

ON THE EXACT EVALUATION OF THE VARIANCES AND COVARIANCES OF ORDER STATISTICS IN SAMPLES FROM THE EXTREME-VALUE DISTRIBUTION¹

BY JULIUS LIEBLEIN

National Bureau of Standards

Summary. This paper develops explicit closed formulas for the covariances of order statistics in samples from the extreme-value distribution which involve only tabulated functions. Such results do not appear to have been given previously. They have been used in an investigation of the estimation of extreme-value parameters by means of order statistics which will be presented in a fuller report to be submitted to the National Advisory Committee for Aeronautics.

1. Problem. We are concerned with random samples of size n from the "extreme-value" distribution whose cdf is

$$(1.1) \quad F(x) = \exp(-e^{-y}), \quad y = \frac{x - u}{\beta}, \quad -\infty < x < \infty.$$

(This distribution was derived as a limiting form of the distribution of the largest value in a sample by Fisher and Tippett [1] and has been extensively studied by Gumbel (e.g. [4], [5]). However, this paper is not concerned with the extremal properties of this distribution.) If the n values after ordering in size are denoted by

$$x_1, x_2, \dots, x_n, \quad x_1 \leq x_2 \leq \dots \leq x_n,$$

then we seek the second-order moments of the x_i, x_j , namely, the variances σ_i^2 and covariances σ_{ij} . The first moments have been tabulated [6] for samples of $n \leq 100$.

The second moments involve integrals which at first sight look more difficult than the corresponding ones for the normal distribution, which latter have required a very extensive amount of numerical integration. In this paper a method is shown for evaluating the extreme-value integrals in closed form ((3.10) below) involving only tabulated functions. Thus, the extreme-value distribution is brought into the select circle, which previously included only the normal (at least for $n \leq 6$ —see [3]), exponential, and rectangular distributions, and perhaps some others, for which the second moments of the order statistics can be evaluated explicitly without quadratures.

2. Theory. The density function of the i th order statistic, x_i , from the distribu-

Received 12/4/52.

¹This paper is based on research sponsored by the National Advisory Committee for Aeronautics.

tion (1.1) is

$$(2.1) \quad p(x) = \frac{n!}{(i-1)!(n-i)!} [F(x)]^{i-1} [1 - F(x)]^{n-i} f(x), \quad -\infty < x < \infty,$$

where $x = x_i$, $f(x) = F'(x)$. The joint d.f. of the i th and j th order statistics x_i, x_j , is

$$(2.2) \quad p(x, y) = \frac{n!}{(i-1)!(j-i-1)!(n-j)!} [F(x)]^{i-1} [F(y) - F(x)]^{j-i-1} \cdot [1 - F(y)]^{n-j} f(x) f(y), \quad -\infty < x \leq y < \infty,$$

where $x = x_i, y = x_j, i < j, i, j = 1, 2, \dots, n$. Without loss of generality, we shall henceforth refer only to the standardized or "reduced" extreme-value distribution, with the parameters $\beta = 1, u = 0$,

$$(2.3) \quad P(y) = \exp(-e^{-y}), \quad -\infty < y < \infty,$$

and denote its variable by y .

From the density functions (2.1) and (2.2) we obtain

$$(2.4) \quad \begin{aligned} E(y_i^k) &= \frac{n!}{(i-1)!(n-i)!} \int_{-\infty}^{\infty} x^k e^{-(i-1)e^{-x}} (1 - e^{-e^{-x}})^{n-i} e^{-x-e^{-x}} dx \\ &= \frac{n!}{(i-1)!(n-1)!} \sum_{r=0}^{n-i} (-1)^r C_r^{n-i} \int_{-\infty}^{\infty} x^k e^{-x-(i+r)e^{-x}} dx, \end{aligned}$$

$$(2.5) \quad \begin{aligned} E(y_i y_j) &= \frac{n!}{(i-1)!(j-i-1)!(n-j)!} \\ &\cdot \int_{-\infty}^{\infty} \int_{-\infty}^y x y e^{-y-e^{-y}} e^{-x-ie^{-x}} (e^{-e^{-y}} - e^{-e^{-x}})^{j-i-1} (1 - e^{-e^{-y}})^{n-j} dx dy \\ &= \frac{n!}{(i-1)!(j-i-1)!(n-j)!} \\ &\cdot \sum_{r=0}^{j-i-1} \sum_{s=0}^{n-j} (-1)^{r+s} C_r^{j-i-1} C_s^{n-j} \phi(i+r, j-i-r+s), \end{aligned}$$

where the function ϕ is the double integral

$$(2.6) \quad \phi(t, u) = \int_{-\infty}^{\infty} \int_{-\infty}^y x y e^{-x-te^{-x}} e^{-y-ue^{-y}} dx dy, \quad t, u > 0$$

whose evaluation is the main point of this paper.

3. Evaluation of the integrals.

3.1 *Variance-type integrals.* These integrals are of the general form

$$g_k(c) = \int_{-\infty}^{\infty} x^k e^{-x-ce^{-x}} dx, \quad c > 0.$$

The evaluation given here is not new, but is presented for completeness.

The change of variable $e^{-x} = v$ gives

$$g_k(c) = \int_0^{\infty} (-\log v)^k e^{-cv} dv,$$

which for k a nonnegative integer¹

$$\begin{aligned} (3.1) \quad &= (-1)^k \frac{d^k}{dt^k} \int_0^{\infty} v^{t-1} e^{-cv} dv \Big|_{t=1} \\ &= (-1)^k \frac{d^k}{dt^k} [\Gamma(t)c^{-t}] \Big|_{t=1} \end{aligned}$$

The needed first two values are

$$(3.2) \quad g_1(c) = -\left[\frac{\Gamma'(1)}{c} - \frac{\Gamma(1)}{c} \log c \right] = \frac{1}{c} (\gamma + \log c),$$

where $\gamma = -\Gamma'(1)$ is Euler's constant, .5772156649 \dots . Likewise,

$$(3.3) \quad g_2(c) = \frac{1}{c} \left[\frac{\pi^2}{6} + (\gamma + \log c)^2 \right].$$

3.2 Covariance integrals. An integration by parts applied to the inner integral in (2.6) with "dv" equal to the exponential factor gives

$$\int_{-\infty}^y x e^{-x-tx} dx = t^{-1} y e^{-te-y} - t^{-1} \int_{-\infty}^y e^{-te-x} dx.$$

Hence from (2.6) and (3.1),

$$(3.4) \quad t\phi(t, u) = g_2(t+u) - \psi(t, u),$$

where

$$(3.5) \quad \psi(t, u) = \int_{-\infty}^{\infty} y e^{-y-ue-y} \left[\int_{-\infty}^y e^{-te-x} dx \right] dy.$$

The function ψ regarded as a simple integral containing a parameter ($t > 0$) may be differentiated under the integral sign, giving, by (3.2),

$$(3.6) \quad \frac{\partial \psi}{\partial t} = \frac{1}{t} g_1(t+u) = -\frac{1}{t(t+u)} [\gamma + \log(t+u)], \quad t, u > 0.$$

Before integrating this equation, it is convenient to make the change of variable $w = 1 + (t/u)$. After the substitution integrate (3.6) with respect to w from $w = 2$ to $w = w$, and replace the upper limit w in the resulting expression by its value in terms of t , noting that the corresponding limits for t are $t = u$ to $t = t$. The result is

$$\begin{aligned} (3.7) \quad &u[\psi(t, u) - \psi(u, u)] = \gamma \log(1 + u/t) + \frac{1}{2} [\log(t+u)]^2 \\ &- \log u \log t/u - \gamma \log 2 - \frac{1}{2} (\log 2u)^2 - \int_2^{1+t/u} \frac{\log w}{w-1} dw. \end{aligned}$$

The integral on the right is immediately expressible in terms of *Spence's integral* (or function)

$$(3.8) \quad L(1+x) = \int_1^{1+x} \frac{\log w}{w-1} dw = \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n^2} x^n.$$

Several tables of this function are cited in [2]. The most extensive of these is given by F. W. Newman [7] to twelve decimal places.

It remains only to evaluate $\psi(u, u)$. From (2.6) and (3.2),

$$\begin{aligned} \phi(u, u) &= \int_{-\infty}^{\infty} ye^{-y-ue^{-y}} \left(\int_{-\infty}^y xe^{-x-ue^{-x}} dx \right) dy \\ &= \int_{-\infty}^{\infty} \frac{1}{2} \frac{d}{dy} \left(\int_{-\infty}^y xe^{-x-ue^{-x}} dx \right)^2 dy \\ &= \frac{1}{2u^2} (\gamma + \log u)^2. \end{aligned}$$

This value when substituted in (3.4) gives $\psi(u, u)$. Combining this result with (3.7) and the easily obtainable value $L(2) = \pi^2/12$ gives, after a little algebra, the following formula:

$$(3.9) \quad \begin{aligned} 2tu \phi(t, u) &= (u-t)g_2(t+u) + t^2[g_1(t)]^2 + 2L\left(1 + \frac{t}{u}\right) \\ &\quad - \left(\log \frac{t}{u}\right)^2 - \frac{\pi^2}{6}, \end{aligned}$$

where the functions $g_1(t)$, $g_2(t)$ are given by (3.2), (3.3). This may be simplified a little by use of the following property of Spence's function:

$$L(1+x) + L\left(1 + \frac{1}{x}\right) = \frac{1}{2}(\log x)^2 + \frac{\pi^2}{6},$$

giving the result

$$(3.10) \quad 2tu \phi(t, u) = (u-t)g_2(t+u) + t^2[g_1(t)]^2 - 2L\left(1 + \frac{u}{t}\right) + \frac{\pi^2}{6}.$$

The above results (3.9), (3.10), together with (2.4), (2.5), make possible the evaluation of all the variances and covariances. This requires the calculation of n values of g_1 and of g_2 , and $\frac{1}{2}n(n-1)$ values of ϕ .

The calculation may be simplified with the aid of the relation

$$(3.11) \quad \phi(t, u) + \phi(u, t) = g_1(t)g_1(u),$$

which may be derived from (2.6) and (3.2) by means of a change in the order of integration. Thus (3.10) need be used only for $t \leq u$, so that (3.11) reduces the number of values of ϕ by almost half, unless n is small, say $n < 10$.

4. Illustration. The above formulas have been used by the author in an investigation of estimation of extreme-value parameters by means of order statistics.

The results of this research, including a table of the first two moments for small samples, will be reported elsewhere.

The following computations for $n = 3$ illustrate the procedure described in this article.

From (3.2),

$$\begin{aligned}g_1(1) &= \gamma = 0.57721\ 57 \\g_1(2) &= 0.63518\ 14 \\g_1(3) &= 0.55860\ 93.\end{aligned}$$

The means are then given by (2.4) and (3.2) as

$$\begin{aligned}E(y_1) &= 3[g_1(1) - 2g_1(2) + g_1(3)] = -0.40361\ 4 \\E(y_2) &= 6[g_1(2) - g_1(3)] = +0.45943\ 3 \\E(y_3) &= 3g_1(3) = +1.67582\ 8.\end{aligned}$$

As a simple check, these three values sum to 3γ to within six decimal places, and also agree with those in [6]. (The notation in the table cited differs from that used here: $E(y_i)$ in this paper corresponds to $E(y_{n-i})$ in the table.) Next, from (3.3), we have

$$\begin{aligned}g_2(1) &= \frac{\pi^2}{6} + \gamma^2 = 1.97811\ 2 \\g_2(2) &= 1.62937\ 8 \\g_2(3) &= 1.48444\ 4.\end{aligned}$$

The mixed function $\phi(t, u)$ is then given by (3.9) and (3.10):

$$\begin{aligned}\phi(1, 1) &= \frac{1}{2}(\gamma^2) = 0.16658\ 9 \\ \phi(1, 2) &= \frac{1}{4}[g_2(3) + \gamma^2 + 2L(1\frac{1}{2}) - (1n2)^2 - \pi^2/6] = 0.14726\ 6 \\ \phi(2, 1) &= \frac{1}{4}\{-g_2(3) + 4[g_1(2)]^2 - 2L(1\frac{1}{2}) + \pi^2/6\} = 0.21937\ 1.\end{aligned}$$

Newman's table [7] provides the value of the function $L(1\frac{1}{2}) = 0.44841\ 42069$. Finally, equations (2.4) and (2.5) give, for the moments about the origin,

$$\| E(y_i y_j) \| = \begin{bmatrix} 0.61140 & 0.11594 & -0.43263 \\ 0.11594 & 0.86960 & 1.31622 \\ -0.43263 & 1.31622 & 4.45333 \end{bmatrix}$$

whence the moments about the mean are given by

$$\| \sigma(y_i y_j) \| = \begin{bmatrix} 0.44850 & 0.30137 & 0.24376 \\ 0.30137 & 0.65852 & 0.54629 \\ 0.24376 & 0.54629 & 1.64493 \end{bmatrix}.$$

The final results are correct to about four decimal places. (One additional place is shown for checking purposes.)

As a check, we should have

$$\sum_{j=1}^3 \sum_{i=1}^3 \sigma(y_i y_j) = \sigma^2 \left(\sum_{i=1}^3 y_i \right) = 9\sigma^2(\bar{y}) = 3\sigma_y^2 = \frac{\pi^2}{2},$$

since σ_y^2 , the variance of the distribution $P(y)$ in (2.3), is known to be $\pi^2/6$. The left side of this equation is found to be 4.93479; the right side, 4.93480. This type of check cannot be considered to be very effective, however, as only gross errors, and not compensating ones, will ordinarily be revealed.

The reader should be cautioned that, unless n is fairly small, it may be necessary to carry out the calculations to a considerably greater number of places than is desired in the results. This results from the presence of binomial coefficients and alternating signs in formulas (2.4) and (2.5), both of which operate to reduce accuracy rapidly as n increases.

REFERENCES

- [1] R. A. FISHER AND L. H. C. TIPPETT, "Limiting forms of the frequency distribution of the largest or smallest member of a sample," *Proc. Cambridge Philos. Soc.*, Vol. 24 (1928), pp. 180-190.
- [2] A. FLETCHER, J. C. P. MILLER, AND L. ROSENHEAD, *An Index of Mathematical Tables*, McGraw-Hill Book Company, Inc., 1946, pp. 343-344.
- [3] H. J. GODWIN, "Some low moments of order statistics," *Ann. Math. Stat.*, Vol. 20 (1949), pp. 279-285.
- [4] E. J. GUMBEL, "Les valeurs extrêmes des distributions statistiques," *Ann. Inst. H. Poincaré*, Vol. 4 (1935), pp. 115-158.
- [5] E. J. GUMBEL, "The return period of flood flows," *Ann. Math. Stat.*, Vol. 12 (1941), pp. 163-190.
- [6] *Table of the First Moment of Ranked Extremes*, National Bureau of Standards Report 1167, September 20, 1951, (special report submitted to the National Advisory Committee for Aeronautics.)
- [7] F. W. NEWMAN, *The Higher Trigonometry, Superrationals of Second Order*, Macmillan and Bowes, Cambridge University Press, (1892), pp. 64-65.