



On a Test Whether Two Samples are from the Same Population

Author(s): A. Wald and J. Wolfowitz

Source: *The Annals of Mathematical Statistics*, Vol. 11, No. 2 (Jun., 1940), pp. 147-162

Published by: Institute of Mathematical Statistics

Stable URL: <https://www.jstor.org/stable/2235872>

Accessed: 09-08-2020 15:59 UTC

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <https://about.jstor.org/terms>



Institute of Mathematical Statistics is collaborating with JSTOR to digitize, preserve and extend access to *The Annals of Mathematical Statistics*

ON A TEST WHETHER TWO SAMPLES ARE FROM THE SAME POPULATION¹

BY A. WALD² AND J. WOLFOWITZ

1. The Problem.³ Let X and Y be two independent stochastic variables about whose cumulative distribution functions nothing is known except that they are continuous. Let x_1, x_2, \dots, x_m be a set of m independent observations on X and let y_1, \dots, y_n be a set of n independent observations on Y . It is desired to test the hypothesis (the null hypothesis) that the distribution functions of X and Y are identical.

An important step in statistical theory was made when "Student" proposed his ratio of mean to standard deviation for a similar purpose. In the problem treated by "Student" the distribution functions were assumed to be of known (normal) form and completely specified by two parameters. It is clear that in the problem to be considered here the distributions cannot be specified by any finite number of parameters.

It might nevertheless be argued that by virtue of the limit theorems of probability theory, "Student's" ratio might be used in our problem for large samples. Such a procedure is open to very serious objections. The population distributions may be of such form (e.g., Cauchy distribution) that the limit theorems do not apply. Furthermore, the distributions of X and Y may be radically different and yet have the same first two moments; clearly "Student's" ratio will not distinguish between two such distributions.

The Pearson contingency coefficient is a useful test specifically designed for the problem we are discussing here, but one which also possesses some disadvantages. The location of the class intervals is to a considerable extent arbitrary. In order to use the χ^2 distribution, the numbers in each class interval must not be small; often this can be done only by having large class intervals, thus entailing a loss of information.

2. Preliminary remarks. Denote by $P\{X < x\}$ the probability of the relation in braces. Let $f(x)$ and $g(x)$ be the distribution functions of X and Y respectively; e.g., $P\{X < x\} = f(x)$. Throughout this paper we shall assume that $f(x)$ and $g(x)$ are continuous.

Let the set of $m + n$ elements x_1, \dots, x_m and y_1, \dots, y_n be arranged in

¹ Presented to the Institute of Mathematical Statistics at Philadelphia, December 27, 1939.

² Research under a grant-in-aid from the Carnegie Corporation of New York.

³ The authors are indebted to Prof. S. S. Wilks for proposing this problem to them.

ascending order of magnitude, and let the sequence be designated by Z , thus: $Z = z_1, z_2, \dots, z_{m+n}$, where $z_1 < z_2 < \dots < z_{m+n}$. ($f(x)$ and $g(x)$ were assumed to be continuous. Hence the probability is 0 that $z_i = z_{i+1}$ and therefore we may exclude this case.) Let $V = v_1, v_2, \dots, v_{m+n}$ be a sequence defined as follows: $v_i = 0$ if z_i is a member of the set x_1, \dots, x_m and $v_i = 1$ if z_i is a member of the set y_1, \dots, y_n . It is easy to show that any statistic S used to test the null hypothesis should be invariant under any continuous, reciprocally one-to-one transformation of the real axis. That is to say, if $t' = \varphi(t)$ is any such transformation, then

$$(1) \quad S(x_1, \dots, x_m, y_1, \dots, y_n) \equiv S(\varphi(x_1), \dots, \varphi(x_m), \varphi(y_1), \dots, \varphi(y_n)).$$

The reason for this requirement on S is the fact that the transformed stochastic variables $X' = \varphi(X)$ and $Y' = \varphi(Y)$ are continuous and have identical distributions if and only if X and Y have identical distributions. Hence S must be a function of V only, with the added restriction that $S(V) = S(V')$, where $V' = v_{m+n}, v_{m+n-1}, \dots, v_1$. For if S were a function of $x_1, \dots, x_m, y_1, \dots, y_n$ which cannot be expressed as a function of V alone, then there exists a continuous reciprocally one-to-one transformation $t' = \varphi(t)$ such that (1) is not true. On the other hand, any continuous reciprocally one-to-one transformation of the entire line into itself is monotonic and hence either leaves V invariant or else transforms it into V' .

3. Previous results. In an interesting paper on this problem W. R. Thompson [1] proceeds as follows: Let the sets x_1, \dots, x_m and y_1, \dots, y_n be ordered in ascending order of magnitude, thus: $x_{p_1}, x_{p_2}, \dots, x_{p_m}$ and $y_{p'_1}, y_{p'_2}, \dots, y_{p'_n}$ where $x_{p_1} < x_{p_2} < \dots < x_{p_m}$ and $y_{p'_1} < y_{p'_2} < \dots < y_{p'_n}$. Let $P\{x_{p_k} < y_{p'_k}\}$ denote the probability of the relation in braces under the null hypothesis ($f(x) \equiv g(x)$). This probability is shown to be independent of $f(x)$ and the relation

$$(2) \quad P\{x_{p_k} < y_{p'_k}\} = \psi(m, n, k, k')$$

holds, where the right member, which is given explicitly by Thompson, is a function only of the arguments exhibited. To make a test of the null hypothesis with, say, a 5% level of significance, this writer proposes to choose k and k' so that $\psi(m, n, k, k') = .05$. The test would then consist of noticing whether $x_{p_k} < y_{p'_k}$ or not. In the former case the null hypothesis is to be considered as disproved.

It is clear that this test cannot be very efficient, ignoring as it does so many of the relations among the observations. Except under certain rather narrow restrictions on the admissible alternatives, for example, that $g(x) \equiv f(x + c)$, where c is an arbitrary constant, the test suffers the further defect of not being "consistent" in a way which will be discussed below. Hence the test suggested by Thompson can scarcely be regarded as a satisfactory solution of the problem. This criticism, of course, does not apply to those sections of Thompson's paper which deal with the question of estimating the so-called normal range.

4. The statistic U . A subsequence $v_{s+1}, v_{s+2}, \dots, v_{s+r}$ of V (where r may also be 1) will be called a "run" if $v_{s+1} = v_{s+2} = \dots = v_{s+r}$ and if $v_s \neq v_{s+1}$ when $s > 0$ and if $v_{s+r} \neq v_{s+r+1}$ when $s + r < m + n$. For example, $V = 1, 0, 0, 1, 1, 0$ contains the following runs: 1; 0, 0; 1, 1; 0. The statistic⁴ U defined as the number of runs in V seems a suitable statistic for testing the hypothesis that $f(x) \equiv g(x)$. In the event that the latter identity holds, the distribution of U is independent of $f(x)$. A difference between $f(x)$ and $g(x)$ tends to decrease U . U is consistent in a sense which will be discussed below.

In order to derive the distribution of U under the null hypothesis, we first note that all the $\frac{(m+n)!}{m!n!}$ ($= {}^{m+n}C_m$) possible sequences V have the same probability $\left(= \frac{m!n!}{(m+n)!} \right)$. To see this, consider the sequence V where $v_i = 0$ ($i = 1, 2, \dots, m$) and $v_i = 1$ ($i = m+1, m+2, \dots, m+n$). Clearly the probability of the sequence is

$$q = \frac{m(m-1) \dots 1 \cdot n(n-1) \dots 1}{(m+n)(m+n-1) \dots (n+1)n(n-1) \dots 1}.$$

Furthermore, the probability of any other sequence is equal to the product of the factors in the numerator of q taken in a different order, divided by the product of the factors in the denominator taken in the same order. The quotient is, of course, $= q$.

Let e_0 be the number of runs in V whose elements are 0 and let e_1 be the number of runs whose elements are 1. Obviously $U = e_0 + e_1$. Let the runs of each kind be arranged in the ascending order of the indices of the v_i . Let r_{0j} be the number of elements 0 in the j^{th} run of that kind ($j = 1, 2, \dots, e_0$) and let $r_{1j'}$ be the number of elements 1 in the j'^{th} run of that kind ($j' = 1, 2, \dots, e_1$). The following relations obviously hold:

$$(3) \quad \sum_{j=1}^{e_0} r_{0j} = m,$$

$$(4) \quad \sum_{j'=1}^{e_1} r_{1j'} = n,$$

$$(5) \quad 1 \leq e_0 \leq m, \quad 1 \leq e_1 \leq n,$$

$$(6) \quad |e_0 - e_1| \leq 1.$$

⁴When this paper was already in proof, our attention was called to a paper by W. L. Stevens, entitled "Distribution of groups in a sequence of alternatives," *Annals of Eugenics*, Vol. 9 (1939). There a statistic, which is essentially the U statistic, is proposed for a problem different from that considered by us and the distribution of U is obtained in a different manner. However, the application of the U statistic for the purpose herein described, the proof of consistency and the other results of our paper are not contained in it.

Hence if $U = 2k$, then $e_0 = e_1 = k$, and if $U = 2k - 1$, then either $e_0 = k$, $e_1 = k - 1$ or $e_0 = k - 1$, $e_1 = k$. The element v_1 of V together with the numbers $r_{01}, r_{02}, \dots, r_{0e_0}, r_{11}, r_{12}, \dots, r_{1e_1}$, completely determines the sequence V whose probability is q .

Without loss of generality we may assume that $m \leq n$. If $U = 2k$, $1 \leq k \leq m$, $v_1 = 0$, any two sequences of k positive numbers each may constitute a sequence of $r_{01}, \dots, r_{0e_0}, r_{11}, \dots, r_{1e_1}$ provided only that (3) and (4) are satisfied. The number of sequences $r_{01}, r_{02}, \dots, r_{0k}$ which satisfy (3) is the coefficient of a^m in the purely formal expansion of

$$(a + a^2 + a^3 + \dots)^k = \left(\frac{a}{1 - a} \right)^k$$

and hence is ${}^{m-1}C_{k-1}$. Similarly the number of sequences $r_{11}, r_{12}, \dots, r_{1k}$ which satisfy (4) is found to be ${}^{n-1}C_{k-1}$. Bearing in mind the case $U = 2k$, $v_1 = 1$, we obtain

$$(7) \quad P\{U = 2k\} = \frac{2({}^{m-1}C_{k-1} \cdot {}^{n-1}C_{k-1})}{m+nC_m}, \quad (k = 1, 2, \dots, m),$$

where the left member denotes the probability of the relation in braces under the null hypothesis. In a similar manner we obtain

$$(8) \quad P\{U = 2k - 1\} = \frac{({}^{m-1}C_{k-1} \cdot {}^{n-1}C_{k-2} + {}^{m-1}C_{k-2} \cdot {}^{n-1}C_{k-1})}{m+nC_m},$$

$$(k = 2, \dots, m + 1),$$

with the proviso that ${}^aC_b = 0$ if $a < b$.

We shall now briefly indicate a method of obtaining the mean $E(U)$ and variance $\sigma^2(U)$ of U . For example, $E(U)$ may be obtained by performing several summations of the type

$$(9) \quad \sum_{i=0}^{m-1} i \cdot {}^{m-1}C_i \cdot {}^{n-1}C_i.$$

It is easy to verify that the expression (9) is the term free of a in the purely formal expansion in a of:

$$(10) \quad (m-1) \cdot (1+a)^{m-2} \cdot a \cdot \left(1 + \frac{1}{a}\right)^{n-1},$$

and hence is

$$(11) \quad (m-1) \cdot {}^{m+n-3}C_{n-2}.$$

The other summations required for the mean and variance can be carried out in a similar manner. We shall omit these tedious calculations. The results are:

$$(12) \quad E(U) = \frac{2mn}{m+n} + 1,$$

$$(13) \quad \sigma^2(U) = \frac{2mn(2mn - m - n)}{(m+n)^2(m+n-1)}.$$

The critical region for testing the null hypothesis on a level of significance β is given by the inequality $U < u_0$, where u_0 is a function of m and n such that $P\{U < u_0\} = \beta$.

5. The asymptotic distribution of U . Let $m/n = \alpha$, a positive constant. Then, as $m \rightarrow \infty$,

$$E(U) \sim \frac{2m}{1+\alpha},$$

$$\sigma^2(U) \sim \frac{4\alpha m}{(1+\alpha)^3}.$$

THEOREM I. *If t is any real number, the probability of the relation $U < \frac{2m}{1+\alpha} + 2\left[\frac{\alpha m}{(1+\alpha)^3}\right]^{\frac{1}{2}}t$ converges uniformly in t to*

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-\frac{1}{2}w^2} dw$$

as $m \rightarrow \infty$.

The proof of this theorem is essentially the same as the classical proof that the binomial law converges to the normal distribution (see, for example, Fréchet [2], p. 89) and it will be unnecessary to give the details. Since the asymptotic distribution of the subpopulation of even U is the same as that of odd U , it will be sufficient to consider only the right member of (7). Let $m' = m - 1$, $n' = n - 1$, and $k' = k - 1$. We make the substitution

$$(14) \quad w = \frac{k' - \frac{m'}{1+\alpha'}}{\sqrt{m'}}, \quad \text{where} \quad \alpha' = \frac{m'}{n'},$$

$$(15) \quad dw = \frac{1}{\sqrt{m'}},$$

and evaluate the factorials by Stirling's formula. We shall give here only the results of successive simplifications. At each step we shall omit the factors free of k or w , since their product may be reconstructed from the final exponential form. Thus instead of the right member of (7) we can consider the expression:

$$(16) \quad m^{-1}C_{k-1} \cdot n^{-1}C_{k-1}.$$

Omitting factors free of k , we get

$$(17) \quad \frac{1}{(k-1)!(m-k)!(k-1)!(n-k)!}$$

and by Stirling's formula, since k and m are both large:

$$(18) \quad \frac{1}{k^{(2k'+1)}(m'-k')^{(m'-k'+\frac{1}{2})}(n'-k')^{(n'-k'+\frac{1}{2})}}.$$

Now apply (14). We obtain

$$(19) \quad \left(\sqrt{m'}w + \frac{m'}{1+\alpha'}\right)^{-2\sqrt{m'}w - \frac{2m'}{1+\alpha'} - 1} \cdot \left(-\sqrt{m'}w + \frac{m'\alpha'}{1+\alpha'}\right)^{\sqrt{m'}w - \frac{m'\alpha'}{1+\alpha'} - \frac{1}{2}} \\ \cdot \left(-\sqrt{m'}w + \frac{m'}{\alpha'(1+\alpha')}\right)^{\sqrt{m'}w - \frac{m'}{\alpha'(1+\alpha')} - \frac{1}{2}}.$$

Dividing inside the parentheses by $\frac{m'}{1+\alpha'}$, $\frac{m'\alpha'}{1+\alpha'}$, $\frac{m'}{\alpha'(1+\alpha')}$, respectively, and again omitting factors free of w , we get

$$(20) \quad \left(1 + \frac{(1+\alpha')w}{\sqrt{m'}}\right)^{-2\sqrt{m'}w - \frac{2m'}{1+\alpha'} - 1} \cdot \left(1 - \frac{(1+\alpha')w}{\alpha'\sqrt{m'}}\right)^{\sqrt{m'}w - \frac{m'\alpha'}{1+\alpha'} - \frac{1}{2}} \\ \cdot \left(1 - \frac{\alpha'(1+\alpha')w}{\sqrt{m'}}\right)^{\sqrt{m'}w - \frac{m'}{\alpha'(1+\alpha')} - \frac{1}{2}}.$$

Taking logarithms, expanding in powers of $\frac{w}{\sqrt{m'}}$ and neglecting terms in $\frac{w^3}{m'}$ and higher orders, the results are

$$(21) \quad -\left(2\sqrt{m'}w + \frac{2m'}{1+\alpha'} + 1\right)\left(\frac{(1+\alpha')w}{\sqrt{m'}} - \frac{(1+\alpha')^2w^2}{2m'}\right) \\ -\left(\sqrt{m'}w - \frac{m'\alpha'}{1+\alpha'} - \frac{1}{2}\right)\left(\frac{(1+\alpha')w}{\alpha'\sqrt{m'}} + \frac{(1+\alpha')^2w^2}{2\alpha'^2m'}\right) \\ -\left(\sqrt{m'}w - \frac{m'}{\alpha'(1+\alpha')} - \frac{1}{2}\right)\left(\frac{\alpha'(1+\alpha')w}{\sqrt{m'}} + \frac{\alpha'^2(1+\alpha')^2w^2}{2m'}\right)$$

which equals

$$(22) \quad -\frac{w^2(1+\alpha')^3}{2\alpha'} + O(m'^{-\frac{1}{2}}).$$

The proof of the fact that the distribution of w converges uniformly to the normal distribution with zero mean and variance $\frac{\alpha'}{(1+\alpha')^3}$ can be carried out in the same way as the classical proof that the binomial law converges to the normal distribution.

It is obvious that

$$w^* = \frac{k - \frac{m}{1 + \alpha}}{\sqrt{m}}$$

has the same distribution as w . From this and from the fact that $U = 2k$ or $2k - 1$ THEOREM I follows.

In using conventional tables of the Gaussian function to make tests of significance on U when m and n are large, the reader is urged not to forget that the critical region of U lies in only one tail of the curve.

6. An example. We give here a simple example illustrating the use of the statistic U and THEOREM I.

Suppose 50 observations were made on X and 50 observations on Y . Suppose further that these observations are arranged in ascending order and that the i^{th} element of this sequence is said to have the rank i . The observations on X occupy the following ranks: 1, 5, 6, 7, 12, 13, 14, 15, 16, 17, 19, 20, 21, 25, 26, 27, 28, 31, 32, 38, 42, 43, 44, 45, 50, 51, 52, 53, 54, 56, 57, 58, 62, 63, 64, 65, 68, 69, 75, 79, 80, 81, 86, 87, 89, 90, 91, 93, 94, 95.

The observations on Y occupy the remaining ranks.

In this case, $U = 34$.

For $m = n = 50$,

$$E(U) = 51,$$

$$\sigma^2(U) = 24.747.$$

The probability of getting 34 runs or less when the distribution functions of X and Y are continuous and identical is therefore less than $5 \cdot 10^{-4}$.

7. Consistency. We shall say that a test is "consistent" if the probability of rejecting the null hypothesis when it is false (i.e., the complement of the probability of a type II error, cf. Neyman and Pearson, [3]) approaches one as the sample number approaches infinity. In the literature of statistics a function of the observations which converges stochastically to a population parameter as the sample number approaches infinity, is called a "consistent" statistic. If a test of a hypothesis about a population parameter is made by a proper use of a consistent (statistic) estimate of the parameter, the test will be consistent also according to our definition, which thus furnishes an extension of the idea of consistency to the case where the alternatives to the null hypothesis cannot be specified by a finite number of parameters.

It is obvious that consistency ought to be a minimal requirement of any good test. It is the purpose of this section to prove that, subject to some slight and from the practical statistical point of view, unimportant, restrictions on the distribution functions, the test furnished by the statistic U is consistent.

We shall say that the distribution functions $f(x)$ and $g(x)$ satisfy the condition A , if, for any arbitrarily small positive δ , there exist a finite number of

closed intervals, such that the probability of the sum I of these intervals is $> 1 - \delta$ according to at least one of the distribution functions $f(x)$ and $g(x)$, and such that $f(x)$ and $g(x)$ have positive continuous derivatives $f'(x)$ and $g'(x)$ in I .

In all that follows, although m and n are considered as variables, their ratio m/n is to be a constant, denoted by α . Let $\beta > 0$ denote the level of significance on which the test is to be made, so that, if $f(x) \equiv g(x)$,

$$(23) \quad P\{U < u_0(m)\} = \beta$$

where the critical region for two samples of size m and n , respectively, is given by

$$U < u_0(m).$$

THEOREM II. *If $f(x)$ and $g(x)$ satisfy condition A, and if*

$$(24) \quad f(x) \not\equiv g(x),$$

then

$$(25) \quad \lim_{m \rightarrow \infty} P\{U < u_0(m)\} = 1.$$

The proof of this theorem will be given in several stages.

Let $E\left(\frac{U}{m}; f; g\right)$ and $\sigma^2\left(\frac{U}{m}; f; g\right)$ denote the mean and variance, respectively, of $\frac{U}{m}$, when X and Y have the distribution functions $f(x)$ and $g(x)$, respectively, and the sample numbers are m and n . Let the set $x_1 \dots x_m; y_1 \dots y_n$ be arranged in ascending order of magnitude, thus:

$$(26) \quad Z = z_1, z_2, \dots, z_{m+n},$$

where $z_1 < z_2 < \dots < z_{m+n}$. The sequence

$$(27) \quad V = v_1, v_2, \dots, v_{m+n}$$

is defined as follows: $v_i = 0$ if z_i is a member of the set $x_1 \dots x_m$ and $v_i = 1$ if z_i is a member of the set $y_1 \dots y_n$.

LEMMA 1. *If the following are fulfilled:*

- | | |
|----|--|
| a) | $f(x) \equiv 0 \quad x < 0,$ |
| | $f(x) \equiv x \quad 0 \leq x \leq 1,$ |
| | $f(x) \equiv 1 \quad x > 1.$ |
| b) | $g(x) \equiv 0 \quad x \leq 0,$ |
| | $g(x) \equiv 1 \quad x \geq 1.$ |

c) *The derivative $g'(x)$ of $g(x)$ exists, is continuous and positive everywhere in the interval $0 \leq x \leq 1$.*

d) k is an arbitrary but fixed positive integer. For every m , $i_{1m} < i_{2m} < \dots < i_{km}$ are a set of k positive integers subject only to the restriction that the least upper bound γ of the sequence $\frac{i_{km}}{m+n}$ is less than 1.

Then the expected value

$$E\left(\prod_{j=1}^k v_{i_{jm}}\right) \text{ of } \prod_{j=1}^k v_{i_{jm}}$$

satisfies the inequality

$$(28) \quad \left| E\left(\prod_{j=1}^k v_{i_{jm}}\right) - \prod_{j=1}^k \frac{g'(a_{\lambda_{jm}})}{\alpha + g'(a_{\lambda_{jm}})} \right| < \varphi(m)$$

where $\lambda_{jm} = \frac{i_{jm}}{m+n}$ and $a_{\lambda_{jm}}$ ($j = 1 \dots k$) is the root of

$$(29) \quad ma_{\lambda_{jm}} + ng(a_{\lambda_{jm}}) = \lambda_{jm}(m+n)$$

and $\varphi(m)$ depends only on m and is such that

$$(30) \quad \lim_{m \rightarrow \infty} \varphi(m) = 0.$$

It is easy to verify that the root $a_{\lambda_{jm}}$ of (29) exists and is unique.

PROOF: It will be sufficient to show that, for any specified set of values of

$$v_{i_{1m}} \dots v_{i_{(r-1)m}}, \quad v_{i_{(r+1)m}} \dots v_{i_{km}} \quad (r = 1 \dots k)$$

the conditional probability $P\{v_{i_{rm}} = 1\}$ of the relation in braces satisfies the inequality

$$(31) \quad \left| \frac{g'(a_{\lambda_{rm}})}{\alpha + g'(a_{\lambda_{rm}})} - P\{v_{i_{rm}} = 1\} \right| < \psi(m),$$

where $\psi(m)$ depends only on m and is such that

$$(32) \quad \lim_{m \rightarrow 0} \psi(m) = 0.$$

For each m let

$$(33) \quad V'_m = v'_{i_{1m}}, v'_{i_{2m}} \dots v'_{i_{(r-1)m}}, v'_{i_{(r+1)m}} \dots v'_{i_{km}}$$

be a fixed sequence whose elements are either 0 or 1. We shall consider the conditional probability $P\{v_{i_{rm}} = s\}$, ($s = 0, 1$) of the relation in braces subject to the condition that

$$(34) \quad v_{i_{jm}} = v'_{i_{jm}}, \quad (j = 1, 2, \dots (r-1), (r+1), (r+2), \dots k).$$

Let a and b be two numbers such that $0 < a < b < 1$, and let m^* be a non-negative integer such that $m^* \leq m$, and $m^* \leq [\gamma(m+n)]$ where $[\gamma(m+n)]$ denotes the largest integer $\leq \gamma(m+n)$. Let $Q_m(a, b, m^*)$ denote the proba-

bility that, if m^* observations are made on X and $[\gamma(m+n)] - m^*$ observations are made on Y , the following conditions will be fulfilled:

- (a) the total number of observations $< a$ is exactly $i_{rm} - 1$
- (b) all observations are $< b$
- (c) if the $[\gamma(m+n)]$ observations are arranged in ascending order and if $v_j^* = 0$ or 1 according as the j^{th} element is an observation on X or on Y , then

$$(35) \quad v_{i_{jm}}^* = v'_{i_{jm}} \quad (j = 1, 2, \dots, r-1),$$

and

$$(36) \quad v_{i_{jm-1}}^* = v'_{i_{jm}} \quad (j = r+1, r+2, \dots, k).$$

It is easy to see that the probability P_0 of the simultaneous fulfillment of the relations (34) and of $v_{i_{rm}} = 0$ is given by

$$(37) \quad P_0 = \int_0^1 \int_0^b \sum_{m^*} R_m(a, b, m^*) m' (1-b)^{m'-1} (1-g(b))^{n'} da db,$$

where

$$(38) \quad R_m(a, b, m^*) = {}^m C_{m^*} {}^n C_{[\gamma(m+n)]-m^*} \frac{dQ_m}{db}(a, b, m^*),$$

$$(39) \quad m' = m - m^*,$$

and

$$(40) \quad n' = n - [\gamma(m+n)] + m^*.$$

Similarly, the probability P_1 of the simultaneous fulfillment of the relations (34) and of $v_{i_{rm}} = 1$ is given by

$$(41) \quad P_1 = \int_0^1 \int_0^b \sum_{m^*} R_m(a, b, m^*) n' g'(a) (1-b)^{m'} (1-g(b))^{n'-1} da db.$$

Then

$$(42) \quad \frac{P\{v_{i_{rm}} = 0\}}{P\{v_{i_{rm}} = 1\}} = \frac{P_0}{P_1}.$$

Let $n_0 = \sum_{i > [\gamma(m+n)]} v_i$ and $m_0 = m + n - [\gamma(m+n)] - n_0$. The variables $(z_{i_{rm}} - a_{\lambda_{rm}})$, $(z_{[\gamma(m+n)]} - a_\gamma)$, $\left(\frac{m_0}{n_0} - \frac{\alpha(1-a_\gamma)}{(1-g(a_\gamma))}\right)$ all converge stochastically to zero.

Let $P_0(\epsilon)$ and $P_1(\epsilon)$ denote the values of the right members of (37) and (41), respectively, if the integration is restricted to the region where $a \leq b$, $|a - a_{\lambda_{rm}}| < \epsilon$, $|b - a_\gamma| < \epsilon$ and the summation is restricted to those values

of m^* for which $\left| \frac{m'}{n'} - \frac{\alpha(1-a_\gamma)}{(1-g(a_\gamma))} \right| < \epsilon$. Hence, because of the aforementioned stochastic convergence, for all sufficiently large m

$$(43) \quad |P_s(\epsilon) - P_s| < \epsilon \quad s = 1, 2.$$

Since $P_s > 0$, for sufficiently large m , also

$$(44) \quad \left| \frac{P_0(\epsilon)}{P_1(\epsilon)} - \frac{P_0}{P_1} \right| < \epsilon.$$

Since $g(x)$ and $g'(x)$ are continuous in the interval $[0, 1]$ and hence uniformly continuous, it is clear that

$$(45) \quad \left| \frac{P_0(\epsilon)}{P_1(\epsilon)} - \frac{\alpha}{g'(a_{\lambda_{rm}})} \right| < c\epsilon,$$

where c is a fixed constant independent of m . From (44) and (45) it follows easily that, for any arbitrarily small ϵ' ,

$$(46) \quad \left| \frac{P_0}{P_1} - \frac{\alpha}{g'(a_{\lambda_{rm}})} \right| < \epsilon'$$

for sufficiently large m .

Since $P\{v_{i,rm} = 1\} = \frac{P_1}{P_0 + P_1}$, the required relation (31) follows. This completes the proof of LEMMA 1.

LEMMA 2. *If conditions a, b, and c of Lemma 1 are satisfied, then*

$$(47) \quad \lim_{m \rightarrow \infty} E \left(\frac{U}{m}; f; g \right) = 2 \int_0^1 \frac{g'(x)}{\alpha + g'(x)} dx$$

and

$$(48) \quad \lim_{m \rightarrow \infty} \sigma^2 \left(\frac{U}{m}; f; g \right) = 0.$$

PROOF: Since

$$(49) \quad \begin{aligned} \frac{U}{m} &= \frac{1}{m} + \frac{1}{m} \sum_{j=2}^{m+n} (v_j - v_{j-1})^2 \\ &= \frac{1 + v_1 + v_{m+n}}{m} + \frac{2}{m} \sum_{j=2}^{m+n-1} v_j - \frac{2}{m} \sum_{j=2}^{m+n} v_{j-1} v_j, \end{aligned}$$

we have from LEMMA 1,

$$(50) \quad \begin{aligned} E \left(\frac{U}{m} \right) &= \frac{2}{m} \left[\sum_i \frac{g'(a_{jm})}{\alpha + g'(a_{jm})} - \sum_i \left(\frac{g'(a_{jm})}{\alpha + g'(a_{jm})} \right)^2 \right] + \eta(m) + \eta^*(\gamma) \\ &= \frac{2}{m} \sum \left[\frac{\alpha g'(a_{jm})}{(\alpha + g'(a_{jm}))^2} \right] + \eta(m) + \eta^*(\gamma), \end{aligned}$$

where

$$(51) \quad \lim_{m \rightarrow \infty} \eta(m) = \lim_{\gamma \rightarrow 1} \eta^*(\gamma) = 0$$

and a_{jm} is the root of the equation

$$(52) \quad ma_{jm} + ng(a_{jm}) = j \quad (j = 2 \dots m + n).$$

From equation (52) it follows that

$$(53) \quad \lim_{m \rightarrow \infty} (a_{jm} - a_{(j-1)m})(m + ng'(a_{jm})) = 1$$

uniformly in j . Since γ may be chosen arbitrarily near to 1, the required result (47) follows easily from (50).

It remains to consider the variance of $\frac{U}{m}$. The expression

$$\frac{1 + v_1 + v_{m+n}}{m} + \frac{2}{m} \sum_{j=2}^{m+n-1} v_j$$

differs from $\frac{2}{\alpha}$ by at most $\frac{1}{m}$, so that its variance converges to zero with $m \rightarrow \infty$.

In order to prove (48), it will be sufficient to show that the variance of

$$(54) \quad W = \frac{1}{m} \sum_{j=2}^{m+n} v_{j-1} v_j$$

goes to zero with increasing m . From LEMMA 1 it follows that

$$(55) \quad -z(m) < [E(v_i v_j v_k v_l) - E(v_i v_j) E(v_k v_l)] < z(m),$$

where $\lim_{m \rightarrow \infty} |z(m)| = 0$, provided only that the integers i, j, k, l are distinct and $< \gamma(m + n)$. The variance of mW is the sum of terms of the type occurring in (55). The number of terms for which i, j, k, l are distinct is of the order m^2 . All other terms are of size at most 2 and their number is of the order m . Since the number γ may be chosen arbitrarily near to 1, the variance of W converges to zero with $m \rightarrow \infty$.

This proves LEMMA 2.

LEMMA 3. *If conditions a, b, and c of Lemma 1 are fulfilled, and if (24) holds, then*

$$(56) \quad T = \int_0^1 \frac{g'(x)}{\alpha + g'(x)} dx < \frac{1}{1 + \alpha}.$$

Let $a_1 < a_3$ be any two real numbers and designate $\frac{a_1 + a_3}{2}$ by a_2 . Let $F(x)$ be defined as follows:

$$(57) \quad \begin{aligned} F(a_1) &= 0, \\ F(x) &= (x - a_i)b_i + F(a_i), \quad (a_i \leq x \leq a_{i+1}; i = 1, 2). \end{aligned}$$

Let c be defined by

$$(58) \quad F(a_3) = c(a_3 - a_1).$$

Then it is easy to verify that the maximum of

$$(59) \quad T^* = \int_{a_1}^{a_3} \frac{F'(x)}{\alpha + F'(x)} dx$$

with respect to b_1 and b_2 , subject to the restrictions that b_1 and b_2 be non-negative, and that a_1 , a_3 and c be fixed ($c > 0$), occurs when and only when

$$(60) \quad b_1 = b_2 = c.$$

Now define

$$(61) \quad \begin{aligned} P_{ij} &= \frac{i}{2^j}, & P_{0j} &= 0, \\ l_{ij} &= \frac{g(P_{ij}) - g(P_{(i-1)j})}{2^j} \end{aligned}$$

and

$$S_j = \frac{1}{2^j} \sum_{i=1}^{2^j} \frac{l_{ij}}{\alpha + l_{ij}}, \quad (i = 1, 2, \dots, 2^j; j = 0, 1, 2, \dots).$$

Repeated application of the result of the previous paragraph easily gives

$$(62) \quad S_j \geq S_{j+1}.$$

From (24) it follows that there exists a positive integer j' such that $S_{j'} > S_{j'+1}$. Obviously

$$(63) \quad S_0 = \frac{1}{1 + \alpha}$$

and

$$(64) \quad \lim_{j \rightarrow \infty} S_j = T.$$

Hence LEMMA 3 is proved.

Proof of Theorem II: Let $\delta_1 > \delta_2 > \dots > \delta_j > \dots$ be an arbitrary but fixed sequence such that $\lim \delta_j = 0$. For $\delta = \delta_j$, let $I_1, \dots, I_{k(j)}$ be a set of closed intervals such that no two intervals have an interior point in common and within which, by condition (A), $f'(x)$ and $g'(x)$ exist, are positive, and continuous. Let I_{0j} be the complementary set (with respect to the whole line). (It is easy to see that, if condition (A) is fulfilled, such a system can be constructed.) Let U_i ($i = 1, 2, \dots, k(j)$) and U_{0j} denote, respectively, the runs caused by the observations which fall in the intervals I_i, I_{0j} . Then

$$(65) \quad \left| U - \sum_{i=1}^{k(j)} U_i - U_{0j} \right| \leq 2(k(j)).$$

From condition (A) it follows that, with a probability arbitrarily close to 1, for sufficiently large m ,

$$(66) \quad U_{0j} < 3pm\delta_j,$$

$$\text{where} \quad p = \max \left[1, \frac{1}{\alpha} \right], \quad (j = 1, 2 \dots).$$

Let $[a_i \leq x < b_i]$, $i = 1, 2 \dots$ denote the interval I_i , and let m_i and n_i denote the number of observations on X and Y , respectively, which fall in the interval I_i . Then $\frac{m_i}{m}$ and $\frac{n_i}{n}$ converge stochastically with increasing m to $[f(b_i) - f(a_i)]$ and $[g(b_i) - g(a_i)]$, respectively.

Within the interval I_i ($i = 1, 2 \dots k$) we perform the transformation

$$(67) \quad X^* = f(X), \quad Y^* = f(Y),$$

which leaves U_i invariant. For fixed m_i , n_i the relative distribution of X^* is uniform and the relative distribution of Y^* fulfills condition (c) of LEMMA 1.

Hence from LEMMA 2 we obtain that $\frac{U_i}{m}$ converges stochastically to

$$(68) \quad \lim_{m \rightarrow \infty} E \left(\frac{U_i}{m}; f; g \right) \leq \frac{2[f(b_i) - f(a_i)][g(b_i) - g(a_i)]}{[g(b_i) - g(a_i)] + \alpha[f(b_i) - f(a_i)]}.$$

It can be verified that the sum of the second members in (68) over all values i is less than or equal to $\frac{2}{1 + \alpha}$.

From (24) and condition (A) we get that, for sufficiently small δ_j , there exists at least one interval for which the first member of (68) is less than the second member. Hence

$$(69) \quad \Sigma < \frac{2}{1 + \alpha},$$

where

$$(70) \quad \Sigma = \sum_{i=1}^{\infty} \lim_{m \rightarrow \infty} E \left(\frac{U_i}{m}; f; g \right).$$

Now take j so large that

$$(71) \quad 3p\delta_j < \epsilon,$$

where

$$(72) \quad 0 < 3\epsilon < \frac{2}{1 + \alpha} - \Sigma.$$

Since $\frac{U_i}{m}$ converges stochastically to its expected value, from (65), (66), (70), (71), and (72), it follows that, with a probability arbitrarily close to 1, for sufficiently large m ,

$$(73) \quad \frac{U}{m} < \frac{2}{1+\alpha} - \epsilon.$$

From (23) and THEOREM I we get

$$(74) \quad \lim_{m \rightarrow \infty} \frac{u_0(m)}{m} = \frac{2}{1+\alpha}.$$

THEOREM II follows easily from (73) and (74).

8. Remarks on a proposed test. We have already remarked in Section 3 that the test proposed by W. R. Thompson is not consistent. To show this, we shall give two distribution functions $f(x)$ and $g(x)$ such that, although these functions will be very different, the probability of rejecting the hypothesis that they are the same will not approach one as the sample number approaches infinity.

Suppose, to simplify the notation, that the observations have been ordered according to size, i.e., that $x_1 < x_2 < \dots < x_m$ and $y_1 < y_2 < \dots < y_n$. Suppose further that $m = n$, and that the test is to be made on a level of significance $\beta > 0$. In the right member of (2) we need not exhibit n and shall replace k and k' by $k(m)$ and $k'(m)$ to show the dependence on m . We have, under the null hypothesis,

$$(75) \quad P\{x_{k(m)} < y_{k'(m)}\} = \psi(m, k(m), k'(m)) = \beta.$$

The sequence $\frac{k(m)}{m}$ is bounded, so that there exists a monotonically increasing subsequence m_1, m_2, \dots of the sequence of integers $1, 2, \dots$ and a number h , $0 \leq h \leq 1$, such that

$$(76) \quad \lim_{i \rightarrow \infty} \frac{k(m_i)}{m_i} = h.$$

It is easy to see that then also

$$(77) \quad \lim_{i \rightarrow \infty} \frac{k'(m_i)}{m_i} = h.$$

We shall now assume that $0 < h < 1$. If $h = 0$ or 1 only a trivial alteration will be needed in the argument to follow. Let ϵ and δ be arbitrarily small positive numbers. We now consider two populations, A and B described as follows:

$$A) \quad f(x) \equiv g(x) \equiv x \quad (0 \leq x \leq 1),$$

$$B) \quad f(x) \equiv x \quad (0 \leq x \leq 1),$$

$$g(x) \equiv g(a_i) + \frac{(x - a_i)(g(a_{i+1}) - g(a_i))}{(a_{i+1} - a_i)} \quad (a_i \leq x \leq a_{i+1}; i = 0, 1, \dots, 4),$$

where

$$\begin{array}{ll}
 a_0 = 0 & g(a_0) = 0 \\
 a_1 = h - 2\delta > 0 & g(a_1) = 0 \\
 a_2 = h - \delta & g(a_2) = a_2 \\
 a_3 = h + \delta < 1 - \delta & g(a_3) = a_3 \\
 a_4 = 1 - \delta & g(a_4) = a_3 \\
 a_5 = 1 & g(a_5) = 1
 \end{array}$$

The definition of $f(x)$ and $g(x)$ outside the interval $0 \leq x \leq 1$ is obvious. It will be shown that even for such different populations as A and B and for samples of size greater than that of any arbitrarily assigned number, the probability of rejecting the null hypothesis if B is true will be at most $\beta + \epsilon$.

Let h_1, h_2, h_3 denote the number of observations on X which fall in the intervals $0 < x \leq a_2, a_2 < x \leq a_3, a_3 < x \leq 1$, respectively (m fixed, of course). Let h'_1, h'_2, h'_3 be the corresponding numbers for Y . For a fixed m , the probability of a set $h_1, h_2, h_3, h'_1, h'_2, h'_3$ is the same whether the sample be drawn from the population A or B. From (76), (77), and multinomial law it follows that for all sufficiently large m_i the probability is at least $1 - \epsilon$ of the occurrence of a set $h_1, h_2, h_3, h'_1, h'_2, h'_3$ for which $x_{k(m_i)}$ and $y_{k'(m_i)}$ will both fall in the interval $a_2 < x \leq a_3$. Furthermore it is obvious that for all samples with fixed h_2, h'_2 the distribution within the interval $a_2 < x \leq a_3$ is the same whether the sample came from the population A or B. Hence even when the sample is drawn from the population B, the first member of (75) is $< \beta + \epsilon$. This completes the proof of the inconsistency of the test based on (75).

This test is consistent if the alternatives to the null hypothesis are limited, for example, to those where $g(x) \equiv f(x + c)$, c a constant.

REFERENCES

- [1] WILLIAM R. THOMPSON, *Annals of Math. Stat.*, Vol. 9, (1938), p. 281.
- [2] MAURICE FRÉCHET, *Généralités sur les Probabilités. Variables aléatoires*, Paris, (1937).
- [3] J. NEYMAN AND E. S. PEARSON. *Statistical Research Memoirs*. University College, London. Vol. 1, (1936).

COLUMBIA UNIVERSITY,
NEW YORK, N. Y.