

The significance probability of the Smirnov two-sample test

By J. L. HODGES, Jr.¹

With 3 figures in the text

1. Introduction

In 1939 N. V. Smirnov proposed the following rank-order test for the two-sample problem. Let x_1, \dots, x_m and y_1, \dots, y_n be samples of independent observations from populations with continuous distribution functions F and G , respectively. Form from the samples the empirical distribution functions F_m and G_n ; that is, $mF_m(u)$ is the number of the observations x_1, \dots, x_m which do not exceed u , with $nG_n(u)$ defined analogously. To test the hypothesis $F = G$ we use the statistic $D = \sup_u |F_m(u) - G_n(u)|$, large values of which are significant. We may without loss of generality assume $m \geq n$.

It is clear that the significance probability $Pr \{D \geq d | F = G\}$, which we shall denote throughout by P_2 , is independent of the common value of $F = G$; that is, the test (like all rank-order tests for the two-sample problem) is similar over the class of all continuous distributions. Further, the fact that $\sup_u |F_m(u) - F(u)|$ tends to 0 in probability as $m \rightarrow \infty$ implies that the test is consistent against all alternatives $F \neq G$. These properties of similarity and consistency, together with a certain mathematical elegance, give the test wide appeal to mathematical statisticians. A considerable literature has developed, the proposer of the test has been awarded a Stalin prize (Kolmogorov and Hinč'in 1951), and the test has begun to appear in applied handbooks. The test is not very powerful against specific alternatives such as shift (van der Waerden 1953), but this could hardly be expected in view of its consistency.

Smirnov's test was suggested by analogy with the earlier test of Kolmogorov (1933) for the one-sample problem. In fact, Smirnov's test generalizes Kolmogorov's, for when $n \rightarrow \infty$ we may replace G_n by G , and D becomes Kolmogorov's statistic for the hypothesis that F equals a completely specified G . Thus general results on the Smirnov test usually give (by the limit passage $m \rightarrow \infty$) results on the Kolmogorov test. We shall not however attempt to discuss the significance problem for Kolmogorov's test, nor shall we take up the many variants of Smirnov's test which have been suggested.

Smirnov's test also appears in a one-sided version. We may use $D^+ = \sup_u [F_m(u) - G_n(u)]$ to test the hypothesis that $F(u) \leq G(u)$ for all u . This form of the test is in

¹ This paper was written while the author was a fellow of the John Simon Guggenheim memorial foundation, and a guest of Stockholms högskola. It is a pleasure to record appreciation for the courtesies extended to me by the högskola and its rector, Professor Harald Cramér.

fact often appropriate when we wish to know whether a new method of treatment (producing the population G) gives larger values than the standard treatment (population F) for some quantile. We shall denote by P_1 the quantity $Pr \{D^+ \geq d | F = G\}$, which is the size of the one-sided test. One can also define a non-symmetrical two-sided test statistic, which includes D and D^+ as special cases. While no essential difficulties appear in doing this, we shall not discuss the general statistic as its usefulness does not appear to justify the considerable notational complication required.

In spite of the fact that the Smirnov test has been in use for nearly twenty years, the computation of P_1 and P_2 has not been satisfactorily dealt with, nor does it even seem to be widely realized that a computational problem exists. We review in Section 2 the methods now available, and give in Section 3 the results of a brief numerical investigation showing the inadequacy of those methods. In Section 4 we present a technique which is useful when $m - n$ is small, and which serves to illuminate the complexity of the problem. Among the remarks in the concluding Section 5 is an interpolation formula which appears to give considerably better values than the currently standard technique.

2. Methods for obtaining values of P_1 and P_2

The user of Smirnov's test, faced with the problem of obtaining a value of P_1 or P_2 corresponding to the observed value of D^+ or D , has available a variety of methods, which we shall now review.

(a) *Direct computation*

As the distribution $F = G$ is continuous, we may assume that the $m + n$ observations are distinct. Since the samples are drawn from the same population, we may imagine that they were obtained by first drawing $m + n$ observations, and then selecting at random m of these to form the first sample. Thus, under the null hypothesis $F = G$, each of the $\binom{m+n}{n}$ possible orderings of the samples with respect to each other has the same probability $1 / \binom{m+n}{n}$. As with all rank-order tests, the problem of computing P_1 and P_2 reduces to a purely combinatorial one.

The solution of the combinatorial problem, and incidentally the calculation of the values of D^+ and D , is aided by a graphical device. We arrange the $m + n$ observations in increasing order on a common sequence, and associate with this arrangement a path in the x, y plane. We begin at the origin, and (reading the observations from smallest to largest) take a unit step to the right for each x -observation, and a unit step up for each y -observation. The path terminates at the point (m, n) , and from it we can reproduce the ranks of the two samples in their common array. Figure 1 illustrates the process for samples in the order indicated by $xyxyxyyx$.

Such a graphical representation is of course a standard method in classical probability, for example in the problem of gambler's ruin. It is usually presented with steps to the right and left instead of to the right and up (Korolyuk and Yaroševs'kii 1951, Gnedenko and Korolyuk 1951) but the present version, which is used by Drion (1952), is more convenient for use with ordinary graph paper.

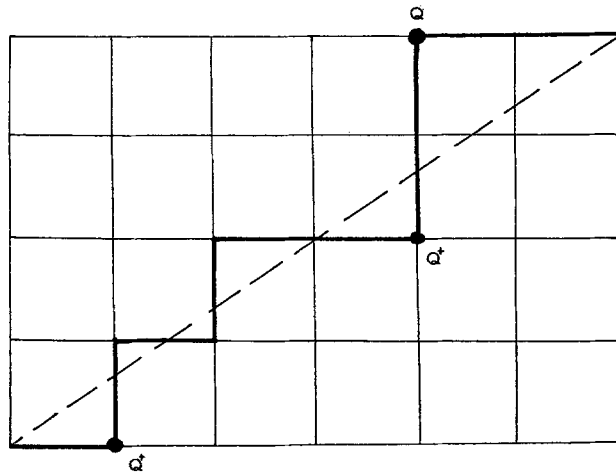


Fig. 1.

Suppose that of the first $x + y$ observations in the common array, just x come from the first sample. Then between the $(x + y)$ th and $(x + y + 1)$ st observations, we shall have $F_m(x) - G_n(x) = x/m - y/n = (nx - my)/mn$, which is proportional to the distance of (x, y) from the diagonal of the rectangle with corners $(0, 0)$ and (m, n) . Therefore to determine $D^+(D)$ we need only locate those points $Q^+(Q)$ on the path which are farthest below (farthest from) the diagonal. In Figure 1, $D^+ = \frac{1}{6}$, corresponding to either of the points labelled Q^+ , while $D = \frac{1}{3}$, corresponding to the point labelled Q .

To compute P_2 for an observed path, we may count those paths which never get as far from the diagonal as the farthest point on the observed path. Construct *boundaries* (Figure 2), parallel to the diagonal and at the same distance from it as Q . The number $A(x, y)$ of ways to go from $(0, 0)$ to (x, y) while staying strictly inside the boundaries satisfies the recursion formula

$$A(x, y) = A(x - 1, y) + A(x, y - 1) \tag{2.1}$$

with starting values $A(0, y) = A(x, 0) = 1$. The desired value $A(m, n)$ can be computed by simple additions, as illustrated on Figure 2, where $P_2 = 1 - A(m, n) / \binom{m+n}{n} = 97/105$. This technique, which will be referred to as the *inside method*, is effective for small sample sizes but becomes rapidly less so as m and n increase. For fixed values of P_2 the number of additions increases as $n^{1/2}$ and the size of the numbers increases exponentially.

The work can be substantially reduced by two devices: (i) Until the boundary is reached, we are simply generating combinatorials by Pascal's triangle, so that an initial part of the work is unnecessary. (ii) By symmetry, we need carry on only until $x + y \geq \frac{1}{2}(m + n)$, and use the fact that the number of paths through (x, y) is the product of $A(x, y)$ and $A(m - x, n - y)$. Finally, when m and n have a large common factor, and P_2 is large, Polya's (1948) method for exact sequential analysis may be useful.

This inside method can also be used for P_1 , but since the number of additions now increases as n^2 , the alternative *outside method* is usually preferable. We count the paths which reach the (lower) boundary, classifying them according to the point (x, y) at which they first reach (or pass) it. If $B(x, y)$ is the number of ways to go from $(0, 0)$ to (x, y) without previously reaching the boundary, then the total number of paths reaching the boundary is

$$\sum_y B(x, y) \binom{m+n-x-y}{n-y} \tag{2.2}$$

We may compute B successively for $y = 0, 1, \dots$, by observing that $B(x, y)$ is $\binom{x+y}{y}$ diminished by the number of those ways of going from $(0, 0)$ to (x, y) which previously reach the boundary. The latter number can easily be found with the aid of earlier B values. This process, proposed by Korolyuk (1955a), is illustrated below for the problem of Figure 2.

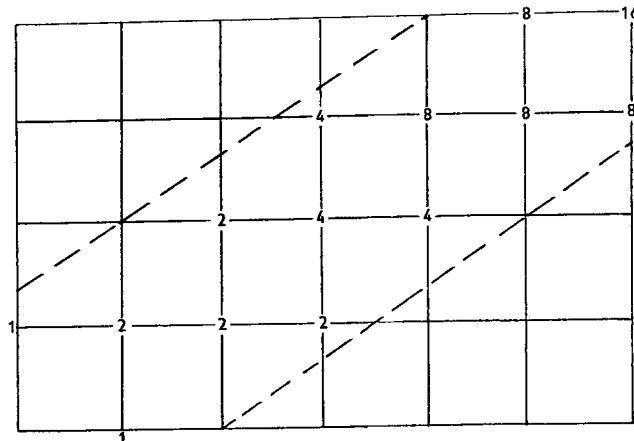


Fig. 2.

y	x	$\binom{x+y}{y}$	$B(x, y)$	$\binom{m+n-x-y}{n-y}$
0	2	1	3	10
1	4	5	2	2
2	5	21	7	3

Thus, $B(4, 1) = 5 - 3$, where the 3 represents the number $\binom{3}{1}$ of ways to go from $(2, 0)$ to $(4, 1)$. The sum of products of the final column is 111, the number of paths reaching the boundary, so that $P_1 = 111/210$.

We note that the outside method can also be used for P_2 when the boundaries are so far apart that a path cannot reach both, in which case $P_2 = 2P_1$. Even if a few paths can reach both boundaries, it may be best to allow for these separately and use the outside method.

(b) Tables

Massey has published two tables of P_2 , both computed by the inside method.

(i) In Massey (1951) we are given $1 - P_2$, to from 2 to 6 significant figures, for $m = n = 1(1)40$ and for $a = 2(1)13$ where $d = a/n$. At $m = n = 40$ the last decimal is not reliable. This table is a useful complement to formula (2.4) below.

(ii) Massey (1952) gives $1 - P_2$ to 5 decimals, for all $n \leq m \leq 10$ and all d , and for selected values of $n < m$ and d for $m > 10$. Unfortunately this table does not seem to be reliable. In checking about 125 values, using the method of Section 4, I found the following 16 instances in which I could not verify Massey's value of P_2 . The first of the given values is taken from Massey (1952), while the second is mine.

Table 1.

n	m	d	$Pr(D \geq d)$		n	m	d	$Pr(D \geq d)$	
6	7	36	00874	00816	7	10	56	00607	00452
6	8	21	00533	00466	8	10	22	09877	09511
7	8	26	34297	31313			23	07683	07043
		27	28205	25221	9	10	54	03436	03027
		28	22238	19254			63	00719	00704
		35	06371	05594	12	16	30	00577	00525
7	10	50	01584	01399			31	00376	00331
		53	01100	00946	15	20	31	01414	01363

It seems unlikely that Massey's second table will be extended in its present form to cover very large sample sizes, because the number of possible values of D increases so rapidly. It can be shown that D has $1 + [\frac{1}{2}rs] + (t-1)rs$ possible values, where $m = rt$, $n = st$, and r and s are relatively prime. While 568 entries suffice for all $n < m \leq 10$, 8707 would be required for $n < m \leq 20$, and one can show (using known facts on the density of relative primes) that the number of entries required to cover $n < m \leq M$ is asymptotically $M^4[2\zeta(3) - \zeta(4)]/16\zeta(2) = 0.0502 M^4$. While P_1 and P_2 could be programmed for efficient electronic computation, the cost of publishing an adequate table would be excessive.

(c) Closed expressions

In a few special cases one can obtain expressions for P_1 and P_2 which are relatively easy to compute.

(i) When $m = n$ our rectangle becomes a square, the boundaries have slope 1, and we can use a reflectional method. Let the boundary be $y = x - a$, where $a = 1, 2, \dots, m$. A path reaching the boundary will first do so at some point Q . We reflect that part of the path from $(0,0)$ to Q about the boundary, generating a path from $(a, -a)$ to (n,n) . As the number of these is $\binom{2n}{n-a}$, and as the correspondence is one-to-one, we have

$$P_1 = \binom{2n}{n-a} / \binom{2n}{n}. \tag{2.3}$$

To compute P_2 , we must allow for paths which reach both boundaries, which can be done by repeated reflections. We find

$$P_2 = 2 \left[\binom{2n}{n-a} - \binom{2n}{n-2a} + \binom{2n}{n-3a} - \dots \right] / \binom{2n}{n}. \quad (2.4)$$

The argument is classical, and is given for example by Gnedenko and Korolyuk (1951) and Drion (1952).

(ii) It has recently been pointed out that similar results hold when $m = np$, where p is a positive integer. Recall the quantity B entering into the inside method. Korolyuk (1955a) noted that of those paths to the point (x, y) on the boundary $py = x - a$ which have previously reached it, exactly $1/(p + 1)$ approach (x, y) through $(x, y - 1)$. From this it follows that

$$B(x, y) = \binom{x+y}{y} - (p+1) \binom{x+y-1}{y-1}.$$

When this is substituted into (2.2) we see that the number of paths reaching the boundary is

$$N(a) = \sum_{y=0}^{\lfloor n-a/p \rfloor} \frac{a}{(p+1)y+a} \binom{(p+1)y+a}{y} \binom{(p+1)(n-y)-a}{n-y}.$$

We then of course have $P_1 = N(a) / \binom{(p+1)n}{n}$.

Two attempts have been made to obtain a corresponding result for P_2 . Korolyuk (1955a, p. 86) produced a formula, but it contains a partitionial sum that would be difficult to compute. The formula given in Blackman (1956) is unfortunately incorrect, and the revised formula (Blackman 1957) is not suited for easy computation.

(iii) A closed formula (4.4) for P_1 when $m = n + 1$ is developed in Section 4 below.

(d) *Large-sample approximations*

(i) In his original papers (1939a, b) Smirnov proved that, as m and $n \rightarrow \infty$ so that $m/n \rightarrow q$ we have, for fixed $z > 0$,

$$Pr \left\{ \sqrt{\frac{mn}{m+n}} D^+ \geq z \right\} \rightarrow e^{-2z^2}, \quad (2.5)$$

$$Pr \left\{ \sqrt{\frac{mn}{m+n}} D \geq z \right\} \rightarrow 1 - K(z) = 2 [e^{-2z^2} - e^{-2(2z)^2} + e^{-2(3z)^2} - \dots]. \quad (2.6)$$

The function K , which also appears in the limit theory of the Kolmogorov test, has been tabled by Smirnov (1939b, 1948) to 6 decimals for its argument at intervals of 0.01 over the entire range. Alternative proofs or heuristic proofs of (2.5) and 2.6) have been given by Feller (1948), and Doob (1949). (See also Donsker 1952.)

Smirnov's theorems suggest the introduction of new random variables $Z = \sqrt{mn/(m+n)} D$ and $z^+ = \sqrt{mn/(m+n)} D^+$. We shall hereafter always restrict z

to possible values of Z and Z^+ , and shall refer to z as the *distance variable* of the boundary. It is clear that, if z' and z'' are consecutive possible values of Z or Z^+ , and if $z' > z \geq z''$, then $Pr\{z^{(+)} \geq Z\} = Pr\{Z^{(+)} \geq z''\}$, so we lose no generality by restricting z . To permit z to assume arbitrary values, as is customary in the literature, leads to considerable and needless notational complication. Similarly, d will always denote a possible value of D or D^+ , with $z = \sqrt{mn/(m+n)}d$, etc.

It is a truism that a limit theorem can be used to justify many different large-sample approximations. For example, we might on the basis of (2.5) approximate $Pr(D^+ \geq d)$ by $\exp\{-2mnd^2/(m+n)\}$, but this quantity could with equal reason be used to approximate $Pr(D^+ > d)$, which in some cases is substantially different.

It has been noted by Drion (1952) that good results are obtained, when $m = n$, if we approximate $P_2 = Pr(D \geq d)$ by the quantity $1 - K(\sqrt{\frac{1}{2}n}d)$. The corresponding observation for P_1 has been made by H. E. Daniels (*J. Roy. Stat. Soc. B* (1956), V. 18, p. 22). These observations suggest the use of the large-sample approximations

$$\tilde{P}_1 = e^{-2z^2} \quad \text{and} \quad \tilde{P}_2 = 1 - K(z),$$

for P_1 and P_2 respectively. This numerical finding can be reinforced by an asymptotic expansion of (2.3) and (2.4). Using Stirling's formula and expanding the logarithms involved, it is easy to show that when $n \rightarrow \infty$ and $a = 0(\sqrt{n})$,

$$P_1 = \exp\left\{-\frac{a^2}{n} + \frac{a^2}{2n^2} - \frac{a^4}{6n^3} + 0\left(\frac{1}{n^2}\right)\right\}.$$

The substitution $a = z\sqrt{2n}$ gives

$$P_1 = \exp\left\{-2z^2 + \frac{z^2}{n}\left(1 - \frac{2}{3}z^2\right) + 0\left(\frac{1}{n^2}\right)\right\}, \tag{2.7}$$

which shows that the use of \tilde{P}_1 as an approximation for P_1 will lead to a relative error of order $1/n$, when $m = n$. Note that in this case, the customary continuity correction would introduce an error of order $1/\sqrt{n}$; but see Section 5(a). A similar analysis shows that

$$P_2 = e^{-2z^2} - e^{-2(2z)^2} + e^{-2(3z)^2} - \dots + 0\left(\frac{1}{n}\right) = P_1 - P_1^2 + P_1^3 - \dots + 0\left(\frac{1}{n}\right). \tag{2.8}$$

Essentially these expansions are to be found in Gnedenko (1952). (The version in Gnedenko (1954) appears to be in error.)

(ii) Korolyuk (1954, 1955b) has recently developed asymptotic expansions for P_1 and P_2 giving the terms of order $1/\sqrt{n}$ and $1/n$ for arbitrary m, n . His formulae would appear to provide means for dealing with the practical problem of determining P_1 and P_2 . However, their correctness has been challenged by Blackman, in his review of Korolyuk's paper [*Mathematical Reviews* 16 (1955) 839], and again in Blackman (1956). He finds that Korolyuk's results are not consistent with earlier results of Smirnov (1944) on the one-sample problem, with some of Gnedenko's results for $m = n$, nor with Blackman's own results for $m = np$.

As a consequence of the theory developed in Section 4, it follows that

$$\sqrt{n} [e^{2z^2_{m,n}} P r \{Z^+ \geq z_{m,n}\} - 1]$$

does not tend to a limit as $m, n \rightarrow \infty, z_{m,n} \rightarrow z, m/n \rightarrow 1$. Thus, not only is the particular expression for P_1 given by Korolyuk wrong, but no general expression of this simple character can be right. Korolyuk uses analytic machinery on a function of two discrete variables in a formal way, without verifying its applicability. In particular he ignores the lattice structure of the boundary of his region, and it follows from the development of Section 4 that the structure of the boundary in some cases influences the term of order $1/\sqrt{n}$.

3. Numerical examination of the Smirnov approximations

The preceding section leads to this main conclusion: except for quite small m , we shall usually either have to carry out a rather heavy computation to obtain the significance probability, or else rely on the approximations \tilde{P}_1 and \tilde{P}_2 based on Smirnov's principal term limit theorem. Although these have been in use for nearly twenty years, I have not been able to find any numerical examination of their accuracy, except when $m = n$. The present section presents the results of a brief numerical study, for $n = 12, m = 13(1)18$, and $0.05 > P_2 > 0.002$.

We shall report our results in terms of the *relative*, not the *absolute*, accuracy of the approximations. This may be motivated from the point of view of the Neyman-Pearson theory, by observing that an error of given percentage in determining the significance level of a test will usually result in a percentage error of the same order of magnitude in the power function generally. Again, from the Bayesian viewpoint, a percentage error in small P results in a *posteriori* odds in error by about the same percentage, without regard to the absolute error in P .

In the range studied, \tilde{P}_2/P_2 and \tilde{P}_1/P_1 will be nearly identical. From (2.6) we see that $\tilde{P}_2 = 2P_1[1 - \tilde{P}_1^3 + \tilde{P}_1^8 - \dots]$ so that the two-tailed approximation is almost exactly double the one-tailed approximation. Correspondingly, the actual value P_2 is either exactly $2P_1$ or else very nearly so throughout our range. The required values of P_1 were computed by the method of Section 4 and by the outside method.

The results are shown as Figure 3, which for each m gives $\log_e(\tilde{P}_2/P_2)$ against z . An examination of this figure reveals several interesting features.

(i) The relative errors are perhaps surprisingly large. For comparison, for $m = n = 12$, the relative errors of the two values of \tilde{P}_2 , in the same range of P_2 , are 9% and 23%. When we add a single observation to one sample, the relative errors are increased to range from 38% to 180%. This example may serve as a useful warning against the common belief that increased sample sizes always favor an asymptotic approximation.

(ii) The relative error is by no means monotone in z ; instead there are wild oscillations at $m = 13$, which gradually subside as m is increased to 16, but which seem to reappear at $m = 17$ only to vanish at $m = 18$.

(iii) If we try to average out the oscillations by some sort of trend line, we see that the "average" relative error increases—perhaps linearly—with z .

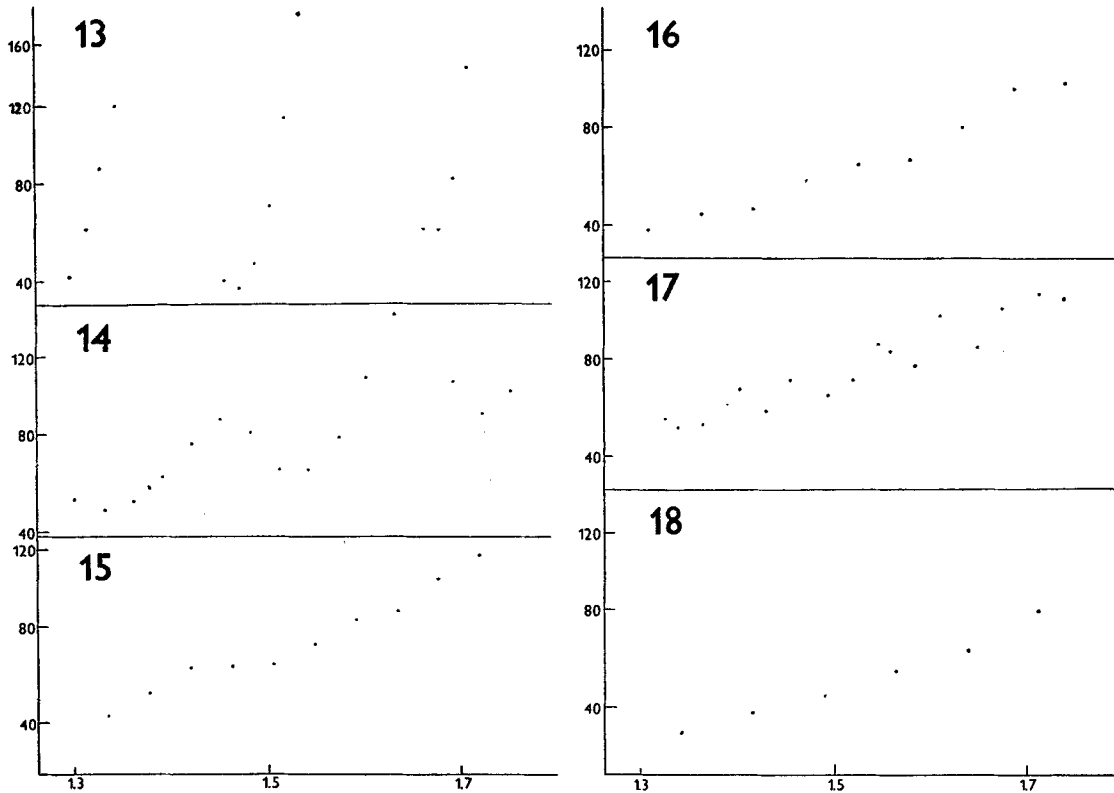


Fig. 3.

(iv) But, for fixed z , the average is not monotone in m . In general it falls as m is increased, but at $m = 17$ it is higher than at $m = 16$.

These numerical results raise both practical and mathematical problems. The Smirnov approximation is seen to be highly inaccurate for values of m and n which are already large enough for direct computations to be arduous. Further, some mathematical explanation is desirable for the phenomena just observed. We shall in the two following sections give partial solutions for these problems.

4. Nearly equal sample sizes

Considerations of efficiency or symmetry usually lead to specifying $m = n$ in the design of comparative experiments, but one or more observations is often lost. As a result, the problem of calculating P_1 and P_2 when m and n are nearly but not exactly equal acquires practical importance, especially since the Smirnov approximation is particularly bad in this case. We shall now develop a theory appropriate when $m - n = c$ is small. We begin by establishing a general lower bound for P_1 .

For brevity we shall refer to the line segments $y = x - a$, for $a = c, c + 1, \dots, m$ and $y \geq 0, x \leq m$, as *mirrors*. Each path reaches at least one mirror; of these, the

one with largest a will be called the *end mirror*. A path touches its end mirror in at least one point; the lowest of these is the *end point*. The random coordinates of the end point will be denoted by (X, Y) and we write $X - Y = A$. We shall use the random variable A to index the mirrors.

The reflection method (Section 2c) gives at once

$$Pr(A \geq a) = \binom{m+n}{n+a} / \binom{m+n}{n}. \quad (4.1)$$

Furthermore, it enables us to calculate the number of ways to go from $(0,0)$ to $(a+y, y)$ without previously touching the mirror $A = a$, which will be denoted by $I(a, y)$. This number is the same as the number of ways to go from $(0,0)$ to $(a+y-1, y)$ without touching $A = a$. There are $\binom{a+2y-1}{y}$ ways to reach $(a+y-1, y)$, but by reflection we see that $\binom{a+2y-1}{y-1}$ of these have touched $A = a$. Hence

$$I(a, y) = \binom{a+2y-1}{y} - \binom{a+2y-1}{y-1}. \quad (4.2)$$

Similarly, the number of ways to go from $(a-y, y)$ to (m, n) without crossing (but possibly touching) the mirror $A = a$ is $I(n+a+1-m, m-a-y)$. Thus

$$\binom{m+n}{n} Pr(A = a \text{ and } Y \leq u) = \sum_{y=0}^u I(a, y) I(n+a+1-m, m-a-y). \quad (4.3)$$

Consider now a lower boundary line L , corresponding to a most distant point (x^+, y^+) . The line $L: m(y - y^+) = n(x - x^+)$ will intersect at least one mirror; let the highest mirror which L intersects be $A = \alpha$. Then α is the smallest integer greater than or equal to $x^+ - m y^+ / n$. It is easy to show that the ordinate of the point of intersection of L and $A = \alpha$, which we shall denote by $v n$, is equal to $n[1 - g(x^+ - m y^+ / n)] / c$, where $g(u)$ denotes the fractional part of u . The variable v thus defined will be called the *phase variable* of L . The ordinate of the point of intersection of L with $A = \alpha + i$ is then $n(v + i/c)$. Finally, for the distance variable z we have the expression

$$z = \sqrt{\frac{m n}{m+n}} \left(\frac{\alpha + v n}{m} - \frac{v n}{n} \right),$$

so that
$$\alpha = v c + z \sqrt{\frac{(n+c)(2n+c)}{c}}.$$

Any path whose end point lies on or below L must of course reach or pass L . Therefore

$$\binom{m+n}{n} P_1 \geq \binom{m+n}{m+\alpha} + \sum_{a=\alpha}^{\alpha+c} \sum_{y=0}^{\lfloor n(v+i/c) \rfloor} I(a, y) I(n+a+1-m, m-a-y). \quad (4.4)$$

We shall denote the right side of (4.4) by $\binom{m+n}{n} P_1^*$. P_1^* is thus a lower bound for P_1 , and equals, P_1 if and only if L cuts a single mirror, which will always be the case when $c = 1$. Thus (4.4) provides an exact closed expression for P_1 when $m = n + 1$. It is intuitively plausible that P_1^* should be close to P_1 when c is small, and it is shown below that $P_1 - P_1^* = O(1/n)$ when $n \rightarrow \infty$ and c is bounded.

The most notable feature of P_1^* is that it depends in a simple way on a function, I , of only two variables, while P_1 itself depends on the four variables m, n, x^+ , and y^+ . We give a short table of I , suitable for computing P_1 when $m \leq 30$. From the 323 values of Table 2, and a table of combinationals, one can easily obtain any of thousands of accurate estimates for P_1 . (As explained in Section 3, $2P_1^*$ may then be used to estimate P_2 if the latter is not too large.)

As an example, we take $m = 20, n = 16$, and seek P_1 corresponding to a path with farthest point (17, 6). The boundary L is $5(y - 6) = 4(x - 17)$, which cuts the mirrors $A = 10, 11$, so that $\alpha = 10$. The ordinates of the points of intersection of L with $A = 10$ and $A = 11$ are respectively 2 and 6. As most of the mirror $A = 11$ lies below L , it is quicker to compute the complement of (4.3) with respect to $\binom{m+n}{n}$

$Pr(A = a)$. We have

$$\begin{aligned} \binom{36}{16} P_1^* &= \binom{36}{9} - [I(11, 7) I(8, 2) + I(11, 8) I(8, 1) + I(11, 9) I(8, 0)] + \\ &\quad + [I(10, 0) I(7, 8) + I(10, 1) I(7, 9) + I(10, 2) I(7, 10)] \\ &= 9.1407 \times 10^7 \end{aligned}$$

This gives $P_1^* = 0.01251$. Massey (1952) gives $P_2 = 0.02511$, and since the probability of reaching both boundaries is negligible, this implies $P_1 = 0.01256$. In the range $P_2 < 0.1$, $2 P_1^*$ never differs from P_2 as given by Massey by as much as 1% of P_2 . The accuracy of P_1^* will of course decrease as c increases, but even at $n = 12, m = 20$ P_1^* is considerably superior to \tilde{P}_1 .

We shall now develop an asymptotic expression for P_1 , correct to terms of order $1/\sqrt{n}$, for what may be called "nearly equal" sample sizes. Consider a sequence (m_k, n_k) of sample sizes, such that $n_k \rightarrow \infty$ and $m_k - n_k = c_k$ is bounded. With each (m_k, n_k) is associated a boundary L_k with distance variable z_k , phase variable v_k , and significance probability P_{1k} . We assume that z_k is bounded and bounded away from 0, so that α_k will be of exact order $\sqrt{n_k}$. To simplify typography we shall omit the subscript k hereafter.

A straightforward application of Stirling's formula to (4.1) shows that

$$Pr(A \geq a) = \exp \left\{ -\frac{a^2}{n} + \frac{ac}{n} + O\left(\frac{1}{n}\right) \right\} \tag{4.5}$$

when $a = O(\sqrt{n})$. From this it follows that

$$Pr(A = a) = \frac{2a}{n} e^{-a^2/n} + O\left(\frac{1}{n}\right).$$

Table 2. Table of $I(a, y)$.

The superscript is the exponent of that power of 10 by which the entry is to be multiplied.
 $I(a, 0) = 1$; $I(a, 1) = a$; $I(1, y + 1) = I(2, y)$.

y	$a=2$	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
2	5	9	14	20	27	35	44	54	65	77	90	104	119	135	152	170	189
3	14	28	48	75	110	154	208	273	350	440	544	663	798	950	1120	1309	1518
4	42	90	165	275	429	637	910	1260	1700	2244	2907	3705	4655	5775	7084	8602	10351
5	132	297	572	1001	1638	2548	3808	5508	7752	10661	14361	19021	24791	31881	40481	50831	63181
6	429	1001	2002	3640	6188	9996	15501	23261	33921	48281	67301	92091	12402	16442	21532	27852	35632
7	1430	3432	7072	13261	23261	38761	62021	95931	14422	21152	30362	42762	59202	80732	10862	14422	18932
8	4862	11931	25191	48451	87211	14922	24522	38942	60092	90452	13322	19242	27312	38172	52592	71522	96122
9	16801	41991	90441	17762	32692	57202	96142	15622	24672	37992	57232	84542	12274	17534	24684	34304	47074
10	58791	14922	32692	65382	12262	21872	37492	62162	10024	15742	24194	36464	53904	78684	11302	16002	22392
11	20802	53492	11892	24142	46022	83512	14574	24584	40324	64514	10102	15502	23362	34662	50672	73062	10402
12	74292	19322	43462	89482	17304	31874	56454	96774	16132	26232	41722	65092	99752	15042	22352	32752	47392
13	26742	70202	15974	33274	65134	12162	21832	37962	64192	10592	17102	27072	42122	64462	97212	14462	
14	96952	25664	58934	12412	24562	46402	84362	14862	25452	42552	69622	11172	17622	27342	41802		
15	35364	94294	21832	46402	92802	17722	32572	58022	10062	17022	28192	45812	73152	11502			
16	12962	34802	81202	17402	35122	67692	12572	22632	39652	67842	11362	18682	30182				
17	47762	12902	30302	65412	13312	25882	48512	88152	15602	26962	45642	75822					
18	17672	47972	11342	24652	50532	99042	18722	34322	61282	10692	18272						
19	65642	17902	42552	93082	19212	37982	72252	13352	24052	42322							
20	24472	67022	16012	35222	73152	14542	27892	51942	94262								
21	91482	25162	60382	13352	27892	55792	10772	20202									
22	34312	94682	22822	50712	10652	21422	41622										
23	12902	35722	86432	19292	40722	82342											
24	48622	13512	32802	73512	15592												
25	18372	51172	12472	28052													
26	69532	19422	47472														
27	26372	73852															
28	10022																

On substituting $a = \sqrt{2n} z + vc + 0(1/\sqrt{n})$ into (4.5), we find

$$Pr(A \geq \alpha + c) = \exp \left\{ -2z^2 - (2vc + c) \frac{\sqrt{2}z}{\sqrt{n}} + 0\left(\frac{1}{n}\right) \right\}. \quad (4.6)$$

Similarly

$$Pr(A = \alpha + i) = \frac{2\sqrt{2}z}{\sqrt{n}} e^{-2z^2} + 0\left(\frac{1}{n}\right) \text{ for } i = 0, 1, \dots, c. \quad (4.7)$$

Since the mirrors cut by L have probability of order $1/\sqrt{n}$, it is enough to determine $Pr(Y \leq un | A = \alpha + i)$ to constant order. From (4.2), it can be shown that

$$I(a, y) = 2^{2y+a-1} a \exp \left\{ -\frac{a^2}{4y} + 0\left(\frac{1}{\sqrt{n}}\right) \right\} / \sqrt{\pi y}^{1/2} \quad (4.8)$$

when $y = 0(n)$. When we substitute this into (4.3) and make the usual integral approximation of a sum, we find

$$Pr(Y \leq un | A = \alpha + i) = \Phi \left[\frac{z(2u-1)}{\sqrt{u(1-u)}} \right] + 0(1), \quad (4.9)$$

where Φ is the normal distribution function. On substituting (4.6), (4.7), and (4.9) into (4.4), we obtain

$$P_1^* = \exp \left\{ -2z^2 - \frac{\sqrt{2}z}{\sqrt{n}} + \frac{2\sqrt{2}z}{\sqrt{n}} \sum_{i=0}^{c-1} \Psi_z \left(v + \frac{i}{c} \right) \right\} \quad (4.10)$$

where
$$\Psi_z(u) = \Phi \left[\frac{z(2u-1)}{\sqrt{u(1-u)}} \right] - u.$$

It remains to show that $P_1 - P_1^* = 0(1/n)$ in order to establish (4.10) as an expression for P_1 with error of order $1/n$. The argument will be given for simplicity only in the case $m = n + 2$. A path counted in P_1 but not in P_1^* must (i) reach but not pass the mirror $A = a$ for $x + y < n$, and then (ii) reach but not pass the mirror $A = a + 1$ for $x + y > n$. Let S be the point in which such a path crosses $x + y = n$. For a given value of S , the behavior of the path after S is conditionally independent of its behavior before S . As a consequence of (4.6), each of the events (i) and (ii) has probability of order $1/\sqrt{n}$, uniformly in S .

The expression (4.10) has a number of interesting features. (i) Each of the three terms in the exponent has an interpretation. Thus $-2z^2$ corresponds to the Smirnov approximation \bar{P}_1 . The second term $-\sqrt{2}z/\sqrt{n}$ represents the "average" error of the Smirnov approximation; we note that it is of order $1/\sqrt{n}$, and that it is proportional to z , as was suggested empirically in section 3. Further, it does not depend on c . The third term represents the oscillatory component. This is also of order $1/\sqrt{n}$, which shows that it is not possible to express P_1 to accuracy $1/\sqrt{n}$ in a formula of the type proposed by Korolyuk.

(ii) The function Ψ was derived by Malmquist (1954), using the heuristic argument of Doob (1949), as the limiting conditional distribution of the ordinate of the furthest point, say Y^+ , given the value of Z^+ . Our argument, which is restricted to the case when c is bounded, permits this limiting distribution to be derived rigorously when Z^+ is restricted to an appropriate interval, and also shows that the result is not correct for arbitrary small intervals for Z^+ . We thus have an example of a problem in which the heuristic argument fails unless properly restricted.

(iii) From symmetry we see that $\Psi_z(u) + \Psi_z(1-u) = 0$. It follows that $\sum_{i=0}^{c-1} \Psi_z\left(v + \frac{i}{c}\right) = 0$ for $2cv$ an integer, and that $\int_0^{1/c} \sum_{i=0}^{c-1} \Psi_z\left(v + \frac{i}{c}\right) dv = 0$, so that the oscillatory part averages out to zero over a cycle of the phase variable v . Further, it can be shown from the Euler-Mclaurin formula that $\max_v \sum_{i=0}^{c-1} \Psi_z\left(v + \frac{i}{c}\right) \rightarrow 0$ as $c \rightarrow \infty$, which "explains" the phenomenon, observed in Section 3, of damped oscillations as m is increased. A few values of the maximum oscillation are given in Table 3.

Table 3. $\max_v \sum_{i=0}^{c-1} \Psi_z\left(x + \frac{i}{c}\right)$.

z	$c = 1$	2	3
1	0.13	0.04	0.006
1.5	0.21	0.03	0.013
2	0.26	0.10	0.011

(iv) The formula shows that the error of the Smirnov approximation may be quite large even for substantial values of n , particularly when $c = 1$. For example, when $m = n + 1$ and $z = 1.5$ corresponding to $\bar{P}_2 = 0.022$, the relative error of P_2 at the least favorable phasing is about $3/\sqrt{n}$, so that the sample sizes must be about 900 to assure a 10% relative error. The same goal is achieved at $n = 20$ when the sample sizes are equal.

We conclude with an extension of these results to the two-sided test. As remarked above, in practice the approximation $P_2 = 2P_1$ will usually serve; and since the results for P_2 are considerably more complicated than those for P_1 , we shall give only an outline of methods. Corresponding to the end mirror A , we now have upper and lower end mirrors, say A_1 and A_2 . Instead of (4.1) we now obtain, analogously to (2.4),

$$\binom{m+n}{n} Pr(A_2 \geq a \text{ or } A_1 \leq a - c) = 2 \binom{m+n}{n+a} - \binom{m+n}{n+2a} - \binom{m+n}{n+2a-c} + \dots \tag{4.11}$$

The first term on the right side of (4.4) is thus replaced by (writing $\alpha + c$ for a in (4.11),

$$2 \binom{m+n}{m+\alpha} - \binom{m+n}{n+2\alpha+2c} - \binom{m+n}{n+2\alpha+c} + \dots$$

If we carry out work analogous to that which gave (4.6), we find

$$Pr(A_1 \leq \alpha \text{ or } A_2 \geq \alpha + c) = 2Pr(A \geq \alpha + c) \{1 - [Pr(A \geq \alpha + c)]^2 + \dots\} \quad (4.12)$$

a relation analogous to the formula $\tilde{P}_2 = 2\tilde{P}_1(1 - \tilde{P}_1^2 + \dots)$ of Section 2. If n or P_2 is very large, the second term in (4.12) may be worth using.

A similar extension is possible for the second term on the right side of (4.4). Instead of the function I , we have J defined by $J(a, y) =$ the number of ways to go from $(0, 0)$ to $(a + y, y)$ without previously touching either $A = a$ or $A = -a$. An application of the reflection method yields

$$J(a, y) = I(a, y) - I(3a, y - a) + I(5a, y - 2a) - I(7a, y - 3a) + \dots$$

which enables one to compute an analog to P_1^* . Corresponding to (4.8), $J(a, y)$ has an expansion consisting of the right side of (4.8) multiplied by $1 - 3e^{-2a^2/y} + \dots$.

If we ignore paths influencing P_2 only to order $1/n$, we can classify the paths reaching one or both boundaries into the following sets. (a) Paths reaching $A = \alpha + c$ or $A = -\alpha$ or both; (b_{*i*}) paths reaching $A = \alpha + i$ below its intersection with the lower boundary, but not reaching $A = \alpha + i + 1$ nor $A = -\alpha$; (c_{*i*}) a set of paths analogous to b_{*i*}, but with the boundaries interchanged. The set (b_{*i*}) will contribute to P_2 a quantity analogous to (4.9), but with (4.9) replaced by a series whose first term is (4.9). The remaining terms will be of order $1/\sqrt{n}$, but in practice they are usually negligibly small.

5. Conclusion

We conclude with three remarks or conjectures on the nature of the general problem of computing P_1 to order $1/\sqrt{n}$.

(a) It seems likely that the oscillatory behaviour demonstrated above for m/n near 1, will also obtain for m/n near any rational number r/s , where $r > s$ and r and s are relatively prime. We shall give an heuristic argument which also yields a quantitative formula for the oscillation. If the latter is correct, the oscillation will be of practical interest in only a few cases.

Suppose m is "near" $(r/s)n$. Analogously to the end mirrors of section 4, we define the *end line* of a path as the line $r(y - y^+) = s(x - x^+)$, and again index these lines by their x -intercept, a . Now sa takes on consecutive integral values, and consecutive z values differ by $1/\sqrt{nr(r+s)}$. This suggests that the limiting probability of an end line is $4ze^{-2z^2/\sqrt{nr(r+s)}}$. This is $\sqrt{2}/r(r+s)$ times its value when $m = n$.

We define the *end point* of a path as the lowest (highest) point which the path reaches on its end line, when $m > (<) rn/s$. Malmquist's heuristic argument suggests that the limiting conditional distribution of the ordinate of the end point, given the end line, is again given by (4.9); in fact, any other form for this distribution would not be reconcilable with the result of Malmquist cited above.

Finally, we notice that the number of end lines cut by the boundary is, in the limit, (i) cs when $m = rn/s \pm c$, or (ii) cr when $n = sm/r \pm c$. Combining, we should find the asymptotic magnitude of the oscillation to be $\sqrt{2}/r(r+s)$ times as great as it is when $m = n + cs$ in case (i), or $m = n + cr$ in case (ii). In particular, the oscilla-

tion at $(n = 12, m = 17)$ should be about $\sqrt{2/15} = 0.37$ times as great as that at $(n = 12, m = 14)$. The empirical results of Section 3 are in reasonable agreement with this prediction.

A comparison of the conclusion of this heuristic argument with Table 3 suggests that the oscillatory aspect of P_1 is likely to be of practical interest only when m and n are nearly equal, and to a lesser extent when m/n is near $\frac{3}{2}$. Thus, for m/n near $\frac{4}{3}$, the oscillations should be at worst only $\sqrt{\frac{2}{4}}(4 + 3) = 0.27$ as great as at $m = n + 3$; or for z near 1.5, they should contribute at worst $0.015/\sqrt{n}$ to the relative error. Even near $\frac{3}{2}$, the maximum oscillation is at worst about $0.05/\sqrt{n}$.

(b) We now ignore the oscillations, and examine the "average" value of P_1 , say \bar{P}_1 , for m and n nearly equal. From (4.10) we have

$$\bar{P}_1 = \exp \left\{ -2z^2 - \frac{\sqrt{2}z}{\sqrt{n}} + o\left(\frac{1}{n}\right) \right\} \quad (5.1)$$

for any $c > 0$. On the other hand, for $c = 0$, we have from (2.7)

$$P_1 = \bar{P}_1 \exp \left\{ -2z^2 + o\left(\frac{1}{n}\right) \right\}.$$

We are confronted with a disconcerting discontinuity in \bar{P}_1 as a function of c . It is of interest to note that the two formulae are reconciled if a continuity correction is used.

In fact, for $c > 0$, the probabilities of individual z -values are of order $1/n^{3/2}$, so that (5.1) continues to hold when a continuity correction is employed. At $c = 0$, on the other hand, $Pr(Z^+ = z) = (2\sqrt{2}z/\sqrt{n})e^{-2z^2}$. The use of a continuity correction implies increasing the estimate by half of this amount, so that (5.1) used with the correction will yield, to order $1/n$, the value \bar{P}_1 . (Incidentally, the oscillations for $m = n + 1$ are made considerably more regular when the continuity correction is used.)

(c) Since (5.1) holds for any fixed c , it is tempting to conjecture that this formula may be used for arbitrary values of m and n . This conjecture, however, is easily disproved by considering the case $m = np$, where p is an integer. By use of Stirling's formula it can be shown that the y th term of the sum given in Section 2 for $N(a)$ is

$$\frac{a\sqrt{n}}{\sqrt{2\pi p(p+1)y^3(n-y)}} \exp \left\{ -\frac{na^2}{2p(p+1)y(n-y)} - \frac{na}{2(p+1)y(n-y)} + \frac{(2y-n)a}{2py(n-y)} + \frac{(2p+1)n(n-2y)a^3}{6[p(p+1)y(n-y)]^2} + o\left(\frac{1}{n}\right) \right\}$$

when y and $n - y$ are both of order n . If we make the substitutions $y = vn$ and $a = \sqrt{p(p+1)}z$, we find

$$P_1 = \frac{z}{\sqrt{2\pi v^3(1-v)}} \int_0^1 \exp \left(-\frac{z^2}{2v(1-v)} \right) \left\{ 1 - \frac{z}{2\sqrt{p(p+1)}n} \left[\frac{p}{v(1-v)} - \frac{(p+1)(2v-1)}{v(1-v)} - \frac{(2p+1)(1-2v)z^2}{3v^2(1-v)^2} \right] \right\} dv.$$

The various integrals may be reduced by substituting $1 - 2v = u$ when $0 < v < \frac{1}{2}$ and $2v - 1 = u$ when $\frac{1}{2} < v < 1$, adding the integrands, and transforming to the gamma integral. We find

$$P_1 = \exp \left\{ -2z^2 - \frac{2(p-1)}{3\sqrt{p(p+1)}} \frac{z}{\sqrt{n}} + o\left(\frac{1}{n}\right) \right\} \quad (5.2)$$

which agrees with the appropriate special case of Corollary 2 in Blackman (1956) after making three changes, apparently misprints, in the later.

When a continuity correction is employed, formula (5.2) becomes

$$P_1 = \exp \left\{ -2z^2 - \frac{2z}{3} \frac{m+2n}{\sqrt{mn(m+n)}} + o\left(\frac{1}{n}\right) \right\}. \quad (5.3)$$

This simple formula is thus shown to be correct for $m = n$, for m an integral multiple of n , and "on the average" for m nearly equal to n . On the basis of a limited numerical investigation, it appears to work reasonably well in other cases, and can perhaps be recommended as a general interpolation formula.

University of California, Berkeley, Cal., U.S.A.

REFERENCES

- BLACKMAN, JEROME (1956), "An extension of the Kolmogorov distribution". *Ann. Math. Statistics* 27, 513-520.
- (1957), "Correction to 'An extension of the Kolmogorov distribution'", unpublished.
- DONSKER, MONROE D. (1952), "Justification and extension of Doob's heuristic approach to the Kolmogorov-Smirnov theorems". *Ann. Math. Statistics* 23, 277-281.
- DOOB, J. L. (1949), "Heuristic approach to the Kolmogorov-Smirnov theorems". *Ann. Math. Statistics* 20, 393-403.
- DRION, E. F. (1952), "Some distribution-free tests for the difference between two empirical cumulative distribution functions". *Ann. Math. Statistics* 23, 563-574.
- FELLER, W. (1948), "On the Kolmogorov-Smirnov limit theorems for empirical distributions". *Ann. Math. Statistics* 19, 177-189.
- GNEDENKO, B. V. (1952), "Some results on the maximum discrepancy between two empirical distributions". *Doklady Akad. Nauk. SSSR* 82, 661-663.
- GNEDENKO, B. (1954), "Kriterien für die Unveränderlichkeit der Wahrscheinlichkeitsverteilung von zwei unabhängigen Stichprobenreihen". *Math. Nachr.* 12, 29-66.
- GNEDENKO, B. V., and KOROLYUK, V. S. (1951), "On the maximum discrepancy between two empirical distributions". *Doklady Akad. Nauk SSSR* 80, 525-528.
- GNEDENKO, B. V., and RVAČEVA, E. L. (1952), "On a problem of comparison of two empirical distributions". *Doklady Akad. Nauk SSSR* 82, 513-516.
- GNEDENKO, B. V., and STUDNEV, YU. P. (1952), "Comparison of the effectiveness of several methods of testing homogeneity of statistical material". *Dopovidi Akad. Nauk Ukrain. RSR*, pp. 359-363.
- KOLMOGOROFF, A. (1933), "Sulla determinazione empirica di una legge di distribuzione". *Giorn. Ist. Ital. Attuari* 4, 83-91.
- KOLMOGOROV, A. N., and HINČIN, A. YA (1951), "The work of N. V. Smirnov on the investigation of properties of variational series and on nonparametric problems of mathematical statistics". *Uspehi Matem. Nauk* 6, no. 4 (44), p. 190-192.
- KOROLYUK, V. S. (1954), "Asymptotic expansions for A. N. Kolmogorov's and N. V. Smirnov's criteria of fit." *Doklady Akad. Nauk SSSR* 95, 443-446.
- (1955 a), "On the discrepancy of empiric distributions for the case of two independent samples". *Izvestiya Akad. Nauk SSSR. Ser. Mat.* 19, 81-96.
- (1955 b), "Asymptotic expansions for the criteria of fit of A. N. Kolmogorov and N. V. Smirnov". *Izvestiya Akad. Nauk SSSR. Ser. Mat.* 19, 103-124.

J. L. HODGES, JR., *The Smirnov two-sample test*

- KOROLYUK, V. S., and YAROŠEVSKII, B. I. (1951), "Study of the maximum discrepancy of two empirical distributions". *Dopovidi Akad. Nauk Ukrain RSR*, pp. 243-247.
- MALMQUIST, STEN (1954), "On certain confidence contours for distribution functions", *Ann. Math. Statistics* 25, 523-533.
- MASSEY, F. J. JR. (1950), "A note on the power of a non-parametric test". *Ann. Math. Statistics* 21, 440-443.
- MASSEY, FRANK J., JR. (1951), "The distribution of the maximum deviation between two sample cumulative step functions". *Ann. Math. Statistics* 22, 125-128.
- (1952), "Distribution table for the deviation between two sample cumulatives". *Ann. Math. Statistics* 23, 435-441.
- POLYA, GEORGE (1948), "Exact formulas in the sequential analysis of attributes". *Univ. California Publ. Math.* 1, 229-239.
- RVAČEVA, E. L. (1952), "On the maximum discrepancy between two empirical distributions". *Ukrain. Mat. Žurnal* 4, 373-392.
- SMIRNOFF, N. (1939), "Sur les écarts de la courbe de distribution empirique". *Mat. Sbornik* 6 (48), p. 3-26.
- SMIRNOV, N. (1939), "On the estimation of the discrepancy between empirical curves of distribution for two independent samples". *Bull. Math. Univ. Moscou* 2:2.
- (1948), "Table for estimating the goodness of fit of empirical distributions". *Ann. Math. Statistics* 19, 279-281.
- VAN DER WAERDEN, B. L. (1953), "Order tests for the two-sample problem. II". *Indagationes Math.* 15, 303-310.

Tryckt den 17 december 1957

Uppsala 1957. Almqvist & Wiksells Boktryckeri AB