

- [4] A. WALD: "The fitting of straight lines if both variables are subject to error," *Annals of Math. Stat.*, Vol. 11 (1940), pp. 284 ff.
- [5] T. HAAVELMO: "The probability approach in econometrics," *Econometrica*, Vol. 12 (1944), Supplement.
- [6] T. KOOPMANS: *Linear Regression Analysis in Economic Time Series*, Haarlem, 1937.
- [7] G. TINTNER: "An application of the variate difference method to multiple regression," *Econometrica*, Vol. 12 (1944), pp. 97 ff.
- [8] H. HOTELLING: "Simplified calculation of principal components," *Psychometrika*, Vol. 1 (1936), pp. 27 ff.
- [9] T. W. ANDERSON AND M. A. GIRSHICK: "Some extensions of the Wishart distribution," *Annals of Math. Stat.*, Vol. 15 (1944), pp. 354 ff.

### NOTE ON THE DISTRIBUTION OF THE SERIAL CORRELATION COEFFICIENT<sup>1</sup>

BY WILLIAM G. MADOW

*Bureau of the Census*

The distribution of the serial correlation coefficient when  $\rho = 0$  has been previously obtained.<sup>2</sup> The purpose of this note is to derive the distribution of the serial correlation coefficient, using the circular definition, when  $\rho \neq 0$ .

Let us assume that the random variables  $x_1, \dots, x_N$  have a joint normal distribution<sup>3</sup>  $p(x_1, \dots, x_N | A, B, \mu)$  where

$$\log p(x_1, \dots, x_N | A, B, \mu) \\ = \log K_1 - \frac{1}{2} \left[ A \sum_i (x_i - \mu)^2 + 2B \sum_i (x_i - \mu)(x_{i+L} - \mu) \right]$$

the term in the bracket is positive definite,  $K_1$  is independent of the  $x_i$  and if  $i + L > N$  then  $x_{i+L} = x_{i+L-N}$ . It is then clear that  $\bar{x}$ ,  $V_N$ , and  ${}_L C_N$ , where  $\bar{x}$  is the arithmetic mean,  $V_N = \sum_i (x_i - \bar{x})^2$  and

$${}_L C_N = \sum_i (x_i - \bar{x})(x_{i+L} - \bar{x})$$

are sufficient statistics with respect to the estimation of  $\mu$ ,  $A$ , and  $B$ .

Let  $V_N {}_L R_N = {}_L C_N$  define  ${}_L R_N$ , the serial correlation coefficient. Then if

<sup>1</sup> Presented at a meeting of the Cowles Commission for Economic Research in Chicago, January 31, 1945.

<sup>2</sup> See R. L. Anderson, "Distribution of the serial correlation coefficient", pp. 1-13 and T. Koopmans, "Serial correlation and quadratic forms in normal variables", pp. 14-33, *Annals of Math. Stat.*, Vol. XIII, No. 1, March, 1942.

<sup>3</sup> The expression  $p(\xi_1, \dots, \xi_m | \theta_1, \dots, \theta_g)$  means the probability density or the distribution of the random variables  $\xi_1, \dots, \xi_m$  for the given values of the parameters  $\theta_1, \dots, \theta_g$ . When used as an index of summation or multiplication, the letter  $i$  will assume all values from 1 through  $N$ .

$A = 1, B = 0$  Anderson has shown<sup>4</sup> that, if  $N$  is odd, the joint distribution of  ${}_1R_N$  and  $V_N$  is given by

$$(1) \quad D(R_N, V_N) = KV_N^{\frac{1}{2}(N-3)} e^{-\frac{1}{2}V_N} \sum_{i=1}^m (\lambda_i - R_N)^{\frac{1}{2}(N-5)/\alpha_i}, \quad \text{for } \lambda_{m+1} \leq R_N \leq \lambda_m$$

where

$$R_N = {}_1R_N, \quad \lambda_k = \cos \frac{2\pi k}{N}, \quad \alpha_i = \prod_{j=1}^{\frac{1}{2}(N-1)} (\lambda_i - \lambda_j), \quad \text{for all } j \neq i$$

and  $K^{-1} = 2^{\frac{1}{2}(N-1)} \Gamma[\frac{1}{2}(N-3)]$ ; while if  $N$  is even, the same formula holds except that

$$\alpha_i = \prod_{j=1}^{\frac{1}{2}(N-2)} (\lambda_i - \lambda_j) \sqrt{(\lambda_i + 1)}, \quad \text{for all } j \neq i.$$

We now extend Anderson's distributions to the case where it is not assumed that  $A = 1$  and  $B = 0$ .

As a means of extending<sup>5</sup> Anderson's distribution let us recall that if  $x_1, \dots, x_N$  have a distribution  $p(x_1, \dots, x_N | \theta_1, \dots, \theta_\sigma)$  depending on several parameters  $\theta_1, \dots, \theta_\sigma$ , and if  $z_1, \dots, z_k$  are a sufficient set of statistics with respect to  $\theta_1, \dots, \theta_\sigma$ , i.e.

$$p(x_1, \dots, x_N | \theta_1, \dots, \theta_\sigma) = h(z_1, \dots, z_k | \theta_1, \dots, \theta_\sigma) m(x_1, \dots, x_N)$$

where  $m(x_1, \dots, x_N)$  is independent of  $\theta_1, \dots, \theta_\sigma$ , then if the distribution of  $z_1, \dots, z_k$  is found, assuming  $\theta_1, \dots, \theta_\sigma$  have specific values  $\theta_1^0, \dots, \theta_\sigma^0$ , then it follows that

$$p(z_1, \dots, z_k | \theta_1, \dots, \theta_\sigma) = p(z_1, \dots, z_k | \theta_1^0, \dots, \theta_\sigma^0) \frac{h(z_1, \dots, z_k | \theta_1, \dots, \theta_\sigma)}{h(z_1, \dots, z_k | \theta_1^0, \dots, \theta_\sigma^0)}$$

We may call Anderson's distribution given in (1),  $p(R_N, V_N | 1, 0)$ , i.e.

$$p(R_N, V_N | 1, 0) = D(R_N, V_N)$$

Furthermore,  $\bar{x}$  is distributed independently of  $R_N$  and  $V_N$  for all values of  $A$  and  $B$  and hence by a simple transformation,<sup>6</sup> we can apply the above theorem.

<sup>4</sup> Anderson loc. cit. p. 3 and p. 5. Although the remainder of the note deals only with the case where  $L = 1$  the procedure is general and may be easily carried through for other lags.

<sup>5</sup> See W. G. Madow Contributions to the "Theory of multivariate statistical analysis", *Trans. of the Amer. Math. Soc.*, Vol. 44, No. 3, November 1938, p. 461.

<sup>6</sup> For a proof that an orthogonal transformation of the variable  $x_i - \mu$  exists such that  $V_N$  and  ${}_L C_N$  are simultaneously reduced to canonical forms involving the same  $N - 1$  of the variables of the transformation, and  $\sqrt{N}(\bar{x} - \mu)$  is the  $N$ th variable of the transformation, see J. von Neumann, "Distribution of the ratio of the mean square successive difference to the variance, *Annals of Math. Stat.*, Vol. XII, No. 4, December 1941, pp. 368, 369. The proof there is given for  $V_N$  and  $\sum_i (x_i - x_{i+1})^2$  but is easily extended to this case.

Then it is easy to show that  $N(\bar{x} - \mu)$  is independently distributed of  $V_N$ , and  ${}_L C_N$  and has distribution  $\log p[\sqrt{N}(\bar{x} - \mu) | A, B] = \log K_2 - \frac{1}{2}[A + 2B]N(\bar{x} - \mu)_2$  where  $K_2 = (2\pi)^{-\frac{1}{2}}(A + 2B)^{\frac{1}{2}}$  and  $K'_1 K_2 = K_1$ .

Then

$$p(R_N, V_N | A, B) = p(R_N, V_N | 1, 0)\Omega$$

where

$$\Omega = \frac{K_1' e^{-\frac{1}{2}(AV_N + 2BR_N V_N)}}{(2\pi)^{-\frac{1}{2}N} e^{-\frac{1}{2}VN}}$$

Hence it follows that,

$$p(R_N, V_N | A, B) = KK_1'(2\pi)^{\frac{1}{2}N} V_N^{\frac{1}{2}(N-3)} e^{-\frac{1}{2}V_N(A+2BR_N)} \sum_{i=1}^m (\lambda_i - R_N)^{\frac{1}{2}(N-5)}/\alpha_i,$$

for  $\lambda_{m+1} \leq R_N \leq \lambda_m$ , where the  $\alpha_i$  have different values according to whether  $N$  is odd or even. In order to evaluate  $p(R_N | A, B)$  we then need only integrate out  $V_N$ . Now

$$\int_0^\infty V_N^{\frac{1}{2}(N-3)} e^{-\frac{1}{2}V_N(A+2BR_N)} dV_N = \Gamma[\frac{1}{2}(N-1)](A/2 + BR_N)^{-\frac{1}{2}(N-1)}.$$

Hence

$$p(R_N | A, B) = KK_1'(2\pi)^{\frac{1}{2}N} \Gamma[\frac{1}{2}(N-1)](A/2 + BR_N)^{-\frac{1}{2}(N-1)} \sum_{i=1}^m (\lambda_i - R_N)^{\frac{1}{2}(N-5)}/\alpha_i.$$

The parameters  $K_1'$ ,  $A$  and  $B$  depend on the different types of assumptions that may be made. In general

$$K_1 = (2\pi)^{-\frac{1}{2}N} \Delta^{1/2}$$

where  $\Delta$  is a circulant  $(a_1, \dots, a_N)$  such that

$$a_1 = A, \quad a_{1+L} = B, \quad a_{1+(N-L)} = B, \quad a_i = 0 \text{ otherwise,}$$

and hence

$$\Delta = \prod_i \left( A + B \cos \frac{2\pi iL}{N} \right) = \prod (A + B\lambda_i).$$

Then, one assumption is

$$A = \frac{1}{\sigma^2}, \quad B = -\rho/\sigma^2$$

where  $\rho$  is the "true" serial correlation coefficient. Other assumptions are possible.<sup>7</sup> However, these vary with the problem under consideration and may be left for further examination.

---

<sup>7</sup> One possible alternative definition is given by W. J. Dixon, "Further contributions to the problem of serial correlation", *Annals of Math. Stat.*, Vol. XV, No. 2, June 1944, p. 120, equation (2.1).