

# NOTES

*This section is devoted to brief research and expository articles, and notes on methodology.*

---

## A NOTE ON THE BEST LINEAR ESTIMATE

BY ALLEN T. CRAIG

*University of Iowa*

**1. Introduction.** Let the chance variable  $x$  be subject to the distribution function  $D(x)$  and as usual let  $E[g(x)]$  denote the mathematical expectation of the function  $g(x)$ . If  $x_1, x_2, \dots, x_n$  constitute a sample of  $n$  independent values of  $x$ , the function  $y = c_1x_1 + c_2x_2 + \dots + c_nx_n$  is frequently called the best linear estimate of  $E(x)$  when the  $c$ 's are so chosen that  $E(y) = E(x)$ , and  $E[y - E(x)]^2 = \sigma_y^2$  is a minimum. It is the purpose of this note to give an example of an estimate  $y$ , best in the sense defined, yet such that,  $y'$  being another estimate,

$$Pr[E(x) - \delta \leq y \leq E(x) + \delta] \leq Pr[E(x) - \delta \leq y' \leq E(x) + \delta],$$

for every  $\delta > 0$ .

**2. The rectangular distribution.** Consider  $D(x) = 1/a, 0 \leq x \leq a$ , and let the  $n$  items of each sample be arranged in ascending order of magnitude so that  $x_1 \leq x_2 \leq \dots \leq x_n, n \geq 2$ . The generating function  $G(t)$  of the moments of the distribution of  $y = c_1x_1 + c_2x_2 + \dots + c_nx_n$  is

$$G(t) = E(e^{ty}) = \frac{n!}{a^n} \int_0^a \int_0^{x_n} \int_0^{x_{n-1}} \dots \int_0^{x_2} e^{t(c_1x_1 + \dots + c_nx_n)} dx_1 dx_2 \dots dx_n.$$

Thus

$$E(y) = G'(0) = \frac{a}{n+1} [c_1 + 2c_2 + 3c_3 + \dots + nc_n],$$

and

$$\begin{aligned} E(y^2) = G''(0) = & \frac{a^2}{(n+1)(n+2)} [1 \cdot 2c_1^2 + 2 \cdot 3c_2^2 + \dots + n(n+1)c_n^2 \\ & + 2\{1 \cdot 3c_1c_2 + 1 \cdot 4c_1c_3 + \dots + 1 \cdot (n+1)c_1c_n \\ & + 2 \cdot 4c_2c_3 + \dots + 2(n+1)c_2c_n \\ & \quad \vdots \\ & + (n-1)(n+1)c_{n-1}c_n\}]. \end{aligned}$$

From  $E(y) = E(x) = a/2$ , we have

$$c_1 = \frac{1}{2}(n+1) - 2c_2 - \dots - nc_n.$$

Thus  $\sigma_y^2 = G''(0) - a^2/4$  with  $c_1$  in  $G''(0)$  replaced by  $\frac{1}{2}(n+1) - 2c_2 - \dots - nc_n$ . From  $\frac{\partial \sigma_y^2}{\partial c_j} = 0, j = 2, 3, \dots, n$ , we obtain the following system of  $n - 1$  non-homogeneous linear equations in  $n - 1$  unknowns:

$$\begin{aligned} 4c_2 + 6c_3 + \dots + 2nc_n &= n + 1 \\ 6c_2 + 12c_3 + \dots + 4nc_n &= 2(n + 1) \\ 8c_2 + 16c_3 + \dots + 6nc_n &= 3(n + 1) \\ \vdots &\vdots \\ 2nc_2 + 4nc_3 + \dots + 2n(n - 1)c_n &= (n - 1)(n + 1). \end{aligned}$$

Since the determinant of the coefficients is not zero, the solution  $c_2 = c_3 = \dots = c_{n-1} = 0, c_n = (n + 1)/2n$ , is unique. Further, we see that  $c_1 = 0$  so the best linear estimate of the mean of the rectangular population is  $y = (n + 1)x_n/2n$ , where  $x_n$  is the largest item in the sample.

The distribution function of  $y$  is readily found to be

$$D(y) = n \left[ \frac{2n}{a(n + 1)} \right]^n y^{n-1}, \quad 0 \leq y \leq \frac{n + 1}{2n} a.$$

From this, it follows that  $\sigma_y^2 = \frac{a^2}{4n(n + 2)}$ .

It has long been known<sup>1</sup> that the sampling distribution of the statistic  $\omega = \frac{1}{2}(x_1 + x_n)$ , where  $x_1$  and  $x_n$  are respectively the smallest and largest items in samples of size  $n$  from a rectangular population, has a smaller variance than does that of the arithmetic mean  $\bar{x}$  of all  $n$  items. The distribution function of  $\omega$  is

$$\begin{aligned} D(\omega) &= \frac{2^{n-1} n \omega^{n-1}}{a^n}, & 0 \leq \omega \leq \frac{1}{2}a, \\ &= \frac{2^{n-1} n}{a^n} (a - \omega)^{n-1}, & \frac{1}{2}a \leq \omega < a, \end{aligned}$$

so that  $E(\omega) = \frac{1}{2}a$  and  $\sigma_\omega^2 = \frac{a^2}{2(n + 1)(n + 2)}$ . Thus  $\sigma_y^2 = \frac{1}{2}\sigma_\omega^2$ , approximately.

Yet Pittman has recently proved that for every  $\delta > 0, Pr[E(x) - \delta \leq \omega \leq E(x) + \delta]$  exceeds the probability that any other estimate, including  $y$ , will fall in this interval of length  $2\delta$  about the mean  $a/2$ .

If we write  $u = \frac{y - a/2}{\sigma_y}$  and  $v = \frac{\omega - a/2}{\sigma_\omega}$ , then the limits of  $D(u)$  and  $D(v)$  as  $n$  approaches infinity are respectively  $e^{u-1}, -\infty \leq u \leq 1$ , and  $\frac{1}{\sqrt{2}} e^{-\sqrt{2}|v|}$ ,

---

<sup>1</sup> R. A. Fisher, "Theoretical foundations of mathematical statistics," *Phil. Trans. Roy. Soc. London*, Series A, Vol. 222 (1921), pp. 309-368.

$-\infty \leq v \leq \infty$ . Thus neither  $y$  nor  $\omega$  has an asymptotic normal distribution. It is, of course, this fact which makes the criterion of minimum variance illusory.

**3. Other polynomial distribution functions.** Let repeated samples of  $n$  independent values of  $x$  be drawn from a population characterized by  $D(x) = \frac{k+1}{a^{k+1}} x^k$ ,  $0 \leq x \leq a$ , and  $k$  a positive integer or zero. It can be shown that the best linear estimate of the mean of the population is  $y = \frac{(k+1)n+1}{n(k+2)} x_n$ , where as before  $x_n$  is the largest item of the sample. The sampling distribution of  $y$  is easily obtained. It follows that

$$\sigma_y^2 = \frac{(k+1)a^2}{(k+2)^2[(k+1)n^2+2n]} = \frac{k+3}{n(k+1)+2} \sigma_x^2,$$

where as usual  $\bar{x}$  is the arithmetic mean of the sample. Again, if we write  $u = \left(y - \frac{k+1}{k+2} a\right) / \sigma_y$ , the limit of the distribution of  $u$  as  $n$  approaches infinity is, as before,  $e^{u-1}$ ,  $-\infty \leq u \leq 1$ .

---

### A NOTE ON TOLERANCE LIMITS

BY EDWARD PAULSON<sup>1</sup>

*Columbia University*

Among various statistical problems arising in the process of controlling quality in mass production, a rather important one appears to be the determination of tolerance limits when the variability of the product is known to be due to random factors. This problem was recently treated in a pioneer article by Wilks. This note will point out a relationship between tolerance limits and confidence limits (used in the sense of Neyman), and will use this concept to establish tolerance limits when the product is described by two qualities, the measurements on which are assumed to have a bivariate normal distribution.

For the case of a single variate, the problem of finding tolerance limits as stated by Wilks is to find a sample size  $n$ , and two functions  $L_1(x_1 \cdots x_n)$  and  $L_2(x_1 x_2 \cdots x_n)$  so that if  $P = \int_{L_1}^{L_2} f(x) dx$  denotes the conditional probability of a future observation falling between the random variates  $L_2$  and  $L_1$ , then

$$E(P) = \alpha, \quad \text{and Prob. } [\alpha - \Delta_1 \leq P \leq \alpha + \Delta_2] \geq \beta.$$

The relationship between confidence limits and tolerance limits will arise if confidence limits are determined, not for a parameter of the distribution, but for

---

<sup>1</sup> Work done under a grant-in-aid from the Carnegie Corporation of New York.