

ON THE ELIMINATION OF SYSTEMATIC ERRORS DUE TO GROUPING

By
JOHN R. ABERNETHY

In the calculation of the moments of a frequency distribution it is often desirable, or even necessary, to consider not the distribution itself but another derived from it by certain groupings. As a first approximation to the moments of the original distribution we take the corresponding moment of the grouped distribution. But this first approximation is not satisfactory, and it is necessary to obtain some method for the elimination of part of the error committed in replacing the moments of the original distribution by the corresponding moments of the grouped distribution.

This problem was first discussed by W. F. Sheppard in a paper : *On the Calculation of the most Probable Value of Frequency-Constants, for Data arranged according to Equidistant Divisions of a Scale.*¹ If we denote the n -th moment of the original distribution by μ_n and the n -th moment of the grouped distribution by ν_n , we will have Sheppard's corrections in the form :

$$\begin{aligned}\mu_1 &= \nu_1 = 0, \\ \mu_2 &= \nu_2 - \frac{1}{12}, \\ \mu_3 &= \nu_3, \\ \mu_4 &= \nu_4 - \frac{1}{2} \nu_2 + \frac{7}{240}, \\ \mu_5 &= \nu_5 - \frac{5}{6} \nu_3, \qquad \text{etc.}\end{aligned}$$

As pointed out by Karl Pearson² the hypotheses under which these formulae have been obtained are: (a) that Taylor's theorem

¹ Proceedings London Mathematical Society, Vol. 29, p. 353-380.

² *On an elementary proof of Sheppard's formulae for correcting raw moments and other allied points*, editorial in *Biometrika*, Vol. 3, p. 308-312.

may be applied to the frequency function throughout the range; (b) $x^i f^{(j)}(x)$ is finite and continuous throughout the range; (c) $f(x)$ and its derivatives vanish at the limits of the range. These hypotheses are not always satisfied by the frequency functions with which the statistician has to work; and as it is impossible to tell before calculating the moments of a distribution whether the corresponding theoretical frequency function satisfies these conditions, it is desirable to study the problem from another standpoint.

A comparison of the title of Sheppard's paper and the paper itself suggests the question, in what sense do Sheppard's formulae give the most probable value of the moments of a distribution? A partial answer is given by B. L. Shook in the *Synopsis of Elementary Mathematical Statistics*.³ Miss Shook presents⁴ the formulae

$\mu_1 = \nu_1 = 0$, $\mu_2 = \nu_2 - \frac{1}{12} \left(1 - \frac{1}{m^2}\right)$, and $\mu_3 = \nu_3$ for a discrete distribution with m values of the variable grouped in each class interval and shows that for a particular distribution these formulae serve to eliminate the systematic errors from M , μ_2 , and μ_3 . Two problems are suggested by the *synopsis*: the derivation of formulae for the class of discrete distributions, as these three formulae are stated without proof;⁵ the proof that this larger set of formulae and those of Sheppard do serve under all conditions to eliminate the systematic error due to grouping, subject only to the existence of the moments involved. When we have solved these two problems, we shall be in position to understand the true nature of the approxi-

³ THE ANNALS OF MATHEMATICAL STATISTICS, Vol. 1; p. 34-40.

⁴ These formulae are only special cases of a more general formula stated by H. C. Carver in an editorial: ANNALS OF MATHEMATICAL STATISTICS, Vol. 1; formula (14), p. 111.

⁵ Two methods of developing this formula suggest themselves: (a) the elimination of the moment of a continuous graduating function expressed in terms of fine groupings of class intervals of $\frac{1}{m}$ on the one hand and in the terms of coarse groupings of unit class intervals on the other; (b) by a process similar to that of Sheppard, employing for example Lubbock's formula instead of the Euler-Maclaurin sum formula. According to a statement made by Professor Carver, the formulae in question were derived by the latter process.

mation involved in employing Sheppard's corrections and corrections similar to them for discrete distributions.

The problem we wish to consider is this. Given the probabilities $\Phi(x_i)$ that a value of the statistical variable x taken at random will fall within the interval $x_i - \frac{1}{2} < x < x_i + \frac{1}{2}$, we wish to find the moments of the distribution. We consider this problem for two classes of distributions: the distribution of a discrete variable; the distribution of a continuous variable. In either case we shall work with the uni-frequency distributional function $f(x)$. For the discrete distribution $\frac{1}{m} f\left(\frac{j}{m}\right)$ represents the probability that a value of x taken at random will be the number $\frac{j}{m}$; m denotes a definite positive integer; ($j = -2, -1, 0, 1, 2, \dots$). For the continuous distribution $\int_a^b f(x) dx$ represents the probability that a value of x taken at random will fall within the interval $a < x < b$. Thus the function $f(x)$ has the value zero outside the range of the distribution and we may for the sake of convenience denote the limits of summation and integration as $\pm \infty$. For the n -th moment about the origin we have:

$$\mu_n' = \sum_{j=-\infty}^{+\infty} \left(\frac{j}{m}\right)^n f\left(\frac{j}{m}\right) \cdot \frac{1}{m},$$

for the discrete distribution; and

$$\mu_n' = \int_{-\infty}^{+\infty} x^n f(x) dx,$$

for the continuous distribution. What we want is the value of μ_n' . What we are able to find is the value of

$$\nu_n' = \sum_{i=-\infty}^{+\infty} (x_i)^n \Phi(x_i).$$

In establishing approximate relations between the set of true moments $\{\mu_n^i\}$ and the set of raw moments $\{\nu_n^i\}$ we shall employ another set of statistical constants $\{\bar{\nu}_n^i\}$. For the discrete distribution there are m distinct sets of groupings that can be made, leading to m values of the raw moment ν_n^i ; $\bar{\nu}_n^i$ is used to represent the average of these. Similarly for the continuous distribution, $\bar{\nu}_n^i$ is used to denote the average of the moments ν_n^i corresponding to $x_i = i+t$ for all values of t satisfying $0 \leq t < 1$. We shall call this intermediate set of statistical constants $\{\bar{\nu}_n^i\}$ the average grouped moments of the distribution. We then divide the problem into two parts. First we seek the expression of μ_n^i in terms of the $\{\bar{\nu}_n^i\}$. Secondly we seek the nature of approximation in replacing $\bar{\nu}_n^i$ by ν_n^i . The first of these can be solved completely without approximation and without any assumption other than the existence of the moments involved. We can best understand the nature of the approximation involved in the second after the first of our two problems has been solved.

The m values of ν_n^i corresponding to the m distinct methods of grouping a discrete distribution are given by

$$\nu_n^i(t) = \sum_{i=-\infty}^{+\infty} \left(i+t + \frac{m-1}{2m}\right)^n \sum_{j=0}^{m-1} \frac{1}{m} f\left(i+t + \frac{j}{m}\right); \quad mt = 0, 1, \dots, m-1.$$

The average of these is

$$\bar{\nu}_n^i = \frac{1}{m} \sum_{k=0}^{m-1} \sum_{i=-\infty}^{+\infty} \left(i + \frac{k}{m} + \frac{m-1}{2m}\right)^n \sum_{j=0}^{m-1} \frac{1}{m} f\left(i + \frac{k+j}{m}\right).$$

We shall first express the average grouped moment $\bar{\nu}_n^i$ in terms of the true moments $\{\mu_n^i\}$, and then solve for the $\{\mu_n^i\}$ in terms of the $\{\bar{\nu}_n^i\}$. We wish to arrange the right hand side of the above equation according to values of the argument x appearing in $f(x)$; we therefore let $s = mi+k+j$. This equation then becomes

$$\bar{\nu}_n^i = \frac{1}{m} \sum_{s=-\infty}^{+\infty} \sum_{j=0}^{m-1} \left(\frac{s}{m} + \frac{m-1-2j}{2m}\right)^n \frac{1}{m} f\left(\frac{s}{m}\right)$$

from which, by means of the binomial theorem, we obtain

$$\bar{\nu}_n^i = \sum_{i=0}^n \binom{n}{i} \left\{ \sum_{j=0}^{m-1} \left(\frac{m-1-2j}{2m}\right)^i \frac{1}{m} \right\} \left\{ \sum_{s=-\infty}^{+\infty} \left(\frac{s}{m}\right)^{n-i} f\left(\frac{s}{m}\right) \frac{1}{m} \right\}.$$

But

$$\mu_{n-i}^i = \sum_{s=-\infty}^{+\infty} \left(\frac{s}{m}\right)^{n-i} f\left(\frac{s}{m}\right) \frac{1}{m}.$$

We therefore have

$$(1) \quad \bar{\nu}_n^i = \sum_{i=0}^n \binom{n}{i} b_i(m) \mu_{n-i}^i,$$

where

$$(2) \quad b_i(m) = \sum_{j=0}^{m-1} \left(\frac{m-1-2j}{2m}\right)^i \frac{1}{m}.$$

We shall sometimes write b_i instead of $b_i(m)$ in order to simplify the expression of an equation. The change of order of summation

is based on the assumption that the m summations $\mathcal{V}_n^1(t)$ converge absolutely, an assumption equivalent to that of the existence of μ_n^1 since $f(x)$ has only positive or zero values.

We see immediately from (2) that $b_{2k+1}^{(m)} = 0$, since the terms of the summation cancel each other in pairs, with the possible addition of a middle term equal to zero. The calculation of $b_{2k}^{(m)}$ may readily be effected by means of the Euler-Maclaurin sum formula⁶

$$\sum_{j=0}^{m-1} g\left(j+\frac{1}{2}\right) = \int_0^m g(t) dt + \sum_{i=1}^{\infty} \left[\frac{D_{2i}}{4^i (2i)!} g^{(2i-1)}(t) \right]_0^m,$$

where

$$\begin{aligned} D_0 &= 1, \\ D_2 &= -\frac{1}{3}, \\ D_4 &= \frac{7}{15}, \\ D_6 &= -\frac{31}{21}, \\ D_8 &= \frac{127}{15}. \end{aligned}$$

We substitute

$$g(t) = \frac{1}{m^{2k+1}} \left(t - \frac{1}{2}m\right)^{2k}$$

⁶ See for example Nörlund's *Differenzenrechnung*, Berlin 1924; especially formulae (39), p. 27; (42), p. 28; and (49), p. 30. Formula (39) is

$$\sum_{i=0}^n \binom{n}{i} D_{n-i} - \sum_{i=0}^n (-1)^i \binom{n}{i} D_{n-i} = 0 \quad \text{for } n > 1, D_0 = 1.$$

From this we may obtain the values of D_{2i} and show that $D_{2i+1} = 0$; also

we obtain $\sum_{i=0}^n \binom{2n+1}{2i} D_{2i} = 0$ for $n > 0$ which we shall employ in the

proof of our formula (7).

obtaining

$$(3) \quad b_{2k}(m) = \frac{1}{4^k(2k+1)} \sum_{i=0}^k \binom{2k+1}{2i} \frac{D_{2i}}{m^{2i}}$$

The first few values are:

$$b_0(m) = 1,$$

$$b_2(m) = \frac{1}{12} \left(1 - \frac{1}{m^2}\right),$$

$$b_4(m) = \frac{1}{240} \left(3 - \frac{10}{m^2} + \frac{7}{m^4}\right),$$

$$b_6(m) = \frac{1}{1344} \left(3 - \frac{21}{m^2} + \frac{49}{m^4} - \frac{31}{m^6}\right),$$

$$b_8(m) = \frac{1}{11520} \left(5 - \frac{60}{m^2} + \frac{294}{m^4} - \frac{620}{m^6} + \frac{381}{m^8}\right).$$

A control on the values of $b_i(m)$ may be obtained by substituting $m=1$; then all except b_0 vanish as $b_i(1) = 0$, for $i > 0$.

Having in (1) and (3) obtained the expression of the average grouped moment of a discrete distribution in terms of the true moments, we wish to solve for the true moments in terms of the average grouped moments. We shall obtain this solution by the method of undetermined coefficients. Let

$$(4) \quad \mu_n' = \sum_{j=0}^n \binom{n}{j} A_{n-j} \bar{v}_j'$$

Substituting this in (1) we shall have

$$\bar{v}_n' = \sum_{i=0}^n \binom{n}{i} b_i \sum_{j=0}^{n-i} \binom{n-i}{j} A_{n-i-j} \bar{v}_j'$$

from which we obtain

$$\bar{v}_n^1 = \sum_{j=0}^n \binom{n}{j} \bar{v}_j^1 \sum_{i=0}^{n-j} \binom{n-j}{i} b_i A_{n-j-i},$$

by a change of order of summation effected by applying the Dirchlet sum formula⁷

$$\sum_{i=0}^n \sum_{j=0}^{n-i} w(i, j) = \sum_{j=0}^n \sum_{i=0}^{n-j} w(i, j).$$

Equating coefficients of \bar{v}_j^1 gives us the recurring formula

$$\sum_{i=0}^k \binom{k}{i} b_i A_{k-i} = 0 \text{ for } k > 0,$$

together with the initial condition $A_0 = 1$. This may also be written

$$(5) \quad A_k = -\sum_{i=0}^{k-1} \binom{k}{i} b_{k-i} A_i, \text{ for } k \geq 1; A_0 = 1.$$

Ordinarily in an expression such as (4) we would have written $A_{n-j}(n)$ instead of A_{n-j} ; had we done so in this case, we would now drop the functional expression as we have shown that the value of $A_{n-j}(n)$ depending only on $n-j$ is completely independent of n . The coefficients A_{n-j} are also independent of the position of the origin since if in

$$\mu'_{n:x+h} = \sum_{i=0}^n \binom{n}{i} h^i \mu'_{n-i:x},$$

⁷ For the method of derivation of this formula see, for example, Steffensen's *Interpolation* (Baltimore, 1927), p. 91-92.

we substitute

$$\mu'_{n-l;x} = \sum_{j=0}^{n-l} \binom{n-l}{j} A_j \bar{v}'_{n-l-j;x},$$

we shall have

$$\mu'_{n;x+h} = \sum_{j=0}^n \binom{n}{j} A_j \sum_{l=0}^{n-j} \binom{n-j}{l} h^l \bar{v}'_{n-j-l;x},$$

and hence

$$\mu'_{n;x+h} = \sum_{j=0}^n \binom{n}{j} A_j \bar{v}'_{n-j;x+h}.$$

If in (5) we substitute $k=1$ we shall obtain $A_1 = -b_1 = 0$.

Moreover in general $A_{2i+1} = 0$, since by induction if

$$A_1 = A_3 = \dots = A_{2i-1} = 0,$$

the terms of the summation (5) will have respectively the zero factors

$$b_{2i+1}, A_1, b_{2i-1}, A_3, \dots, b_3, A_{2i-1}, b_1.$$

Also from (5) we obtain:

$$A_0 = 1,$$

$$A_2 = -b_2$$

$$A_4 = -b_4 + 6(b_2)^2,$$

$$A_6 = -b_6 + 30b_2b_4 - 90(b_2)^3,$$

$$A_8 = -b_8 + 56b_2b_6 + 70(b_4)^2 - 1260(b_2)^2b_4 + 2520(b_2)^4.$$

An observation of these expressions of the A_i 's suggests the formula:

$$A_{2i} = \sum \frac{(-1)^k (2i)! k! (b_2)^{a_1} (b_4)^{a_2} \dots (b_{2j})^{a_j}}{(2!)^{a_1} (4!)^{a_2} \dots ((2j)!)^{a_j} (a_1!) (a_2!) \dots (a_j!)}$$

where $k = a_1 + a_2 + \dots + a_j$, the summation extending over all positive integral or zero values of a_1, a_2, \dots, a_j satisfying $a_1 + 2a_2 + \dots + ja_j = i$. That this formula holds in general may be proved by induction: assume it true for $i = 0, 1, \dots, j-1$ and substitute in (5) for $k = 2j$. Upon collecting terms according to products of the b 's we shall have established this formula also for $i = j$, and hence for every positive integral value.

If in the expressions of the A 's in terms of the b 's we substitute the values of the b 's in terms of m , we shall obtain the expression of the A 's in terms of m . Thus we have

$$\begin{aligned} A_0 &= 1, \\ A_2 &= -\frac{1}{12} \left(1 - \frac{1}{m^2}\right), \\ (6) \quad A_4 &= \frac{1}{240} \left(7 - \frac{10}{m^2} + \frac{3}{m^4}\right), \\ A_6 &= -\frac{1}{1344} \left(31 - \frac{49}{m^2} + \frac{21}{m^4} - \frac{3}{m^6}\right), \\ A_8 &= \frac{1}{11520} \left(381 - \frac{620}{m^2} + \frac{294}{m^4} - \frac{60}{m^6} + \frac{5}{m^8}\right). \end{aligned}$$

A comparison of the values of A_0 with b_0 , of A_2 with b_2 , of A_4 with b_4 , of A_6 with b_6 , and of A_8 with b_8 shows a remarkable similarity between the coefficients in A_{2i} and those in b_{2i} ; in fact

we observe that $A_{2i} = \frac{1}{m^{2i}} b_{2i} \left(\frac{1}{m}\right)$. Substituting $\frac{1}{m}$ for m in (3) and dividing by m^{2i} , we obtain

$$(7) \quad A_{2k}(m) = \frac{1}{4^k(2k+1)} \sum_{i=0}^k \binom{2k+1}{2i+1} \frac{D_{2k-2i}}{m^{2i}}$$

In order to prove that (7) is true in general, we assume it true up to a certain point and prove it true for the next highest value of k .

That is we assume

$$A_{2i} = \frac{1}{m^{2i}} b_{2i} \left(\frac{1}{m}\right) \text{ for } i = 0, 1, \dots, k-1$$

and substitute in

$$A_{2k} = - \sum_{i=0}^{k-1} \binom{2k}{2i} b_{2k-2i} A_{2i},$$

another form of (5) since $b_{2k-2i-1} = 0$. From (3) we have

$$b_{2k-2i} = \frac{1}{4^{k-i}(2k-2i+1)} \sum_{j=0}^{k-i} \binom{2k-2i+1}{2j} \frac{D_{2j}}{m^{2j}},$$

and

$$A_{2i} = \frac{1}{4^i(2i+1)} \sum_{r=0}^i \binom{2i+1}{2r} \frac{D_{2r}}{m^{2i-2r}}$$

After this substitution we arrange the terms of A_{2k} according to powers of $\frac{1}{m}$ obtaining

$$A_{2k} = \frac{1}{4^k(2k+1)} \sum_{s=0}^k \binom{2k+1}{2s+1} \frac{1}{m^{2s}} \left\{ \sum_j \binom{2s+1}{2j} D_{2j} \right\} \\ \cdot \left\{ - \frac{1}{(2k-2s+1)} \sum_r \binom{2k-2s+1}{2r} D_{2r} \right\},$$

where $s = i - r + j$. When $s = 0, 1, \dots, k-1$ the summation extends from $r=0$ to $r=k-s$ for $j \neq 0$, but from $r=0$ to

$r = k - s - 1$ for $j = 0$. Since

$$\sum_{r=0}^{k-s} \binom{2k-2s+1}{2r} D_{2r} = 0,$$

for $s < k$, the summation as to r gives zero for $j \neq 0$ but

$$-(2k-2s+1)D_{2k-2s}, \text{ for } j = 0.$$

At the same time for $j = 0$, the factor $\binom{2s+1}{2j} D_{2j}$ equals unity and we have the desired terms for $s = 0, 1, \dots, k-1$.

For $s = k$ we have constantly $r = 0$, the summation as to j being from $j = 1$ to $j = k$; we therefore have

$$\sum_{j=1}^s \binom{2s+1}{2j} D_{2j} = -D_0,$$

at the same time that

$$-\frac{1}{(2k-2s+1)} \sum_{r=0}^0 \binom{2k-2s+1}{2r} D_{2r} = -1.$$

We therefore come again to formula (7) with i replaced by s . Hence formula (7) is true for every positive integral value of k . The first few values of the A 's have been calculated in (6), others may be easily obtained by substituting the value of the Eulerian numbers from some table of D_{2i} .⁸

Formulae (4) and (7) give us the expression of the true moment

$$\mu'_n = \sum_{j=-\infty}^{+\infty} -\left(\frac{j}{m}\right)^n f\left(\frac{j}{m}\right),$$

of a discreet distribution in terms of the set of average grouped moments

⁸ Nörlund, loc. cit., Tafel 4, p. 458, gives the value up to D_{20}

$$\bar{v}_i' = \sum_{j=-\infty}^{+\infty} \frac{1}{m} \left(\frac{j}{m} + \frac{m-1}{2m}\right)^i \phi\left(\frac{j}{m} + \frac{m-1}{2m}\right).$$

Employing the particular values given in (6), we have the formula

$$\begin{aligned} \mu_n' = & \bar{v}_n' - \frac{1}{12} \binom{n}{2} \left(1 - \frac{1}{m^2}\right) \bar{v}_{n-2}' + \frac{1}{240} \binom{n}{4} \left(7 - \frac{10}{m^2} + \frac{3}{m^4}\right) \bar{v}_{n-4}' \\ (8) \quad & - \frac{1}{1344} \binom{n}{6} \left(31 - \frac{49}{m^2} + \frac{21}{m^4} - \frac{3}{m^6}\right) \bar{v}_{n-6}' \\ & + \frac{1}{11520} \binom{n}{8} \left(381 - \frac{620}{m^2} + \frac{294}{m^4} - \frac{60}{m^6} + \frac{5}{m^8}\right) \bar{v}_{n-8}' \end{aligned}$$

For any particular value of n , this series terminates and we may therefore apply the ordinary theory of limits to (8). Thus we obtain

$$\begin{aligned} \mu_n' = & \bar{v}_n' - \frac{1}{12} \binom{n}{2} \bar{v}_{n-2}' + \frac{7}{240} \binom{n}{4} \bar{v}_{n-4}' \\ (9) \quad & - \frac{31}{1344} \binom{n}{6} \bar{v}_{n-6}' + \frac{127}{3840} \binom{n}{8} \bar{v}_{n-8}' - \dots, \end{aligned}$$

the expression of the true moment $\mu_n' = \int_{-\infty}^{+\infty} x^n f(x) dx$ of a continuous distribution in terms of the set of average grouped moments

$$\bar{v}_i' = \int_{-\infty}^{+\infty} \left(x + \frac{1}{2}\right)^i \phi\left(x + \frac{1}{2}\right) dx = \int_{-\infty}^{+\infty} x^i \phi(x) dx.$$

We have thus completely solved the first of our two problems; we have obtained the expression of the true moments in terms of the average grouped moments without any assumption other than the existence of μ_n' . The existence of μ_n' requires the convergence

of the summation or integration as the lower limit approaches $-\infty$ and as the upper limit approaches $+\infty$ independently.

If in (8) we replace

$$\bar{v}_i' = \sum_{j=-\infty}^{+\infty} \frac{1}{m} \left(\frac{j}{m} + \frac{m-1}{2m} \right)^i \phi \left(\frac{j}{m} + \frac{m-1}{2m} \right)$$

by

$$v_i' = \sum_{j=-\infty}^{+\infty} (x_j)^i \phi(x_j),$$

we will have the general Sheppard-Carver formula. Since there is no approximation involved in (8), any error in the Sheppard-Carver formulae must be a result of the error involved in replacing the average grouped moments \bar{v}_i' by the raw moments v_i' . By definition \bar{v}_i' is the average of the

$$\bar{v}_i'(t) = \sum_{j=-\infty}^{+\infty} (j+t)^i \phi(j+t),$$

and, therefore, if we take any particular grouping at random,

$$v_i' = \sum_{j=-\infty}^{+\infty} (x_j)^i \phi(x_j),$$

is the mean of a random sample of one from the parent distribution $\bar{v}_i'(t)$ and hence the most probable value of \bar{v}_i' . The Sheppard-Carver formula, therefore, gives the most probable value of the true moment μ_n' of a discrete distribution in the sense that these formulae eliminate the systematic errors due to grouping.

Similarly, we shall obtain Sheppard's corrections if in (9) we

replace

$$\bar{v}_i' = \int_{-\infty}^{+\infty} x^i \phi(x) dx,$$

by

$$v_i' = \sum_{j=-\infty}^{+\infty} (x_j)^i \phi(x_j).$$

These formulae give the most probable value of the true moments μ_n^i for a continuous distribution in the same sense as do the Sheppard-Carver formulae for a discrete distribution.

The Sheppard corrections for continuous distributions and the Sheppard-Carver corrections for discrete distributions give the most probable value of the true moments $\{\mu_n^i\}$ of a distribution $f(x)$ in the sense that they give an approximate value for μ_n^i which is correct *on the average*. That is these formulae eliminate the systematic errors due to grouping whatever the distributional function $f(x)$ so long as the moments under consideration exist. While it is true that the accidental errors not accounted for in these corrections may not be negligible, these formulae do give the most probable value of μ_n^i for a particular grouping and hence have a basis for universal application.

John R. Abernethy