

SAMPLING IN THE CASE OF CORRELATED OBSERVATIONS

By

CECIL C. CRAIG

National Research Fellow

Dr. E. C. Rhodes, in a paper in the *Journal of the Royal Statistical Society*,¹ has considered the distribution of characteristics of samples of N when the individual observations are not assumed to be independent. As he points out, there are many important cases in which the usual assumption of independence or randomness in the observations is not justifiable. In the present paper will be explained a method based on the semi-invariants of Thiele for the calculation of the characteristics of the sought distributions in this case which is especially to be preferred to the method based on moments when it is supposed that the observations are normally correlated. In the case it is further assumed that only consecutive observations are correlated, in addition to Dr. Rhodes' results, the third semi-invariant (which is the same as the third moment about the mean) of the variance and the mean and the variance of the third and fourth moments about the mean are given.

Let the N observations composing a sample be given by values of x_1, x_2, \dots, x_N respectively and let $F_N(x_1, x_2, \dots, x_N)$ be the n -way probability function of x_1, x_2, \dots , and x_N .

¹The Precision of Means and Standard Deviations When the Individual Errors Are Correlated, Vol. 90 (1927), pp. 135-143.

Then the semi-invariants, $\lambda_{rst} \dots$ of x_1, x_2, \dots, x_N are defined by

$$\begin{aligned}
 & e^{(\sum_1^N \lambda_i t_i) + \frac{1}{2} (\sum_1^N \lambda_i t_i)^{(2)} + \frac{1}{3!} (\sum_1^N \lambda_i t_i)^{(3)} + \dots} \\
 (1) \quad & = \int_{-\infty, \dots, -\infty}^{\infty, \dots, \infty} dF_N(x_1, x_2, \dots, x_N) e^{(\sum_1^N x_i t_i)} \quad 1
 \end{aligned}$$

which is to be regarded as a formal identity in t_1, t_2, \dots, t_N . $(\sum_1^N \lambda_i t_i)^{(k)}$ is first expanded by the multinomial law and then each term $\lambda_1^r, \lambda_2^s, \lambda_3^t \dots$ in the result is replaced by $\lambda_{rst} \dots$.

We shall pass over the characteristics of distributions of means, since the method of semi-invariants is equivalent to that of moments in this case, and take up the distribution of moments about the mean in samples of N . Following the method previously used by the author in the case of independent observations,² let

$$\begin{aligned}
 (2) \quad & \delta_i = x_i - \sum_1^N \frac{x_j}{N} \\
 & = \sum_{j=1}^N a_{ij} x_j \quad \text{with} \quad \begin{cases} a_{ij} = -\frac{1}{N} \\ a_{ii} = \frac{N-1}{N} \end{cases}
 \end{aligned}$$

Then let $V(\delta_1, \delta_2, \dots, \delta_{N-1})$ be the probability function of $\delta_1, \delta_2, \dots, \delta_{N-1}$, ($\sum_1^N \delta_i = 0$). The semi-invariants $\lambda'_{rst} \dots$ of $\delta_1, \delta_2, \dots, \delta_{N-1}$ are defined by

¹Following Cramér, I distinguish between probability and frequency functions. $F_N(x_1, x_2, \dots, x_N)$ is the "cumulative" frequency function and thus the integral is an n -way Stieltjes integral.

²An Application of Thiele's Semi-invariants to the Sampling Problem; *Metron*, Vol. 7, No. 4 (1928), pp. 3-75.

$$\begin{aligned}
 & e^{\left(\sum_1^{N-1} \lambda'_i t_i\right) + \frac{1}{2} \left(\sum_1^{N-1} \lambda'_i t_i\right)^{(2)} + \frac{1}{3!} \left(\sum_1^{N-1} \lambda'_i t_i\right)^{(3)} + \dots} \\
 (3) \quad & = \int_{-\infty, \dots, -\infty}^{\infty, \dots, \infty} dV(\delta_1, \dots, \delta_{N-1}) e^{\left(\sum_1^{N-1} \delta_i t_i\right)} \\
 & = \int_{-\infty, \dots, -\infty}^{\infty, \dots, \infty} dF_N(x_1, x_2, \dots, x_N) e^{\left(\sum_{i=1}^{N-1} \sum_{j=1}^N a_{ij} x_j t_i\right)}
 \end{aligned}$$

We have at once,

$$\left(\sum_{i=1}^{N-1} t_i \sum_{j=1}^N \lambda_j a_{ij}\right)^{(K)} = \left(\sum_{i=1}^{N-1} \lambda'_i t_i\right)^{(K)}$$

and as the author has previously remarked,¹ we can also write

$$(4) \quad \left(\sum_1^N t_i \sum_{j=1}^N \lambda_j a_{ij}\right)^{(K)} = \left(\sum_1^N \lambda'_i t_i\right)^{(K)}$$

so long as the relation is only used to find the values of $\lambda'_{rst\dots}$'s in which at least one of the subscripts is zero.

Then $S_K(V_n)$, the k 'th semi-invariant of the n 'th moment about the mean in samples of N , is given by the formula

$$\begin{aligned}
 S_K(V_n) = & \\
 & \frac{1}{N^K \Sigma \Sigma \dots} \frac{(-1)^{(r+s+t+\dots)-k} [(r+s+t+\dots)-k]! K! \int_{a_1, a_2, \dots}^r \int_{b_1, b_2, \dots}^s \int_{c_1, c_2, \dots}^t \dots}{[a_1! a_2! \dots]^r [b_1! b_2! \dots]^s [c_1! c_2! \dots]^t r! s! t! \dots}
 \end{aligned}$$

¹loc. cit., pp. 18, 19.

the notation V'_{uvw} referring to moments of $\delta_1, \dots, \delta_{N-1}, \delta_N$, the summation including all terms for which

$$r(a_1 + a_2 + \dots) + s(b_1 + b_2 + \dots) + t(c_1 + c_2 + \dots) + \dots = k,$$

$$a_1 \geq a_2 \geq \dots$$

$$b_1 \geq b_2 \geq \dots$$

$$c_1 \geq c_2 \geq \dots$$

.....

$$(a_1 + a_2 + \dots) > (b_1 + b_2 + \dots) > (c_1 + c_2 + \dots) > \dots$$

In particular:

$$S_1(V_N) = \frac{1}{N} \sum V'_{n,0}, \quad (\sum V'_{n,0} = V'_{n,0\dots 0} + V'_{0,n,0\dots 0} + V'_{00,n,0\dots 0}, \dots)$$

$$S_2(V_N) = \frac{1}{N^2} [\sum V'_{2n,0} + 2\sum V'_{n,n,0} - (\sum V'_{n,0})^2],$$

(5)

$$S_3(V_N) = \frac{1}{N^3} [\sum V'_{3n,0} + 3\sum V'_{2n,n,0} + 6\sum V'_{n,n,n,0} - 3(\sum V'_{2n,0})(\sum V'_{n,0}) - 6(\sum V'_{n,n,0})(\sum V'_{n,0}) + 2(\sum V'_{n,0})^3],$$

On writing out the moments V'_{uvw} in terms of the semi-invariants λ'_{rst} ² and then using (4) the sought semi-invariants are obtained.

In the case that the N observations are normally correlated and $F_N(x_1, x_2, \dots, x_N)$ is the N -dimensional normal probability function, the left-hand member of (4) vanishes for $k \geq 3$.

If we suppose that the standard deviations of x_1, x_2, \dots, x_N are all equal (which we shall always do) and take as the simplest case that x_1, x_2, \dots, x_N are normally correlated and that

¹See the author's paper cited, p. 21, formula (25).

²For a detailed explanation of this kind of calculation see the author's paper cited, pp. 23-27.

the correlation as measured by the Pearsonian coefficient, $r_{x_i x_j}$, is the same for each pair, x_i, x_j , of the set of N observations, we get

$$\lambda'_{20} = \lambda'_{020} = \lambda'_{0020} = \dots = \frac{N-1}{N} (\lambda_{20} - \lambda_{11}) = \frac{N-1}{N} (1-r) \lambda_{20},$$

$$\lambda'_{110} = \lambda'_{1010} = \lambda'_{0110} = \dots = -\frac{1}{N} (\lambda_{20} - \lambda_{11}) = -\frac{1}{N} (1-r) \lambda_{20},$$

if the common value of $r_{x_i x_j}$ be denoted simply by r . But if the observations are independent and the parent population is normal we have

$$\lambda'_{20} = \lambda'_{020} = \lambda'_{0020} = \dots = \frac{N-1}{N} \lambda_{20},$$

$$\lambda'_{110} = \lambda'_{1010} = \lambda'_{0110} = \dots = -\frac{1}{N} \lambda_{20}.$$

Thus it follows that the distributions of the characteristics of samples of N in this particular case of dependent observations are the same as if the observations were independent and taken from a normal population of variance $(1-r) \lambda_{20}$.

In case $F_N(x_1, x_2, \dots, x_N)$ is normal it is convenient to express the right hand members of (5) directly in terms of the semi-invariants $\lambda'_{rst\dots}$ for $n=2, 3, 4$. For that purpose we shall adopt the following notation. Let the linear form $\sum_{j=1}^N a_{ij} \lambda_j$ be denoted by A_i . Then (4) becomes

$$(6) \quad \left(\sum_i^N A_i t_i \right)^{(k)} = \left(\sum_i^N \lambda'_i t_i \right)^{(k)}.$$

Thus in a symbolic sense A_i 's and λ'_i 's are equivalent. But with regard to the subscripts of the A terms in the expansion of the left member of (6) we use a different convention than for the subscripts of the λ 's. We set

¹See the author, loc. cit., p. 19.

$$\lambda'_{20} = A_{11}, \lambda'_{020} = A_{22}, \dots$$

$$\lambda'_{110} = A_{12}, \lambda'_{1010} = A_{13}, \dots$$

We get

$$S_1(V_2) = \frac{1}{N} \sum A_{ii}$$

$$S_2(V_2) = \frac{1}{N^2} [3\sum A_{ii}^2 + 2\sum A_{ii} A_{jj} + 4\sum A_{ij}^2 - (\sum A_{ii})^2], i \neq j,$$

the summations, of course, running over all values of i and j from 1 to N . But since

$$\sum A_{ii}^2 + 2\sum A_{ii} A_{jj} = (\sum A_{ii})^2$$

the second relation reduces to

$$S_2(V_2) = \frac{2}{N^2} (\sum A_{ii}^2 + 2\sum A_{ij}^2).$$

Similarly

$$S_3(V_2) = \frac{8}{N^3} (\sum A_{ii}^3 + 3\sum A_{ii} A_{ij}^2 + 6\sum A_{ij} A_{ik} A_{jk}),$$

$$S_4(V_2) = \frac{48}{N^4} (\sum A_{ii}^4 + 4\sum A_{ii}^2 A_{ij}^2 + 4\sum A_{ii} A_{jj} A_{ij}^2 + 2\sum A_{ij}^4$$

(7) $+ 8\sum A_{ii} A_{ij} A_{ik} A_{jk} + 4\sum A_{ij}^2 A_{ik}^2 + 8\sum A_{ij} A_{ik} A_{jl} A_{kl}),$
 $\dots \dots \dots$

$$S_1(V_3) = 0,$$

$$S_2(V_3) = \frac{3}{N^2} (5\sum A_{ii}^3 + 6\sum A_{ii} A_{jj} A_{ij} + 4\sum A_{ij}^3),$$

$$S_3(V_3) = 0,$$

$\dots \dots \dots$

$$S_1(V_4) = \frac{3}{N} \sum A_{ii}^2,$$

$$S_2(V_4) = \frac{48}{N^2} (2\sum A_{ii}^4 + 3\sum A_{ii} A_{jj} A_{ij}^2 + \sum A_{ij}^2).$$

To illustrate the use of these formulas and to give some results in a case of practical interest, let us suppose that the set of N observations composing a sample may be assigned an order in which only consecutive observations are correlated and in a constant degree. Thus our observations might be prices or indices taken at the ends of consecutive time intervals. We suppose, then, that

$$\lambda_{110} = \lambda_{0110} = \lambda_{00110} = \dots = r \lambda_{20},$$

$$\lambda_{101} = \lambda_{1001} = \lambda_{0101} = \dots = 0$$

The first step in the calculation is to obtain the values of the various A 's which enter into the formulas (7). A_{11} is found from A_1^2 , A_{12} from $A_1 A_2$ and so on. We get

$$\begin{aligned} A_{11} &= A_{N,N} = \left(1 - \frac{1}{N} - \frac{2r}{N^2}\right) \lambda_{20}, \\ A_{22} &= A_{33} = \dots = A_{N-1,N-1} = \left(1 - \frac{1+2r}{N} - \frac{2r}{N^2}\right) \lambda_{20}, \\ A_{12} &= A_{N-1,N} = \left(r - \frac{1+r}{N} - \frac{2r}{N^2}\right) \lambda_{20}, \\ (8) \quad A_{23} &= A_{34} = \dots = A_{N-2,N-1} = \left(r - \frac{1+2r}{N} - \frac{2r}{N^2}\right) \lambda_{20}, \\ A_{13} &= A_{14} = \dots = A_{1,N-1}, \\ &= A_{2,N} = A_{3,N} = \dots = A_{N-2,N} = \left(-\frac{1+r}{N} - \frac{2r}{N^2}\right) \lambda_{20}, \\ A_{i,N} &= \left(-\frac{1}{N} - \frac{2r}{N^2}\right) \lambda_{20}, \\ A_{ij} &= \left(-\frac{1+2r}{N} - \frac{2r}{N^2}\right) \lambda_{20} \quad \begin{cases} 1 < i < N-1 \\ 1 < j < N-1 \\ |i-j| > 1 \end{cases} \end{aligned}$$

Then, on substitution in (7), we have finally

$$S_1(V_2) = \left(1 - \frac{1}{N}\right) \left(1 - \frac{2r}{N}\right) \lambda_{20},$$

$$S_2(V_2) = \frac{2}{N} \left[\left(1 - \frac{1}{N}\right) \left(1 - \frac{4r}{N}\right) + 2r^2 \left(1 - \frac{3}{N} + \frac{2}{N^2} + \frac{2}{N^3}\right) \right] \lambda_{20}^2.$$

These two results are given by Dr. Rhodes, loc. cit., though there is a slight misprint in the second one as given there. The remainder of the results given here are believed to be new.

$$S_3(V_2) = \frac{8}{N^2} \left[\left(1 - \frac{1}{N}\right) \left(1 - \frac{6r}{N}\right) + 6r^2 \left(1 - \frac{3}{N} + \frac{2}{N^2} + \frac{2}{N^3}\right) - \frac{4r^3}{N} \left(2 - \frac{3}{N} - \frac{3}{N^2} - \frac{2}{N^3}\right) \right] \lambda_{20}^3,$$

$$S_1(V_3) = 0,$$

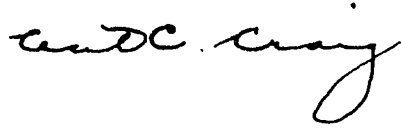
$$S_2(V_3) = \frac{6}{N} \left[\left(1 - \frac{1}{N}\right) \left(1 - \frac{2}{N}\right) \left(1 - \frac{6r}{N}\right) - \frac{6r^2}{N} \left(1 - \frac{5}{N} + \frac{12}{N^3}\right) + \frac{2r^3}{N^2} \left(1 - \frac{7}{N} + \frac{14}{N^2} + \frac{2}{N^3} - \frac{24}{N^4} - \frac{40}{N^5}\right) \right] \lambda_{20}^4,$$

$$S_3(V_3) = 0,$$

$$S_1(V_4) = 3 \left[\left(1 - \frac{1}{N}\right)^2 \left(1 - \frac{4r}{N}\right) + \frac{4r^2}{N^2} \left(1 - \frac{3}{N^2}\right) \right] \lambda_{20}^2,$$

$$S_2(V_4) = \frac{24}{N} \left[\left(1 - \frac{1}{N}\right) \left(4 - \frac{9}{N} + \frac{6}{N^2}\right) \left(1 - \frac{8r}{N}\right) + 6r^2 \left(1 - \frac{3}{N} + \frac{25}{N^2} - \frac{23}{N^3} - \frac{74}{N^4} + \frac{68}{N^5}\right) - \frac{8r^3}{N} \left(4 - \frac{19}{N} + \frac{33}{N^2} + \frac{30}{N^3} - \frac{54}{N^4} - \frac{108}{N^5}\right) + 2r^4 \left(1 - \frac{9}{N} + \frac{44}{N^2} - \frac{64}{N^3} - \frac{114}{N^4} + \frac{192}{N^5} + \frac{360}{N^6} + \frac{288}{N^7}\right) \right] \lambda_{20}^4.$$

It should be observed that the expressions for $S_1(V_n)$ for $N < 3$ and for $S_k(V_n)$, $k \geq 2$ for $N < 5$ are in general not valid, since it can be seen by reference to (8) that all the types of A 's used in the formulas (7) do not exist for values of N so small. But for these small values of N , the values of the characteristics for which expressions are given above can be readily computed directly.



Stanford University.