

**ON THE ANALYSIS OF VARIANCE IN CASE OF MULTIPLE
CLASSIFICATIONS WITH UNEQUAL CLASS FREQUENCIES**

BY ABRAHAM WALD¹

Columbia University

In a previous paper² the author considered the case of a single criterion of classification with unequal class frequencies and derived confidence limits for σ'^2/σ^2 where σ'^2 denotes the variance associated with the classification, and σ^2 denotes the residual variance. The scope of the present paper is to extend those results to the case of multiple classifications with unequal class frequencies.

For the sake of simplicity of notations we will derive the required confidence limits in the case of a two-way classification, the extension to multiple classifications being obvious.

Consider a two-way classification with p rows and q columns. Let y be the observed variable, and let n_{ij} be the number of observations in the i th row and j th column. Denote by $y_{ij}^{(k)}$ the k th observation on y in the i th row and j th column ($k = 1, \dots, n_{ij}$). Let the total number of observations be N . We order the N observations and let y_α be the α th observation on y in that order. Consider the variables:

$$t, t_1, \dots, t_p, v_1, \dots, v_q,$$

and denote by t_α the α th observation on t , by $t_{i\alpha}$ the α th observation on t_i and by $v_{j\alpha}$ the α th observation on v_j . The values of t_α , $t_{i\alpha}$ and $v_{j\alpha}$ are defined as follows:

$$t_\alpha = 1 \quad (\alpha = 1, \dots, N),$$

$$t_{i\alpha} = 1 \text{ if } y_\alpha \text{ lies in the } i\text{th row,}$$

$$t_{i\alpha} = 0 \text{ if } y_\alpha \text{ does not lie in the } i\text{th row,}$$

$$v_{j\alpha} = 1 \text{ if } y_\alpha \text{ lies in the } j\text{th column,}$$

$$v_{j\alpha} = 0 \text{ if } y_\alpha \text{ does not lie in the } j\text{th column.}$$

We make the assumptions

$$y_{ij}^{(k)} = x_{ij}^{(k)} + \epsilon_i + \eta_j,$$

where the variates $x_{ij}^{(k)}$, ϵ_i , η_j ($i = 1, \dots, p$; $j = 1, \dots, q$; $k = 1, \dots, n_{ij}$) are independently and normally distributed, the variance of $x_{ij}^{(k)}$ is σ^2 , the variance of ϵ_i is σ'^2 , the variance of η_j is σ''^2 , and the mean values of ϵ_i and η_j are zero.

¹ Research under a grant-in-aid from the Carnegie Corporation of New York.

² "A note on the analysis of variance with unequal class frequencies," *Annals of Math. Stat.*, Vol. 11 (1940).

Let the sample regression of y on $t, t_1, \dots, t_{p-1}, v_1, \dots, v_{q-1}$ be

$$Y = at + b_1t_1 + \dots + b_{p-1}t_{p-1} + d_1v_1 + \dots + d_{q-1}v_{q-1}.$$

We want to derive confidence limits for

$$\sigma'^2 / \sigma^2 = \lambda^2.$$

Let us introduce the notations:

$$\begin{aligned} \sum_{\alpha} t_{\alpha}t_{i\alpha} &= a_{0i} && (i = 1, \dots, p - 1), \\ \sum t_{\alpha}v_{j\alpha} &= a_{0p-1+j} && (j = 1, \dots, q - 1), \\ \sum t_{i\alpha}t_{j\alpha} &= a_{ij} && (i, j = 1, \dots, p - 1), \\ \sum t_{i\alpha}v_{j\alpha} &= a_{ip-1+j} && (i = 1, \dots, p - 1; j = 1, \dots, q - 1), \\ \sum v_{i\alpha}v_{j\alpha} &= a_{p-1+i, p-1+j} && (i, j = 1, \dots, q - 1), \\ \|c_{ij}\| &= \|a_{ij}\|^{-1} && (i, j = 0, 1, \dots, p + q - 2). \end{aligned}$$

Let the regression of $x_{ij}^{(k)}$ on $t, t_1, \dots, t_{p-1}, v_1, \dots, v_{q-1}$ be

$$X = a^*t + b_1^*t_1 + \dots + b_{p-1}^*t_{p-1} + d_1^*v_1 + \dots + d_{q-1}^*v_{q-1}.$$

The regression of $\epsilon_i + \eta_j$ on the same independent variables is evidently equal to

$$\begin{aligned} \epsilon_1t_1 + \dots + \epsilon_p t_p + \eta_1v_1 + \dots + \eta_q v_q \\ = (\eta_q + \epsilon_p)t + (\epsilon_1 - \epsilon_p)t_1 + \dots + (\epsilon_{p-1} - \epsilon_p)t_{p-1} \\ + (\eta_1 - \eta_q)v_1 + \dots + (\eta_{q-1} - \eta_q)v_{q-1}, \end{aligned}$$

since $t_p = t - t_1 - \dots - t_{p-1}$ and $v_q = t - v_1 - \dots - v_{q-1}$. Hence

$$(1) \quad b_i = b_i^* + (\epsilon_i - \epsilon_p), \quad (i = 1, \dots, p - 1),$$

and therefore

$$(2) \quad \begin{aligned} \sigma_{b_i, b_j} &= \sigma_{b_i^*, b_j^*} + \sigma_{(\epsilon_i - \epsilon_p)(\epsilon_j - \epsilon_p)} = c_{ij}\sigma^2 + \sigma_{\epsilon_i \epsilon_j} + \sigma_{\epsilon_p \epsilon_p} \\ &= [c_{ij} + (1 + \delta_{ij})\lambda^2]\sigma^2, \quad (i, j = 1, \dots, p - 1), \end{aligned}$$

where δ_{ij} is the Kronecker delta, i.e. $\delta_{ij} = 0$ for $i \neq j$ and $\delta_{ii} = 1$. Denote $c_{ij} + (1 + \delta_{ij})\lambda^2$ by c'_{ij} . Since the expected value of b_i^* is equal to zero, on account of (1) also the expected value of b_i is equal to zero. Let

$$\|g_{ij}\| = \|c'_{ij}\|^{-1}, \quad (i, j = 1, \dots, p - 1).$$

Then

$$(3) \quad \frac{1}{\sigma^2} \sum_{j=1}^{p-1} \sum_{i=1}^{p-1} g_{ij} b_i b_j$$

has the χ^2 -distribution with $p - 1$ degrees of freedom. The expression

$$(4) \quad \frac{1}{\sigma^2} \sum_{\alpha=1}^N (y_\alpha - Y_\alpha)^2,$$

has the χ^2 -distribution $N - p - q + 1$ degrees of freedom. The expressions (3) and (4) are independently distributed. Hence

$$(5) \quad \frac{N - p - q + 1}{p - 1} \frac{\sum \sum g_{ij} b_i b_j}{\sum (y_\alpha - Y_\alpha)^2},$$

has the F -distribution (analysis of variance distribution). We will now show that (5) is a monotonic function of λ^2 . It is known that $\sum \sum g_{ij} b_i b_j$ is invariant under linear transformations, i.e.

$$\sum \sum g_{ij} b_i b_j = \sum \sum g'_{ij} b'_i b'_j,$$

where b'_i is an arbitrary linear function, say $\mu_{i1} b_1 + \dots + \mu_{ip-1} b_{p-1}$ of b_1, \dots, b_{p-1} ($i = 1, \dots, p - 1$) and

$$\|g'_{ij}\| = \|\sigma_{b_i b_j}\|^{-1}.$$

We can choose the matrix $\|\mu_{ij}\|$ such that

$$\epsilon'_i = \mu_{i1}(\epsilon_1 - \epsilon_p) + \dots + \mu_{ip-1}(\epsilon_{p-1} - \epsilon_p), \quad (i = 1, \dots, p - 1),$$

are independently distributed and $\sigma_{\epsilon'_i}^2 = \sigma^2$. The coefficients μ_{ij} of course do not depend on σ' . We have

$$\sigma_{b'_i b'_j} = \sigma_{b_i^* b_j^*} + \delta_{ij} \sigma'^2, \quad (\delta_{ij} = \text{Kronecker delta}).$$

Now let

$$b''_i = \nu_{i1} b'_1 + \dots + \nu_{ip-1} b'_{p-1}, \quad (\nu = 1, \dots, p - 1),$$

where $\|\nu_{ij}\|$ is an orthogonal matrix and is chosen such that b''_1, \dots, b''_{p-1} are independently distributed. On account of the orthogonality of $\|\nu_{ij}\|$ we obviously have

$$\sigma_{b''_i}^2 = \sigma_{b_i^*}^2 + \sigma'^2; \quad \sigma_{b'_i b'_j} = 0 \quad \text{for } i \neq j.$$

Hence

$$(6) \quad \sum \sum g_{ij} b_i b_j = \sum_{i=1}^{p-1} \frac{b_i''^2}{\sigma_{b_i^*}^2 + \lambda^2 \sigma^2}.$$

The right hand side of (6) is evidently a monotonic function of λ^2 which proves our statement. The endpoints of the confidence interval for λ^2 are the roots in λ^2 of the equations

$$(7) \quad \frac{N - p - q + 1}{p - 1} \frac{\sum \sum g_{ij} b_i b_j}{\sum (y_\alpha - Y_\alpha)^2} = F_2; \quad \frac{N - p - q + 1}{p - 1} \frac{\sum \sum g_{ij} b_i b_j}{\sum (y_\alpha - Y_\alpha)^2} = F_1,$$

where F_2 denotes the upper, and F_1 the lower critical value of F .

The derivation of the required confidence limits in case of classifications in more than two ways can be carried out in the same way and I shall merely state here the results.

Consider r criterions of classifications and denote by p_u the number of classes in the u th classification ($u = 1, \dots, r$). Denote by $n_{i_1 \dots i_r}$ the number of observations which belong to the i_1 th class of the first classification, i_2 th class of the second classification, \dots , and to the i_r th class of the r th classification. Let $y_{i_1 \dots i_r}^{(k)}$ be the k th observation on y in the set of observations belonging to the classes mentioned above ($k = 1, \dots, n_{i_1 \dots i_r}$). We make the assumption

$$y_{i_1 \dots i_r}^{(k)} = x_{i_1 \dots i_r}^{(k)} + \epsilon_{i_1}^{(1)} + \dots + \epsilon_{i_r}^{(r)},$$

where the variates

$x_{i_1 \dots i_r}^{(k)}, \epsilon_{i_1}^{(1)}, \dots, \epsilon_{i_r}^{(r)}$ ($i_u = 1, \dots, p_u; u = 1, \dots, r; k = 1, \dots, n_{i_1 \dots i_r}$), are independently and normally distributed, the variance of $x_{i_1 \dots i_r}^{(k)}$ is σ^2 , the variance of $\epsilon_{i_u}^{(u)}$ is σ_u^2 and the mean value of $\epsilon_{i_u}^{(u)}$ is zero ($i_u = 1, \dots, p_u; u = 1, \dots, r$).

Let N be the total number of observations. We order the observations in a certain order and denote by y_α the α th observation in that order ($\alpha = 1, \dots, N$). Consider the variables:

$$t, t_{i_u}^{(u)}, \quad (u = 1, \dots, r; i_u = 1, \dots, p_u),$$

and denote by t_α the α th observation on t and by $t_{i_u \alpha}^{(u)}$ the α th observation on $t_{i_u}^{(u)}$. The values of t_α and $t_{i_u \alpha}^{(u)}$ are given as follows:

$$t_\alpha = 1 \quad (\alpha = 1, \dots, N),$$

$$t_{i_u \alpha}^{(u)} = 1 \text{ if } y_\alpha \text{ lies in the } i_u \text{th class of the } u \text{th classification,}$$

$$t_{i_u \alpha}^{(u)} = 0 \text{ if } y_\alpha \text{ does not lie in the } i_u \text{th class of the } u \text{th classification.}$$

Let the sample regression of y on $t, t_{i_u}^{(u)}$ be given by

$$Y = at + \sum_{u=1}^r \sum_{i_u=1}^{p_u-1} b_{i_u}^{(u)} t_{i_u}^{(u)}.$$

Let the covariance of $b_{i_u}^{(u)}$ and $b_{j_u}^{(u)}$ be given by $C_{i_u j_u}^{(u)} \sigma^2$ under the assumption that $\sigma_1 = \sigma_2 = \dots = \sigma_r = 0$. The matrix $\|C_{i_u j_u}^{(u)}\|$ ($i_u, j_u = 1, \dots, p_u - 1$) can be calculated by known methods of the theory of least squares. Let

$$\|g_{i_u j_u}^{(u)}\| = \|C_{i_u j_u}^{(u)} + (1 + \delta_{i_u j_u}) \lambda_u^2\|^{-1} \quad (i_u, j_u = 1, \dots, p_u - 1),$$

where $\delta_{i_u j_u}$ is the Kronecker delta and $\lambda_u^2 = \sigma_u^2 / \sigma^2$. Then the lower and upper confidence limits for λ_u^2 are given by the roots in λ_u^2 of the equations

$$(8) \quad \frac{N - \sum_{u=1}^r p_u + r - 1}{p_u - 1} \frac{\sum_{j_u=1}^{p_u-1} \sum_{i_u=1}^{p_u-1} g_{i_u j_u}^{(u)} b_{i_u}^{(u)} b_{j_u}^{(u)}}{\sum_{\alpha=1}^N (y_\alpha - Y_\alpha)^2} = F_i \quad (i = 1, 2),$$

where F_2 is the upper and F_1 the lower critical value of the analysis of variance distribution with $p_u - 1$ and $N - \sum_{u=1}^r p_u + r - 1$ degrees of freedom. In case of a single criterion of classification the confidence limits (8) are identical with those given in my previous paper.

THE FREQUENCY DISTRIBUTION OF A GENERAL MATCHING PROBLEM

BY T. N. E. GREVILLE

Bureau of the Census

1. Introduction. This paper considers the matching of two decks of cards of arbitrary composition, and the complete frequency distribution of correct matchings is obtained, thus solving a problem proposed by Stevens.¹ It is also shown that the results can be interpreted in terms of a contingency table.

Generalizing a problem considered by Greenwood,² let us consider the matching of two decks of cards consisting of t distinct kinds, all the cards of each kind being identical. The first or "call" deck will be composed of i_1 cards of the first kind, i_2 of the second, etc., such that

$$i_1 + i_2 + i_3 + \dots + i_t = n;$$

and the second or "target" deck will contain j_1 cards of the first kind, j_2 of the second, etc., such that

$$j_1 + j_2 + \dots + j_t = n.$$

Any of the i 's or j 's may be zero. It is desired to calculate, for a given arrangement of the "call" deck, the number of possible arrangements of the "target" deck which will produce exactly r matchings between them ($r = 0, 1, 2, \dots, n$). It is clear that these frequencies are independent of the arrangement of the call deck. For convenience the call deck may be thought of as arranged so that all the cards of the first kind come first, followed by all those of the second kind, and so on.

2. Formulae for the frequencies. Let us consider the number of arrangements of the target deck which will match the cards in the k_1 th, k_2 th, \dots , k_s th positions in the call deck, regardless of whether or not matchings occur elsewhere. Let the cards in these s positions in the call deck consist of c_1 of the first kind, c_2 of the second, etc. Then:

$$c_1 + c_2 + \dots + c_t = s.$$

The number of such arrangements of the target deck is

$$(1) \quad \frac{(n-s)!}{\prod_{h=1}^t (j_h - c_h)!}.$$

¹ W. L. STEVENS, *Annals of Eugenics*, Vol. 8 (1937), pp. 238-244.

² J. A. GREENWOOD, *Annals of Math. Stat.*, Vol. 9 (1938), pp. 56-59.