

AN EXACT TEST FOR RANDOMNESS IN THE NON-PARAMETRIC CASE BASED ON SERIAL CORRELATION¹

BY A. WALD AND J. WOLFOWITZ

Columbia University

1. Introduction. A sequence of variates x_1, \dots, x_N is said to be a random series, or to satisfy the condition of randomness, if x_1, \dots, x_N are independently distributed with the same distribution; i.e., if the joint cumulative distribution function (c.d.f.) of x_1, \dots, x_N is given by the product $F(x_1) \cdots F(x_N)$ where $F(x)$ may be any c.d.f.

The problem of testing randomness arises frequently in quality control of manufactured products. Suppose that x in some quality character of a product and that x_1, x_2, \dots, x_N are the values of x for N consecutive units of the product arranged in some order (usually in the order they were produced). The production process is said to be in a state of statistical control if the sequence (x_1, \dots, x_N) satisfies the condition of randomness. A number of tests of randomness have been devised for purposes of quality control, all having the following features in common: 1) They are based on runs in the sequence x_1, \dots, x_N . 2) The test procedure is invariant under topologic transformation of the x -axis, i.e., the test procedure leads to the same result if the original variates x_1, \dots, x_N are replaced by x'_1, \dots, x'_N where $x'_\alpha = f(x_\alpha)$ and $f(t)$ is any continuous and strictly monotonic function of t . 3) The size of the critical region, i.e., the probability of rejecting the hypothesis of randomness when it is true, does not depend on the common c.d.f. $F(x)$ of the variates x_1, \dots, x_N . Condition (3) is *a fortiori* fulfilled if condition (2) is satisfied and if $F(x)$ is continuous. The fulfillment of condition (3) is very desirable, since in many practical applications the form of the c.d.f. $F(x)$ is unknown.

Tests of randomness are of importance also in the analysis of time series (particularly of economic time series) where they are frequently based on the so-called serial correlation. The serial correlation coefficient with lag h is defined by the expression² (see, for instance, Anderson [1])

$$(1) \quad R_h = \frac{\sum_{\alpha=1}^N x_\alpha x_{h+\alpha} - \left(\sum_{\alpha=1}^N x_\alpha\right)^2 / N}{\sum_{\alpha=1}^N x_\alpha^2 - \left(\sum_{\alpha=1}^N x_\alpha\right)^2 / N}$$

where $x_{h+\alpha}$ is to be replaced by $x_{h+\alpha-N}$ for all values of α for which $h + \alpha > N$. The distribution of R_h has recently been studied by R. L. Anderson [1], T. Koopmans [2], L. C. Young [3], J. v. Neumann [4, 5], B. I. Hart and J. v. Neu-

¹ Presented to the Institute of Mathematical Statistics and the American Mathematical Society at a joint meeting at New Brunswick, New Jersey, on September 13, 1943.

² Some authors (see, for instance, [2] p. 27, equation (61)) use a non-circular definition.

mann [6], and J. D. Williams [7], under the assumption that x_1, \dots, x_N are independently distributed with the same normal distribution. Thus, in addition to the randomness of the series (x_1, \dots, x_N) it is assumed that the common c.d.f. of the variates x_1, \dots, x_N is normal. This is a restrictive assumption since frequently the form of the common c.d.f. $F(x)$ of the variates x_1, \dots, x_N is unknown.

The purpose of this paper is to develop a test procedure based on R_h such that (a) if $F(x)$ is continuous the size of the critical region does not depend on the common c.d.f. $F(x)$ of the variates x_1, \dots, x_N , thus making an exact test of significance possible also when nothing is known about $F(x)$ except its continuity; (b) if $F(x)$ is not continuous, but all its moments are finite and its variance is positive, the size of the critical region approaches, as $N \rightarrow \infty$, the value it would have if $F(x)$ were continuous. Thus in the limit an exact test is possible in this case as well. We will refer to the case where the form of $F(x)$ is unknown as the non-parametric case, in contrast to the case when it is known that $F(x)$ is a member of a finite parameter family of c.d.f.'s.

The test based on the serial correlation seems to be suitable if the alternative to randomness is the existence of a trend³ or of some regular cyclical movement in the data. In the analysis of time series it is frequently assumed that this is the case and this is perhaps the reason why tests based on serial correlation are widely used in the analysis of time series. In quality control of manufactured products the existence of a trend is often considered as the alternative to randomness, caused perhaps by the steady deterioration of a machine in the production process. Thus, tests of randomness based on serial correlation could also be used in quality control.

2. An exact test procedure based on R_h . Let a_α be the observed value of x_α ($\alpha = 1, \dots, N$). Consider the subpopulation where the set (x_1, \dots, x_N) is restricted to permutations of a_1, \dots, a_N . In this subpopulation the probability that (x_1, \dots, x_N) is any particular permutation (a'_1, \dots, a'_N) of (a_1, \dots, a_N) is equal to $1/N!$ if the hypothesis to be tested, i.e., that of randomness, is true. (If two of the a_i ($i = 1, 2, \dots, N$) are identical we assume that some distinguishing index is attached to each so that they can then be regarded as distinct and so that there still are $N!$ permutations of the elements a_1, \dots, a_N .)

The probability distribution of R_h in this subpopulation can be determined as follows: Consider the set of $N!$ values of R_h which are obtained by substituting for (x_1, \dots, x_N) all possible permutations of (a_1, \dots, a_N) . (A value which occurs more than once is counted as many times as it occurs.) Each of these values of R_h has the probability $1/N!$. On the basis of this distribution of R_h an exact test of significance can be carried out. Suppose that α is the level of significance, i.e., the size of the critical region. We choose as critical region a subset of M values out of the set of $N!$ values of R_h where $M/N! = \alpha$. The sub-

³ If the existence of a trend is feared it may be preferable to use the non-circular statistic discussed, for example, in [2].

set of M values which constitute the critical region will depend in each particular problem on the possible alternatives to randomness. For example, if a linear trend is the only possible alternative to randomness, then the critical region will consist of the M largest values⁴ of R_h . The value of the lag h will also be chosen on the basis of the alternatives under consideration. For instance, if some cyclical movement in the data is suspected the choice of h will depend on the form of these cycles. The general idea underlying the choice of the subset of M values and of the lag is to make the power of the test with respect to the alternatives which are particularly feared as high as possible.

If R_h has the same value for several permutations of (a_1, \dots, a_N) , it may be impossible to have a critical region consisting of exactly M values of R_h . For example, if $a_1 = a_2 = \dots = a_N$, then all the $N!$ values of R_h are equal, and the number of values of R_h included in the critical region must be either 0 or $N!$. If $F(x)$ is continuous the probability that two values of R_h be equal is zero. This explains why an exact test is always possible when $F(x)$ is continuous. On the other hand, if $F(x)$ is not continuous, the probability that several values of R_h be equal is positive. However, the theorem we shall prove in Section 4 shows that in the limit an exact test is possible even when $F(x)$ is not continuous, but has finite moments and a positive variance. For if the latter is true, the probability is one that the weaker conditions for the validity of our theorem (given at the end of Section 4) will be fulfilled.

Consider the statistic

$$(2) \quad \bar{R}_h = \sum_{\alpha=1}^N x_\alpha x_{h+\alpha}$$

where $x_{h+\alpha}$ is to be replaced by $x_{h+\alpha-N}$ for all values of α for which $h + \alpha > N$. Since in the subpopulation under consideration $\sum_{\alpha=1}^N x_\alpha$ and $\sum_{\alpha=1}^N x_\alpha^2$ are constants, the statistic \bar{R}_h is a linear function of R_h in this subpopulation. Hence, the test based on \bar{R}_h is equivalent to the test based on R_h . Since \bar{R}_h is simpler than R_h , in what follows we shall restrict ourselves to the statistic \bar{R}_h .

We shall now show that, if h is prime to N , the totality T_h of the $N!$ values taken by \bar{R}_h is the same as T_1 , the totality of the $N!$ values taken by \bar{R}_1 .

In the argument which follows it is to be understood that, whenever a positive integer is greater than N , it is to be replaced by that positive integer less than or equal to N which differs from it by an integral multiple of N .

Clearly it will be sufficient to show the existence of a permutation p_1, p_2, \dots, p_N of the first N integers such that

$$p_i + 1 = p_{i+h} \quad (i = 1, 2, \dots, N).$$

Such a permutation is given by

$$j = p_{(j-1)h+1} \quad (j = 1, 2, \dots, N).$$

For if $j \neq j'$ then $(j-1)h+1 \neq (j'-1)h+1$ because h is prime to N . Hence to every positive integer i there is a unique positive integer j , ($i, j \leq N$) such

⁴ See footnote 3.

that

$$i = (j - 1)h + 1$$

Now

$$p_i + 1 = p_{(j-1)h+1} + 1 = j + 1 = p_{jh+1} = p_{i+h},$$

which is the required result.

In what follows we shall restrict ourselves to the case when h is prime to N . This is not a very restrictive assumption since in practice h will be small as compared with N and by omitting a few observations we can always make N prime to h . Since T_h is the same as T_1 we shall deal with the statistic \bar{R}_1 only. To simplify the notation we shall write R instead of \bar{R}_1 . Thus, the test procedure will be based on the statistic

$$(4) \quad R = \sum_{\alpha=1}^{N-1} x_{\alpha} x_{\alpha+1} + x_N x_1.$$

If N is very small an exact test of significance can be carried out by actually calculating the $N!$ possible values of R . However, this procedure is practically impossible if N is not small. In Section 3 the exact mean value and variance of R will be calculated, and in section 4 the normality of the limiting distribution of R will be proved. Thus, if N is sufficiently large so that the limiting distribution of R can be used, a test of significance can easily be carried out. Difficulties in carrying out the test arise if N is neither sufficiently small to make the computation of the $N!$ values of R practically possible, nor sufficiently large to permit the use of the limiting distribution. In such cases it may be helpful to determine the third and fourth, and perhaps higher, moments of R , on the basis of which upper and lower limits for the cumulative distribution of R can be derived. (For a description of the Tchebycheff inequalities by which this can be done see, for example, Uspensky, [8], pp. 373-380.) Since the limiting distribution is normal it may be useful to approximate the distribution by a Gram-Charlier series or to employ similar methods.

3. Mean value and variance of R .⁵ It is clear that

$$(5) \quad \begin{aligned} E(R) &= NE(x_1 x_2) = \frac{N}{N(N-1)} \sum_{\alpha \neq \beta} a_{\alpha} a_{\beta} \\ &= \frac{1}{N-1} [(a_1 + \dots + a_N)^2 - (a_1^2 + \dots + a_N^2)]. \end{aligned}$$

To calculate the variance of R we first calculate the second moment of R about the origin. We have

$$(6) \quad \begin{aligned} E(R^2) &= E(x_1 x_2 + \dots + x_{N-1} x_N + x_N x_1)^2 \\ &= NE x_1^2 x_2^2 + 2NE x_1 x_2^2 x_3 + (N^2 - 3N) E x_1 x_2 x_3 x_4. \end{aligned}$$

⁵ The first four moments of a similar statistic have been obtained by Young [3].

To express the expected values $Ex_1^2x_2^2$, $Ex_1x_2^2x_3$, and $Ex_1x_2x_3x_4$ we shall introduce the following notations for the symmetric functions of a_1, \dots, a_N : For any set of positive integers i_1, i_2, \dots, i_k the symbol $S_{i_1i_2\dots i_k}$ denotes the symmetric function $\sum_{\alpha_k} \dots \sum_{\alpha_1} a_{\alpha_1}^{i_1} \dots a_{\alpha_k}^{i_k}$ where the summation is to be taken over all possible sets of k positive integers $\alpha_1, \dots, \alpha_k$ subject to the restriction that $\alpha_u \leq N$ and $\alpha_u \neq \alpha_v$ ($u, v = 1, \dots, k$).

From (6) we easily obtain

$$\begin{aligned} E(R^2) &= \frac{N}{N(N-1)} S_{22} + \frac{2N}{N(N-1)(N-2)} S_{121} \\ &+ \frac{N^2 - 3N}{N(N-1)(N-2)(N-3)} S_{1111} \\ (7) \quad &= \frac{S_{22}}{(N-1)} + \frac{2S_{121}}{(N-1)(N-2)} + \frac{S_{1111}}{(N-1)(N-2)}. \end{aligned}$$

It will probably facilitate computation to express each of the symmetric functions in the right member of (7) by a sum of terms, each a product of factors S_r ($r = 1, 2, \dots$). One can easily verify the relationships

$$(8) \quad S_{11} = S_1^2 - S_2$$

$$(9) \quad S_{12} = S_{21} = S_1S_2 - S_3$$

$$(10) \quad S_{13} = S_{31} = S_1S_3 - S_4$$

$$(11) \quad S_{22} = S_2^2 - S_4$$

$$\begin{aligned} (12) \quad S_{111} &= S_{11}S_1 - 2S_{12} = (S_1^2 - S_2)S_1 - 2(S_1S_2 - S_3) \\ &= S_1^3 - 3S_1S_2 + 2S_3 \end{aligned}$$

$$\begin{aligned} (13) \quad S_{112} &= S_{121} = S_{211} = S_{11}S_2 - 2S_{13} \\ &= (S_1^2 - S_2)S_2 - 2(S_1S_3 - S_4) \\ &= S_1^2S_2 - S_2^2 - 2S_1S_3 + 2S_4 \end{aligned}$$

$$\begin{aligned} (14) \quad S_{1111} &= S_{111}S_1 - 3S_{112} \\ &= S_1^4 - 3S_1^2S_2 + 2S_1S_3 - 3S_1^2S_2 + 3S_2^2 + 6S_1S_3 - 6S_4 \\ &= S_1^4 - 6S_1^2S_2 + 8S_1S_3 + 3S_2^2 - 6S_4. \end{aligned}$$

It follows from (5) that

$$(15) \quad E(R) = \frac{1}{N-1} (S_1^2 - S_2),$$

and from (7), (11), (13), (14), and (15) that the variance of R is given by

$$\begin{aligned} \sigma^2(R) &= E(R^2) - [E(R)]^2 \\ (16) \quad &= \frac{S_2^2 - S_4}{N-1} + \frac{S_1^4 - 4S_1^2S_2 + 4S_1S_3 + S_2^2 - 2S_4}{(N-1)(N-2)} - \frac{1}{(N-1)^2} (S_1^2 - S_2)^2. \end{aligned}$$

The mean value and variance of R can easily be computed from (15) and (16) as soon as the values of S_1, S_2, S_3 , and S_4 have been determined.

The formulas (15) and (16) are considerably simplified if $S_1 = 0$. In the special case that $S_1 = 0$ we have

$$(15') \quad E(R) = -\frac{S_2}{N-1}$$

and

$$(16') \quad \sigma^2(R) = \frac{S_2^2 - S_4}{N-1} + \frac{S_2^2 - 2S_4}{(N-1)(N-2)} - \frac{S_2^2}{(N-1)^2}.$$

We can always make S_1 equal to zero by replacing a_α by $b_\alpha = a_\alpha - N^{-1} \Sigma a_\alpha$. This substitution is permissible, since it changes the statistic R only by an additive constant and consequently leaves the test procedure unaffected. Thus, in practical applications it may be convenient to replace a_α by b_α and to use formulas (15') and (16').

4. Limiting distribution of R . Let $\{a_\alpha\}$ ($\alpha = 1, 2, \dots$ ad inf.) be a sequence of real numbers with the following properties:

a) There exists a sequence of numbers $A_1, A_2, \dots, A_r, \dots$ such that

$$(17) \quad \frac{1}{N} \left| \sum_{\alpha=1}^N a_\alpha^r \right| \leq A_r \quad (r = 1, 2, \dots \text{ ad inf.})$$

for all N . (This condition means that the moments about the origin of the sequence a_1, a_2, \dots, a_N are bounded functions of N .)

b) If

$$\delta(N) = \frac{1}{N} \left[\sum_{\alpha=1}^N a_\alpha^2 - \frac{1}{N} \left(\sum_{\alpha=1}^N a_\alpha \right)^2 \right],$$

then

$$(18) \quad \liminf_N \delta(N) > 0.$$

(This condition means that the dispersion of the N values a_1, a_2, \dots, a_N is eventually bounded below.)

Let $R(N)$ be the serial correlation coefficient R as defined in (4), where x_1, \dots, x_N is a random permutation of a_1, a_2, \dots, a_N . We shall prove the following

THEOREM: As $N \rightarrow \infty$, the probability that

$$\frac{R(N) - E(R(N))}{\sigma(R(N))} < t$$

approaches the limit

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-\frac{1}{2}x^2} dx,$$

For any function $f(N)$ and any positive function $\phi(N)$ let

$$f(N) = O(\phi(N))$$

mean that $|f(N)/\phi(N)|$ is bounded from above for all N , and let

$$f(N) = \Omega(\phi(N))$$

mean that

$$f(N) = O(\phi(N))$$

and that $\liminf_N |f(N)/\phi(N)| > 0$. Also let

$$f(N) = o(\phi(N))$$

mean that

$$\lim_{N \rightarrow \infty} \frac{f(N)}{\phi(N)} = 0.$$

Let $[\rho]$ denote the largest integer less than or equal to ρ .

To simplify the proof we shall temporarily assume:

c) There exists a positive constant K such that, for every positive integral N ,

$$(19) \quad -K \leq S_1 = \sum_{\alpha=1}^N a_\alpha \leq K.$$

This restriction will be removed later.

LEMMA 1:

$$\sum_{\alpha_1 < \dots < \alpha_k} a_{\alpha_1} a_{\alpha_2} \dots a_{\alpha_k} = O(N^{\lfloor k/2 \rfloor}).$$

PROOF: $\sum_{\alpha_1 < \dots < \alpha_k} a_{\alpha_1} \dots a_{\alpha_k}$ can be written as the sum of a finite number of terms where each term is a product of factors S_r ($r = 1, 2, \dots$). This representation will be called the normal representation of $\sum \dots \sum a_{\alpha_1} \dots a_{\alpha_k}$. Since $S_1 = O(1)$ by (19) and $S_r = O(N)$ by (17) and since the number of factors S_r ($r > 1$) in a single term of the normal representation of $\sum \dots \sum a_{\alpha_1} \dots a_{\alpha_k}$ is at most $\lfloor \frac{1}{2}k \rfloor$, the equation $\sum \dots \sum a_{\alpha_1} \dots a_{\alpha_k} = O(N^{\lfloor k/2 \rfloor})$ must hold.

LEMMA 2: Let $y = x_1 \dots x_k z$, where $z = x_{k+1}^{i_1} \dots x_{k+r}^{i_r}$ and $i_j > 1$ ($j = 1, \dots, r$). If (x_1, \dots, x_N) is a random permutation of a_1, \dots, a_N , and if k, r, i_1, \dots, i_r are fixed values independent of N , then $E(y) = O(N^{\lfloor k/2 \rfloor - k})$.

PROOF: Let $E(y | x_{k+1}, \dots, x_{k+r})$ be the conditional expected value of y when x_{k+1}, \dots, x_{k+r} are fixed. It follows easily from Lemma 1 that

$$E(y | x_{k+1}, \dots, x_{k+r}) = O(N^{\lfloor k/2 \rfloor - k}).$$

Hence also $E(y) = O(N^{\lfloor k/2 \rfloor - k})$ and Lemma 2 is proved.

Denote $x_\alpha x_{\alpha+1}$ by $y_\alpha (\alpha = 1, \dots, N - 1)$ and $x_N x_1$ by y_N , and consider the expansion of $(y_1 + \dots + y_N)^r$. Let y be a term of this expansion, i.e., $y = \frac{N!}{i_1! \dots i_u!} y_{\alpha_1}^{i_1} \dots y_{\alpha_u}^{i_u} (\alpha_1 < \alpha_2 < \dots < \alpha_u)$. We will say that two factors y_α and y_β are neighbors if $|\alpha - \beta + 1|$ or $|\alpha - \beta - 1|$ is either 0 or N . The set of u factors $y_{\alpha_1}, \dots, y_{\alpha_u}$ can be subdivided into cycles as follows: The first cycle contains y_{α_1} and all those y_α which can be reached from y_{α_1} by a succession of neighboring y_α . The second cycle contains the first y_α of the remaining sequence and all those which can be reached from the first y_α by a succession of neighboring y_α . The third cycle is similarly constructed from the remaining sequence, etc. After a finite number of cycles have been withdrawn the sequence will be exhausted. If m is the number of such cycles we will say that y has m cycles.

LEMMA 3: Let y be a term of the expansion $(x_1 x_2 + \dots + x_N x_1)^r = (y_1 + \dots + y_N)^r$ (r fixed). Let m be the number of cycles in y and k be the number of linear factors in y if y is written as a function of x_1, \dots, x_N (i.e., if we replace y_α by $x_\alpha x_{\alpha+1}$). Then the maximum value of $m + [\frac{1}{2}k] - k$ is equal to $[\frac{1}{2}r]$.

PROOF: First we maximize $m + [\frac{1}{2}k] - k$ with respect to k when m is fixed. If $m \leq [\frac{1}{2}r]$, then the minimum value of k is obviously zero. Let $m = [\frac{1}{2}r] + r'$ ($r' > 0$). The minimum value of k is reached if each cycle consists of a single factor y_α and if each factor y_α in y is either linear or squared. If r is even, then the minimum value of k is $4r'$ and if r is odd then the minimum value of k is $4r' - 2$. Hence for $m = [\frac{1}{2}r] + r'$ we have

$$\max_k (m + [\frac{1}{2}k] - k) = [\frac{1}{2}r] - r' \quad \text{if } r \text{ is even}$$

and

$$= [\frac{1}{2}r] - r' + 1 \text{ if } r \text{ is odd.}$$

Hence maximizing with respect to m and k we obtain

$$\max (m + [\frac{1}{2}k] - k) = [\frac{1}{2}r],$$

and Lemma 3 is proved.

LEMMA 4: The expected value of the sum of all those terms in the expansion of $(x_1 x_2 + \dots + x_N x_1)^r$ for which m is the number of cycles and k the number of linear factors (if y is expressed in terms of x_1, \dots, x_N) is equal to $O(N^{m + [\frac{1}{2}k] - k})$.

This Lemma follows from Lemma 2 and the fact that the number of terms y with the required properties is $O(N^m)$.

LEMMA 5:

$$E(x_1 x_2 + \dots + x_N x_1)^r = O(N^{[\frac{1}{2}r]}).$$

This follows from Lemmas 3 and 4.

LEMMA 6: If r is even then

$$E(x_1 x_2 + \dots + x_N x_1)^r = \left(C_{\frac{1}{2}r}^N \left(\frac{r!}{2^{\frac{1}{2}r}} \right) E(x_1^2 x_2^2 \dots x_{\frac{1}{2}r}^2) \right) + o(N^{[\frac{1}{2}r]}).$$

PROOF: It follows easily from our considerations in proving Lemma 3 that $m + [\frac{1}{2}k] - k < \frac{1}{2}r$ for all terms in the expansion of $(x_1x_2 + \dots + x_Nx_1)^r$ which are not of the type $x_1^2 \dots x_r^2$. Hence it follows from Lemma 4 that the expected value of the sum of all those terms in the expansion of $[x_1x_2 + \dots + x_Nx_1]^r$ which are not of the type $x_1^2 \dots x_r^2$ is equal to $o(N^{\frac{1}{2}r})$. Lemma 6 follows from the fact that $2^{-\frac{1}{2}r}r!$ is the coefficient of the terms of the type $x_1^2 \dots x_r^2$ in the expansion of $(x_1x_2 + \dots + x_Nx_1)^r$ and that the number of terms of such type is equal to $C_{\frac{1}{2}r}^N$.

LEMMA 7. $\lim_{N \rightarrow \infty} \frac{E(x_1x_2 + \dots + x_Nx_1)^r}{\{E(x_1x_2 + \dots + x_Nx_1)^2\}^{\frac{1}{2}r}} = 0$ if r is odd and $= 2^{-\frac{1}{2}r}r!/(\frac{1}{2}r)!$ if r is even.

PROOF: From Lemma 6 it follows that

$$(20) \quad E(x_1x_2 + \dots + x_Nx_1)^2 = NE(x_1^2x_2^2) + o(N) = \Omega(N).$$

The first half of Lemma 7 follows from Lemma 5 and equation (20). If r is even then it follows from (20) that

$$(21) \quad \lim_{N \rightarrow \infty} \frac{E(x_1x_2 + \dots + x_Nx_1)^r}{\{E(x_1x_2 + \dots + x_Nx_1)^2\}^{\frac{1}{2}r}} = \lim_{N \rightarrow \infty} \frac{2^{-\frac{1}{2}r}C_{\frac{1}{2}r}^N r! E(x_1^2 \dots x_r^2)}{N^{\frac{1}{2}r}(E(x_1^2x_2^2))^{\frac{1}{2}r}} \\ = \lim_{N \rightarrow \infty} \frac{r!}{2^{\frac{1}{2}r}(\frac{1}{2}r)!} \frac{E(x_1^2 \dots x_r^2)}{(E(x_1^2x_2^2))^{\frac{1}{2}r}}.$$

It follows from (17), (19), and the normal representation of symmetric functions that

$$k! \sum_{a_{\alpha_1} < a_{\alpha_2} < \dots < a_{\alpha_k}} a_{\alpha_1}^2 \dots a_{\alpha_k}^2 = S_2^k + O(N^{k-1}).$$

From (17) and (18) we have $S_2 = \Omega(N)$. Since

$$E(x_1^2 \dots x_r^2) = r! \left(\sum_{a_{\alpha_1} < a_{\alpha_2} < \dots < a_{\alpha_r}} a_{\alpha_1}^2 \dots a_{\alpha_r}^2 \right) [N(N-1) \dots (N-r+1)]^{-1},$$

we obtain

$$(22) \quad \lim_{N \rightarrow \infty} \frac{E(x_1^2 \dots x_r^2)}{(E(x_1^2x_2^2))^{\frac{1}{2}r}} = 1.$$

The second half of Lemma 7 follows from (21) and (22).

LEMMA 8:

$$(23) \quad \lim_{N \rightarrow \infty} \frac{E(R(N))}{\sigma(R(N))} = 0,$$

$$(24) \quad \lim_{N \rightarrow \infty} \frac{E(R^2(N))}{\sigma^2(R(N))} = 1.$$

PROOF: Equation (24) is a trivial consequence of (23). From (15) $E(R) = O(1)$ and from (16) $\sigma(R) = \Omega(N^{\frac{1}{2}})$. The lemma follows easily from these relations.

PROOF OF THE THEOREM: According to Lemma 7 the r -th moment of $R[E(R^2)]^{-\frac{1}{2}}$ approaches the r -th moment of the normal distribution as $N \rightarrow \infty$. From this and Lemma 8 the required result follows if condition (c) holds. It remains therefore merely to remove condition (c). Assume now only that $a_1, a_2, \dots, a_\alpha, \dots$ satisfy conditions (a) and (b).

$R(N)$ is formed from the population of values a_1, a_2, \dots, a_N . Addition of a constant q to a_1, \dots, a_N adds the same constant to all the values of $R(N)$ and hence leaves $[R(N) - E(R(N))]/\sigma(R(N))$ unaltered. Let $q^{(N)}$ be $-\sum_{\alpha=1}^N a_\alpha/N$ and write $b_\alpha^{(N)} = a_\alpha + q^{(N)}$. Consider the sequences

$$B^{(i)} = b_1^{(i)}, b_2^{(i)}, \dots, b_i^{(i)} \quad (i = 1, 2, \dots, \text{ad inf.}).$$

From (17) it follows that the $|q^{(N)}|$ are bounded for all N . Hence the sequences $B^{(i)}$ satisfy condition (a). They obviously satisfy condition (c). Since $\delta(j)$ is invariant under addition of a constant we have

$$\liminf_j \frac{1}{j} \left(\sum_{\alpha=1}^j (b_\alpha^{(j)})^2 - \frac{1}{j} \left(\sum_{\alpha=1}^j b_\alpha^{(j)} \right)^2 \right) > 0,$$

so that the $B^{(i)}$ satisfy condition (b). Since $[R(N) - E(R(N))]/\sigma(R(N))$ has the same distribution in the sequence a_1, a_2, \dots, a_N as in the sequence $B^{(N)}$, the theorem follows.

It should be remarked that the theorem remains valid if conditions (a) and (b) are replaced by the weaker condition

$$\mu_r/\mu_2^{\frac{1}{2}r} = O(1) \quad (r = 3, 4, \dots, \text{ad inf.})$$

where

$$\mu_r = \frac{1}{N} \sum_{\alpha=1}^N \left(a_\alpha - \frac{1}{N} \sum_{\alpha=1}^N a_\alpha \right)^r.$$

This follows easily from the fact that $[R(N) - E(R(N))]/\sigma(R(N))$ remains unaltered if we replace the sequence a_1, \dots, a_N by the sequence $c_1^N, c_2^N, \dots, c_N^N$ where

$$c_\alpha^N = \left(a_\alpha - \frac{1}{N} \sum_1^N a_\alpha \right) / \left[\frac{1}{N} \sum \left(a_\alpha - \frac{1}{N} \sum a_\alpha \right)^2 \right]^{\frac{1}{2}}.$$

Conditions (a) and (b) are obviously satisfied by the sequence c_1^N, \dots, c_N^N .

5. Transformation of the original observations.

Let $f(t)$ be a continuous and strictly monotonic function of t ($-\infty < t < +\infty$). Suppose we replace the original observations a_1, \dots, a_N by d_1, \dots, d_N , where $d_\alpha = f(a_\alpha)$ ($\alpha = 1, \dots, N$). We obtain a valid test of significance if we carry out the test procedure as if d_1, \dots, d_N were the observed values instead of a_1, \dots, a_N . We could also replace the observed values a_1, \dots, a_N by their ranks. The question arises whether there is any advantage in making the test on the transformed values instead of on the original observations. It may well

be that by certain transformations we could considerably increase the power of the test with respect to alternatives under consideration. This problem needs further study.

6. Summary. A test procedure based on serial correlation is given for testing the hypothesis that x_1, \dots, x_N are independent observations from the same population, i.e., that x_1, \dots, x_N is a random series. By considering the distribution of the serial correlation coefficient in the subpopulation consisting of all permutations of the actually observed values a test procedure is obtained such that

a) if the common c.d.f. $F(x)$ is continuous, the size of the critical region, i.e., the probability of rejecting the hypothesis of randomness when it is true, does not depend upon $F(x)$,

b) if $F(x)$ is not continuous but all its moments are finite and its variance is positive, the size of the critical region approaches, as $N \rightarrow \infty$, the value it would have if $F(x)$ were continuous. Thus in the limit an exact test is possible in this case as well.

It is shown that the test based on the serial correlation with lag h is equivalent to the test based on the statistic⁶

$$\sum_{\alpha=1}^N x_{\alpha} x_{h+\alpha}$$

where $x_{h+\alpha}$ is to be replaced by $x_{h+\alpha-N}$ for all values of α for which $h + \alpha > N$. If h is prime to N , the distribution of $\sum_1^N x_{\alpha} x_{h+\alpha}$ is exactly the same as the distribution of $R = \sum_1^N x_{\alpha} x_{1+\alpha}$.

The mean value and variance of R are given by the following expressions:

$$E(R) = (S_1^2 - S_2)/(N - 1)$$

and

$$\sigma^2(R) = \frac{S_2^2 - S_4}{N - 1} + \frac{S_1^4 - 4S_1^2 S_2 + 4S_1 S_3 + S_2^2 - 2S_4}{(N - 1)(N - 2)} - \frac{(S_1^2 - S_2)^2}{(N - 1)^2}$$

where $S_r = x_1^r + \dots + x_N^r$.

It is shown that under some mild restrictions the limiting distribution of R is normal. The test procedure can therefore be easily carried out when N is sufficiently large to permit the use of the limiting distribution of R .

REFERENCES

- [1] R. L. ANDERSON, *Annals of Math. Stat.*, Vol. 13 (1942), p. 1.
- [2] T. KOOPMANS, *ibid.*, p. 14.
- [3] L. C. YOUNG, *Annals of Math. Stat.*, Vol. 12 (1941), p. 293.
- [4] J. v. NEUMANN, *Annals of Math. Stat.*, Vol. 12 (1941), p. 367.
- [5] J. v. NEUMANN, *Annals of Math. Stat.*, Vol. 13 (1942), p. 86.
- [6] B. I. HART and J. v. NEUMANN, *Annals of Math. Stat.*, Vol. 13 (1942), p. 207.
- [7] J. D. WILLIAMS, *Annals of Math. Stat.*, Vol. 12 (1943), p. 239.
- [8] J. USPENSKY, *Introduction to Mathematical Probability*, New York, 1937.

⁶ If the non-circular definition of the serial correlation coefficient is used, the term $x_N x_{N+h}$ should be omitted.