

## THE DISTRIBUTION OF GALTON'S STATISTIC<sup>1</sup>

BY SHULAMITH GROSS AND PAUL W. HOLLAND

*Harvard University*

**0. Summary.** Let  $X_{(1)} < \dots < X_{(n)}$  and  $Y_{(1)} < \dots < Y_{(n)}$  be the order statistics of two independent random samples from the absolutely continuous distribution functions  $F(x)$  and  $G(y)$ , respectively. Let  $T_n$  be the proportion of pairs,  $(X_{(i)}, Y_{(i)})$ , for which  $X_{(i)} \geq Y_{(i)}$ . Tests of the equality of  $F$  and  $G$  based on  $T_n$  are among the oldest nonparametric procedures in the literature, going back at least to Galton's analysis of Darwin's data [3]. Hodges [5] showed the null distribution of  $nT_n$  to be uniform over  $0, 1, \dots, n$ . Bickel and Hodges [1] treated the asymptotic distribution of the Lehmann estimate based on the one-sample version of  $T_n$ . In this note we use very elementary methods to derive expressions for the distribution and moments of  $T_n$  from which conditions for the consistency of tests based on  $T_n$  follow immediately. More generally we can show that (unnormalized)  $T_n$  always has an asymptotic distribution for any pair  $(F, G)$ . This distribution is degenerate at zero if  $Y$  happens to be stochastically larger than  $X$ . We give informative expressions for the first two moments of this asymptotic distribution. Our technique is to express the distribution of  $T_n$  in terms of integrals of certain multinomial probabilities.

**1. Results.** Define  $I(x, y) = 1$  if  $x \geq y$  and zero otherwise. Then  $T_n = n^{-1} \sum_{k=1}^n I(X_{(k)}, Y_{(k)})$ . Let  $f$  and  $g$  be the densities of  $F$  and  $G$ . We assume  $F^{-1}$  exists and define  $h(u) = G(F^{-1}(u))$  for  $0 \leq u \leq 1$ . Then  $h(u)$  is increasing on  $[0, 1]$ . For any pair  $1 \leq j \leq n$  and any  $0 \leq u_1 < \dots < u_n \leq 1$  let  $(M_1, \dots, M_{j+1})$  and  $(N_1, \dots, N_{j+1})$  be independent multinomial random vectors with parameters  $(n - j; u_1, u_2 - u_1, \dots, u_j - u_{j-1}, 1 - u_j)$  and  $(n; h(u_1), h(u_2) - h(u_1), \dots, h(u_j) - h(u_{j-1}), 1 - h(u_j))$  respectively. Finally we define

$$(1) \quad p_{n,j}(u_1, \dots, u_j) = P\{\sum_{i=1}^k (N_i - M_i) \geq k; k = 1, \dots, j\}$$

for  $1 \leq j \leq n$ .

LEMMA 1. For every  $1 \leq m \leq n$ ,

$$P\{nT = m\} = \sum_{j=m}^n (-1)^{j-m} \binom{j}{m} n! ((n-j)!)^{-1} \int \dots \int_{0 \leq u_1 < \dots < u_j \leq 1} p_{n,j}(u_1, \dots, u_j) du_1, \dots, du_j.$$

PROOF. For every  $1 \leq m \leq n$  we have

$$(2) \quad P\{nT_n = m\} = \sum_{j=m}^n (-1)^{j-m} \binom{j}{m} \cdot \sum_{1 \leq i_1 < \dots < i_j \leq n} P\{X_{(i_1)} \geq Y_{(i_1)}, \dots, X_{(i_j)} \geq Y_{(i_j)}\}.$$

Received 6 October 1967.

<sup>1</sup> This work was facilitated by a grant from the National Science Foundation (GS-341).

Now for  $1 \leq i_1 < i_2 < \dots < i_j \leq n$

$$(3) \quad P\{X_{(i_1)} \geq Y_{(i_1)}, \dots, X_{(i_j)} \geq Y_{(i_j)}\} \\ = \int \dots \int_{-\infty < y_1 < \dots < y_j < \infty} P\{Y_{(i_1)} \leq y_1, \dots, Y_{(i_j)} \leq y_j\} \\ \cdot f_{i_1 \dots i_j}(y_1 \dots y_j) dy_1 \dots dy_j,$$

where  $f_{i_1 \dots i_j}(y_1, \dots, y_j)$  is the joint density of  $X_{(i_1)}, \dots, X_{(i_j)}$  and is given by

$$(4) \quad = n! f(y_1) \dots f(y_j) [(i_1 - 1)! (i_2 - i_1 - 1)! \dots (n - i_j)!]^{-1} \\ \cdot [F(y_1)]^{i_1 - 1} [F(y_2) - F(y_1)]^{i_2 - i_1 - 1} \dots [1 - F(y_j)]^{n - i_j}$$

for  $y_1 < y_2 < \dots < y_j$ , and zero otherwise. Another easy calculation shows that for  $1 \leq i_1 < i_2 < \dots < i_j \leq n$  and  $y_1 < y_2 < \dots < y_n$

$$(5) \quad P\{Y_{(i_1)} \leq y_1, \dots, Y_{(i_j)} \leq y_j\} \\ = \sum_{l_j=i_j}^n \dots \sum_{l_2=i_2}^{l_3} \sum_{l_1=i_1}^{l_2} (l_1, l_2 - l_1, \dots, n - l_j) \\ \cdot [G(y_1)]^{l_1} [G(y_2) - G(y_1)]^{l_2 - l_1} \\ \dots [G(y_j) - G(y_{j-1})]^{l_j - l_{j-1}} [1 - G(y_j)]^{n - l_j}.$$

Now putting (2), (3), (4), and (5) together and making the transformation  $u_i = F(y_i)$ , one obtains

$$(6) \quad P\{nT_n = m\} = \sum_{j=m}^n (-1)^{j-m} \binom{j}{m} \\ \int \dots \int_{0 < u_1 < \dots < u_j < 1} p_{n,j}^*(u_1, \dots, u_j) du_1 \dots du_j$$

where

$$(7) \quad p_{n,j}^*(u_1, \dots, u_j) \\ = \sum_{1 \leq i_1 < \dots < i_j \leq n} \sum_{l_j=i_j}^n \\ \dots \sum_{l_1=i_1}^{l_2} n! [(i_1 - 1)! \dots (n - i_j)!]^{-1} u_1^{i_1 - 1} (u_2 - u_1)^{i_2 - i_1 - 1} \\ \dots (1 - u_1)^{n - i_j} (l_1, l_2 - l_1, \dots, n - l_j) [h(u_1)]^{l_1} [h(u_2) - h(u_1)]^{l_2 - l_1} \\ \dots [1 - h(u_j)]^{n - l_j}.$$

If, in (7) we replace  $i_k$  by  $i_k - k$ , then  $p_{n,j}^*(u_1, \dots, u_j)$  becomes

$$(8) \quad n! [(n - j)!]^{-1} \sum_{i_1=0}^{i_2} \sum_{i_2=i_1}^{i_3} \dots \sum_{i_j=i_{j-1}}^{n-j} \sum_{l_1=i_1+1}^{l_2} \sum_{l_2=i_2+2}^{l_3} \dots \sum_{l_j=i_j+j}^n \\ \cdot (i_1, i_2 - i_1, \dots, n - j - i_j) u_1^{i_1} (u_2 - u_1)^{i_2 - i_1} \dots (1 - u_j)^{n - j - i_j} \\ \cdot (l_1, l_2 - l_1, \dots, n - l_j) [h(u_1)]^{l_1} [h(u_2) - h(u_1)]^{l_2 - l_1} \dots [1 - h(u_j)]^{n - l_j}.$$

Inspection of the summation in (8) reveals it to be  $p_{n,j}(u_1, \dots, u_j)$ .  $\square$

Using Lemma 1 to compute the expectation of  $T_n^k$  we easily obtain

COROLLARY 1. *If  $k \geq 1$  then*

$$(9) \quad E(T_n^k) = \sum_{j=1}^{min\{k, n\}} \Delta^j(0^k) n! [n^k(n-j)]^{-1} \int \cdots \int_{0 \leq u_1 < \cdots < u_j \leq 1} p_{n,j}(u_1 \cdots u_j) du_1 \cdots du_j.$$

The numbers,  $\Delta^j(0^k)$  are shorthand for  $\Delta^j(X^k) |_{x=0}$  where  $\Delta$  is the difference operator. Two useful properties of these numbers needed in the proofs of Corollaries 1 and 2 are:  $\Delta^j(0^k) = 0$  if  $j > k$  and  $\Delta^k(0^k) = k!$  (see [6] page 36–50).

To show that the moments of  $T_n$  all converge we apply dominated convergence to the integrals in (9). To show  $\lim_{n \rightarrow \infty} p_{n,j}(u_1, \dots, u_j) = p_j(u_1, \dots, u_j)$  exists we first define these subsets of  $[0, 1]$ .

$$(10) \quad S_+ = \{u: h(u) > u\}, S_- = \{u: h(u) < u\}, S_0 = \{u: h(u) = u\}.$$

Observe that the vector  $n^{-\frac{1}{2}}((N_1 - M_1) - n(h(u_1) - u_1), \dots, \sum_{i=1}^j (N_i - M_i) - n(h(u_j) - u_j))$  has a limiting multivariate normal distribution with zero mean vector and covariance matrix,  $\sum (u_i, \dots, u_j)$ , whose elements are given by

$$(11) \quad \sigma_{l,m} = h(u_m)(1 - h(u_l)) + u_m(1 - u_l), \quad 1 \leq l \leq m \leq j.$$

If  $u_i \in S_0$  for all  $i = 1, \dots, j$ ,  $\sigma_{l,m}$  reduces to

$$(12) \quad \sigma_{l,m} = 2u_m(1 - u_l).$$

There are several possibilities for  $p_j(u_1, \dots, u_j)$ . (i) If  $u_i \in S_-$  for some  $i = 1, \dots, j$  then  $p_j(u_1, \dots, u_j) = 0$ . (ii) If  $u_i \in S_+$  for all  $i = 1, \dots, j$ , then  $p_j(u_1, \dots, u_j) = 1$ . (iii) If  $u_{i_1}, \dots, u_{i_l} \in S_0$  while the remaining  $u$ 's are in  $S_+$ , then  $p_j(u_1, \dots, u_j)$  is the probability content of the positive orthant in  $l$ -dimensional space given by the multivariate normal distribution with mean vector zero and covariances  $\sigma_{i_\alpha i_\beta} = 2u_{i_\alpha}(1 - u_{i_\beta})$ ;  $i_\alpha \leq i_\beta$ . This shows that  $p_j(u_1, \dots, u_j)$  exists and we obtain

COROLLARY 2. *If  $k \geq 1$  then*

$$(13) \quad \lim_{n \rightarrow \infty} E(T_n^k) = k! \int \cdots \int_{0 \leq u_1 < \cdots < u_k \leq 1} p_k(u_1, \dots, u_k) du_1 \cdots du_k.$$

Since  $T_n$  is bounded and all of its moments converge,  $T_n$  has an asymptotic distribution for any choice of  $F$  and  $G$  and the limiting moments are the moments of this limiting distribution. For  $k > 2$ , (13) does not readily simplify. For  $k = 1, 2$ , it may be applied to give interesting results. Let  $\lambda(S)$  denote the Lebesgue measure of a set  $S \subseteq [0, 1]$  and  $I_S(x)$  denote the indicator function of  $S$ . Recall the fact (see Cramér [2], page 290, for example) that the probability content of the first quadrant of a central bivariate normal distribution with correlation  $\rho$  is given by  $4^{-1} + (2\pi)^{-1} \sin^{-1}(\rho)$ . We have

COROLLARY 3.

$$(a) \quad \lim_{n \rightarrow \infty} E(T_n) = \frac{1}{2}\lambda(S_0) + \lambda(S_+),$$

$$(b) \quad \lim_{n \rightarrow \infty} \text{Var}(T_n)$$

$$= \pi^{-1} \int_0^1 \int_0^v \sin^{-1}(uv^{-1}(1-v)(1-u)^{-1})^{\frac{1}{2}} I_{S_0}(u) I_{S_0}(v) du dv.$$

Applying the last corollary we obtain

**THEOREM 1.**  $T_n$  has a degenerate distribution if and only if  $\lambda(S_0) = 0$ . If  $\lambda(S_0) = 0$ , then  $T_n$  converges in probability to  $\lambda(S_+)$ .

We observe that if  $G(x) < F(x)$  for all  $x$ , then  $\lambda(S_0) = \lambda(S_+) = 0$  so that rejection for small values of  $T_n$  yields a consistent test of the hypothesis  $F = G$ . By applying the fact that if  $F = G$ ,  $T_n$  has an asymptotic uniform distribution on  $[0, 1]$  we obtain a simple evaluation of the definite integral

$$\int_0^1 \int_0^v \sin^{-1}(uv^{-1}(1-v)(1-u)^{-1})^{\frac{1}{2}} du dv = \frac{1}{2}\pi.$$

**2. Remarks.** The assumption of absolute continuity may be reduced to simple continuity by the following device, essentially given in [4]. (a) Find a distribution function  $H$  such that  $H \gg F$  and  $H \gg G$  (use  $H = \frac{1}{2}(F + G)$  for example); (b) If  $T$  is a measurable transformation of the line into itself and  $H^*$ ,  $F^*$ , and  $G^*$  are the induced distributions then  $H^* \gg F^*$  and  $G^*$ ; (c) Use  $H$  itself as  $T$ . Then if  $F$  and  $G$  are continuous so is  $H$  and therefore  $H^*(x) \equiv x$  a.e. in  $[0, 1]$ . Hence  $F^*$  and  $G^*$  both possess densities with respect to Lebesgue measure in  $[0, 1]$ ; (d) Since  $H \gg F$  and  $G$ , is continuous and monotone (though not necessarily strictly monotone)  $T_n$  is unchanged under the transformation of the  $X$ 's and  $Y$ 's by  $H$ , with probability one.

Strict monotonicity of  $F$  was only used in the transformation of the integrals in (6). If  $F^{-1}(y)$  is defined as  $\inf \{x: F(x) \geq y\}$  for  $0 \leq y \leq 1$ , this restriction may be removed. Because  $h(u) = G(F^{-1}(u)) = G^*(F^{*-1}(u))$  it is unnecessary to actually transform the problem from  $(F, G)$  to  $(F^*, G^*)$  to use the results of this paper.

It may be possible to use the fact that  $p_{n,j}(u_1, \dots, u_j)$  is the probability of a large deviation when  $\lambda(S_-) = 1$ , and Lemma 1 to obtain more useful expressions for the distribution of  $T_n$  under the hypothesis of stochastic ordering. We have no results in this direction yet.

Finally, we would like to thank Dr. F. Mosteller for helpful discussions of this research.

#### REFERENCES

- [1] BICKEL, P. J. and HODGES, J. L., JR. (1967). The asymptotic theory of Galton's test and a related simple estimate of location. *Ann. Math. Statist.* **38** 73-89.
- [2] CRAMÉR, H. (1946). *Mathematical Methods of Statistics*. Princeton Univ. Press.
- [3] DARWIN, C. (1878). *The Effect of Cross- and Self-fertilization in the Vegetable Kingdom* (2nd edition). John Murray, London.
- [4] GOVINDARAJULU, Z., LECAM, L. and RAGHAVACHARI, M. (1967). Generalizations of theorems of Chernoff and Savage on the asymptotic normality of test statistics. *Proc. Fifth Berkeley Symp. Math. Statist. Prob.* **1** 609-638. Univ. of California Press.
- [5] HODGES, J. L., JR. (1955). Galton's rank order test. *Biometrika* **42** 261-262.
- [6] MILNE-THOMPSON, L. M. (1951). *The Calculus of Finite Differences*. Macmillan, New York.