

Dana Vorlíčková

Remark on the rank tests in the case of censored samples

*Aplikace matematiky*, Vol. 20 (1975), No. 5, 372--377

Persistent URL: <http://dml.cz/dmlcz/103602>

## Terms of use:

© Institute of Mathematics AS CR, 1975

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://project.dml.cz>

REMARK ON THE RANK TESTS IN THE CASE OF CENSORED SAMPLES

DANA VORLÍČKOVÁ

(Received June 10, 1974)

1. INTRODUCTION

The problem is to test the hypothesis of randomness of the observations  $X_1, \dots, X_N$  against a regression alternative by means of a rank test procedure. We know exact values of some observations but about the other ones we can only say that they lie in a known open interval.

Let us consider  $k$  open intervals  $(y_0, y_1), \dots, (y_{k-1}, y_k)$ ,  $y_0 < y_1 < \dots < y_k$ , and let  $N_j$  observations lie in the  $j$ -th interval,  $1 \leq j \leq k$ , where  $N_1, \dots, N_k$  are random variables,  $0 \leq N_j \leq N$ ,  $\sum_{j=1}^k N_j \leq N$ . The exact values of the observations lying in the intervals  $(y_{j-1}, y_j)$ ,  $1 \leq j \leq k$ , are unknown, the exact values of the other ones are known. It is possible to have  $y_0 = -\infty$  or  $y_k = \infty$ . If  $k = 1$  we have the situation discussed in [3] and [2]. In the former paper some asymptotically efficient rank test procedures are derived for the hypothesis of randomness against two- and  $k$ -sample alternatives. In the latter paper the locally most powerful rank test of the hypothesis of randomness against the alternative of two samples is derived and the asymptotic distribution of the two sample statistic is established.

We try to find a rank test based on the censored sample described above in another way than it is done in [2] and [3]. We see that it is irrelevant that we cannot distinguish the values of observations lying in the same interval  $(y_{j-1}, y_j)$ ,  $1 \leq j \leq k$ , but it is essential that we cannot assign the exact ranks to them. We only know the set of possible values of ranks of these observations. Thus, we are in a similar situation as in the case of ties when a sample is taken from a noncontinuous distribution.

Without loss of generality we may suppose  $k = 1$ . Further, let  $X_1, \dots, X_N$  have a continuous distribution function  $F$ . We observe the following order statistics:

$$(1.1) \quad X^{(1)} < \dots < X^{(R_j)}, \quad X^{(R_j+N_1+1)} < \dots < X^{(N)},$$

where  $R_i$  stands for the rank of  $X_i$ , with probability 1.  $N_1$  observations lie in the

interval  $(y_0, y_1) = I$  and we can say about their ranks  $R_i$  only that  $R_j < R_i \leq R_j + N_1$ . We shall treat the observations from  $I$  as if all of them take on the same value  $y \in I$  with the probability  $F(y_1) - F(y_0)$  and form a tie. Then the censored sample under consideration behaves as a sample of random variables with the following common distribution function  $H$ :

$$(1.2) \quad \begin{aligned} H(x) &= F(x), & x < y_0, x > y_1, \\ &= F(y_0), & y_0 \leq x \leq y, \\ &= F(y_1), & y < x \leq y_1. \end{aligned}$$

The number  $y \in I$  is arbitrary but fixed.

For this situation Conover's work [1] about rank tests under noncontinuous distributions is applicable.

Let  $c_1, \dots, c_N, a_1, \dots, a_N$  be arbitrary real constants. Applying the method of randomization we get for the observations from  $I$  the ranks  $R_i^*$ ,  $R_i^* = R_j + 1, R_j + 2, \dots, R_j + N_1$ , respectively, and for the test of the hypothesis of randomness  $H_0$ :

$$P(X_1 \leq x_1, \dots, X_N \leq x_N) = \prod_{i=1}^N P(X_i \leq x_i) = \prod_{i=1}^N F(x_i)$$

we have the linear rank statistic

$$(1.3) \quad S^* = \sum_{i=1}^N c_i a(R_i^*) = \sum_{X_i \notin I} c_i a(R_i) + \sum_{X_i \in I} c_i a(R_i^*).$$

Applying the method of averaged scores we obtain for testing  $H_0$  the statistic

$$(1.4) \quad \bar{S} = \sum_{i=1}^N c_i a(R_i, R),$$

where

$$(1.5) \quad \begin{aligned} R_i &= \text{number of } X_j < X_i, \quad 1 \leq j \leq N, \quad X_i \notin I, \\ &= R_j + N_1, \quad X_i \in I, \end{aligned}$$

and

$$(1.6) \quad \begin{aligned} a(i, R) &= \frac{1}{N_1} \sum_{k=R_j+1}^{R_j+N_1} a(k), \quad i = R_j + N_1, \\ &= a(i), \quad i \neq R_j + N_1, \end{aligned}$$

so that

$$(1.7) \quad \bar{S} = \sum_{X_i \notin I} c_i a(R_i) + a(R_j + N_1, R) \sum_{X_i \in I} c_i.$$

The statistic (1.7) coincides with the statistic (3.2) of [2] in the case of two samples with the respective scores and with  $c_i = 1$  whenever  $X_i$  is from the first sample,  $c_i = 0$  otherwise,  $1 \leq i \leq N$ .

## 2. ASYMPTOTIC DISTRIBUTION OF TEST STATISTICS UNDER $H_0$

Let  $\varphi(u)$ ,  $0 < u < 1$ , (throughout the paper) be an arbitrary nonconstant square-integrable function. Let the scores  $a_i$ ,  $1 \leq i \leq N$ , satisfy

$$(2.1) \quad \int_0^1 (a(1 + [uN]) - \varphi(u))^2 du \rightarrow 0, \quad \text{as } N \rightarrow \infty,$$

and let for the regression constants  $c_i$ ,  $1 \leq i \leq N$ ,

$$(2.2) \quad \sum_{i=1}^N (c_i - \bar{c})^2 / \max_{1 \leq i \leq N} (c_i - \bar{c})^2 \rightarrow \infty \quad \text{as } N \rightarrow \infty$$

hold.

**Theorem 2.1.** Under  $H_0$ , (2.1), and (2.2), the statistic (1.3) is asymptotically normal  $((1/N) \sum_{i=1}^N c_i \sum_{i=1}^N a_i, (1/(N-1)) \cdot \sum_{i=1}^N (c_i - \bar{c})^2 \sum_{i=1}^N (a_i - \bar{a})^2)$ .

*Proof.* The assertion follows from [1], Theorem 4.4.

Denote by  $G$  the common distribution function of  $Y_i = H(X_i)$ ,  $1 \leq i \leq N$ , where  $X_i$ 's have the distribution function  $H$  given by (1.2), under  $H_0$ . Let  $G(\{y\}) = P(Y = y)$  in discontinuity points of  $G$  and equal zero elsewhere. Let us define

$$(2.3) \quad \begin{aligned} \varphi_a(u) &= \varphi(u) \quad \text{if } G(\{G^{-1}(u)\}) = 0, \\ &= \frac{1}{b(u) - a(u)} \int_{a(u)}^{b(u)} \varphi(t) dt, \quad u \in (a(u), b(u)), \end{aligned}$$

where  $a(u) = G^{-1}(u) - G(\{G^{-1}(u)\})$  and  $b(u) = G^{-1}(u)$ .

Put  $\bar{\varphi} = \int_0^1 \varphi(t) dt$ .

**Theorem 2.2.** Under  $H_0$ , (2.1), (2.2), and  $\int_0^1 (\varphi_a(u) - \bar{\varphi})^2 du > 0$ , the conditional distribution of the statistic (1.7) given  $N_1$  observations lie in  $I$  is asymptotically normal  $(N\bar{a}\bar{c}, \text{var}(\bar{S} | R))$ , where

$$(2.4) \quad \text{var}(\bar{S} | R) = \frac{1}{N-1} \sum_{i=1}^N (c_i - \bar{c})^2 \sum_{j=1}^N (a(j, R) - \bar{a})^2,$$

with  $a(j, R)$  given by (1.5).

*Proof.* The assertion follows from [1], Theorem 4.2.

### 3. ASYMPTOTIC DISTRIBUTION OF TEST STATISTICS UNDER ALTERNATIVES

Let us consider the alternative

$$(3.1) \quad P(X_1 \leq x_1, \dots, X_N \leq x_N) = \prod_{i=1}^N F(x_i, \theta_i),$$

where  $\theta_i$ ,  $1 \leq i \leq N$ , are real parameters, satisfying

$$(3.2) \quad \max_{1 \leq i \leq N} \theta_i^2 \rightarrow 0 \quad \text{as } N \rightarrow \infty,$$

and  $F$  is a continuous distribution function. The censored sample, which is observed, will be considered as a sample of random variables with the distribution functions  $H(\cdot, \theta_i)$  given by (1.2) with  $F(x) = F(x, \theta_i)$ ,  $1 \leq i \leq N$ .

Put

$$(3.3) \quad h(x, \theta) = \frac{dH(x, \theta)}{dH(x, 0)}$$

and define the generalized Fisher's information as

$$(3.4) \quad I(H, \theta) = \int_{-\infty}^{\infty} \left[ \frac{(\partial/\partial\theta) h(x, \theta)}{h(x, \theta)} \right]^2 dH(x, \theta).$$

Suppose

$$(3.5) \quad 0 < \lim_{\theta \rightarrow 0} I(H, \theta) = I(H) < \infty,$$

and

$$(3.6) \quad \lim_{N \rightarrow \infty} I(H) \sum_{i=1}^N \theta_i^2 = b^2, \quad 0 < b^2 < \infty.$$

Let  $J$  be an open interval containing zero. We shall consider a family of  $h$  given by (3.3) satisfying the following conditions for a.a.  $\theta \in J$  with respect to  $H(x, 0)$ :

- (3.7)      a)  $h(x, \theta)$  exists ,  
               b)  $\frac{\partial}{\partial\theta} h(x, \theta)|_{\theta=0} = \lim_{\theta \rightarrow 0} \frac{h(x, \theta) - h(x, 0)}{\theta}$  exists ,  
               c)  $h(x, 0) = \lim_{\theta \rightarrow 0} h(x, \theta)$  exists .

A family of population distribution functions  $F$  is determined by the conditions (3.7), too, through the relation (1.2) for  $F$  and  $H$ .

**Theorem 3.1.** *Let (3.2), (3.5), (3.6), (3.7) hold. Then, under (3.1), (2.1), and (2.2), the statistic (1.3) is asymptotically normal with parameters*

$$(3.8) \quad N\bar{c}\bar{a} + \sum_{i=1}^N c_i \theta_i \int_0^1 \varphi(u) \varphi(u, H, 0) du ,$$

and

$$(3.9) \quad \sum_{i=1}^N (c_i - \bar{c})^2 \int_0^1 (\varphi(u) - \bar{\varphi})^2 du ,$$

where

$$\varphi(u, H, 0) = \frac{(\partial/\partial\theta)h(H^{-1}(u, \theta), \theta)|_{\theta=0}}{h(H^{-1}(u, 0), 0)} .$$

Proof. The assertion follows from [1], Theorem 8.3 and the proof of Theorem 4.4.

Let  $\varphi_a$  be defined by (2.3) with  $G$  corresponding to  $H(x, 0)$  and let

$$(3.10) \quad \int_0^1 (\varphi_a(u) - \bar{\varphi})^2 du > 0 .$$

**Theorem 3.2.** Let (3.2), (3.5), (3.6), (3.7) hold. Then, under (3.1), (2.1), (2.2), and (3.10), the statistic (1.7) is asymptotically normal with parameters (3.8) and

$$\sum_{i=1}^N (c_i - \bar{c})^2 \int_0^1 (\varphi_a(u) - \bar{\varphi})^2 du .$$

Proof. The assertion follows from [1], Theorem 8.3 and the proof of Theorem 4.2.

**Remark.** Conover's conclusions about the efficiency of the tests continue to hold even for rank tests in the case of censored samples.

The considerations introduced above are obviously applicable also in the case of noncontinuous distribution functions  $F$  of  $X_i$ ,  $1 \leq i \leq N$ . It is only necessary to take into account, in addition, that ties of exactly measurable observations may occur.

A method analogical to the described one may be also used for censored samples in the case of testing symmetry but we omit it because of its similarity to the case of the hypothesis of randomness.

#### References

- [1] Conover W. J.: Rank tests for one sample, two samples and k samples without the assumption of a continuous distribution function. Ann. Statist. 1 (1973), 1105—1125.
- [2] Johnson R. A., Mehrotra K. G.: Locally most powerful rank tests for the two sample problem with censored data. Ann. Math. Statist. 43 (1972), 823—831.
- [3] Peto R., Peto J.: Asymptotically efficient rank invariant test procedures. J. Roy. Statist. Soc. Ser. A, 135 (1972), 185—206.

## Souhrn

### POZNÁMKA K POŘADOVÝM TESTŮM V PŘÍPADĚ CENSOROVANÝCH VÝBĚRŮ

DANA VORLÍČKOVÁ

Uvažujme náhodný výběr, o jehož některých pozorováních víme pouze, že leží v intervalu  $(y_{j-1}, y_j)$ ,  $1 \leq j \leq k$ ,  $y_0 < \dots < y_k$ , přičemž může být  $y_0 = -\infty$  nebo  $y_k = \infty$ . U ostatních pozorování známe přesné hodnoty. U pozorování, která padnou do téhož intervalu, nelze rozlišit pořadí. Proto při hledání pořadového testu postupujeme tak, jakoby pozorování padnuvší do téhož intervalu nabyla téže hodnoty s pravděpodobností rovnou rozdílu hodnot distribuční funkce v koncových bodech intervalu. S cenzorovaným výběrem pak zacházíme jako s výběrem z nespojitého rozdělení, v němž nastaly shody, a při konstrukci pořadových statistik užijeme buď metody znáhodnění nebo metody průměrných skóre. S využitím teorie z [1] pak dostáváme asymptotické rozdělení těchto statistik za hypotézy náhodnosti i při kontiguitních alternativách.

Analogického postupu by se dalo použít i pro testování hypotézy symetrie.

*Author's address:* RNDr. Dana Vorličková, CSc., katedra matematické statistiky MFF KU, Sokolovská 83, 186 00 Praha 8.