

# GOODNESS OF FIT OF AN ASSIGNED SET OF SCORES FOR THE ANALYSIS OF ASSOCIATION IN A CONTINGENCY TABLE

A. M. KSHIRSAGAR\*

(Received Sept. 8, 1969)

## Abstract

The problem of association between two attributes in a  $p \times q$  contingency table can be looked upon as the problem of relationship between two vector variables  $\mathbf{x}$  and  $\mathbf{y}$ . If there is only one true non-zero canonical correlation between  $\mathbf{x}$  and  $\mathbf{y}$ , the association between the two attributes is of rank 1 and in this case, one set of scores is adequate to describe the association completely; these scores are nothing but the coefficients in the canonical variates corresponding to the true non-zero canonical correlation. Given a set of hypothetical scores  $\alpha_1, \alpha_2, \dots, \alpha_p$  for the rows, one is interested in testing their goodness of fit. Tests for this are suggested in this paper. For obtaining these tests, a preliminary result about direction and collinearity factors in discriminant analysis, when  $S$  irrelevant variables are eliminated, is needed. This is derived in part one of this paper.

## 1. Relationship between two vectors

The problem of association between two vectors  $\mathbf{x}(p \times 1)$  and  $\mathbf{y}(q \times 1)$  arises in regression analysis, multivariate analysis of variance, discriminant analysis and in contingency table analysis. This relationship has different interpretations and implications in these fields, but in each case it can be expressed in terms of canonical correlations and canonical variables. The canonical correlations  $r_1, r_2, \dots, r_p$  ( $p \leq q$ ) in a sample, are the roots of the equation

$$(1.1) \quad \begin{vmatrix} -r^2 C_{xx} & C_{xy} \\ C_{yx} & -r^2 C_{yy} \end{vmatrix} = 0$$

---

\* This research was sponsored by the Office of Naval Research, Contract No. N00014-68-A-0515, Project No. NR042-260. Reproduction in whole or in part is permitted for any purpose of the United States Government.

and the canonical variables corresponding to  $r_i^2$  are  $l_{(i)}'x$  and  $m_{(i)}'y$  ( $i=1, 2, \dots, p$ ), where the column vectors  $l_{(i)}$ ,  $m_{(i)}$  satisfy the equation

$$(1.2) \quad \left[ \begin{array}{c|c} -r_i^2 C_{xx} & C_{xy} \\ \hline C_{yx} & -r_i^2 C_{yy} \end{array} \right] \left[ \begin{array}{c} l_{(i)} \\ m_{(i)} \end{array} \right] = 0.$$

Here

$$(1.3) \quad C = \left[ \begin{array}{c|c} C_{xx} & C_{xy} \\ \hline C_{yx} & C_{yy} \end{array} \right]$$

is the matrix of the corrected sum of squares (s.s.) and sum of products (s.p.) of observations on  $x$  and  $y$  and is based on  $n$  degrees of freedom (d.f.). The true or population canonical correlations are denoted by  $\rho_1, \rho_2, \dots, \rho_p$ . If all the  $\rho$ 's are null, there is no association between  $x$  and  $y$  and under the assumption of normality, this is tested by using any one of the following criteria:

$$(1.4) \quad \text{Wilks's [9] } A \text{ criterion; } A = |A|/|A+B|$$

or

$$(1.5) \quad \text{Pillai's [7] criterion; } t_r(A+B)^{-1}B$$

where

$$(1.6) \quad B = C_{xy}C_{yy}^{-1}C_{yx}, \quad A = C_{xx} - C_{xy}C_{yy}^{-1}C_{yx}, \quad A+B = C_{xx}.$$

If only  $\rho_1 \neq 0$  but  $\rho_2 = \dots = \rho_p = 0$ , we say that the association between  $x$  and  $y$  is of rank 1. In this case, the entire association can be adequately described by the canonical variates corresponding to  $\rho_1$ . In discriminant analysis, this means that the means of  $q+1$  groups to be discriminated are collinear and a single discriminant function is adequate. Testing the goodness of fit of a single discriminant function  $\alpha'x = \alpha_1x_1 + \dots + \alpha_px_p$ , in this context, means that one wishes to test (1) whether  $\alpha'x$  agrees with the true canonical variate corresponding to  $\rho_1$  and (2) whether one linear function is adequate at all to describe completely the relationship between  $x$  and  $y$ . (1) is called the 'direction' aspect and (2) is called the collinearity aspect of the goodness of fit test. Bartlett [1] and Williams [11] derived tests for this purpose by factorizing Wilks'  $A$  as

$$(1.7) \quad A = A_1 \cdot A_2 \cdot A_3.$$

where (see Kshirsagar [4])

$$(1.8) \quad A_1 = \alpha'A\alpha/\alpha(A+B)\alpha$$

$$(1.9) \quad A_2 = 1 - \frac{\alpha'B(A+B)^{-1}B\alpha/\alpha'B\alpha}{\alpha'A\alpha/\alpha'(A+B)\alpha}$$

$$(1.10) \quad A_3 = A/A_1A_2$$

$A_2$  is the direction factor and  $A_3$  is the 'partial' collinearity factor. Bartlett has given an alternative factorization also viz,

$$(1.11) \quad A = A_1A_4A_5,$$

where

$$(1.12) \quad A_4 = A \left\{ 1 + \frac{\alpha'BA^{-1}Ba}{\alpha'Ba} \right\}$$

$$(1.13) \quad A_5 = A/A_1A_4$$

$A_4$  is the collinearity factor and  $A_5$  is the 'partial' direction factor. A statistic  $t$  is said to have a  $A(n, p, q)$  distribution, if it is distributed as  $\prod_{i=1}^p U_i$  where  $U_i$ 's are independent and  $U_i$  has the distribution

$$(1.14) \quad \text{Const. } U_i^{(n-q-i-1)/2} (1-U_i)^{(q-2)/2} dU_i$$

Bartlett [1] has shown that, in this case,

$$(1.15) \quad - \left\{ n - \frac{1}{2}(p+q+1) \right\} \log_e t$$

has a  $\chi^2$  distribution with  $pq$  d.f. in large samples. If the null-hypothesis of goodness of fit of  $\alpha'x$  is true, he shows that  $A_2$  is a  $A(n-1, 1, p-1)$  and  $A_3$  is an independent  $A(n-2, q-1, p-1)$ . Alternatively  $A_4$  is  $A(n-1, q-1, p-1)$  and  $A_5$  is an independent  $A(n-q, 1, p-1)$ . Briefly,  $A_2$  is based on  $p-1$  d.f.,  $A_3$  on  $(p-1)(q-1)$  d.f.,  $A_4$  on  $(p-1)(q-1)$  d.f. and  $A_5$  on  $(p-1)$  d.f.

The author [5] has shown, that the other criterion  $t_r B(A+B)^{-1}$  can also be partitioned, analogous to this factorization of  $A$ , as

$$(1.16) \quad nt_r B(A+B)^{-1} = \gamma_1 + \gamma_2 + \gamma_3$$

where

$$(1.17) \quad \begin{aligned} \frac{1}{n} \gamma_1 &= \frac{\alpha'Ba}{\alpha'(A+B)\alpha}, & \frac{1}{n} \gamma_2 &= \frac{\alpha'B(A+B)^{-1}Ba}{\alpha'Ba} - \frac{1}{n} \gamma_1, \\ \frac{1}{n} \gamma_3 &= t_r B(A+B)^{-1} - \frac{1}{n} \gamma_1 - \frac{1}{n} \gamma_2. \end{aligned}$$

Here  $\gamma_2$  is the 'direction' part and  $\gamma_3$  is the 'collinearity' part and under the null hypothesis of goodness of fit of  $\alpha'x$ , they are distributed independently as  $\chi^2$  with  $p-1$  d.f. and  $\chi^2$  with  $(p-1)(q-1)$  d.f. respectively, in large samples.

## 2. Elimination of irrelevant variables

In some situations, it so happens that one is interested in studying the relationship between—not  $x$  and  $y$ —but between *residual* variates  $z$  and  $w$ , where the latter are obtained from  $x$  and  $y$  by eliminating the first  $S$  sample canonical variables. These first  $S$  sample canonical variables are known apriori to be irrelevant and are therefore to be excluded. Let  $L_1x$  and  $M_1y$ , where

$$(2.1) \quad L_1 = [l_{(1)} | l_{(2)} | \cdots | l_{(S)}]'$$

$S \times p$

and

$$(2.2) \quad M_1 = [m_{(1)} | m_{(2)} | \cdots | m_{(S)}]'$$

$S \times q$

be the first  $S$  canonical variables. On account of (1.2), we find

$$(2.3) \quad C_{xx}L_1'R = C_{xy}M_1'$$

where  $R$  is the  $S \times S$  diagonal matrix of  $r_i^2$  ( $i=1, 2, \dots, S$ ). One can also show from (1.2) that

$$(2.4) \quad BL_1' = (A+B)L_1'R$$

and

$$(2.5) \quad AL_1' = (A+B)L_1'(I-R).$$

Let  $L_2$  be a  $(p-S) \times p$  matrix and  $M_2$  a  $(q-S) \times q$  matrix such that

$$(2.6) \quad L_2C_{xx}L_1' = 0, \quad M_2C_{yy}M_1' = 0$$

i.e.,  $L_2x$  and  $L_1x$  are uncorrelated and so also are  $M_2y$  and  $M_1y$ . From (2.4), (2.5), (2.6) it can be seen easily that

$$(2.7) \quad L_2B_1' = 0, \quad L_2AL_1' = 0.$$

We can now take

$$(2.8) \quad z = L_2x, \quad w = M_2y$$

as our residual variables, after eliminating  $L_1x$  and  $M_1y$ . We now want to test the goodness of fit of an assigned function  $\alpha'x$  for the relationship between  $z$  and  $w$ . It is obvious that this assigned function must be so chosen that it is uncorrelated with the eliminated variables  $L_1x$ ; in other words, it must be a linear function, say  $k'z$  of  $z$  alone. If so,  $k$  will satisfy

$$(2.9) \quad \alpha = L_2'k.$$

We define  $C_{zz}$ ,  $C_{zw}$ ,  $C_{ww}$  in the same way as in (1.3) and then  $A_z$ ,  $B_z$  and  $A_z + B_z$  as in (1.6). We can, then easily write down the new direction and collinearity factors  $A_{2z}$ ,  $A_{3z}$ ,  $A_{4z}$ ,  $A_{5z}$  or  $\gamma_{2z}$ ,  $\gamma_{3z}$  etc. by using  $k'z$  instead of  $\alpha'x$  and  $A_z$ ,  $B_z$  for  $A$  and  $B$  in (1.9), (1.10), (1.12), (1.13) and (1.17). We must also replace  $n$  by  $n-S$ ,  $p$  by  $p-S$  and  $q$  by  $q-S$  as  $S$  variables have been eliminated from  $x$  and from  $y$ . We, however, wish to express these test statistics in terms of our old matrices  $A$ ,  $B$  and the assigned vector  $\alpha$ . This can be done as below:

From (2.3) and (2.6),

$$(2.10) \quad L_2 C_{xy} M'_1 = 0.$$

Hence

$$(2.11) \quad \begin{aligned} B_z &= C_{zw} C_{ww}^{-1} C_{wz} = L_2 C_{zw} C_{ww}^{-1} C_{wz} L'_2 \\ &= L_2 [C_{xy} C_{yy}^{-1} C_{yx} - C_{xy} M'_1 (M_1 C_{yy} M'_1)^{-1} M_1 C_{yx}] L'_2 \\ &= L_2 B L'_2, \quad \text{on account of (2.10).} \end{aligned}$$

Also

$$(2.12) \quad \begin{aligned} A_z + B_z &= C_{zz} = L_2 C_{xx} L'_2 \\ &= L_2 (A + B) L'_2. \end{aligned}$$

Let

$$(2.13) \quad L = \begin{bmatrix} L_1 \\ L_2 \end{bmatrix}_{p-s}^s.$$

Then

$$(2.14) \quad \begin{aligned} (A + B)^{-1} &= L' (L C_{xx} L')^{-1} L \\ &= L' \left[ \begin{array}{c|c} L_1 C_{xx} L'_1 & 0 \\ \hline 0 & L_2 C_{xx} L'_2 \end{array} \right]^{-1} L \quad \text{on account of (2.6)} \\ &= \sum_{i=1}^2 L_i (L_i C_{xx} L'_i)^{-1} L_i. \end{aligned}$$

Hence

$$(2.15) \quad \begin{aligned} k' B_z (A_z + B_z)^{-1} B_z k &= k' L_2 B L'_2 (L_2 C_{xx} L'_2)^{-1} L_2 B L'_2 k \\ &= \alpha' B \{ (A + B)^{-1} - L'_1 (L_1 C_{xx} L'_1)^{-1} L_1 \} B \alpha \\ &= \alpha' B (A + B)^{-1} B \alpha \end{aligned}$$

on account of (2.7). In exactly the same way, it can be shown that

$$(2.16) \quad k' B_z A_z^{-1} B_z k = \alpha' B A^{-1} B \alpha.$$

Note also that

$$\begin{aligned}
 (2.17) \quad A &= \frac{|A|}{|A+B|} = \frac{|LAL'|}{|L(A+B)L'|} = \frac{|L_1AL'_1||L_2AL'_2|}{|L_1(A+B)L'_1||L_2(A+B)L'_2|} \\
 &= \prod_{i=1}^s (1-r_i^2) \frac{|A_s|}{|A_s+B_s|} \\
 &= A_s \prod_{i=1}^s (1-r_i^2),
 \end{aligned}$$

on account of (2.4) and (2.5). Also

$$(2.18) \quad k'B_s k = k'L_2 B L'_2 k = a' B a$$

and

$$(2.19) \quad k'A_s k = k'L_2 A L'_2 k = a' A a.$$

Substituting (2.15), (2.16), (2.17), (2.18) and (2.19) in  $A_{2s}$ ,  $A_{3s}$ ,  $A_{4s}$ ,  $A_{5s}$ ,  $\gamma_{2s}$  and  $\gamma_{3s}$ , we find that these 'new' direction and collinearity factors or parts are exactly the same as the old ones vis.  $A_2$ ,  $A_3$ ,  $A_4$ ,  $A_5$ ,  $\gamma_2$ ,  $\gamma_3$  for  $x$  and  $y$ , except that  $A$  must be changed to  $A / \prod_{i=1}^s (1-r_i^2)$ ,  $n$  to  $n-S$ ,  $p$  to  $p-S$  and  $q$  to  $q-S$ .

We are now in a position to apply these results to the analysis of a contingency table, which we do in Section 3 of this paper.

### 3. Association between two attributes

Consider a  $p \times q$  contingency table with the rows corresponding to  $p$  categories  $a_1, a_2, \dots, a_p$  of an attribute 'a' and columns to  $q$  categories  $b_1, b_2, \dots, b_q$  of another attribute 'b'. Let  $n_{ij}$  ( $i=1, \dots, p$ ,  $j=1, \dots, q$ ) be the frequency in the  $(i, j)$ th cell. Let  $n_{i.}$  ( $i=1, \dots, p$ ) be the row totals and  $n_{.j}$  ( $j=1, \dots, q$ ) be the column totals. Let  $n = \sum_i n_{i.} = \sum_j n_{.j}$  be the total frequency. We define

$$(3.1) \quad N = [n_{ij}] \quad (i=1, \dots, p, j=1, \dots, q)$$

$$(3.2) \quad D_1 = \begin{bmatrix} n_{1.} & & & \\ & n_{2.} & & \\ & & \ddots & \\ & & & n_{p.} \end{bmatrix}$$

and

$$(3.3) \quad D_2 = \begin{bmatrix} n_{.1} & & & \\ & n_{.2} & & \\ & & \ddots & \\ & & & n_{.q} \end{bmatrix}.$$

The problem of assigning optimum scores to the rows and columns has received considerable attention in the literature (Yates [13], Fisher [2], [3], Maung [6], Bartlett [1], Williams [10]). It has been shown that the vectors of optimum scores  $\xi$  and  $\eta$  corresponding to the  $a$ 's and  $b$ 's are obtainable from the equations

$$(3.4) \quad \left[ \begin{array}{c|c} -r^2 D_1 & N \\ \hline N' & -r^2 D_2 \end{array} \right] \begin{bmatrix} \xi \\ \eta \end{bmatrix} = 0.$$

If we, therefore, consider two vector variables  $x(p \times 1)$  and  $y(q \times 1)$ , with the variance-covariance matrix,

$$(3.5) \quad \left[ \begin{array}{c|c} D_1 & N \\ \hline N' & D_2 \end{array} \right]$$

it is evident from (3.4) that  $\xi'x$  and  $\eta'y$  are nothing but the canonical variates corresponding to the canonical correlation  $r^2$ . In other words, the association between two sets of categories in a contingency table can also be looked upon as a problem of relationship between two vector variables. In general, one set of scores will not be adequate to describe the association between ' $a$ ' and ' $b$ ' completely. We shall need as many sets of scores, as there are significant canonical correlations between  $x$  and  $y$ . If, however, only one canonical correlation is significant, one set of scores will be adequate. We say, in this case, that the association is 'linear' or of rank 1.

In the notation of Section 1,  $C_{xx}=D_1$ ,  $C_{xy}=N$  and  $C_{yy}=D_2$  and hence

$$(3.6) \quad B = ND_2^{-1}N' = \left[ \sum_{j=1}^q \frac{n_{ij}n_{hj}}{n_{.j}} \right] \quad (i, h=1, \dots, p)$$

$$(3.7) \quad A = D_1 - ND_2^{-1}N'$$

$$(3.8) \quad A + B = D_1.$$

We shall denote by  $A_0$ ,  $B_0$  and  $D_1^0$ , the matrices obtained from  $A$ ,  $B$  and  $D_1$  respectively, by deleting the last row and the last column. It is readily observed from (1.1) that  $r^2=1$  is a canonical correlation between  $x$  and  $y$ , the corresponding canonical variates being  $x_1 + \dots + x_p$  and  $y_1 + \dots + y_q$ . Obviously, these are irrelevant to our present problem

of assigning scores to the  $a$ 's and  $b$ 's. We must therefore eliminate these variables and study the residual variates  $z$  and  $w$  as in Section 2. By taking regression on  $\sum_1^p x_i$  and  $\sum_1^q y_j$ , we can take the new variables to be

$$(3.9) \quad z_i = x_i - \frac{n_{i.}}{n} (x_1 + \cdots + x_p); \quad i=1, 2, \dots, p-1$$

and

$$(3.10) \quad w_j = y_j - \frac{n_{.j}}{n} (y_1 + \cdots + y_q); \quad j=1, 2, \dots, q-1.$$

We can easily calculate  $C_{zz}$ ,  $C_{zw}$ ,  $C_{ww}$  and hence  $A_z$ ,  $B_z$  from these. They turn out to be

$$(3.11) \quad A_z = A_0, \quad B_z = B_0 - \frac{1}{n} d_0 d'_0,$$

where,

$$(3.12) \quad d_0 = \begin{bmatrix} n_{1.} \\ n_{2.} \\ \vdots \\ n_{p-1.} \end{bmatrix}, \quad d = \begin{bmatrix} d_0 \\ n_p \end{bmatrix}.$$

Note that

$$(3.13) \quad |A_z + B_z| = \left| D_1^0 - \frac{1}{n} d_0 d'_0 \right| = n_{1.} n_{2.} \cdots n_{p.} / n.$$

Consider now the problem of testing the goodness of fit of a set of hypothetical scores  $\alpha_1, \alpha_2, \dots, \alpha_p$  for the rows. The null hypothesis here comprises of two aspects (i) the association between  $a$ 's and  $b$ 's is linear and (ii) the true scores corresponding to this linear association are  $\alpha_1, \alpha_2, \dots, \alpha_p$ . (i) is the collinearity part and (ii) is the direction part of the null hypothesis.

Since we have eliminated  $\sum_1^p x_i$ , the assigned function  $\alpha'x$ , where  $\alpha' = [\alpha_1, \dots, \alpha_p]$ , must—as we noticed in Section 2—be uncorrelated with  $\sum x_i$  i.e.

$$(3.14) \quad d'\alpha = 0.$$

On account of this,  $\alpha'x$  can be written, in terms of the residual variables  $z$  as  $k'z$ , where

$$(3.15) \quad k' = [\alpha_1 - \alpha_p, \alpha_2 - \alpha_p, \dots, \alpha_{p-1} - \alpha_p].$$



We cannot obtain the 'direction' and 'collinearity' factors straightway from Section 2, in this case, because they involve  $|A|$ ,  $A^{-1}$  and these do not exist in the present case, as

$$(3.16) \quad Ae=0.$$

where

$$(3.17) \quad e' = [1, 1, \dots, 1] = [e'_0 | 1]_{1 \times p}$$

and thus  $A$  is singular. We must, therefore find the direction and collinearity factors by directly working with  $A_z$  and  $B_z$ , especially for  $A_3$ ,  $A_4$  and  $A_5$ .  $A_2$  does not involve  $A^{-1}$  and can be written directly.

Partition  $A$ ,  $B$  and  $\alpha$  as

$$(3.18) \quad B = \left[ \begin{array}{c|c} B_0 & t \\ \hline t' & b_{pp} \end{array} \right]_{p-1}^{p-1}, \quad A = \left[ \begin{array}{c|c} A_z & -t \\ \hline -t' & a_{pp} \end{array} \right], \quad \alpha = \left[ \begin{array}{c} \alpha_0 \\ \hline \alpha_p \end{array} \right]_{p-1}^{p-1}.$$

From (3.16) and (3.18)

$$(3.19) \quad A_z e_0 = t.$$

Let

$$(3.20) \quad Ba = f = \left[ \begin{array}{c} f_1 \\ \vdots \\ f_p \end{array} \right] = \left[ \begin{array}{c} f_0 \\ \hline f_p \end{array} \right]$$

so that

$$f_i = \sum_{h=1}^p \sum_{j=1}^q n_{ij} n_{hj} \alpha_h / n_{.j}.$$

Then  $e'f = e'Ba = d'\alpha = 0$  on account of (3.14). The equations

$$(3.21) \quad Ag = f$$

in the  $p$  unknowns  $g' = [g_1, \dots, g_p] = [g'_0 | g_p]$  are soluble. A solution is

$$(3.22) \quad g = A^- f$$

where  $A^-$  is a pseudo inverse of  $A$  (see Rao [8]). But (3.21) and (3.18) yield

$$A_z g_0 - g_p t = f_0$$

or

$$g_0 - g_p A_z^{-1} t = A_z^{-1} f_0$$

or

$$(3.23) \quad g_0 - g_p e_0 = A_z^{-1} f_0 \quad (\text{on account of (3.19)}).$$

Also observe that

$$\begin{aligned} B_z k &= \left( B_0 - \frac{1}{n} d_0 d_0' \right) (\alpha_0 - \alpha_p e_0) \\ &= B_0 \alpha_0 + \alpha_p t = f_0, \quad \text{on account of (3.20).} \end{aligned}$$

Hence

$$\begin{aligned} (3.24) \quad k' B_z A_z^{-1} B_z k &= f_0' A_z^{-1} f_0 \\ &= f_0' (g_0 - g_p e_0), \quad \text{from (3.23)} \\ &= f' g \\ &= f' A^{-1} f \\ &= \alpha' B A^{-1} B \alpha. \end{aligned}$$

Hence  $\Lambda_{4z}$  and  $\Lambda_{5z}$  are the same as  $\Lambda_4$  and  $\Lambda_5$ , even if  $A^{-1}$  does not exist, provided we use  $A^{-}$  for  $A^{-1}$ . Hence the direction and collinearity factors are

$$\begin{aligned} (3.25) \quad \Lambda_{2z} = \Lambda_2 &= 1 - \frac{\alpha' B (A+B)^{-1} B \alpha / \alpha' B \alpha}{\alpha' A \alpha / \alpha' (A+B) \alpha} \\ &= \frac{1 - \sum_{i=1}^p \frac{1}{n_i} f_i^2 / \sum_{i=1}^p \alpha_i f_i}{1 - \sum_{i=1}^p f_i \alpha_i / \sum_{i=1}^p n_i \alpha_i^2}. \end{aligned}$$

But

$$\begin{aligned} (3.26) \quad \Lambda_{3z} &= \Lambda_z / \Lambda_{1z} \Lambda_{2z} \\ &= |A_z| / |A_z + B_z| \Lambda_{1z} \Lambda_{2z} \\ &= \frac{n |A_z|}{(n_1 n_2 \cdots n_p) \left( 1 - \sum_{i=1}^p \frac{1}{n_i} f_i^2 / \sum_{i=1}^p \alpha_i f_i \right)} \end{aligned}$$

$$(3.27) \quad \Lambda_{4z} = \frac{n |A_z|}{(n_1 n_2 \cdots n_p)} \left\{ 1 + \frac{\sum_{i=1}^p f_i g_i}{\sum_{i=1}^p f_i \alpha_i} \right\}$$

and

$$\Lambda_{5z} = \frac{\sum_{i=1}^p n_i \alpha_i^2}{\sum n_i \alpha_i^2 - \sum f_i \alpha_i} \cdot \frac{\sum \alpha_i f_i}{\sum n_i \alpha_i^2 - \sum f_i \alpha_i + \sum f_i g_i}$$

$\Lambda_{2z}$  is  $\Lambda(n-2, 1, p-2)$ .  $\Lambda_{3z}$  is  $\Lambda(n-3, q-2, p-2)$ ,  $\Lambda_{4z}$  is  $\Lambda(n-2, q-2, p-2)$

and  $A_{3z}$  is  $A(n-q, 1, p-2)$ . Under the null hypothesis, therefore, from (1.15)

$$-\left\{(n-2)-\frac{1}{2}(1+p-2+1)\right\}\log_e A_{2z} \quad \text{is } \chi^2 \text{ with } p-2 \text{ d.f.}$$

and

$$-\left\{(n-3)-\frac{1}{2}(q-2+p-2+1)\right\}\log_3 A_{3z} \quad \text{is } \chi^2 \text{ with } (p-2)(q-2) \text{ d.f.}$$

They pertain to the direction and collinearity aspects respectively of the goodness of fit test. We can write down similar results for  $A_{4z}$  and  $A_{5z}$  of the alternative factorization.

The validity of such tests based on the assumption of normality of  $\mathbf{x}$ , for application to discrete data of contingency tables is questionable. Williams [10] justifies this by an appeal to asymptotic normality and also by the result that elementary symmetric functions of  $r_i^2$  have the same expected values in contingency tables, as for normally distributed  $\mathbf{x}$ . The above tests therefore are approximate but, as pointed out by Williams [10], adequate for practical purposes, especially when  $n$  is large.

In the above analysis, we have used Wilks'  $A$  as the over-all criterion for testing the association between the two attributes 'a' and 'b'. However, the usual practice, while dealing with contingency tables, is to use the  $\chi^2$  test viz., if there is no association

$$(3.28) \quad \gamma = n \left( \sum_{i=1}^p \sum_{j=1}^q n_{ij}^2 / (n_{i.} n_{.j}) - 1 \right)$$

has a  $\chi^2$  distribution with  $(p-1)(q-1)$  d.f. But (3.28) is nothing but

$$(3.29) \quad nt_r B_r (A_r + B_r)^{-1}$$

or Pillai's criterion. This can be written, more simply as

$$(3.30) \quad n[t_r B(A+B)^{-1} - 1].$$

The quantity subtracted in the larger bracket of (3.30) is the eliminated root  $r^2=1$ , corresponding to  $\sum_1^p x_i$ .

The 'direction' and 'collinearity' parts,  $\gamma_{2z}$  and  $\gamma_{3z}$  of this over-all  $\chi^2$  of (3.29) are easily seen, from (1.16), (1.17) and (3.30), to be

$$(3.31) \quad \gamma_{2z} = n \left[ \frac{\sum_1^p \frac{1}{n_{i.}} f_i^2}{\sum_1^p \alpha_i f_i} - \frac{\sum \alpha_i f_i}{\sum n_{i.} \alpha_i^2} \right], \quad \text{d.f. } (p-2)$$

and

$$(3.32) \quad \gamma_{32} = \gamma - \frac{n \sum f_i^2 / n_i}{\sum \alpha_i f_i}, \quad \text{d.f. } (p-2)(q-2).$$

Under the null hypothesis, they have  $\chi^2$  distributions, for large  $n$ .

Williams [10] has given the test of goodness of fit of a set of hypothetical scores, only for the particular cases  $q=2, 3$ . We have here the tests for any  $p$  and  $q$ . Further, we have also given the tests, based on the alternative criterion (Pillai), which in this case is the usual  $\chi^2$  of a contingency table and is thus more in tune with the classical method of partitioning an over-all  $\chi^2$ , corresponding to suspected sources of association.

SOUTHERN METHODIST UNIVERSITY

### REFERENCES

- [1] Bartlett, M. S. (1951). The goodness of fit of a single hypothetical discriminant function in the case of several groups, *Ann. Eugen.*, **16**, 199.
- [2] Fisher, R. A. (1940). The precision of discriminant functions, *Ann. Eugen.*, **10**, 422.
- [3] Fisher, R. A. (1950). *Statistical Methods for Research Workers*, 11th ed., Edinburgh, Oliver and Boyd.
- [4] Kshirsagar, A. M. (1964). Distribution of the direction and collinearity factors in discriminant analysis, *Proc. Camb. Philos. Soc.*, **60**, 217.
- [5] Kshirsagar, A. M. (1969). Correlation between two vector variables, *Jour. Roy. Statist. Soc.*, B, **31**, 477-485.
- [6] Muang, K. (1941). Measurement of associations in a contingency table with special reference to the pigmentation of hair and eye colour of Scottish school children, *Ann. Eugen.*, **11**, 189.
- [7] Pillai, K. C. S. (1955). Some new test criteria in multivariate analysis, *Ann. Math. Statist.*, **26**, 117.
- [8] Rao, C. Radhakrishna (1962). A note on a generalized universe of a matrix with applications to problems in mathematical statistics, *Jour. Roy. Statist. Soc.*, B, **24**, 152-158.
- [9] Wilks, S. S. (1932). Certain generalizations in the analysis of variance, *Biometrika*, **24**, 471.
- [10] Williams, E. J. (1952). Use of scores for the analysis of association in contingency tables, *Biometrika*, **39**, 274.
- [11] Williams, E. J. (1955). Significance tests for discriminant functions and linear functional relationships, *Biometrika*, **42**, 360.
- [12] Williams, E. J. (1967). The analysis of association among many variates, *Jour. Roy. Statist. Soc.*, B, **29**, 199.
- [13] Yates, F. (1948). The analysis of contingency tables with groupings based on quantitative characters, *Biometrika*, **35**, 176.