

SP16
8258

MATHEMATICS

THE ASYMPTOTIC NORMALITY AND CONSISTENCY OF
KENDALL'S TEST AGAINST TREND, WHEN TIES ARE
PRESENT IN ONE RANKING

BY

T. J. TERPSTRA

(Communicated by Prof. D. VAN DANTZIG at the meeting of March 29, 1952)

1. Introduction

If $\mathbf{x}_1, \dots, \mathbf{x}_l$ ¹⁾ are l random variables under observation, H. B. MANN [7] defines an upward trend by means of the inequality

$$(1.1) \quad \sum_{i < j} \varepsilon_{j,i} > 0,$$

where

$$\varepsilon_{j,i} = 2 P[\mathbf{x}_i < \mathbf{x}_j] - 1, \quad -1 \leq \varepsilon_{j,i} \leq +1. \quad 2)$$

If only one observation x_i of \mathbf{x}_i is given for each $i = 1, \dots, l$, all observations being different from each other, the test of significance of rank correlation, given by M. G. KENDALL [4] and based on his ranking coefficient \mathbf{S} can be considered as a test against trend. MANN has shown the consistency and unbiasedness of the test, using a statistic \mathbf{T} , which is defined as the number of pairs (i, j) with $i < j$ and $\mathbf{x}_i < \mathbf{x}_j$, and which, but for a constant, is identical with KENDALL's \mathbf{S} .

In this paper, the asymptotic normality and consistency of KENDALL's test is shown for the case, that any number of observations of each random variable is available, all observations being different from each other. Instead of KENDALL's \mathbf{S} , another statistic is defined, which again, but for a constant term, proves to be identical with \mathbf{S} , this time for the case where ties are present in one ranking. In the definition of this statistic, use is made of the statistic \mathbf{U} of WILCOXON's two-sample test [9], as defined by MANN and WHITNEY [8], for comparing all samples, taken from the variables $\mathbf{x}_1, \dots, \mathbf{x}_l$.

2. The problem

Let $\mathbf{x}_{i,h}$, $i = 1, \dots, l$, $h = 1, \dots, n_i$, $\sum_i n_i = N$, be N independent, continuously distributed random variables, all $\mathbf{x}_{i,h}$, $h \leq n_i$, having the

¹⁾ Random variables will be distinguished from numbers (e.g. from the values, they take in an experiment), by printing them in bold type.

²⁾ Unless explicitly stated, i, j , take the values $1, \dots, l$.

same distribution function as a random variable \mathbf{x}_i . Let $x_{i,h}$ denote an observation of $\mathbf{x}_{i,h}$. We shall call the values $x_{i,h}$, $h \leq n_i$, the sample taken from \mathbf{x}_i .

We want a test, based on the observations $x_{i,h}$ for the hypothesis H_0 , stating that the random variables \mathbf{x}_i all have the same continuous distribution function, against the alternative hypothesis H_1 , stating that the variables $\mathbf{x}_1, \dots, \mathbf{x}_l$ possess an upward trend, as defined by (1.1).

Because of the continuity, all observations may be assumed to be different from each other, this being true with probability 1.

3. The statistic used and its connection with KENDALL'S \mathbf{S}

Let $\mathbf{U}_{i,j}$ be the number of pairs (h, k) , $h \leq n_i$, $k \leq n_j$, with $i < j$ and $\mathbf{x}_{i,h} < \mathbf{x}_{j,k}$ ³⁾. Then we define

$$(3.1) \quad T = \sum_{i < j} \sum \mathbf{U}_{i,j}.$$

For the case, that equal observations occur, this definition may be extended by increasing $\mathbf{U}_{i,j}$ with one half for each pair $(\mathbf{x}_{i,h}, \mathbf{x}_{j,k})$ of equal observations.

J. HEMELRIJK [3] has remarked that the statistic \mathbf{U} of WILCOXON'S two-sample test is connected with KENDALL'S \mathbf{S} for two rankings. In the same way it can be shown that T also is connected with \mathbf{S} . For this purpose we arrange all observations in order of increasing magnitude. A second ranking with l ties of the sizes n_i ⁴⁾ is obtained, by attributing to the n_i observations $x_{i,h}$ the same rank

$$\frac{1}{n_i} \sum_{h=1}^{n_i} \left(\sum_{k=1}^{i-1} n_k + h \right) = \sum_{k=1}^{i-1} n_k + \frac{1}{2} (n_i + 1).$$

If we compute KENDALL'S \mathbf{S} for these two rankings, a pair $(\mathbf{x}_{i,h}, \mathbf{x}_{j,k})$ ($i \leq j$) contributes to \mathbf{S} a score $\text{sgn}(j-i)(\mathbf{x}_{j,k} - \mathbf{x}_{i,h})$, whereas only pairs $(\mathbf{x}_{i,h}, \mathbf{x}_{j,k})$ ($i < j$) are considered for computation of T , contributing to T a score 1, $\frac{1}{2}$ or 0, according to $(\mathbf{x}_{j,k} - \mathbf{x}_{i,h})$ being > 0 , $= 0$, or < 0 .

Consequently

$$(3.2) \quad 2T - \mathbf{S} = \sum_{i < j} \sum n_i n_j.$$

To test the hypothesis H_0 , against the alternative hypothesis H_1 , a critical region $\{T \geq T_\alpha\}$ is defined, where T_α is the smallest integer with

$$P[T \geq T_\alpha | H_0] \leq \alpha.$$

For $T \geq T_\alpha$, H_0 is rejected, the probability of a wrong conclusion being at most α .

³⁾ MANN and WHITNEY define $\mathbf{U}_{i,j}$ as the number of pairs

$(\mathbf{x}_{i,h}, \mathbf{x}_{j,k})$ with $i < j$ and $\mathbf{x}_{i,h} > \mathbf{x}_{j,k}$.

⁴⁾ The case $n_i = 1$ is included.

4. The probability distribution of T under H_0

Lemma 1: Taking for some j all \mathbf{x}_h with $i < j$ and $h \leq n_i$ together, we define $\mathbf{Z}_1 = 0$ and for $j \geq 2$: $\mathbf{Z}_j = \mathbf{U}_{(1, \dots, j-1), j}$ as the number of pairs (h, k) with $h \leq n_i$, $i < j$, $k \leq n_j$ and $\mathbf{x}_{i,h} < \mathbf{x}_{j,k}$. Then $T = \sum_j \mathbf{Z}_j$.

Proof: From the definition of $\mathbf{U}_{i,j}$ and $\mathbf{U}_{(1, \dots, j-1), j}$ follows:

$$(4.1) \quad \begin{cases} \mathbf{U}_{(1, \dots, j-1), j} = \mathbf{U}_{1,j} + \dots + \mathbf{U}_{j-1,j}, \\ \mathbf{T} = \sum_{i < j} \sum \mathbf{U}_{i,j} = \sum_{j=2} \mathbf{U}_{(1, \dots, j-1), j} = \sum_j \mathbf{Z}_j. \end{cases}$$

Lemma 2: If H_0 is true, the random variables \mathbf{Z}_j are (completely) independent.

This lemma follows immediately from Theorem I (below) and the definition of \mathbf{Z}_j .

Let $\mathbf{y}_1, \dots, \mathbf{y}_n$ be n (completely) independent random variables with the same continuous probability distribution. Let R_n be the n -dimensional fundamental probability set and let for every point $E = (y_1, \dots, y_n) \in R_n$ the ranks r_1, \dots, r_n be defined by $r_i = \frac{1}{2} \sum_j \text{sgn}(y_i - y_j) + \frac{1}{2}(n+1)$. Thus r_1, \dots, r_n have a simultaneous probability distribution on R_n .

Theorem I: If the set of random variables $\mathbf{y}_1, \dots, \mathbf{y}_n$ is split into two sub-sets $\mathbf{y}_1, \dots, \mathbf{y}_m$ and $\mathbf{y}_{m+1}, \dots, \mathbf{y}_n$ and if $\{\mathbf{U}\}$ denotes a set of statistics depending on the permutation of the ranks r_1, \dots, r_m only, when arranged according to increasing magnitude, and $\{\mathbf{V}\}$ a set of statistics depending on the ranks r_{m+1}, \dots, r_n only, then $\{\mathbf{U}\}$ is stochastically independent of $\{\mathbf{V}\}$.

Proof: Let $K = K(r_{m+1}, \dots, r_n)$ be the subset of R_n , consisting of all points (y_1, \dots, y_n) for which y_{m+1}, \dots, y_n have the ranks r_{m+1}, \dots, r_n , then $P[\{\mathbf{U}\} = \{\mathbf{U}\} | K] = P[\{\mathbf{U}\} = \{\mathbf{U}\}]$. On K , the statistics $\{\mathbf{V}\}$ have constant values $\{\mathbf{V}\}$. If $A = A\{\mathbf{V}\}$ is the subset of R_n , consisting of all disjoint subsets K , corresponding with the set $\{\mathbf{V}\}$, we have

$$\begin{aligned} P[\{\mathbf{U}\} = \{\mathbf{U}\} | A] &= \frac{\sum_{(K)} P[\{\mathbf{U}\} = \{\mathbf{U}\} \wedge E \in K]}{\sum_{(K)} P[E \in K]} = \\ &= \frac{\sum_{(K)} P[\{\mathbf{U}\} = \{\mathbf{U}\} | K] \cdot P[E \in K]}{\sum_{(K)} P[E \in K]} = P[\{\mathbf{U}\} = \{\mathbf{U}\}]. \end{aligned}$$

Theorem II: If H_0 is true, the mean and the variance of T are

$$(4.2) \quad \begin{cases} \mu_0 = \mathcal{E}(T | H_0) = \frac{1}{4} (N^2 - \sum_i n_i^2), \\ \sigma_0^2 = \text{var}(T | H_0) = \frac{1}{72} \{N(N+1)(2N+1) - \sum_i n_i(n_i+1)(2n_i+1)\} \end{cases}$$

where

$$N = \sum_i n_i.$$

Proof: The proof follows immediately from (3.2) and KENDALL's result (cf. M. G. KENDALL, 1938 and Rank Correlation Methods, p. 43 (4.4)) by remarking that $(n_i - 1)(2n_i + 5)$ and $(n_i + 1)(2n_i + 1)$ have a

constant difference, cancelling after summation over i against the corresponding part of the first term.

Theorem III: *If H_0 is true, T is asymptotically normally distributed with mean and variance given by (4. 2).*

Proof: Because of (4. 1) and Lemma 2, the characteristic function $K(t)$ of the reduced variable $\tilde{T} = T - \mathcal{E}(T | H_0)$ is the product of the characteristic functions $K_j(t)$ of the variables $\tilde{Z}_j = Z_j - \mathcal{E}(Z_j | H_0)$, thus (cf [2]), putting

$$s_n(t) = \prod_{k=1}^n (\sin kt)/kt,$$

$$K(t) = s_N(\frac{1}{2}t) / \prod_i s_{n_i}(\frac{1}{2}t).$$

The characteristic function of the variable $T' = T/\sigma_0$ evidently is $K(t/\sigma_0)$.⁵⁾

Now we assume that n and (or) l depend on a parameter ν and tend to infinity as $\nu \rightarrow \infty$, whereas constants A and B exist with

$$(4. 3) \quad 0 < A \leq n^{-1} n_i \leq B < \infty \text{ for all } i = 1, \dots, l \text{ and all } \nu.$$

Then we have for large ν

$$\log K\left(\frac{t}{\sigma_0}\right) = -\frac{1}{6} \sum_{k=1}^N \left(k \frac{t}{2\sigma_0}\right)^2 + \frac{1}{6} \sum_i \sum_{k=1}^{n_i} \left(k \frac{t}{2\sigma_0}\right)^2 + O(l^{-1} n^{-1}),$$

consequently

$$\lim \log K\left(\frac{t}{\sigma_0}\right) = \lim \frac{1}{24} \frac{t^2}{\sigma_0^2} \left(-\sum_{k=1}^N k^2 + \sum_i \sum_{k=1}^{n_i} k^2\right) = -\frac{1}{2} t^2,$$

which is the logarithm of the characteristic function of a normally distributed variable with mean 0 and variance 1.

Because of (3. 2) we have the following

Corollary: *KENDALL'S rank-correlation coefficient S for two rankings, one ranking consisting of l ties of the sizes n_i ($n_i = 1$ included), whereas in the other ranking all ranks are different from each other, is asymptotically normally distributed with mean 0 and variance $4 \sigma_0^2$.*

5. The consistency of the test

To derive the class of alternative hypotheses H , against which the test is consistent, we define, analogous to D. VAN DANTZIG [1]

$$\tau_{j,k;i,h} = \iota(\mathbf{x}_{j,k} - \mathbf{x}_{i,h}),$$

where

$$\iota(z) = \begin{cases} 1, & \text{if } z \geq 0, \\ 0, & \text{if } z < 0. \end{cases}$$

⁵⁾ I owe this result to Prof. Dr D. VAN DANTZIG.

The statistic T can then be rewritten as

$$(5.1) \quad T = \sum_{i < j} \sum_{h=1}^{n_i} \sum_{k=1}^{n_j} \tau_{j,k;i,h}.$$

We first derive the mean of T and an upper bound for the variance of T under the hypothesis

$$(5.2) \quad H: P[\mathbf{x}_{i,h} < \mathbf{x}_{j,k}] = \frac{1}{2}(1 + \varepsilon_{j,i}), \quad -1 \leq \varepsilon_{j,i} \leq +1,$$

all $\mathbf{x}_{i,h}$ being (completely) independent and continuously distributed.

Because of the continuity

$$P[\mathbf{x}_{i,h} < \mathbf{x}_{j,k}] = P[\mathbf{x}_{i,h} \leq \mathbf{x}_{j,k}] = \frac{1}{2}(1 + \varepsilon_{j,i}),$$

hence $\mathcal{G}\tau_{j,k;i,h} = \frac{1}{2}(1 + \varepsilon_{j,i})$ and

$$(5.3) \quad \mu = \mathcal{G}(T|H) = \frac{1}{2} \sum_{i < j} n_i n_j (1 + \varepsilon_{j,i}).$$

As $\iota(z)^2 = \iota(z)$, $\mathcal{G}(\tau_{j,k;i,h})^2 = \mathcal{G}\tau_{j,k;i,h} = \frac{1}{2}(1 + \varepsilon_{j,i})$ and

$$\text{var } \tau_{j,k;i,h} = \frac{1}{4}(1 - \varepsilon_{j,i}^2).$$

Further we remark that $\tau_{j,k;i,h}$ and $\tau_{j',k';i',h'}$ are only correlated, if the pairs $(\mathbf{x}_{i,h}, \mathbf{x}_{j,k})$ and $(\mathbf{x}_{i',h'}, \mathbf{x}_{j',k'})$ have a variable in common. Then:

$$\begin{aligned} \text{cov}(\tau_{j',k';i',h'}, \tau_{j,k;i,h}) &= P[\mathbf{x}_{i,h} < \mathbf{x}_{j,k} \wedge \mathbf{x}_{i',h'} < \mathbf{x}_{j',k'}] - \frac{1}{4}(1 + \varepsilon_{j,i})(1 + \varepsilon_{j',i'}) \\ &\leq \min\left\{\frac{1}{2}(1 + \varepsilon_{j,i}), \frac{1}{2}(1 + \varepsilon_{j',i'})\right\} - \frac{1}{4}(1 + \varepsilon_{j,i})(1 + \varepsilon_{j',i'}) \leq \frac{1}{4}(1 - \varepsilon_{j,i} \varepsilon_{j',i'}), \end{aligned}$$

and in the same way:

$$\text{cov}(\tau_{j,k;i,h}, \tau_{j',k';i',h'}) \leq \frac{1}{4}(1 - \varepsilon_{j,i} \varepsilon_{j',i'}),$$

$$\text{cov}(\tau_{j',k';i',h'}, \tau_{j,k;i,h}) = P[\mathbf{x}_{i,h} < \mathbf{x}_{j,k} < \mathbf{x}_{j',k'}] - P[\mathbf{x}_{i,h} < \mathbf{x}_{j,k}] \cdot P[\mathbf{x}_{j,k} < \mathbf{x}_{j',k'}] \leq 0. \quad \text{e)}$$

Because of (5.1) and the inequalities, derived above, we obtain

$$(5.4) \quad \sigma^2 = \text{var}(T|H) \leq \frac{1}{4} \sum_{i < j} \sum_{i' < j'} n_i n_j (n_i + n_j - 1) + 2 \sum_{i < j} \sum_{i' < j'} n_i n_j n_{i'}$$

From (4.2) and (5.4) follows:

$$(5.5) \quad \sigma^2 \leq 12 \sigma_0^2.$$

We now proceed to state the following

Theorem IV: *Rejection of the hypothesis H_0 , if and only if $T \geq T_{\alpha}$, is a consistent test of H_0 against the class of alternative hypotheses (5.2), satisfying the conditions*

$$a) \quad \lambda = \frac{1}{n^2 l(l-1)} \sum_{i < j} n_i n_j \varepsilon_{j,i} > 0,$$

(so that $\frac{1}{2} \lambda n^2 l(l-1) = \mu - \mu_0$), and

$$b) \quad \lambda^{-1} = o((ln)^{1/2}),$$

e) Cf. H. B. MANN [1], p. 251.

and for sufficiently small α against no other alternatives, belonging to the class of alternative hypotheses (5. 2).

Proof: The probability that H_0 is not rejected under the hypothesis H is $P[\mathbf{T} < T_\alpha | H]$, where T_α is the smallest integer with

$$\beta = P[\mathbf{T} \geq T_\alpha | H_0] \leq \alpha.$$

We can write $T = T_\alpha$ in the form $T = \mu_0 + c\sigma_0$, where $\mu_0 = \mathcal{E}(\mathbf{T} | H_0)$ and $\sigma_0^2 = \text{var}(\mathbf{T} | H_0)$.

Because of the asymptotic normality of \mathbf{T} , $\lim \beta = \alpha$ and $\lim c = \xi_\alpha$ ⁷⁾, where

$$\frac{1}{\sqrt{2\pi} \xi_\alpha} \int_0^\infty e^{-x^2} dx = \alpha.$$

We then obtain

$$\begin{aligned} P[\mathbf{T} < T_\alpha | H] &= P[\mathbf{T} - \mu < -\frac{1}{2} \lambda l(l-1) n^2 + c\sigma_0 | H] \\ &\leq \frac{\sigma^2}{\sigma_0^2} \{-\frac{1}{2} \lambda l(l-1) n^2 \sigma_0^{-1} + c\}^{-2} \leq 12 \{\frac{1}{2} \lambda l(l-1) n^2 \sigma_0^{-1} - c\}^{-2}, \end{aligned}$$

because of (5. 5) and BIENAYMÉ's inequality, $\mu_0 - \mu + c\sigma_0$ being negative for sufficiently large ν , as c is bounded, $\lambda > 0$. The expression between curved brackets tends to infinity with ν as $\lambda^{-1} = O((ln)^{\frac{1}{2}})$, consequently $\lim P[\mathbf{T} < T_\alpha | H] = 0$ and the probability, that H_0 will be rejected under the hypothesis H tends to 1.

We now prove, as stated in Theorem IV, that the class of alternative hypotheses cannot be extended without loss of the property of consistency.

If condition (a) is not satisfied, then

$$P[\mathbf{T} \geq T_\alpha | H] \leq 12 \{-\frac{1}{2} \lambda l(l-1) n^2 \sigma_0^{-1} + c\}^{-2} \leq 12 c^{-2},$$

thus for sufficiently small α (e.g. sufficiently large c) there remains a positive probability $> 1 - 12c^{-2}$ that H_0 will not be rejected ⁸⁾.

If condition (b) is not satisfied, there is a constant c_1 and a sub-sequence of values ν , with $|\lambda(ln)^{\frac{1}{2}}| < c_1$. For this sub-sequence, we have for sufficiently small α

$$P[\mathbf{T} \geq T_\alpha | H] \leq 12 \{-\frac{1}{2} \lambda l(l-1) n^2 \sigma_0^{-1} + c\}^{-2} \leq 12 c_2^{-2},$$

where c_2 is a small constant. For sufficiently small α there again remains a positive probability $> 1 - 12 c_2^{-2}$ that H_0 will not be rejected.

6. The mean and the variance of \mathbf{T} , when equal observations are present

In section 3, \mathbf{T} was also defined for the case, that equal observations occur and it was shown that (3. 2) also holds.

⁷⁾ All limits are taken under condition (4. 3).

⁸⁾ If condition (b) is satisfied, this probability tends to 1, for each $\alpha < 1$.

The mean and the variance of T can thus be derived from the mean and the variance of KENDALL'S S for two rankings, one ranking consisting of l ties of the sizes n_i and the other ranking consisting of k ties of the sizes m_j , representing the groups of equal observations.

We obtain (cf. M. G. KENDALL [5], p. 43 (4.3))

$$\mathcal{E}(T|H_0) = \frac{1}{4}(N^2 - \sum_i n_i^2) = \frac{1}{4}(N^{12} - \sum_i n_i^{12})$$

and

$$\begin{aligned} \text{var}(T|H_0) = & (36 N^{13})^{-1} (N^{13} - \sum_i n_i^{13}) (N^{13} - \sum_j m_j^{13}) + \\ & + (8 N^{12})^{-1} (N^{12} - \sum_i n_i^{12}) (N^{12} - \sum_j m_j^{12}), \end{aligned}$$

where

$$N = \sum_i n_i = \sum_j m_j,$$

whereas x^{1k} for any natural k is defined by $x^{1k} = x(x-1)\dots(x-k+1)$.

In this case the asymptotic normality of T and S under condition (4.3) has not yet been shown.

Finally I want to thank Prof. Dr D. VAN DANTZIG, whose suggestions helped me to give the paper its final form.

*Publication of the Statistical Department
of the "Mathematisch Centrum",
Amsterdam.*

REFERENCES

1. DANTZIG, D. VAN, On the consistency and the power of WILCOXON'S two-sample test, Proc. Kon. Ned. Akad. van Wetensch. 54, 1-8 (1951); also Indagationes Mathematicae 13, 1-8 (1951).
2. ———, Kadercursus Mathematische Statistiek 1947-1950, Hoofdstuk VI, § 3. (Mathematisch Centrum, Amsterdam).
3. HEMELRIJK, J., Note on WILCOXON'S two-sample test, when ties are present, Ann. Math. Stat. 23, no. 2 (1952).
4. KENDALL, M. G., A new measure of rank correlation, Biometrika 30, 81 (1938).
5. ———, Rank correlation methods (London, 1948).
6. LEHMANN, E. L., Consistency and unbiasedness of certain non-parametric tests, Ann. Math. Stat. 22, 165-180 (1951).
7. MANN, H. B., Non-parametric tests against trend, Econometrica 13, 254-259 (1945).
8. ———, and D. R. WHITNEY, On a test of whether one of two random variables is stochastically larger than the other, Ann. Math. Stat. 18, 50-60 (1947).
9. WILCOXON, F., Individual comparisons by ranking methods, Biometrics Bull. 1, 80-83 (1945).

8250.