

EDITORIAL: A. L. O'TOOLE

ON A BEST VALUE OF R IN SAMPLES OF R FROM A FINITE POPULATION OF N .

In recent years the problem of finding the moment coefficients for samples of n drawn from a finite population of N has been of interest to so many writers¹ that it seems worthwhile to make a few further observations² concerning these moment coefficients—particularly with respect to their dependence on n . In many instances the value of n to be used is at the discretion of the investigator and he would like to know if there is one value of n which is better than any other. An answer to that question will be given here.

¹ H. C. Carver, On the fundamentals of the theory of sampling, *Annals of Mathematical Statistics*, Vol. I, No. 1, pp. 101-121; Vol. I, No. 3, pp. 260-274.

C. C. Craig, An application of Thiele's semi-invariants to the sampling problem, *Metron*, Vol. VII, No. 4, 1928, pp. 3-74.

R. A. Fisher, Moments and product moments of sampling distributions, *Proc. London Math. Soc.*, Series 2, xxix, 1929, pp. 309-321; xxx, 1929, pp. 199-238.

L. Isserlis, On a formula for the product moment coefficients of any order of a normal frequency distribution in any number of variables. *Biometrika*, xii, 1918-19, pp. 134-139.

P. R. Rider, Moments of moments, *Proc. of the National Academy of Sciences*, Vol. 15, 1929, pp. 430-434.

H. E. Soper, Sampling moments of samples of n units each drawn from an unchanging sampled population, from the point of view of semi-invariants, *Journal of the Royal Statistical Soc.*, Vol. 93, 1930, pp. 104-114.

A. A. Tchouproff, On the mathematical expectation of the moments of frequency distributions, *Biometrika*, xii, 1918-19, pp. 140-169 and 184-210; xiii, 1920-21, pp. 283-295.

A. L. O'Toole, On symmetric functions and symmetric functions of symmetric functions, *Annals of Mathematical Statistics*, Vol. II, No. 2, May 1931, pp. 102-149. See Chapter III.

² These observations arose as a result of some very far-reaching suggestions on the theory of sampling made by Professor Carver, during recent conversations with the writer.

The differential operator method developed by this writer³ for finding the moment coefficients not only was a very simple method but had the added advantage of leading directly to some theorems whose generality had not been established previously.

Using the notation of the previous paper let the finite parent population of N be composed of the N variates $x_1, x_2, x_3, \dots, x_N$. From this population draw all of the ${}_N C_r$ different samples and let $z_i = \sum_{\lambda:i}^{x_i} x$, $i = 1, 2, 3, \dots, {}_N C_r$, where $\sum_{\lambda:i} x$ designates the sum of the r values of x which appear in the i^{th} sample. With this notation it has been shown in the paper cited that

$$(1) \quad S_{t; z} = t! \sum \frac{P_i^I \cdot P_j^J \cdot P_k^K \cdots S_{i;x}^I \cdot S_{j;x}^J \cdot S_{k;x}^K \cdots}{(i!)^I \cdot (j!)^J \cdot (k!)^K \cdots (I!)(J!)(K!) \cdots}$$

where $S_{t; z} = \sum_{i=1}^{{}_N C_r} z_i^t$, $t = 1, 2, 3, \dots$

and $S_{w; x} = \sum_{i=1}^N x_i^w$, $w = 1, 2, 3, \dots$

The summation in (1) is to be taken over terms such that $Ii + Jj + Kk + \dots = t$ where $I, J, K, \dots, i, j, k, \dots$ are positive integers, and where P_m is obtained from the m^{th} sampling polynomial $P_m(p)$ by replacing the exponents of the polynomial by corresponding subscripts.

$$(2) \quad P_m(p) = \sum_{i=0}^{m-1} (-1)^i \binom{m-1}{i} p^{i+1}, \text{ OR}$$

$$(3) \quad P_m(p) = \left. \frac{d^m}{dx^m} \log (pe^x + 1 - p) \right]_{x=0}$$

³ Loc. cit.

In particular $P_1 = p_1$, $P_2 = p_1 - p_2$, $P_3 = p_1 - 3p_2 + 2p_3$

$$P_4 = p_1 - 7p_2 + 12p_3 - 6p_4$$

$$P_5 = p_1 - 15p_2 + 50p_3 - 60p_4 + 24p_5$$

$$P_6 = p_1 - 31p_2 + 180p_3 - 390p_4 + 360p_5 - 120p_6$$

$$P_7 = p_1 - 63p_2 + 602p_3 - 2100p_4 + 3360p_5 - 2520p_6 + 720p_7$$

where $p_k = n \cdot {}_n C_{n-k}$, $k \leq n$.

It must be kept in mind that the multiplication of these operators is symbolic. For example, to find $P_i^I P_j^J$ first multiply the polynomials $P_i^I(p)$ and $P_j^J(p)$ by ordinary multiplication and then the result when the exponents in this product are replaced by corresponding subscripts is $P_i^I P_j^J$.

Since in this paper it is desired to consider moments rather than power sums, replace $S_{t:z}$ by $({}_n C_n) \mu'_{t:z}$ and $S_{w:x}$ by $N \mu'_{w:x}$. Then (1) becomes, after dividing by ${}_n C_n$,

$$(4) \quad \mu'_{t:z} = \frac{t!}{n C_n} \sum \frac{P_i^I P_j^J P_k^K \cdots N^{I+J+K} \mu'_{i:x} \mu'_{j:x} \mu'_{k:x} \cdots}{(i!)^I (j!)^J (k!)^K \cdots I! \cdot J! \cdot K!}$$

Now $p_k = n \cdot {}_n C_{n-k}$, $k \leq n$, hence,

$$\frac{p_k}{n C_n} = \frac{n(n-1)(n-2) \cdots (n-k+1)}{n(n-1)(n-2) \cdots (n-k+1)}$$

Substituting this value for each $p_k/n C_n$ in (4) the result is

Equations 5.

$$\mu'_{1:z} = n \mu'_{1:x}$$

$$\mu'_{2:z} = \frac{n}{n-1} [n(n-1) \mu'_{1:x}{}^2 + (n-n) \mu'_{2:x}]$$

$$\mu'_{3:z} = \frac{n}{(n-1)(n-2)} [n^2(n-1)(n-2) \mu'_{1:x}{}^3 + 3n(n-1)(n-n) \mu'_{1:x} \mu'_{2:x} + (n-n)(n-2n) \mu'_{3:x}]$$

$$\mu'_{4;\bar{z}} = \frac{\nu}{(N-1)(N-2)(N-3)} \left[\begin{aligned} & N^3 (\nu-1)(\nu-2)(\nu-3) \mu'_{1;x}{}^4 \\ & + 6 N^2 (\nu-1)(\nu-2)(N-\nu) \mu'_{1;x}{}^2 \mu'_{2;x} \\ & + 4 N (\nu-1)(N-\nu)(N-2\nu+1) \mu'_{1;x} \mu'_{3;x} \\ & + 3 N (\nu-1)(N-\nu)(N-\nu-1) \mu'_{2;x}{}^2 \\ & + (N-\nu)(N^2+6N\nu+6\nu^2+N) \mu'_{4;x} \end{aligned} \right]$$

etc.

Now let $N = a\nu$. Then

Equations 6.

$$\mu'_{1;\bar{z}} = \nu \mu'_{1;x}$$

$$\mu'_{2;\bar{z}} = \frac{\nu^2}{a\nu-1} \left[a(\nu-1) \mu'_{1;x}{}^2 + (a-1) \mu'_{2;x} \right]$$

$$\mu'_{3;\bar{z}} = \frac{\nu^3}{(a\nu-1)(a\nu-2)} \left[\begin{aligned} & a^2(\nu-1)(\nu-2) \mu'_{1;x}{}^3 + 3a(\nu-1)(a-1) \mu'_{1;x} \mu'_{2;x} \\ & + (a-1)(a-2) \mu'_{3;x} \end{aligned} \right]$$

$$\mu'_{4;\bar{z}} = \frac{\nu^3}{(a\nu-1)(a\nu-2)(a\nu-3)} \left[\begin{aligned} & a^3 \nu (\nu-1)(\nu-2)(\nu-3) \mu'_{1;x}{}^4 \\ & + 6 a^2 \nu (\nu-1)(\nu-2)(a-1) \mu'_{1;x}{}^2 \mu'_{2;x} \\ & + 4 a (\nu-1)(a-1) \{ \nu(a-2)+1 \} \mu'_{1;x} \mu'_{3;x} \\ & + 3 a (\nu-1)(a-1) \{ \nu(a-1)-1 \} \mu'_{2;x}{}^2 \\ & + (a-1) \{ \nu(a^2-6a+6)+a \} \mu'_{4;x} \end{aligned} \right]$$

etc.

A partial check at this point is to note that for $a = 1$ only the first term of each of these moment coefficients remains.

Let $a = 2$. Then the above moment coefficients become

$$(7) \left\{ \begin{aligned} \mu'_{1;\bar{x}} &= n \mu'_{1;x} \\ \mu'_{2;\bar{x}} &= \frac{n^2}{2n-1} \left[2(n-1) \mu'^2_{1;x} + \mu'_{2;x} \right] \\ \mu'_{3;\bar{x}} &= \frac{n^3}{2n-1} \left[2(n-2) \mu'^3_{1;x} + 3 \mu'_{1;x} \mu'_{2;x} \right] \\ \mu'_{4;\bar{x}} &= \frac{n^4}{(2n-1)(2n-3)} \left[4n(n-2)(n-3) \mu'^4_{1;x} + 12n(n-2) \mu'^2_{1;x} \mu'_{2;x} \right. \\ &\quad \left. + 4 \mu'_{1;x} \mu'_{3;x} + 3(n-1) \mu'^2_{2;x} - \mu'_{4;x} \right] \\ \mu'_{5;\bar{x}} &= \frac{n^5}{(2n-1)(2n-3)} \left[n(n-3)(n-4) \mu'^5_{1;x} + 20n(n-3) \mu'^3_{1;x} \mu'_{2;x} \right. \\ &\quad \left. - 20 \mu'^2_{1;x} \mu'_{3;x} + 15(n-1) \mu'_{1;x} \mu'^2_{2;x} \right. \\ &\quad \left. - 5 \mu'_{1;x} \mu'_{4;x} \right] \end{aligned} \right.$$

etc.

It is observed that when $a = 2$, i.e. when $n = 2r$, the moment coefficient $\mu'_{3;\bar{x}}$ is independent of the moment coefficient $\mu'_{3;x}$. Also $\mu'_{5;\bar{x}}$ is independent of $\mu'_{5;x}$. But one must not assume that all the odd moment coefficients of \bar{x} are independent of the corresponding odd moment coefficients of x . For $\mu'_{7;\bar{x}}$ is not independent of $\mu'_{7;x}$ as is seen by evaluating F_7 which is the coefficient of $\mu'_{7;x}$ in the expression for $\mu'_{7;\bar{x}}$.

So far the moments considered have been the moments of x with respect to the origin from which x is measured and the moments of \bar{x} with respect to the origin from which \bar{x} is measured. Consider now the moments of x about the mean value of x and the moments of \bar{x} about the mean value of \bar{x} . That is let

$$\begin{aligned} \bar{z}_i &= z_i - M_{\bar{x}}, & i &= 1, 2, 3, \dots, n C_2; \\ \bar{x}_i &= x_i - M_x, & i &= 1, 2, 3, \dots, N. \end{aligned}$$

Then

$$\begin{aligned}\bar{z}_i &= z_i - M_z = \sum_{x_i} x_i - n M_x \text{ since } M_z = n M_x \text{ by (5),} \\ &= \sum_{x_i} (x_i - M_x) = \sum_{x_i} \bar{x}.\end{aligned}$$

$$\begin{aligned}\text{e.g. } \bar{z}_1 &= z_1 - M_z = x_1 + x_2 + x_3 + \dots + x_n - n M_x \\ &= (x_1 - M_x) + (x_2 - M_x) + (x_3 - M_x) + \dots + (x_n - M_x) \\ &= \bar{x}_1 + \bar{x}_2 + \bar{x}_3 + \dots + \bar{x}_n \\ &= \sum_{x_i} \bar{x}.\end{aligned}$$

Hence it is clear that \bar{z} is the same function of \bar{x} as z is of x . In other words $\mu_{z;\bar{z}}$ — (the moment of z about the mean of z) — is the same function of $\mu_{1;x}, \mu_{2;x}, \mu_{3;x}, \dots$ — (the moments of x about the mean of x) — as $\mu_{z;z}$ was of $\mu'_{1;x}, \mu'_{2;x}, \dots$. There is one important simplification however due to the fact that $\mu_{1;x} = 0$ and hence all terms which involve $\mu_{1;x}$ vanish. With this in mind (6) becomes

$$(8) \left\{ \begin{aligned}\mu_{1;\bar{z}} &= 0, & \mu_{2;\bar{z}} &= \frac{n^2(a-1)}{a n - 1} \mu_{2;x}, \\ \mu_{3;\bar{z}} &= \frac{n^3(a-1)(a-2)}{(a n - 1)(a n - 2)} \mu_{3;x} \\ \mu_{4;\bar{z}} &= \frac{3 a n^3 (n-1)(a-1)[n(a-1)-1]}{(a n - 1)(a n - 2)(a n - 3)} \mu_{2;x}^2 \\ &+ \frac{n^4(a^3 - 7a^2 + 12a - 6) + n^3 a(a-1)}{(a n - 1)(a n - 2)(a n - 3)} \mu_{4;x} \\ \mu_{5;\bar{z}} &= \frac{10 a (a-1)(a-2)(n-1)(a n - n - 1)}{(a n - 1)(a n - 2)(a n - 3)(a n - 4)} \mu_{2;x} \mu_{3;x} \\ &+ \frac{n^4(a-1)(a-2)(a^2 n - 12 a n + 12 n + 5 a)}{(a n - 1)(a n - 2)(a n - 3)(a n - 4)} \mu_{3;x}^2\end{aligned}\right.$$

etc.

Here again it is noticed that for $a = 2$, i.e. for $n = 2h$, $\mu_{3;\bar{z}}$ is independent of $\mu_{3;x}$. In other words the skewness of the distribution of z is independent of the skewness of the parent

population of x . Similarly $\mu_{5;\bar{z}}$ is independent⁴ of $\mu_{5;x}$ and also independent of $\mu_{3;x}$. But since P_7 is not zero for $a = 2$, $\mu_{7;\bar{z}}$ is not independent of $\mu_{7;x}$.

Now consider the variance of \bar{z} ,

$$\begin{aligned}\mu_{2;\bar{z}} &= \frac{n^2(a-1)}{a^2-1} \mu_{2;x} \\ &= \frac{Nn - n^2}{N-1} \mu_{2;x} \quad (\text{since } N = an).\end{aligned}$$

Obviously it would be very desirable to have the variance (squared standard deviation) a minimum. Since the variance is a function of n differentiate $\mu_{2;\bar{z}}$ with respect to n .

$$\frac{d}{dn} \mu_{2;\bar{z}} = \frac{N-2n}{N-1} \mu_{2;x}.$$

To make $\mu_{2;\bar{z}}$ a minimum $\frac{N-2n}{N-1} \mu_{2;x} = 0$ and hence $N = 2n$ or, that is, $a = 2$.

$$\text{When } a = 2, \quad \mu_{2;\bar{z}} = \frac{N^2}{4(N-1)} \mu_{2;x}, \quad \sigma_{\bar{z}} = \frac{N}{2\sqrt{N-1}} \sigma_x.$$

In conclusion it may be said that there would seem to be good reason to suggest that, when possible, the investigator arrange to have twice as many variates in the control group or parent population as in each of the samples to be analyzed. Taking $n = \frac{N}{2}$ will insure that the skewness of the samples will be independent of the skewness of the parent population and also that the fifth moment of the samples will be independent of the fifth moment of the sampled population. In addition, taking $n = \frac{N}{2}$ will cause the variance (squared standard deviation) of the samples to be a minimum. Choosing $n = \frac{N}{2}$ presumes, of course, that N is an even number. But in most instances it should be possible to arrange that N be even. For if an odd number of observations are given either another observation may be added or one of the given observations deleted to make N even.

⁴ P_5 vanishes with P_3 because $P_5 = P_3(1 - 12P_2)$. But P_3 is not a factor of P_7 .