

## ON THE SAMPLING DISTRIBUTION OF THE MULTIPLE CORRELATION COEFFICIENT

By S. S. WILKS\*

The problem of finding the distribution of the multiple correlation coefficient in samples from a normal population with a non-zero multiple correlation coefficient was solved in 1928 by Fisher<sup>1</sup> by the application of geometrical methods. In his derivation he used the facts that the population value  $\rho$  of the multiple correlation coefficient is invariant under linear transformations of the independent variates, and that the distribution of the multiple correlation coefficient is independent of all population parameters except  $\rho$ .

In this paper it will be shown that the distribution of the multiple correlation coefficient can be derived directly from Wishart's<sup>2</sup> generalized product moment distribution without making use of geometrical notions and the property of the invariance of  $\rho$  under linear transformations of the independent variates. Furthermore, it will not be necessary to show that the distribution will be independent of all population parameters except  $\rho$ .

The population value of the multiple correlation coefficient between a variate  $x_1$  and a set of variates  $x_2, x_3, \dots, x_n$  is the ordinary correlation coefficient between  $x_1$  and that linear function of the variates  $x_2, x_3, \dots, x_n$  which will make this correlation a maximum. It can be expressed as  $\rho^2 = 1 - \frac{\Delta}{\Delta_1}$ , where  $\Delta$  is the determinant of the correlations among all of the

---

\*National Research Fellow in Mathematics.

<sup>1</sup>R. A. Fisher, The general sampling distribution of the multiple correlation coefficient, Proceedings of the Royal Society of London, series A, vol. 121 (1928), pp. 654-73.

<sup>2</sup>John Wishart, The generalized product moment distribution in samples from a normal multivariate population, Biometrika, vol. 20A (1928) pp. 32-52.

variates  $x_1, x_2, \dots, x_n$  and  $\Delta_1$  is the determinant of correlations among the independent variates  $x_2, x_3, \dots, x_n$ . Denoting the sample value of  $\rho^2$  by  $R^2$  it is well known that  $R^2 = 1 - \frac{D}{D_1}$ , where  $D$  and  $D_1$  are the determinants of sample correlations among the sets of variates  $x_1, x_2, \dots, x_n$  and  $x_2, x_3, \dots, x_n$  respectively.

Let us suppose a sample of  $N$  items to be drawn at random from the normal  $n$ -variate population whose distribution is

$$(1) \quad \frac{\sqrt{A}}{(2\pi)^{\frac{n}{2}}} e^{-\frac{1}{2} \sum_{i,j=1}^n A_{ij} (x_i - m_i)(x_j - m_j)}$$

where  $A_{ij} = \frac{\Delta_{ij}}{\sigma_i \sigma_j \Delta}$ ,  $\Delta = |\rho_{ij}|$  the determinant of correlations among the  $n$  variates,  $\Delta_{ij}$  is the cofactor of  $\rho_{ij}$  in  $\Delta$ ,  $\sigma_i$  is the standard deviation of  $x_i$  and  $A = |A_{ij}|$ .

In the sample, let

$$\bar{x}_i = \frac{1}{N} \sum_{\alpha=1}^N x_{i\alpha},$$

and

$$a_{ij} = \frac{1}{N} \sum_{\alpha=1}^N (x_{i\alpha} - \bar{x}_i)(x_{j\alpha} - \bar{x}_j),$$

where  $x_{i\alpha}$  is the value of  $x_i$  for the  $\alpha$ -th individual of the sample. Wishart<sup>3</sup> has proved that the simultaneous distribution function of the set  $\{a_{ij}\}$ , ( $i, j = 1, 2, \dots, n$ ) is

$$(2) \quad f(\bar{a}) = \frac{\left(\frac{N}{2}\right)^{\frac{n(N-1)}{2}} A^{\frac{N-1}{2}}}{\pi^{\frac{n(n-1)}{4}} \Gamma\left(\frac{N-1}{2}\right) \Gamma\left(\frac{N-2}{2}\right) \dots \Gamma\left(\frac{N-n}{2}\right)} e^{-\frac{N}{2} \sum_{i,j=1}^n A_{ij} a_{ij}} |a_{ij}|^{\frac{N-n-2}{2}}$$

<sup>3</sup>J. Wishart, loc. cit.

where  $|a_{ij}|$  is the determinant of the  $a$ 's.

We shall define a moment-generating function  $\phi(\alpha, k)$  as

$$(3) \quad \phi(\alpha, k) = \int e^{\alpha a_{11}} |a_{ij}|^h |a_{pq}|^{-k} f(\bar{a}) d\bar{a},$$

where the integration is to be taken over the field of all possible values of the  $a$ 's and  $|a_{pq}|$  is the cofactor of  $a_{11}$  in  $|a_{ij}|$ .

From this definition of  $\phi(\alpha, k)$ , it is clear that  $\frac{\partial^h}{\partial \alpha^h} \phi(\alpha, k) \Big|_{\alpha=0}$

is the product moment  $E \left[ a_{11}^{h+k} (1-R^2)^k \right]$ . It will be shown that this

expectation exists for  $h=-k$  which will yield the  $k$ -th moment of  $(1-R^2)$ , from which the distribution of  $R^2$  can be found.

To find  $\phi(\alpha, k)$  we observe that since (2) is a probability function, its value over the field of all possible values of the  $a$ 's is unity. Hence, we must have

$$(4) \quad \int e^{-\frac{N}{2} \sum_{i,j=1}^n A_{ij} a_{ij}} |a_{ij}|^{\frac{N-n-2}{2}} d\bar{a} = G,$$

$$\text{where } G = \frac{\pi^{\frac{n(n-1)}{2}} \Gamma(\frac{N-1}{2}) \Gamma(\frac{N-2}{2}) \cdots \Gamma(\frac{N-n}{2})}{\binom{N}{2}^{\frac{n(N-1)}{2}} A^{\frac{N-1}{2}}}. \text{ This relation}$$

holds for all positive values of  $N > n$  and for all values of  $A_{ij}$  which will make the matrix  $\| A_{ij} \|$  positive definite.

If  $f(\bar{a})$  be integrated with respect to  $a_{11}, a_{12}, \dots, a_{1n}$ , the resulting form will clearly be the distribution of the set of  $a$ 's contained in  $|a_{pq}|$  and will be

$$(5) \frac{\binom{N}{2} \frac{(n-1)(N-1)}{2} \frac{N-1}{2}}{\pi \frac{(n-1)(n-2)}{4} \dots \Gamma(\frac{N-n+1}{2})} e^{-\frac{N}{2} \sum_{p,q=2}^n B_{pq} a_{pq}} |a_{pq}|^{\frac{N-n-1}{2}}$$

where  $B_{pq}$  is the element in the  $p$ -th row and  $q$ -th column of the reciprocal form<sup>4</sup> of the determinant which is the cofactor of the term in the first row and first column of the reciprocal form of  $|A_{ij}|$ . The value of  $B_{pq}$  in terms of correlation coefficients and standard deviations is  $\frac{\Delta''_{pq}}{\sigma_p \sigma_q \Delta''}$ , where

$\Delta'' = \Delta_{11}$ , and  $\Delta''_{pq}$  is the cofactor of  $\rho_{pq}$  in  $\Delta_{11}$ . Furthermore,  $B = |B_{pq}|$ . Hence

$$(6) \int e^{-\frac{N}{2} \sum_{i,j=1}^n A_{ij} a_{ij}} |a_{ij}|^{\frac{N-n-2}{2}} da_1 \dots da_n$$

$$= \pi^{\frac{n-1}{2}} \binom{N}{2}^{-\frac{N-1}{2}} \left(\frac{B}{\lambda}\right)^{\frac{N-1}{2}} \Gamma(\frac{N-n}{2}) e^{-\frac{N}{2} \sum_{p,q=2}^n B_{pq} a_{pq}} |a_{pq}|^{\frac{N-n-1}{2}} d[a-a_1]$$

where  $da_1 = da_{11} da_{12} \dots da_{1n}$  and  $d[a-a_1]$  is the product of the differentials of all  $a$ 's in  $|a_{pq}|$  ( $p, q = 2, 3, \dots, n$ ).

Now, it is clear that (6) is an identity for all values of  $N$  and the population parameters  $\sigma_i$  and  $\rho_{ij}$  ( $i, j = 1, 2, \dots, n$ ;  $i \neq j$ ), for which both sides of (6) exist. Thus, we can perform the following operations on (6):

- (a). Replace  $N$  by  $N+2k$ .
- (b). Replace  $\sigma_i$  by  $\sigma_i \sqrt{\frac{N+2k}{N}}$ , ( $i = 1, 2, \dots, n$ )

---

<sup>4</sup>By the reciprocal form of a determinant  $|c_{ij}|$  we mean the determinant formed by replacing each element  $c_{ij}$  by the ratio  $\frac{C_{ij}}{C}$  where  $C_{ij}$  is the cofactor of  $c_{ij}$  and  $C = |c_{ij}|$ .

(c). Replace  $A_{11}$  by  $A_{11} - \frac{2\alpha}{N}$ .

(d). Multiply both sides of the identity by  $\frac{1}{G}$ .

(e). Multiply both sides by  $|a_{pq}|^{-k}$ .

Accordingly, we find that the integral of the left side of (6) over all possible values of the  $a$ 's is the definition of  $\phi(\alpha, k)$ , which must be equal to the integral of the right side over the field of all possible values of the  $a$ 's in  $|a_{pq}|$ . But the value of the integral of the right side can be deduced at once from (4). Hence, we finally obtain,

$$(7) \quad \phi(\alpha, k) = \left(\frac{N}{2}\right)^{-k} A^{\frac{N-1}{2}} A_{\alpha}^{-\frac{N-1}{2}-k} B_{\alpha}^k \frac{\Gamma\left(\frac{N-n}{2} + k\right)}{\Gamma\left(\frac{N-n}{2}\right)},$$

where  $A_{\alpha}$  is the determinant  $A$  with  $A_{11}$  replaced by  $A_{11} - \frac{2\alpha}{N}$ , and  $B_{\alpha}$  is the reciprocal of the cofactor of the element in the first row and first column of the reciprocal form of  $A_{\alpha}$ .

That is,

$$(8) \quad B_{\alpha} = \frac{A_{\alpha}^{n-1}}{|\bar{A}_{\alpha, pq}|},$$

where  $\bar{A}_{\alpha, pq}$  is the cofactor of the element in the  $p$ -th row and  $q$ -th column of  $A_{\alpha}$ , ( $p, q = 2, 3, \dots, n$ ). The value of  $|\bar{A}_{\alpha, pq}|$  can be readily found by writing

$$|\bar{A}_{\alpha, pq}| = \frac{|\bar{A}_{\alpha, pq}| \cdot |\bar{A}_{\alpha, ij}|}{A_{\alpha}},$$

where  $A_{\alpha, ij} \equiv A_{ij}$  except for  $i=j=1$  and  $A_{\alpha, 11} = A_{11} - \frac{2\alpha}{N}$ .

Increasing  $|\bar{A}_{\alpha, pq}|$  to an  $n$ -th order determinant by inserting, as first row and first column, an additional row and column which will not change the value of the determinant, and multiplying it by  $|A_{\alpha, 1q}|$  we find

$$|\bar{A}_{\alpha, pq}| = A_{\alpha}^{n-2} \left(A_{11} - \frac{2\alpha}{N}\right).$$

Therefore,

$$B_{\alpha} = \frac{A_{\alpha}}{(A_{11} - \frac{2\alpha}{N})}$$

Substituting this for  $B_{\alpha}$  in (7) and using the fact that

$$A_{\alpha} = A - \frac{2\alpha}{N} \bar{A}_{11}, \text{ we finally obtain}$$

$$(9) \quad \varphi(\alpha, k) = \left(\frac{NA}{2\bar{A}_{11}}\right)^{\frac{N-1}{2}} \left(\frac{NA}{2\bar{A}_{11}} - \alpha\right)^{-\frac{N-1}{2}} \left(\frac{NA_{11}}{2} - \alpha\right)^{-k} \frac{\Gamma(\frac{N-n}{2} + k)}{\Gamma(\frac{N-n}{2})}$$

Thus, it is evident that  $\varphi(\alpha, k)$  exists for sufficiently small values of  $\alpha$ . Let us write

$$\left(\frac{NA}{2\bar{A}_{11}} - \alpha\right)^{-\frac{N-1}{2}} = \left(\frac{NA_{11}}{2} - \alpha\right)^{-\frac{N-1}{2}} \left[1 - \frac{\frac{N}{2}(A_{11} - \bar{A}_{11})}{(\frac{NA_{11}}{2} - \alpha)}\right]^{-\frac{N-1}{2}}$$

and expand the second factor on the right into a Taylor series. Substituting in (9), we have the convergent series

$$(10) \quad \varphi(\alpha, k) = \left(\frac{NA}{2\bar{A}_{11}}\right)^{\frac{N-1}{2}} \frac{\Gamma(\frac{N-n}{2} + k)}{\Gamma(\frac{N-n}{2})} \times \sum_{i=0}^{\infty} \frac{\left(\frac{NA_{11}}{2} - \alpha\right)^{-k - \frac{N-1}{2} - i} \left(\frac{N}{2}\right)^i (A_{11} - \bar{A}_{11})^i \Gamma(\frac{N-1}{2} + i)}{i! \Gamma(\frac{N-1}{2})}$$

For the coefficient of  $\frac{\alpha^h}{h!}$  in the expansion of the right side of (10) in powers of  $\alpha$ , we find

$$(11) \quad \left(\frac{NA_{11}}{2}\right)^{-k-h} \left(\frac{A}{\bar{A}_{11} A_{11}}\right)^{\frac{N-1}{2}} \frac{\Gamma(\frac{N-n}{2} + k)}{\Gamma(\frac{N-n}{2})} \times \sum_{i=0}^{\infty} \frac{\left(1 - \frac{A}{A_{11} \bar{A}_{11}}\right)^i \Gamma(\frac{N-1}{2} + i) \Gamma(\frac{N-1}{2} + k + h + i)}{i! \Gamma(\frac{N-1}{2}) \Gamma(\frac{N-1}{2} + k + i)}$$

which is the definition of  $E [a_{ii}^{h+k} (1-R^2)^k]$ . We observe that (11) exists for all values of  $k$  and  $h$  for which

$$\frac{N-n}{2} + k > 0 \quad \text{and} \quad \frac{N-1}{2} + h + k > 0. \quad \text{Placing } h = -k$$

and pointing out that  $\frac{A}{A_{ii} \bar{A}_{ii}} = 1 - \rho^2$ , we have as the  $k$ -th moment of  $1-R^2$ ,

$$M_k [(1-R^2)] = E [(1-R^2)^k] = \frac{(1-\rho^2)^{\frac{N-1}{2}}}{\Gamma(\frac{N-n}{2}) \Gamma(\frac{N-1}{2})} \sum_{i=0}^{\infty} \frac{\rho^{2i} \Gamma^2(\frac{N-1}{2} + i) \Gamma(\frac{N-n}{2} + k)}{i! \Gamma(\frac{N-1}{2} + k + i)}. \quad (12)$$

By using the relation

$$\frac{\Gamma(\frac{N-n}{2} + k)}{\Gamma(\frac{N-1}{2} + k + i)} = \frac{1}{\Gamma(\frac{n-1}{2} + i)} \int_0^1 (1-\theta)^{\frac{N-n}{2} + k - 1} \theta^{\frac{n-1}{2} + i - 1} d\theta$$

we can write (12) in the form

$$E [(1-R^2)^k] = \frac{(1-\rho^2)^{\frac{N-1}{2}}}{\Gamma(\frac{N-n}{2}) \Gamma(\frac{N-1}{2})} \times \sum_{i=0}^{\infty} \int_0^1 \frac{e^{2i} (1-\theta)^{\frac{N-n}{2} + k - 1} \theta^{\frac{n-1}{2} + i - 1} \Gamma^2(\frac{N-1}{2} + i)}{i! \Gamma(\frac{n-1}{2} + i)} d\theta. \quad (13)$$

The series in (13) is uniformly convergent in  $\theta$  for  $0 \leq \theta \leq 1$  and therefore, we can interchange the order of summation and integration and write

$$E [(1-R^2)^k] = \int_0^1 (1-\theta)^k \phi(\theta) d\theta, \quad (14)$$

where

$$\phi(\theta) = \frac{(1-\rho^2)^{\frac{N-1}{2}} (1-\theta)^{\frac{N-n-1}{2}} \theta^{\frac{n-1}{2}-1}}{\Gamma(\frac{N-1}{2}) \Gamma(\frac{N-n}{2})}$$

$$\times \sum_{i=0}^{\infty} \frac{\rho^{2i} \theta^i \Gamma^2(\frac{N-1}{2} + i)}{i! \Gamma(\frac{n-1}{2} + i)}.$$

(15)

Thus, we have a distribution function of a variable  $\theta$  such that the  $k$ -th moment of  $\theta$  is identical with the  $k$ -th moment of  $R^2$  for all positive values of  $k$ . It follows from Stekloff's<sup>5</sup> theory of closure that  $\phi(\theta)$  must be the only continuous solution of (14), where  $E[(1-R^2)^k]$  is defined as (12). Therefore, the distribution of  $R^2$  is identical with that of  $\theta$  and can be written finally as

$$df = \frac{\Gamma(\frac{N-1}{2})}{\Gamma(\frac{N-n}{2}) \Gamma(\frac{n-1}{2})}$$

(16)

$$\times (1-\rho^2)^{\frac{N-1}{2}} (1-R^2)^{\frac{N-n-1}{2}} (R^2)^{\frac{n-3}{2}} F[\frac{N-1}{2}, \frac{N-1}{2}, \frac{n-1}{2}, \rho^2 R^2] d(R^2).$$

which is the distribution found by Fisher except that he uses the notation  $n_1 = n-1$ , the number of independent variates, and  $n_1 + n_2 + 1 = N$ , the sample number.

---

<sup>5</sup>W. Stekloff: Quelques applications nouvelles de la théorie de fermeture au problème de représentation approchée des fonctions et au problème des moments, Memoire de l'Académie Impériale des Sciences de St. Petersburg, vol. 32, no. 4, (1914).

*S. S. Wilks*