{11}^{-1}T_{12} \quad \text{with d.f.} = n_1 + n_2 - 1 - q$$

for the total.

Now applying the Wilks Λ criterion which is the likelihood ratio test applied on the conditional distributions of X_2 given X_1 we have

$$(18) \quad \Lambda = \frac{|S_{22} - S_{21}S_{11}^{-1}S_{12}|}{|T_{22} - T_{21}T_{11}^{-1}T_{12}|} = \frac{|S|}{|S_{11}|} \frac{T}{|T_{11}|}$$

$$= \frac{|S|}{|T|} \cdot \frac{|S_{11}|}{|T_{11}|}$$

It is well known that in the special case when the matrix B , has 1 d.f., the statistic

$$(19) \quad \frac{n_1 + n_2 - p - 1}{p - q} \left(\frac{1}{\Lambda} - 1 \right)$$

has an F distribution on $(p - q)$ and $(n_1 + n_2 - p - 1)$ d.f.

Defining

$$(20) \quad D_p^2 = (n_1 + n_2 - 2)d'S^{-1}d, \quad D_q^2 = (n_1 + n_2 - 2)d_1'S_{11}^{-1}d_1$$

which are estimates of Δ_p^2 and Δ_q^2 and $c = n_1 n_2 / (n_1 + n_2)(n_1 + n_2 - 2)$ we can write the F statistic (19) as

$$(21) \quad \frac{n_1 + n_2 - p - 1}{p - q} \cdot \frac{c(D_p^2 - D_q^2)}{1 + cD_q^2}$$

the form in which the test was originally expressed (Rao, 1949; 1952).

Optimum properties of the test (21) have been recently investigated by Giri (1964).

4. INFERENCE ON DISCRIMINANT FUNCTION COEFFICIENTS

The test criterion (21) was developed for testing the hypothesis that the coefficients of specified variables in the linear discriminant function are zero. We can apply the same test by a suitable transformation of the variables in drawing inferences of various types on the coefficients of a linear discriminant function.

4.1. Test for a proposed discriminant function (Fisher, 1940)

Let $\lambda'X$ be the assigned discriminant function. Then the null hypotheses under test can be written as

$$(22) \quad \lambda'X \propto \delta' \Sigma^{-1} X \rightarrow \delta \propto \Sigma \lambda, (\lambda \text{ given}).$$

We make the transformation (assuming $\lambda_1 \neq 0$ without loss of generality)

$$(23) \quad Y_1 = \lambda'X, Y_2 = X_2, \dots, Y_p = X_p.$$

Then the null hypothesis says that in the discriminant based on Y_1, \dots, Y_p , the coefficients of Y_2, \dots, Y_p are zero. Hence the test (21) applies with $q = 1$. We need the values of D_1^2 and D_p^2 based on the Y values. Since D^2 is invariant under a linear transformation

$$(24) \quad D_p^2(Y_1, \dots, Y_p) = D_p^2(X_1, \dots, X_p) = (n_1 + n_2 - 2)d' S^{-1} d$$

$$(25) \quad D_1^2(Y_1) = (n_1 + n_2 - 2)(\lambda' d)^2 / \lambda' S \lambda$$

so that D_1^2 and D_p^2 are expressed in terms of the statistics d and S defined in (10) and (11) in terms of the original variables (X_1 ,

..., X_p). Substituting the values (24), (25) for D_p^2 and D_1^2 we obtain the statistic

$$(26) \quad \frac{n_1+n_2-p-1}{p-1} \frac{cD_p^2-cD_1^2}{1+cD_1^2}$$

which is an F statistic on $(p-1)$ and (n_1+n_2-p-1) d.f.

If we are using α probability level of significance, the test (26) can be written as

$$(27) \quad \frac{(\lambda'd)^2}{\lambda'S\lambda} \geq \frac{c(n_1+n_2-p-1)D_p^2-(p-1)F_\alpha}{c(p-1)F_\alpha+c(n_1+n_2-p-1)}$$

where F_α is the upper α probability value of F on $(p-1)$ and (n_1+n_2-p-1) d.f. The inequality (27) provides a cone with vertex at the origin, within which the direction vector λ of the coefficients of the true discriminant function will lie with probability $1-\alpha$.

4.2. Test for a given ratio of the coefficients of two variables

Let the ratio of the coefficients of X_1 and X_2 be ρ , i.e. $\lambda_1/\lambda_2 = \rho$. For a given ρ , the discriminant function can be written

$$(28) \quad \lambda_2(\rho X_1 + X_2) + \lambda_3 X_3 + \dots + \lambda_p X_p$$

where $\lambda_2, \dots, \lambda_p$ are unknown. Equating (28) to $\delta \Sigma^{-1} X$ the null hypothesis can be written as

$$(29) \quad \delta \propto \Sigma \begin{pmatrix} \rho \\ \dots \\ b \end{pmatrix}$$

where $\rho' = (\rho, 1)$ and the vector b is unknown.

Let us make the transformation

$$(30) \quad Y_1 = \rho X_1 + X_2, \quad Y_2 = X_3, \dots, Y_{p-1} = X_p, \quad Y_p = X_1.$$

Then the null hypothesis (29) says that in the discriminant function based on Y_1, \dots, Y_p , the coefficient of Y_p is zero. Hence the test (21) applies with $q = p-1$. We need the values of $D_p^2(Y_1, \dots, Y_p) = D_p^2(X_1, \dots, X_p)$ and $D_{p-1}^2(Y_1, \dots, Y_{p-1})$. To compute D_{p-1}^2 , consider the partition of the original variable $\mathbf{X}' = (\mathbf{X}'_1; \mathbf{X}'_2)$ where $\mathbf{X}'_1 = (X_1, X_2)$ and $\mathbf{X}'_2 = (X_3, \dots, X_p)$ with the corresponding partition

$$(31) \quad \mathbf{d}' = (\mathbf{d}'_1; \mathbf{d}'_2)$$

of the sample mean differences and the partition

$$(32) \quad \mathbf{S} = \begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix}$$

of the within pooled dispersion matrix. Then

$$(33) \quad (n_1 + n_2 - 2)^{-1} D_{p-1}^2(Y_1, \dots, Y_{p-1}) \\ = \mathbf{d}'_2 \mathbf{S}_{22}^{-1} \mathbf{d}_2 + \frac{[\boldsymbol{\rho}'(\mathbf{d}_1 - \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{d}_2)]^2}{\boldsymbol{\rho}'(\mathbf{S}_{11} - \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{21}) \boldsymbol{\rho}}$$

and the variance ratio test on 1 and $(n_1 + n_2 - p - 1)$ d.f. for an assigned ρ is

$$(34) \quad \frac{n_1 + n_2 - p - 1}{1} \frac{c(D_p^2 - D_{p-1}^2)}{1 + c D_{p-1}^2}$$

where $(n_1 + n_2 - 2)^{-1} D_p^2 = \mathbf{d}' \mathbf{S}^{-1} \mathbf{d}$ and D_{p-1}^2 is as given in (33). To obtain a confidence interval for ρ , we consider the equation

$$(35) \quad \frac{n_1 + n_2 - p - 1}{1} \frac{cD_p^2 - cD_{p-1}^2}{1 + cD_{p-1}^2} = F_\alpha$$

which is quadratic in ρ . The confidence bounds for ρ are obtained by determining the roots of the equation (35).

4.3. Test for assigned ratios of the coefficients of several variables

Let ρ_1, \dots, ρ_k be the assigned ratios of the coefficients of the first k variables X_1, \dots, X_k , where some of the ρ_i may be zero. In such a case the discriminant function is of the form

$$(36) \quad \lambda_k(\rho_1 X_1 + \dots + \rho_k X_k) + \lambda_{k+1} X_{k+1} + \dots + \lambda_p X_p$$

where $\lambda_k, \dots, \lambda_p$ are unknown. Equating (36) to $\delta \Sigma^{-1} \mathbf{X}$, the null hypothesis can be written as

$$(37) \quad \delta \propto \Sigma \begin{pmatrix} \rho \\ \dots \\ \mathbf{b} \end{pmatrix}$$

where $\rho' = (\rho_1, \dots, \rho_k)$ and \mathbf{b} is unknown.

Let $\rho_k \neq 0$ without loss of generality. Then consider the transformation

$$(38) \quad \begin{aligned} Y_1 &= \rho_1 X_1 + \dots + \rho_k X_k, \\ Y_2 &= X_{k+1}, \dots, Y_{p-k+1} = X_p, \\ Y_{p-k+2} &= X_1, \dots, Y_p = X_{k-1}. \end{aligned}$$

The hypothesis (37) says that in the discriminant function based on Y_1, \dots, Y_p , the coefficients of the last $(k-1)$ variables are all zero. Hence the test (21) is applicable with the value of $q = p - k + 1$. As in the earlier case, we shall express the values of D_p^2 and D_q^2 in terms of the original variables. Consider the partition $(\mathbf{X}_1, \mathbf{X}_2)$ of \mathbf{X}' where \mathbf{X}_1 consists of the components X_1, \dots, X_k and \mathbf{X}_2 of the rest. Corresponding to such a partition of the random variable, we have the following partition of the sample mean difference and the within S.P. matrix.

$$(39) \quad \mathbf{d}' = (\mathbf{d}'_1 : \mathbf{d}'_2), \quad \mathbf{S} = \begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix}.$$

Then

$$(40) \quad (n_1 + n_2 - 2)^{-1} D_p^2 = \mathbf{d}' \mathbf{S}^{-1} \mathbf{d}$$

$$(41) \quad (n_1+n_2-2)^{-1} D_q^2 = \mathbf{d}'_2 S_{22}^{-1} \mathbf{d}_2 + \frac{[\boldsymbol{\rho}'(\mathbf{d} - S_{12} S_{22}^{-1} \mathbf{d}_2)]^2}{\boldsymbol{\rho}'(S_{11} - S_{12} S_{22}^{-1} S_{21}) \boldsymbol{\rho}}$$

substituting the expressions (40), (41) in (21), we find that the test criterion is an explicit function of $\boldsymbol{\rho}$. Hence we can test for any assigned value of $\boldsymbol{\rho}$ or determine the confidence zone of $\boldsymbol{\rho}$.

4.4. Test whether λ belongs to a given linear manifold

Let A be a $(p \times k)$ matrix providing a basis of the given manifold. Then the hypothesis $\lambda \in M(A) \implies \boldsymbol{\delta} = \Sigma A \boldsymbol{\theta}$ where $\boldsymbol{\theta}$ is a $k \times 1$ vector. The values of $\boldsymbol{\delta}$, Σ and $\boldsymbol{\theta}$ are unknown but the hypothesis only specifies a relationship among them through the known matrix A . Consider the transformation

$$(42) \quad \mathbf{Y}_1 = A'X, \quad \mathbf{Y}_2 = B'X$$

where B is chosen such that $B'\Sigma A = \mathbf{0}$, i.e., the linear functions in $B'X$ are uncorrelated with those in $A'X$. If $\boldsymbol{\delta} = \Sigma A \boldsymbol{\theta}$, then

$$(43) \quad E(B'X | H_1) - E(B'X | H_2) = B'\boldsymbol{\delta} = B'\Sigma A \boldsymbol{\theta} = \mathbf{0}.$$

Thus, according to statement (e) in Section 2, $\mathbf{Y}_1 = A'X$ is sufficient for discrimination between the hypotheses H_1 and H_2 or the coefficients of the components of \mathbf{Y}_2 in the discriminant function expressed in terms of $\mathbf{Y}_1, \mathbf{Y}_2$ are all zero. Hence the test (21) applies with $q = k$. The matrix B in (42) is arbitrary subject to the condition $B'\Sigma A = \mathbf{0}$, but we do not need to know B in order to evaluate the test criterion (21). We observe that

$$(44) \quad D_p^2(\mathbf{Y}_1, \mathbf{Y}_2) = D_p^2(\mathbf{X}) = (n_1+n_2-2) \mathbf{d}' S^{-1} \mathbf{d}$$

$$(45) \quad D_k^2(\mathbf{Y}_1) = D_k^2(A'X) = (n_1+n_2-2) \mathbf{d}' A (A'SA)^{-1} A \mathbf{d}$$

which depend only on A, \mathbf{d} and S , where \mathbf{d} and S are as defined in (10) and (11).

5. DISCRIMINANT FUNCTION FOR GENETIC SELECTION

Consider an observable variable X which has the decomposition

$$(46) \quad X = \gamma + \epsilon$$

in terms of two unobservable variables γ and ϵ which are uncorrelated. The variable γ denotes the conceptual genotypic measurements and ϵ denotes the environmental effects so that X may be considered as representing phenotypic measurements. Under the set up (46), Smith (1936) considered the problem of predicting a linear function $a'\gamma$ representing the genetic worth of an individual on the basis of phenotypic measurements X . Let $D(X) = \Sigma$, $D(\gamma) = \Gamma$ and $D(\epsilon) = E$. Since γ and ϵ are uncorrelated,

$$(47) \quad \Sigma = \Gamma + E$$

$$\text{Cov}(X, \gamma) = \Gamma, \quad \text{Cov}(X, a'\gamma) = \Gamma a$$

giving the regression of $a'\gamma$ on X as $\lambda'X$ where

$$(48) \quad \Sigma \lambda = \Gamma a \quad \text{or} \quad \lambda = \Sigma^{-1} \Gamma a.$$

If Σ and Γ are known, then the best predictor of $a'\gamma$ or the best selection index for $a'\gamma$ is the regression function of $a'\gamma$ on X .

As an alternative to the regression function $a'\Gamma \Sigma^{-1}X$ we may consider the *straight selection function* $a'X$, which is simpler to compute and which does not involve Γ and Σ . It is, therefore, of interest to find the conditions under which $a'\Gamma \Sigma^{-1}X$ and $a'X$ are equivalent. Now, $a'\Gamma \Sigma^{-1}X \propto a'X$ implies that there exists a constant μ such that

$$(49) \quad \Sigma^{-1} \Gamma a = \mu a \quad \text{or} \quad \Gamma a = \mu \Sigma a,$$

or a is an eigen vector of the determinantal equation

$$(50) \quad |\Gamma - \mu \Sigma| = 0,$$

The condition (49) is the same as,

$$(51) \quad E\mathbf{a} = \frac{\mathbf{a}'E\mathbf{a}}{\mathbf{a}'\Sigma\mathbf{a}}\Sigma\mathbf{a},$$

writing

$$\Gamma = \Sigma - E.$$

Let a_1 , the first component of \mathbf{a} , be non-zero. Then we consider the transformation from X to Y ,

$$(52) \quad Y_1 = \mathbf{a}'X, Y_2 = X_2, \dots, Y_p = X_p$$

with the corresponding decomposition $Y = \gamma^* + \epsilon^*$ and the dispersion matrices Σ^* and E^* . In terms of Σ^* and E^* , the condition (49) is equivalent to

$$(53) \quad \frac{\Sigma_{j1}^*}{\Sigma_{11}^*} = \frac{E_{j1}^*}{E_{11}^*}, \quad j = 2, \dots, p.$$

i.e., the regression coefficient of ϵ_j^* on ϵ_1^* is the same as that of Y_j on Y_1 . We shall consider a test of the hypothesis (53) in Section 5.1 on the basis of independent estimators of Σ^* and E^* having Wishart distributions.

A more general hypothesis of interest is the assignment of the ratios of the coefficients of a subset of the phenotypic measurements X_i in the selection index (regression function). In such a case the relationship between Σ and Γ can be written as

$$(54) \quad \Gamma\mathbf{a} = \mu \Sigma \begin{pmatrix} \mathbf{p} \\ \dots \\ \mathbf{b} \end{pmatrix}$$

where \mathbf{p} is the vector of assigned ratios of the coefficients of the (say) first q variables and \mathbf{b} is unknown. There does not seem to be a simple test of the hypothesis (54). However, one can determine the likelihood ratio test based on independent estimators of Σ and E and thus provide an asymptotic test with the usual chi-square approximation.

5.1. A test for the hypothesis (53)

Let T and W be independent random matrices such that

$$(55) \quad T \sim \omega_p(\Sigma^*, k)$$

$$(56) \quad W \sim \omega_p(E^*, m)$$

where $\omega(A, b)$ represents Wishart distribution of a $p \times p$ random matrix with the hypothetical matrix A and degrees of freedom b . Given T_{11} , the conditional distributions

$$(57) \quad \left(\frac{T_{12}}{T_{11}}, \dots, \frac{T_{1p}}{T_{11}} \right) \sim \mathcal{N}_p \left[\nu_1, \frac{1}{T_{11}} (\Sigma_{ij \cdot 1}^*) \right]$$

$$(58) \quad (T_{ij \cdot 1}) = \left(T_{ij} - \frac{T_{i1} T_{j1}}{T_{11}} \right) \sim \omega_p[(\Sigma_{ij \cdot 1}^*), k-1]$$

are independent. Similarly

$$(59) \quad \left(\frac{E_{12}}{E_{11}}, \dots, \frac{E_{1p}}{E_{11}} \right) \sim \mathcal{N}_p \left[\nu_2, \frac{1}{E_{11}} (E_{ij \cdot 1}^*) \right]$$

$$(60) \quad (E_{ij \cdot 1}) = \left(E_{ij} - \frac{E_{i1} E_{j1}}{E_{11}} \right) \sim \omega_p[(E_{ij \cdot 1}^*), m-1]$$

are independent. The hypothesis (53) under test is the same as the hypothesis

$$(61) \quad \nu_1 = \nu_2$$

i.e., the means of the two normal distributions (57) and (59) are equal. Since estimates of the dispersion matrices are available, it is possible to provide an appropriate test of the hypothesis (61).

If $(\Sigma_{ij \cdot 1}^*), (E_{ij \cdot 1}^*)$, the true residual dispersion matrices, are equal, then we have the standard test based on the likelihood ratio

$$(62) \quad \Lambda = \frac{|(T_{ij \cdot 1}) + (W_{ij \cdot 1})| |(T_{11} + W_{11})|}{|(T_{ij}) + (W_{ij})|}$$

where

$$(63) \quad \frac{m+k-p}{p-1} \quad \frac{1-\Lambda}{\Lambda}$$

has an F distribution on $(p-1)$ and $(m+k-p)$ d.f.

If the residual dispersion matrices $(\Sigma_{ij \cdot 1}^*)$ and $(E_{ij \cdot 1}^*)$ are not equal, then we have the following approximate test. The differences

$$(64) \quad d_1 = \frac{T_{12}}{T_{11}} - \frac{W_{12}}{T_{11}}, \dots, d_{p-1} = \frac{T_{1p}}{T_{11}} - \frac{W_{1p}}{W_{11}}$$

have the estimated dispersion matrix

$$(65) \quad (C_{ij}) = \frac{1}{(k-1)T_{11}}(T_{ij \cdot 1}) + \frac{1}{(m-1)W_{11}}(W_{ij \cdot 1}).$$

If (C^{ij}) is the reciprocal of (C_{ij}) , then the statistic for testing the significance of the differences d_1, \dots, d_{p-1} is

$$(66) \quad \Sigma \Sigma C^{ij} d_i d_j$$

which can be used approximately as a chi square on $(p-1)$ d.f.

Note : The choice of a test of the hypothesis (53) or (61) depends on whether the residual dispersion matrices $(\Sigma_{ij \cdot 1}^*)$ and $(E_{ij \cdot 1}^*)$ are equal or not. Since we have the estimates $(T_{ij \cdot 1})$ and $(W_{ij \cdot 1})$ the hypothesis of equality $(\Sigma_{ij \cdot 1}^*) = (E_{ij \cdot 1}^*)$ can be subjected to a suitable test.

In problems of genetic selection we have certain families (or individual lines) out of which some have to be selected on the basis of performance of individuals within family. We obtain observations on a certain number of characteristics from each of n individuals in a family. These observations (after transformation to the Y_t variables as in (51)) provide an analysis of dispersion (S.P. matrices) as between and within families, which correspond to the T and W matrices used in the test. For a numerical application of the tests (63) and (66) and further discussion, the reader is referred to a paper by the author (Rao, 1953).

References

- Fisher, R. A. (1940). "The Precision of Discriminant Functions," *Ann. Eugenics*, **10**, 422-429.
- Giri, N. (1964). "On the Likelihood Ratio Test of a Normal Multivariate Testing Problem," *Ann. Math. Statist.*, **35**, 181-189.
- Rao, C. R. (1946). "Tests with Discriminant Functions in Multivariate Analysis," *Sankhyā*, **7**, 407-414.
- Rao, C. R. (1948). "Tests of Significance in Multivariate Analysis," *Biometrika*, **35**, 58-79.
- Rao, C. R. (1949). "On Some Problems Arising Out of Discrimination with Multiple Characters," *Sankhyā*, **9**, 343-364.
- Rao, C. R. (1952). *Advanced Statistical Methods in Biometric Research*, John Wiley and Sons, New York.
- Rao, C. R. (1953). "Discrimination Functions for Genetic Differentiation and Selection," *Sankhyā*, **12**, 229-246.
- Rao, C. R. (1965). *Linear Statistical Inference with Applications*, John Wiley and Sons, New York.
- Smith, Fairfield H. (1936). "A Discriminant Function for Plant Selection," *Ann. Eugenics*, **7**, 240.

(Received Jan. 18, 1966.)