

SHEFFER POLYNOMIALS FOR COMPUTING TAKÁCS'S GOODNESS-OF-FIT DISTRIBUTIONS¹

BY HEINRICH NIEDERHAUSEN

University of Toronto

How often does the empirical distribution pass over an acceptance band around the hypothetical distribution function? Sheffer polynomials are applied to derive the exact distribution of this goodness-of-fit distribution and its two-sample analogue. This method leads to new proofs and extensions of Takacs's (1971) results.

1. Introduction and results. Let X_1, \dots, X_M be i.i.d. random variables with unknown continuous distribution function and empirical distribution function $F_X(x) = \sum_{i=1}^M 1_{(-\infty, X_i]}(x)$. The well-known one-sided Kolmogorov-Smirnov test rejects the null hypothesis that F is the true underlying distribution function, if

$$(1) \quad K^+ = \sup_x \{F(x) - F_X(x)\}$$

is larger than a certain critical value α/M , depending on the significance probability α . In other words, the null hypothesis is rejected if $F_X + \alpha/M$ intersects F .

In 1939, N. Smirnov introduced the test statistic $\sigma(M, \alpha, 1)$ which counts the number of intersections between F and $F_X + \alpha/M$. He also derived the asymptotic distribution of $\sigma(M, \alpha, 1)$. In 1971, L. Takács generalized $\sigma(M, \alpha, 1)$ to $\sigma(M, \alpha, c)$ by counting the number of intersections between cF and $F_X + \alpha/M$. Hence,

$$\sigma(M, \alpha, c) = \sum_{i=0}^M 1_{[cF(X_{i:M}), cF(X_{i+1:M})]}((i + \alpha)/M),$$

where $F(X_{0:M}) = 0$ and $F(X_{M+1:M}) > 1$, such that the last term ($i = M$) always counts an intersection if $c = 1$ and $\alpha = 0$.

In goodness-of-fit testing, $\sigma(m, \alpha, c)$ can be used in two ways. One method means choosing a (small) positive integer s and rejecting the null hypothesis if $K_s^+ = \sup \{\alpha \mid \sigma(M, \alpha, c) = s\}$ is large. For $c = 1$, K_1^+ equals K^+ (see (1)). The other method rejects the null hypothesis if $\sigma(M, 0, 1)$ is too small. M. M. Siddiqui (1982) called this test a matching test. Under any continuous alternative G , the power Π of the matching test depends only on GF^{-1}

$$\begin{aligned} \Pi = P_G(\sigma(M, 0, 1) \leq s_\alpha) &= 1 - P_G\{F(X_{i:M}) \leq i/M < F(X_{i+1:M}) \text{ for more than } s_\alpha \text{ of the} \\ &\quad \text{subscripts } i = 0, \dots, M\} \\ &= 1 - P\{U_{i:M} \leq GF^{-1}(i/M) \leq U_{i+1:M} \text{ for more than } s_\alpha \text{ of the} \\ &\quad \text{subscripts } i = 0, \dots, M\}, \end{aligned}$$

where U_1, \dots, U_M is a random sample from $U(0, 1)$, $U_{0:M} = 0$ and $U_{M+1:M} = 1$.

In general, there is no "closed form" for Π , but if

$$GF^{-1}(i/M) = \begin{cases} 0 & \text{for } 0 \leq i < \max(0, -a), \\ (i + a)/(cM) & \text{for } \max(0, -a) \leq i \leq \min(M, cM - a), \\ 1 & \text{for } \min(M, cM - a) < i \leq M, \end{cases}$$

then $\Pi = P_F(\sigma(M, \alpha, c) \leq s_\alpha)$, the null distribution of σ . Siddiqui (1982) found that $\lim_{M \rightarrow \infty} E_G\{\sigma(M, 0, 1)/\sqrt{M}\}$ equals the integral of $f(x) = \{2\pi x(1-x)\}^{-1/2}$ over the set $\{u \in [0, 1] \mid GF^{-1}(u) = u\}$.

Received June 1982; revised December 1982.

¹ Research partially supported by Grant A5583 of the NSERC Canada.

AMS 1980 subject classifications. Primary 62G10; secondary 05A15.

Key words and phrases. Goodness-of-fit, Sheffer sequences, distribution of crossings.

Thus, the matching test is consistent if this set has Lebesgue measure 0. The function $f(x)$ gives most of its weight to the tails, and therefore the matching test has only small power if F and G are equal (or very close) on the non-zero parts of the tails (cf. Niederhausen, 1982). Consequently, we proceed analogously to Rényi-type statistics and place a window on the order statistics, defining

$$\sigma(M, a, c; b, B) = \sum_{i=b}^B 1_{[cF(X_{(M)}, cF(X_{(i+1), M}))]}((i + a)/M)$$

for any pair of bounds b and B such that

$$(2) \quad \max(0, -a) \leq b < B \leq \min(M, cM - a).$$

If the bounds b and B fall onto the extreme values in (2), then $\sigma(M, a, c; b, B) = \sigma(M, a, c)$. Takács (1971a, Theorem 2) proved that for these extreme bounds and for $a \geq 0$ the null distribution of σ can be obtained from

$$(3) \quad P\{\sigma(M, a, c; b, B) > s\} = \frac{M!}{(cM)^M} \sum_{j=b+s}^B \frac{(cM - a - j)^{M-j}}{(M - j)!} \sum_{i=s}^{j-b} \frac{(a + j - i)^{j-i}}{(j - i)!(i - s)!} [i^{i-s-2} \{s(s + 1) - i\}]$$

for all $0 \leq s \leq M$. The expression in brackets [] equals one if $i = 0$. In the same theorem he also showed that for $a \geq 0$ and extreme bounds b and B the double sum (3) simplifies to

$$(4) \quad P\{\sigma(M, a, c; 0, B) > s\} = \frac{M!(s + a)}{(cM)^M} \sum_{i=s}^B \frac{(cM - i - a)^{M-i}}{(M - i)!} \frac{(i + a)^{i-s-1}}{(i - s)!} = \frac{M!}{(cM)^s} \left\{ \frac{1}{(M - s)!} - (cM)^{s-M}(s + a) \sum_{i=B+1}^M \frac{(cM - i - a)^{M-i}}{(M - i)!} \frac{(i + a)^{i-s-1}}{(i - s)!} \right\}.$$

The only reason for this simplification is that $b = 0$. A similar result holds for the case $B = M$ as can be seen from the following symmetry relations

$$(5) \quad P\{\sigma(M, a, c; b, B) > s\} = P\{\sigma(M, (c - 1)M - a, c; M - B, M - b) > s\}.$$

More references for special cases of (4) are given in Takács (1971a).

We want to apply a recursive method to derive these distributions, using only very simple properties of Sheffer polynomials. The proofs are just another example for an approach which was applied in Niederhausen (1981) to Rényi type distributions. (For brevity, I shall frequently refer to this paper, blending Takács's and my own notations for the convenience of the reader.) As usual, each approach yields its own generalizations which can be stated as follows.

THEOREM 1. Equation (3) holds for any a, b and B which satisfy the condition (2). Equations (4) and (4') hold if $b = 0$.

The advantages of the recursive method are even more apparent in the two sample case. Let Y_1, \dots, Y_N be a second sample of i.i.d. random variables with empirical distribution function F_Y . Takács (1971b) investigated a test statistic which counts the number of subscripts j where

$$(6) \quad F_Y(Y_{j-1:N}) \leq F_X(Y_{j:N}) + a/N < F_Y(Y_{j:N}).$$

If $N = pM$ for any positive integer p , then (6) is equivalent to

$$(7) \quad j - 1 = pMF_X(Y_{j:N}) + [a],$$

and we can assume that a is an integer. A subscript j can be counted in (7) only if it is of

the form $j = pk + 1 + a$. Denote by $\eta(M, N, a + 1; b, B)$ the number of subscripts $b \leq k \leq B$ where

$$(8) \quad k = MF_X(Y_{pk+a+1:N}).$$

Thus $\eta(M, pM, a + 1; b, B) = \eta_a(M, pM)$ in Takács's notation, if the window bounds fall on the extreme values

$$(9) \quad b = \max(0, \lceil -a/p \rceil) \text{ and } B = \min(M, \lfloor (N - a - 1)/p \rfloor).$$

Takács (1971b, Theorem 1) derived, for $a \geq 0$, $N = pM$ and B as in (9), the distribution of η under the null hypothesis that both samples have the same continuous distribution function

$$(10) \quad \begin{aligned} P\{\eta(M, N, a + 1; 0, B) > s\} \\ = p^s \binom{M + N}{M}^{-1} \sum_{j=M-B}^{M-s} \frac{s(p + 1) + a + 1}{(M - j)(p + 1) + a + 1} \\ \cdot \binom{j + jp + N - pM - a - 1}{j} \binom{(M - j)(p + 1) + a + 1}{M - s - j} \end{aligned}$$

$$(10') \quad = p^s \binom{M + N}{M}^{-1} \left\{ \binom{M + N}{M - s} - \sum_{j=0}^{M-B-1} \dots \right\}$$

for all $0 \leq s \leq B$. More can be shown:

THEOREM 2. *Let a and p be integers, $p > 0$ and $a \geq 0$. (10) and (10') hold for all $N \geq M$ and $0 < B \leq \min(M, \lfloor (N - a - 1)/p \rfloor)$. If b and B lie between the extreme bounds (9), then, for any integer a and $0 \leq s \leq B - b$*

$$(11) \quad \begin{aligned} P\{\eta(M, N, a + 1; b, B) > s\} \\ = p^s \binom{M + N}{M}^{-1} \sum_{j=b+s}^B \binom{M + N - j - pj - a - 1}{M - j} \\ \cdot \sum_{i=s}^{j-b} \binom{(p + 1)(j - i) + a}{j - i} \binom{i + ip}{i - s} \frac{s(ip + 2s + 2) - is - i}{i(ip + s + 1)}, \end{aligned}$$

where the fraction equals 1 if $i = 0$.

Use the symmetry relation (see (2.7))

$$P\{\eta(M, N, a + 1; b, B) > s\} = P\{\eta(M, N, N - pM - a; M - B, M - b) > s\}$$

to obtain the analogous formulas to (10) and (10') if $B = M$.

A more balanced treatment of the two samples can be achieved by counting the number of subscripts k where either $MF_M(Y_{pk+a+1:N}) = k$ as in (8) or $NF_N(X_{k:M}) = pk + a + 1$. The exact distribution is even easier to derive than for η , and the asymptotics are equally simple (Niederhausen, 1982).

The existence of closed expressions for the distribution of σ and η is due to the linearity of the sequences $((i + a)/(cM))$ and $(pk + a + 1)$. Exact distributions can be obtained for general sequences by simple algorithms. Such algorithms are used in Niederhausen (1982) to study the exact efficiency of the one-sample matching test versus the one-sample two-sided Kolmogorov-Smirnov test.

References to earlier work on this distribution (for equal sample sizes) are given in Takács (1971b). In addition, we refer to the papers of Anega and Sen (1972), and Saran and Sen (1983).

1. Proof of Theorem 1. Denote the Lebesgue measure on \mathbb{R}^k by λ^k . We write $\lambda^k\{x; s\}$ short for $\lambda^k\{0 \leq u_1 \leq \dots \leq u_k \leq x \mid u_i \leq (i + a)/cM \leq u_{i+1} \text{ for more than } s \text{ of the}$

subscripts $i = b, \dots, B$), where $u_0 = 0, u_{k+1} = \infty, 0 \leq s \leq B - b$ and $k \geq b + s$. Thus,

$$(1.1) \quad P\{\sigma(M, a, c; b, B) > s\} = M! \lambda^M \{1; s\}.$$

It will be convenient to set $a_i = (i + a)/(cM)$.

The following recursions are obvious. For all $x > \min(a_k, a_B)$

$$(1.2) \quad \begin{aligned} \lambda^k \{x; s\} &= \sum_{j=b+s}^{\min(k,B)} \lambda^j \{0 \leq u_1 \leq \dots \leq u_j \leq a_j \mid u_i \leq a_i \leq u_{i+1} \text{ for} \\ &\text{exactly } s + 1 \text{ of the subscripts } i = b, \dots, j\} (x - a_j)^{k-j} / (k - j)! \\ &= \sum_{j=b+s}^{\min(k,B)} [\lambda^j \{a_j; s\} - \lambda^j \{a_j; s + 1\}] (x - a_j)^{k-j} / (k - j)!, \end{aligned}$$

and for $s \geq 1$ and $k \leq B$, using (1.2),

$$(1.3) \quad \begin{aligned} \lambda^k \{a_k; s\} &= \sum_{j=b+s-1}^{k-1} [\lambda^j \{a_j; s - 1\} - \lambda^j \{a_j; s\}] [(k - j) / (cM)]^{k-j} / (k - j)! \\ &= \frac{1}{cM} \lambda^{k-1} \{a_k; s - 1\}. \end{aligned}$$

Equation (1.2) shows why we are mainly interested in those $\lambda^k \{a_k; s\}$ where $b + s \leq k \leq B$. Define a double sequence $(t_n^s)_{n,s \geq 0}$ of polynomials by

$$(1.4) \quad t_{k-b-s}^s(x) = \sum_{j=b+s}^k [\lambda^j \{a_j; s\} - \lambda^j \{a_j; s + 1\}] (x - a_j)^{k-j} / (k - j)!$$

for all $k \geq b + s$ and for all x . The following properties can be easily verified.

- (i) $t_{k-b-s}^s(x) = \lambda^k \{x; s\}$ for all $x > a_k$ and $b + s \leq k \leq B$,
- (ii) $\deg(t_n^s) = n$,
- (iii) $\frac{d}{dx} t_n^s(x) = t_{n-1}^s(x)$,
- (iv) $t_0^s(x) \neq 0$.

The last two properties imply that $(t_n^s)_{n \geq 0}$ is a Sheffer sequence for the differential operator D (Niederhausen, 1982, page 940).

It is essential for our proof that $t_{k-b-s}^s(x)$ and $\lambda^k \{x; s\}$ also agree at $x = a_k$ for all $b + s \leq k \leq B$. To see this, first note that accordant to (1.3) for $s \geq 0$

$$(1.5) \quad \lambda^k \{a_k; s + 1\} = \frac{1}{cM} t_{k-1-b-s}^s(a_k).$$

Next, compute $t_{k-b-s}^s(a_k)$ from (1.4)

$$\begin{aligned} t_{k-b-s}^s(a_k) &= \sum_{j=b+s}^{k-1} [\lambda^j \{a_j; s\} - \lambda^j \{a_j; s + 1\}] [(k - j) / (cM)]^{k-j} / (k - j)! \\ &\quad + \lambda^k \{a_k; s\} - \lambda^k \{a_k; s + 1\}. \end{aligned}$$

The sum equals $t_{k-1-b-s}^s(a_k) / (cM)$ and cancels against the last term on the right hand side because of (1.5). Hence

$$(1.6) \quad t_{k-b-s}^s(a_k) = \lambda^k \{a_k; s\} = \frac{1}{cM} t_{k-1-b-s}^{s-1}(a_k).$$

Equation (1.6) also shows that the Sheffer sequences $(t_n^s)_{n \geq 0}$ and $(t_n^{s-1} / (cM))_{n \geq 0}$ agree at one point for each n . Therefore, they are equal everywhere (Niederhausen, 1982, Lemma A.2), and we get for all $b + s \leq k \leq B$ and $x \geq a_k$

$$(1.7) \quad \begin{aligned} \lambda^k \{x; s\} &= t_{k-b-s}^s(x) = (cM)^{-s} t_{k-b-s}^0(x) = (cM)^{-s} \lambda^{k-s} \{x; 0\} \\ &= (cM)^{-s} \lambda^{k-s} \{0 \leq u_1 \leq \dots \leq u_{k-s} \leq x \mid u_j \leq a_j \leq u_{j+1} \text{ for} \\ &\quad \text{at least one } b \leq j \leq k - s\} \\ &= (cM)^{-s} \left\{ \frac{x^{k-s}}{(k-s)!} - s_{k-s}(x) \right\}, \end{aligned}$$

where $s_{k-s}(x) = \lambda^{k-s} \{0 \leq u_1 \leq \dots \leq u_{k-s} \leq x \mid u_j > a_j \text{ for all } b \leq j \leq k - s\}$. $(s_n)_{n \geq 0}$ is the

Sheffer sequence for D with roots in

$$v_i = \begin{cases} 0 & i = 0, \dots, b-1 \\ (i + a)/(cM) & i \geq b \end{cases}$$

which frequently occurs in exact Rényi-type distributions. It is well known (cf. Niederhausen, 1982, Theorem A.1) that

$$(1.8) \quad \begin{aligned} s_{k-s}(x) &= \sum_{i=0}^{b-1} \frac{((i + a)/(cM))^i}{i!} \{x - (k - s + a)/(cM)\} \frac{\{x - (i + a)/(cM)\}^{k-s-i-1}}{(k - s - i)!} \\ &= \frac{x^{k-s}}{(k - s)!} - \sum_{i=b}^{k-s} \dots \end{aligned}$$

From (1.7) and (1.8) we get

$$\lambda^s\{a_j; s\} = s(cM)^{-j} \sum_{i=b}^{j-s} \frac{(i + a)^i (j - i)^{j-s-i-1}}{i! (j - s - i)!}.$$

Summing up the differences accordant to (1.2) yields (3) from (1.1). If $B = M$, we obtain $\lambda^M\{1; s\}$ directly from (1.7) and (1.8), proving (4) and (4') via symmetry relations (5).

2. Proof of Theorem 2. Denote the combined sample $X_1, \dots, X_M, Y_1, \dots, Y_N$ by V_1, \dots, V_{M+N} . It is a standard technique to represent two-sample rank statistics by lattice path problems (cf. Mohanty, 1979). We construct a path from the ordered combined sample $V_{1 \cdot M+N}, \dots, V_{M+N \cdot M+N}$ by going to the right in the ℓ th step, if $V_{\ell \cdot M+N}$ belongs to the X -sample, and going upwards, if $V_{\ell \cdot M+N}$ belongs to the Y 's. This path ends at the point (M, N) , and we denote the set of all such paths by $\mathcal{T}(M, N)$. Under the null hypothesis each path occurs with probability $1 / \binom{M+N}{M}$.

The condition (8) $k = MF_X(Y_{pk+a+1N})$ means that the corresponding path reaches the point $(k, pk + a + 1)$ in its $(pk + a + 1)$ th vertical step. Denote by $D(k, m; s)$ the number of paths in $\mathcal{T}(k, m)$; which (i) go upwards during their last step, and (ii) reach the line $y = px + a + 1$, where $b \leq x \leq B$, more than s times from below. Hence

$$(2.1) \quad P\{\eta(M, N, a + 1; b, B) > s\} = D(M, N + 1; s) / \binom{M+N}{M}.$$

It will be convenient to set $a_i = pi + a + 1$.

The following recursions are obvious. For all $m > \min(a_k, a_B)$

$$(2.2) \quad D(k, m; s) = \sum_{j=b+s}^{\min(k, B)} [D(j, a_j; s) - D(j, a_j; s + 1)] \binom{k-j+m-1-a_j}{k-j},$$

and for $s \geq 1$ and $k \leq B$, using (2.2),

$$(2.3) \quad \begin{aligned} D(k, a_k; s) &= \sum_{j=b+s-1}^{k-1} [D(j, a_j; s - 1) - D(j, a_j; s)] \binom{(p+1)(k-j)-1}{k-j} \\ &= pD(k-1, a_k+1; s-1). \end{aligned}$$

From now on, everything is completely analogous to the proof of Theorem 1. A Sheffer sequence—but now for the backwards difference operator ∇ —can be constructed from the right hand side of (2.2) which is equal to $D(k, m + 1; s)$ for all $m > a_k$ and $b + s \leq k \leq B$. It follows from (2.3), that equality holds also for $m = a_k$. The uniqueness property of Sheffer polynomials shows again that (2.3) must hold for all $m \geq a_k$ if $b + s \leq k \leq B$.

Hence,

$$\begin{aligned}
 (2.4) \quad D(k, m; s) &= p^s D(k - s, m + s; 0) \\
 &= p^s \#\{\text{paths in } \mathcal{T}(k - s, m + s - 1) \text{ which reach the line} \\
 &\quad y = px + a + 1 \quad (x \geq b) \text{ from below}\} \\
 &= p^s \left\{ \binom{k + m - 1}{k - s} - s_{k-s}(m + s - 1) \right\},
 \end{aligned}$$

where $s_{k-s}(m + s - 1) = \#\{\text{paths in } \mathcal{T}(k - s, m + s - 1) \text{ which stay strictly over the line } y = px + a \text{ (} x \geq b)\}$.

$(s_n)_{n \geq 0}$ is the Sheffer sequence for ∇ with roots in

$$v_i = \begin{cases} -1 & \text{for all } i = 0, \dots, b - 1, \\ pi + a & \text{for all } i \geq b \end{cases},$$

which also occurs in Rényi-type distributions. It is well known that (cf. Niederhausen, 1981, (3.8); Mohanty, 1979)

$$\begin{aligned}
 (2.5) \quad s_{k-s}(m + s - 1) &= \sum_{i=0}^{b-1} \binom{i + pi + a}{i} \\
 &\quad \cdot \frac{m + s - 1 - p(k - s) - a}{m + s - 1 - pi - a} \binom{k + m - 2 - i - pi - a}{k - s - i} \\
 &= \binom{k + m - 1}{k - s} - \sum_{i=b}^{k-s} \dots.
 \end{aligned}$$

From (2.4) and (2.5) we get

$$D(j, a_j; s) = p^s \sum_{i=b}^{j-s} \binom{i + pi + a}{i} \frac{s + ps}{p(j - i) + s} \binom{(p + 1)(j - i) - 1}{j - s - i}.$$

Summing up the differences accordant to (2.2) yields (11) from (2.1). If $B = M$,

$$\begin{aligned}
 (2.6) \quad P\{\eta(M, N, a + 1; b, M) > s\} &= p^s \binom{M + N}{M}^{-1} \left\{ \binom{M + N}{M - s} - s_{M-s}(N + s) \right\} \\
 &= p^s \binom{M + N}{M}^{-1} \left\{ \binom{M + N}{M - s} - \sum_{i=0}^{b-1} \binom{i + pi + a}{i} \right\} \\
 &\quad \cdot \frac{N - pM + s + ps - a}{N + s - pi - a} \binom{M + N - i - pi - a - 1}{M - s - i} \Big\} \\
 &= p^s \binom{M + N}{M}^{-1} \sum_{i=b}^{M-s} \dots.
 \end{aligned}$$

From $B = M$ follows that $pM + a + 1 \leq N$. Therefore, we could not obtain Takács's (1971b) original result, Theorem 1, because there $N = pM$ and $a \geq 0$ is assumed. The standard trick in such a situation is to turn the paths upside down and to run through them backwards. Each point (i, j) is now transformed into $(M - i, N - j)$, and we are counting the number of paths in $\mathcal{T}(M, N)$ which *leave* the line $y = px + N - pM - a - 1$ ($M - b \leq x \leq M - b$) in a *vertical* step more than s times. Hence, they *reach* the line $y = px + N - pM - a$ ($M - B \leq x \leq M - b$) more than s times *from below*, and can therefore be counted as before

$$(2.7) \quad P\{\eta(M, N, a + 1; b, B) > s\} = P\{\eta(M, N, N - pM - a; M - B, M - b) > s\}.$$

Equation (10) is obtained from (2.6) in this way.

REFERENCES

- ANEGA, K. G. and SEN, K. (1972). Random walk and distribution of rank order statistics. *SIAM J. Appl. Math.* **23** 276-287.
- MOHANTY, S. G. (1979). *Lattice Path Counting and Appliances*. Academic Press, New York.
- NIEDERHAUSEN, H. (1981). Sheffer polynomials for computing exact Kolmogorov-Smirnov and Rényi type distributions. *Ann. Math. Statist.* **9** 923-944.
- NIEDERHAUSEN, H. (1982). The exact distribution of the matching test. Tech. Report, Univ. of Toronto, Dept. of Statistics.
- SARAN, G. and SEN, K. (1983). On the distribution of crossings in a generalized random walk. *J. Statist. Plann. Inference*, to appear.
- SIDDIQUI, M. M. (1982). The consistency of a matching test. *J. Statist. Plann. Inference* **6** 227-233.
- SMIRNOV, N. V. (1939). On deviations of the empirical distribution function. (In Russian). *Mat. Sbornik* **6** 3-26.
- TAKÁCS, L. (1971a). On the comparison of a theoretical and an empirical distribution function. *J. Appl. Probab.* **8** 321-330.
- TAKÁCS, L. (1971b). On the comparison of two empirical distribution functions. *Ann. Math. Statist.* **42** 1157-1166.

DEPARTMENT OF STATISTICS
UNIVERSITY OF TORONTO
TORONTO, ONTARIO, CANADA M5S 1A1