

# THE FREQUENCY DISTRIBUTION OF DEVIATES FROM MEANS AND REGRESSION LINES IN SAMPLES FROM A MULTIVARIATE NORMAL POPULATION

BY D. J. FINNEY

*Oxford University, England*

**1. Summary.** The joint frequency distribution has been found for any set of the  $(n - k)$  deviates from their sample mean of each of the  $t$  variates in a sample from a multivariate normal population. Expressions for the variance of any single deviate in this distribution, the correlation coefficient between any pair of deviates, and certain partial correlation coefficients between any pair have also been obtained.

These results have been generalized so as to include the corresponding properties of deviates from a set of  $t$  multiple linear regression equations estimated from the sample, the  $m$  independent variates being the same for each of the  $t$  dependent.

**2. Introduction.** Some years ago, Irwin published results relating to the frequency distribution of the deviations of individual observations from the mean of a sample drawn from a normal population (see [1]). He derived an expression for the joint distribution of any number of these deviates, which distribution is always of the normal multivariate form, and thence obtained the total and partial correlation coefficients between any pair of the deviates.

The purpose of this paper is to discuss the generalization of Irwin's problem, firstly to the properties of the deviates of individual observations from the mean in a sample from a multivariate normal population and secondly to the properties of deviates from a regression equation instead of from a mean. So far as is known to the writer Irwin's results are of little practical importance, and these generalizations are probably of no practical value whatsoever. Nevertheless, they have some interest as additions to the knowledge of the mathematical properties of the normal frequency function, and for that reason alone they are put on record here.

**3. Deviations from the sample mean.** Irwin based his discussion on a normal population with mean  $m$  and variance  $\sigma^2$ , but the algebra is simplified a little, without any real loss of generality in the final results, if, by means of a preliminary transformation, these parameters of position and scale are made zero and unity respectively. The multivariate normal distribution in the  $t$  variates  $y_i$ , ( $i = 1, 2, \dots, t$ ), each with mean zero and variance unity, has the frequency function

$$(1) \quad \frac{1}{(2\pi)^{\frac{1}{2}t} R^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2R} \rho^{ij} y_i y_j \right\},$$

where  $i, j = 1, 2, \dots, t$ ;  $\rho^{ij}$  is the cofactor of the element  $\rho_{ij}$  in the determinant of population correlation coefficients

$$(2) \quad R = \begin{vmatrix} 1 & \rho_{12} & \rho_{13} & \cdots & \rho_{1t} \\ \rho_{12} & 1 & \rho_{23} & \cdots & \rho_{2t} \\ \dots & \dots & \dots & \dots & \dots \\ \rho_{1t} & \rho_{2t} & \rho_{3t} & \cdots & 1 \end{vmatrix}$$

A summation convention for the affixes  $i, j$  is understood throughout this paper, except when the contrary is explicitly stated.

Let  $(y_p)$  represent a sample of  $n$  independent sets of values of the  $t$  variates randomly selected from the population, ( $p = 1, 2, \dots, n$ ). Then the element of probability for the sample is

$$(3) \quad \frac{1}{(2\pi)^{\frac{1}{2}nt} R^{\frac{1}{2}n}} \exp \left\{ -\frac{1}{2R} \rho^{ij} \sum_p y_p y_p \right\} \prod_{\substack{i=1, \dots, t \\ p=1, \dots, n}} \{d(y_p)\}.$$

If  $\bar{y}$  is the mean of the  $n$  sample values of  $y$ , the deviates from the mean are  $(iY_p)$ , where

$$(4) \quad iY_p = y_p - \bar{y} = \sum_q \left( \delta_{pq} - \frac{1}{n} \right) y_q,$$

the summation being taken over  $q = 1, 2, \dots, n$  with

$$\delta_{pq} = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q. \end{cases}$$

Now the  $iY$  are linear combinations of normally distributed variates, and are therefore themselves normally distributed. Clearly

$$(5) \quad E(iY_p) = 0$$

and, from an expansion by means of equation (4) using

$$(6) \quad \begin{aligned} E(y_p y_q) &= \delta_{pq} \rho_{ij}, \\ E(iY_p iY_q) &= \left( \delta_{pq} - \frac{1}{n} \right) \rho_{ij}, \end{aligned}$$

where  $\rho_{ii} = 1$  (not summed). Consequently the variance of any one deviate is

$$(7) \quad \sigma^2(iY_p) = \frac{n-1}{n},$$

and the correlation coefficient between any pair is

$$(8) \quad \rho(iY_p, iY_q) = \frac{n\delta_{pq} - 1}{n-1} \rho_{ij}.$$

Equation (7) and equation (8) for the particular case of  $i = j$  agree with the well-known results that Irwin has already given as equations (10) of his paper.

For any  $i$ , only  $(n - 1)$  of the deviates  ${}_iY_p$  are functionally independent. The joint distribution of these for  $p = 1, 2, \dots, (n - k)$  may be obtained from an inversion of the matrix of correlation coefficients. If  $\Delta$  is the determinant of this matrix and  $\Delta({}_iY_p, {}_jY_q)$  the cofactor corresponding to the two elements specified, this inversion shows that

$$(9) \quad \frac{\Delta({}_iY_p, {}_jY_q)}{\Delta} = \left( \delta_{pq} + \frac{1}{k} \right) \frac{\rho_{ij}}{R} \cdot \frac{n - 1}{n}.$$

The joint distribution is therefore

$$(10) \quad \text{const.} \times \exp \left\{ -\frac{1}{2R} \rho^{ij} \sum_{p \leq q \leq n-k} \left( \delta_{pq} + \frac{1}{k} \right) {}_iY_p {}_jY_q \right\} \prod (dY).$$

Now  $\Delta$  may be evaluated as

$$\Delta = \left( \frac{n}{n - 1} \right)^{t(n-k-1)} \left( \frac{k}{n - 1} \right)^t R^{n-k},$$

and the constant multiplier in equation (10) is therefore

$$(11) \quad \frac{\left( \frac{n}{k} \right)^{\frac{1}{2}t}}{\{ (2\pi)^t R \}^{\frac{1}{2}(n-k)}}.$$

From equation (9), the partial correlation coefficient between any two of the variates in the distribution (10), the remaining  $t(n - k) - 2$  being held constant, is written down as

$$(12) \quad \text{partial correlation coefficient between } {}_iY_p, {}_jY_q = -\frac{k\delta_{pq} + 1}{k + 1} \cdot \frac{\rho^{ij}}{(\rho^{ii}\rho^{jj})^{\frac{1}{2}}};$$

the summation convention is suspended for this equation.

**4. Deviations from regression equations.** The results obtained in section three may be generalized so as to relate to the frequency distribution of deviates from linear or polynomial regression equations instead of to deviates from means. Suppose that there are  $m$  independent variates  $x^\alpha$ , ( $\alpha = 1, 2, \dots, m$ ), which take values  $x_p^\alpha$  corresponding to the sample observations  ${}_iy_p$ ; polynomial regressions may be included by taking powers of an  $x$  as separate variates. If a conventional variate  $x^0$ , whose value is always unity, be introduced, the regression equation of  ${}_iy$  on  $x^\alpha$ , ( $\alpha = 0, 1, 2 \dots, m$ ), may be written

$$(13) \quad {}_i\eta = {}_ib^\alpha x^\alpha,$$

where a summation convention is understood for  $\alpha = 0, 1, \dots, m$  and the regression coefficients are the solutions of the normal equations.

$$(14) \quad {}_ib^\alpha \sum_p x_p^\alpha x_p^\beta = \sum_p {}_iy_p x_p^\beta.$$

Write

$$(15) \quad B^{\alpha\beta} = \sum_p x_p^\alpha x_p^\beta$$

and let  $(B_{\alpha\beta})$  be the inverse matrix of  $(B^{\alpha\beta})$ .

Then the solutions of equations (14) are

$$(16) \quad {}_i b^\alpha = B_{\alpha\beta} \sum_p {}_i y_p x_p^\beta.$$

If the deviation of  ${}_i y_p$  from the regression equation (13) is  ${}_i Z_p$ , then

$$(17) \quad \begin{aligned} {}_i Z_p &= {}_i y_p - {}_i \eta_p \\ &= \sum_q (\delta_{pq} - B_{\alpha\beta} x_p^\alpha x_q^\beta) {}_i y_q, \end{aligned}$$

the summation for  $q$  being over  $q = 1, 2, \dots, n$ . As for equation (5),

$$(18) \quad E({}_i Z_p) = 0.$$

Also

$$(19) \quad E({}_i Z_p {}_j Z_q) = (\delta_{pq} - B_{\alpha\beta} x_p^\alpha x_q^\beta) \rho_{ij}.$$

since by definition

$$B^{\alpha\beta} B_{\alpha\gamma} = \delta_{\beta\gamma}.$$

Write now  $\theta$  for the square matrix of  $(m + 1)$  rows and columns whose elements are the  $B^{\alpha\beta}$ , and  $X_p$  for the single column matrix of values  $x$  corresponding to the  $p$ th observation; i.e.

$$(20) \quad \theta = (B^{\alpha\beta})$$

and

$$(21) \quad X_p = \begin{pmatrix} x_p^0 \\ x_p^1 \\ x_p^2 \\ \cdot \\ \cdot \\ \cdot \\ x_p^m \end{pmatrix}$$

Write also

$$(22) \quad \theta_{p,q,r,\dots} = \theta - X_p X_p' - X_q X_q' - X_r X_r' - \dots$$

Then

$$|\theta_p| = |\theta| \cdot (1 - B_{\alpha\beta} x_p^\alpha x_p^\beta),$$

and

$$|\theta_{pq}| - |\theta_{pq} + X_p X_q'| = -|\theta| \cdot B_{\alpha\beta} x_p^\alpha x_q^\beta.$$

Hence, from equation (19), the variance of a deviate may be written

$$(23) \quad \sigma^2(iZ_p) = \frac{|\theta_p|}{|\theta|},$$

and the correlation coefficient between any pair of deviates is

$$(24) \quad \rho(iZ_p, iZ_q) = \begin{cases} \rho_{ij} & (p = q) \\ \rho_{ij} \frac{|\theta_{pq}| - |\theta_{pq} + X_p X'_q|}{\{|\theta_p| \cdot |\theta_q|\}^{\frac{1}{2}}} & (p \neq q) \end{cases}$$

For any  $i$ , only  $(n - m - 1)$  of the deviates  $iZ_p$  are functionally independent. The joint distribution of these for  $p = 1, 2, \dots, (n - k)$  and any  $k \geq m + 1$  may be found by inversion of the matrix of correlation coefficients obtained from equation (24). The multiplier of the exponential in this distribution of  $t(n - k)$  variables is

$$\frac{|\theta|^{\frac{1}{2}t(n-k)}}{(2\pi)^{\frac{1}{2}t(n-k)} R^{\frac{1}{2}t(n-k)} D^{\frac{1}{2}t}},$$

where

$$D = \begin{vmatrix} |\theta_1| & |\theta_{12}| - |\theta_{12} + X_1 X'_2| & \cdots & |\theta_{1,n-k}| - |\theta_{1,n-k} + X_1 X'_{n-k}| \\ |\theta_{12}| - |\theta_{12} + X_1 X'_2| & |\theta_2| & \cdots & |\theta_{2,n-k}| - |\theta_{2,n-k} + X_2 X'_{n-k}| \\ \dots & \dots & \dots & \dots \\ |\theta_{1,n-k}| - |\theta_{1,n-k} + X_1 X'_{n-k}| & |\theta_{2,n-k}| - |\theta_{2,n-k} + X_2 X'_{n-k}| & \cdots & |\theta_{n-k}| \end{vmatrix}$$

Since  $\theta$  is positive definite, there exists a non-singular matrix  $K$  such that

$$K\theta K' = I.$$

Then the  $X_p$  may be transformed to new column matrices  $W_p$  by

$$KX_p = W_p = \begin{pmatrix} w_p^0 \\ w_p^1 \\ w_p^2 \\ \vdots \\ w_p^m \end{pmatrix}$$

and consequently

$$X_p = K^{-1}W_p.$$

It follows that

$$|\theta_p| = |\theta| \cdot |I - W_p W'_p|,$$

which may be reduced to the form

$$|\theta_p| = |\theta| \cdot (1 - w_p^\alpha w_p^\alpha).$$

Similarly

$$|\theta_{pq}| - |\theta_{pq} + X_p X'_q| = -|\theta| w_p^\alpha w_q^\alpha.$$

Hence

$$D = |\theta|^{n-k} \begin{vmatrix} 1 - w_1^\alpha w_1^\alpha & -w_1^\alpha w_2^\alpha & \cdots & -w_1^\alpha w_{n-k}^\alpha \\ -w_1^\alpha w_2^\alpha & 1 - w_2^\alpha w_2^\alpha & \cdots & -w_2^\alpha w_{n-k}^\alpha \\ \cdots & \cdots & \cdots & \cdots \\ -w_1^\alpha w_{n-k}^\alpha & -w_2^\alpha w_{n-k}^\alpha & \cdots & 1 - w_{n-k}^\alpha w_{n-k}^\alpha \end{vmatrix}$$

This may be transformed into

$$\begin{aligned} D &= |\theta|^{n-k} \begin{vmatrix} & & & & W'_1 \\ & & & & W'_2 \\ & & & & \vdots \\ & & I_{n-k} & & \vdots \\ & & & & W'_{n-k} \\ W_1 & W_2 & \cdots & W_{n-k} & I_{m+1} \end{vmatrix} \\ &= |\theta|^{n-k} \cdot |I - W_1 W'_1 - W_2 W'_2 - \cdots - W_{n-k} W'_{n-k}| \\ &= |\theta|^{n-k-1} \cdot |\theta_{1,2,\dots,(n-k)}|. \end{aligned}$$

Thus, finally, the constant in the distribution is found to be

$$(25) \quad \frac{1}{\{(2\pi)^t R\}^{\frac{1}{2}(n-k)}} \left\{ \frac{|\theta|}{|\Omega_k|} \right\}^{\frac{1}{2}t},$$

in which  $\Omega_k$  has been written for  $\theta_{1,2,\dots,(n-k)}$ , a matrix of the same form as  $\theta$  but calculated from the last  $k$  sets of observations only.

The cofactors of the matrix of correlation coefficients, required for the coefficients of the quadratic form in the distribution, can be derived in a similar manner. The distribution may be written

$$(26) \quad \frac{1}{\{(2\pi)^t R\}^{\frac{1}{2}(n-k)}} \left\{ \frac{|\theta|}{|\Omega_k|} \right\}^{\frac{1}{2}t} \cdot \exp. \left\{ -\frac{1}{2R} \rho^{ij} \sum_{p \leq q \leq n-k} \left( \delta_{pq} - 1 + \frac{|\Omega_k + X_p X'_q|}{|\Omega_k|} \right) Z_p Z_q \right\} \Pi(dZ),$$

of which the distribution (10) is easily seen to be the particular case for  $m = 0$ .

From (26), the partial correlation coefficient between any pair of deviates,  $Z_p$  and  $Z_q$ , may be written down as

$$(27) \quad -\frac{|\Omega_k + X_p X'_q| + (\delta_{pq} - 1) |\Omega_k|}{\{|\Omega_k + X_p X'_p| \cdot |\Omega_k + X_q X'_q|\}^{\frac{1}{2}}} \frac{\rho^{ij}}{(\rho^{ii} \rho^{jj})^{\frac{1}{2}}};$$

in this expression the summation convention is again suspended.

REFERENCE

[1] J. O. IRWIN, "On the frequency distribution of any number of deviates from the mean of a sample from a normal population and the partial correlations between them," *Roy. Stat. Soc. Jour.*, Vol. 92 (1929), pp. 580-584.