

V. DUPAČ and J. HÁJEK (Praha)

ASYMPTOTIC NORMALITY OF THE WILCOXON STATISTIC UNDER DIVERGENT ALTERNATIVES

1. Introduction. Throughout this paper, X_{v1}, \dots, X_{vm_v} and $X_{v, m_v+1}, \dots, X_{v, m_v+n_v}$ denote two samples from populations with continuous distribution functions F_v and G_v , respectively. We write N_v for $m_v + n_v$, R_{vi} for the rank of X_{vi} in the joint sample X_{v1}, \dots, X_{vN_v} , and S_v for the Wilcoxon statistic $\sum_{i=1}^{m_v} R_{vi}$; all symbols being defined for $v = 1, 2, \dots$

We shall be interested in the problem of asymptotic normality of the Wilcoxon statistic: Under which conditions concerning m_v, n_v, F_v, G_v is the distribution of S_v asymptotically normal? A sufficient condition given by the Chernoff-Savage theorem [1] reads: N_v tends to $+\infty$ and each of the sequences

$$(1) \quad \frac{m_v}{N_v}, \quad \frac{n_v}{N_v}, \quad \text{var } G_v(X_{v1}), \quad \text{var } F_v(X_{vN_v})$$

is bounded away from zero.

In a special case (location alternatives), when $F_v \equiv F$, $G_v(x) = F(x - \Delta_v)$, F' exists and its carrier is an interval of length $\Delta \leq +\infty$, the condition on variances reduces to the boundedness of $|\Delta_v|$ away from Δ .

In the present paper we try to get closer to necessary conditions. It will be shown that

(i) all the boundedness-from-zero conditions imposed on sequences (1) can be relaxed simultaneously (so that, especially, for the location alternatives, the asymptotic normality holds even for $|\Delta_v|$ tending — sufficiently slowly — to Δ);

(ii) more detailed statements concerning location alternatives are available for normal or rectangular F .

As to the methodology, the main tools are an inequality on projections of linear rank order statistics (Hájek [2]) and a lemma presenting sufficient conditions for asymptotic normality in terms of conditional distributions.

2. General alternatives.

THEOREM 1. *If*

$$(2) \quad \lim_{\nu \rightarrow +\infty} \frac{m_\nu n_\nu}{N_\nu} \left(\frac{n_\nu}{N_\nu} \text{var } G_\nu(X_{\nu 1}) + \frac{m_\nu}{N_\nu} \text{var } F_\nu(X_{\nu N_\nu}) \right) = +\infty,$$

then S_ν is asymptotically normal $(ES_\nu, \text{var } S_\nu)$ and

$$(3) \quad \lim_{\nu \rightarrow +\infty} \text{var } S_\nu / (m_\nu n_\nu^2 \text{var } G_\nu(X_{\nu 1}) + m_\nu^2 n_\nu \text{var } F_\nu(X_{\nu N_\nu})) = 1.$$

Proof. The essential points of the proof have been already prepared in [2]. We first approximate S (we shall drop the index ν if there is no danger of confusion) by the statistic $\hat{S} = ES + \sum_{i=1}^N Y_i$, where $Y_i = E(S|X_i) - ES$. The Y_i 's are independent, have zero expectations, and can be expressed as

$$(4) \quad Y_i = \sum_{j=1}^m (E(R_j|X_i) - ER_j) = \begin{cases} -nG(X_i) + \kappa_i, & i = 1, \dots, m, \\ mF(X_i) + \kappa'_i, & i = m+1, \dots, N \end{cases}$$

with some constants κ_i, κ'_i . It easily follows that $\text{var } \hat{S}$ is given by the denominator in (3) and that $|Y_i| \leq N$, $i = 1, \dots, N$, which together with (2) entails that $\text{ess. sup } |Y_i| / (\text{var } \hat{S})^{\frac{1}{2}}$ tends to zero. Thus the Lindeberg condition for the asymptotic normality of \hat{S} with parameters $(E\hat{S}, \text{var } \hat{S})$ is trivially satisfied. The asymptotic normality of S with parameters $(ES, \text{var } S)$ and the relation $\text{var } S / \text{var } \hat{S} \rightarrow 1$ will be entailed by $E(S - \hat{S})^2 / \text{var } \hat{S} \rightarrow 0$. This however, follows from an inequality (Th. 4.1. in [2]) which — specialized to our case — gives $E(S - \hat{S})^2 \leq \frac{1}{2} mn$, and from the relation (2).

THEOREM 2. *If for some $A_\nu, \nu \geq 1$, the sequence*

$$(5) \quad P(X_{\nu i} < A_\nu < X_{\nu j}, 1 \leq i \leq m_\nu, m_\nu + 1 \leq j \leq N_\nu)$$

is bounded away from zero, then S_ν is not asymptotically normal (μ_ν, σ_ν^2) for any $\mu_\nu, \sigma_\nu > 0$.

Proof is evident, since (5) implies that $P(S_\nu = \frac{1}{2} m_\nu (m_\nu + 1)), \nu \geq 1$, are bounded away from zero, which makes the asymptotic normality of S_ν impossible.

3. Location alternatives. From now on, we shall confine ourselves to the case

$$(6) \quad F_\nu(x) = F(x), \quad G_\nu(x) = F(x - \Delta_\nu) \quad \text{for all } x \text{ and } \nu.$$

THEOREM 3. *Let the distribution F be symmetric (about the median). Let*

$$(7) \quad \lim_{\nu \rightarrow \infty} \text{limmin}(m_\nu, n_\nu) \text{var } F(X_{\nu 1} - \Delta_\nu) = +\infty.$$

Then S_ν is asymptotically normal ($ES_\nu, \text{var } S_\nu$) and

$$\lim_{\nu \rightarrow +\infty} \text{var } S_\nu / (m_\nu n_\nu N_\nu \text{var } F(X_{\nu 1} - \Delta_\nu)) = 1.$$

Proof. For all x we have $F(x) = 1 - F(2c - x)$, where c is the median. Hence

$$\int_{G_\nu(u) < t} dF(u) = \int_{1 - F(2c - u + \Delta_\nu) < t} dF(u) = \int_{1 - F(v) < t} dG_\nu(t),$$

i.e., $G(X_{\nu 1})$ and $1 - F(X_{\nu N_\nu})$ have the same distributions, which implies that

$$\text{var } F(X_{\nu N_\nu}) = \text{var } G_\nu(X_{\nu 1}) = \text{var } F(X_{\nu 1} - \Delta_\nu).$$

Inserting this into (2) and (3), we get the statement of the theorem.

4. Normal location alternatives. In addition to (6), assume that

$$F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt.$$

Without loss of generality suppose that $m_\nu \leq n_\nu, \nu \geq 1$; further, let $m_\nu \rightarrow +\infty$.

THEOREM 4.

(i) Let

$$\limsup_{\nu \rightarrow +\infty} |\Delta_\nu| / (3 \log m_\nu)^{\frac{1}{2}} < 1;$$

then S_ν is asymptotically normal ($ES_\nu, \text{var } S_\nu$). If, moreover, $|\Delta_\nu| \rightarrow +\infty$, then

$$\lim_{\nu \rightarrow +\infty} \text{var } S_\nu / (3^{3/2} (2\pi)^{-1} \Delta_\nu^{-2} e^{-\Delta_\nu^2/3} m_\nu n_\nu N_\nu) = 1.$$

(ii) Let

$$\liminf_{\nu \rightarrow +\infty} |\Delta_\nu| / (8 \log m_\nu)^{\frac{1}{2}} > 1,$$

let m_ν / N_ν be bounded away from zero; then S_ν is not asymptotically normal (μ_ν, σ_ν^2) for any $\mu_\nu, \sigma_\nu > 0$.

Proof. We can restrict ourselves to $\Delta_\nu \rightarrow +\infty$, since other cases are either known (Theorem 3) or can be obtained by combining known results with results for $\Delta_\nu \rightarrow +\infty$. We write $F' = f$ and drop indices ν .

(i) The main part of the proof consists in finding bounds for $\text{var } G(X_1)$, $G(x) = F(x - \Delta)$. We have

$$\begin{aligned} \text{var } G(X_1) &= \\ &= \int_{-\infty}^{+\infty} G^2(x) f(x) dx - \left(\int_{-\infty}^{+\infty} G(x) f(x) dx \right)^2 = \int_{-\infty}^{q_1 \Delta} + \int_{q_1 \Delta}^{q_2 \Delta} + \int_{q_2 \Delta}^{+\infty} - \left(\int_{-\infty}^{+\infty} \right)^2 = \\ &= I_1 + I_2 + I_3 - (I_4)^2, \end{aligned}$$

where we shall assume that

$$0 < q_1 < \frac{2}{3} < q_2 < 1.$$

Using the formula

$$1 - F(x) = \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{x} \left(1 - \frac{\delta}{x^2}\right) e^{-\frac{x^2}{2}},$$

holding true for $x > 0$ with some $0 < \delta = \delta_x < 1$, we get

$$\begin{aligned} G^2(x)f(x)|_{x=q\Delta} &= [1 - F((1-q)\Delta)]^2 f(q\Delta) = \\ &= (2\pi)^{-\frac{3}{2}} (1-q)^{-2} \Delta^{-2} \left(1 - \frac{\delta}{(1-q)^2 \Delta^2}\right)^2 \exp[-\Delta^2\{(1-q)^2 + \frac{1}{2}q^2\}]. \end{aligned}$$

Since $(1-q)^2 + \frac{1}{2}q^2 = \frac{3}{2}(q - \frac{2}{3})^2 + \frac{1}{3}$,

$$I_1 < G^2(q_1\Delta) \int_{-\infty}^{+\infty} f(x) dx = o(e^{-\Delta^2(1-q_1)^2}),$$

$$I_3 < \int_{q_2\Delta}^{+\infty} f(x) dx = o(e^{-\Delta^2 \cdot q_2^2/2}),$$

$$\begin{aligned} I_2 &= \Delta \int_{q_1}^{q_2} G^2(q\Delta) f(q\Delta) dq = (2\pi)^{-\frac{3}{2}} \Delta^{-1} \int_{q_1}^{q_2} (1-q)^{-2} \times \\ &\quad \times \left(1 - \frac{\delta}{(1-q)^2 \Delta^2}\right)^2 e^{-\Delta^2\{\frac{3}{2}(q - \frac{2}{3})^2 + \frac{1}{3}\}} dq \\ &= 3^{-1/2} (2\pi)^{-1} \Delta^{-2} e^{-\frac{\Delta^2}{3}} \int_{q_1 - 2/3}^{q_2 - 2/3} \left(\frac{1}{3} - u\right)^{-2} \times \\ &\quad \times \left(1 - \frac{\delta}{(\frac{1}{3} - u)^2 \Delta^2}\right)^2 \frac{1}{(2\pi)^{1/2} 3^{-1/2} \Delta^{-1}} e^{-u^2/2 \cdot 3^{-1} \Delta^{-2}} du \\ &= 3^{3/2} (2\pi)^{-1} \Delta^{-2} e^{-\frac{\Delta^2}{3}} \{1 + o(1)\}. \end{aligned}$$

By similar calculations we get $(I_4)^2 = o(e^{-4\Delta^2/9})$.

Choosing $q_1 < 1 - 3^{-1/2}$, $q_2 > 2^{1/2} 3^{-1/2}$, we have ($\varepsilon > 0$) $I_1 + I_3 - (I_4)^2 = o(\exp[-\Delta^2(1 + \varepsilon)/3])$, hence

$$\text{var}G(X_1) / (3^{3/2} (2\pi)^{-1} \Delta^{-2} e^{-\Delta^2/3}) \rightarrow 1.$$

Now, if $\Delta \leq \alpha \sqrt{3 \log m}$ ($\alpha < 1$), then

$$m \Delta^{-2} e^{-\Delta^2/3} \geq \text{const} \frac{m^{1-\alpha^2}}{\log m} \rightarrow +\infty;$$

hence, the condition (7) is satisfied and the assertion (i) is thus proved.

(ii) As m/N is bounded away from zero and as $n/N \geq \frac{1}{2}$, we have $\alpha'\sqrt{2\log m} > \sqrt{2\log n}$ for each $\alpha' > 1$ and for m, n sufficiently large, and, further, $\Delta \geq \alpha'\sqrt{8\log m}$ implies $\Delta > \sqrt{2\log m} + \sqrt{2\log n}$. Put $A = \sqrt{2\log m}$ and apply Theorem 2. We get

$$\begin{aligned} P(X_i < A < X_j, 1 \leq i \leq m, m+1 \leq j \leq N) &= F^m(A)(1-G(A))^n \\ &= F^m(A)F^n(\Delta - A) \geq F^m(\sqrt{2\log m})F^n(\sqrt{2\log n}); \end{aligned}$$

but

$$F(\sqrt{2\log k}) = 1 - \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\sqrt{2\log k}} \left(1 - \frac{\delta}{2\log k}\right) e^{-\log k} = 1 - \frac{o(1)}{k},$$

hence $F^k(\sqrt{2\log k}) \rightarrow 1$ for $k \rightarrow +\infty$, which completes the proof.

5. Rectangular location alternatives. In this section we shall assume, in addition to (6), that

$$F(x) = \begin{cases} 0 & \text{for } x \leq 0, \\ x & \text{for } 0 < x \leq 1, \\ 1 & \text{for } 1 < x. \end{cases}$$

Again, we suppose $m_v \leq n_v, m_v \rightarrow +\infty$.

THEOREM 5.

(i) Let $m_v(1 - |\Delta_v|)$ tend to $+\infty$. Then S_v is asymptotically normal ($ES_v, \text{var } S_v$) with $ES_v = \frac{1}{2}(m_v(m_v+1) + m_v n_v(1 - |\Delta_v|)^2)$ and

$$\lim_{v \rightarrow +\infty} \text{var } S_v / \left(\frac{1}{12} m_v n_v N_v (1 + 3|\Delta_v|)(1 - |\Delta_v|)^3 \right) = 1.$$

(If $|\Delta_v| \rightarrow 1$, the denominator may be reduced to $\frac{1}{3} m_v n_v N_v (1 - |\Delta_v|)^3$.)

(ii) Let $m_v(1 - |\Delta_v|)$ be bounded (from above). Then S_v is not asymptotically normal (μ_v, σ_v^2) for any $\mu_v, \sigma_v > 0$.

The proof makes use of the following lemma.

LEMMA 1. Let T_v be random variables, let Z_v be k -dimensional random vectors, let $a_v, s_v > 0$ be functions of k variables, $v \geq 1$. Suppose that

(A) for each $\varepsilon > 0$, and each real $a < b$,

$$\lim_{v \rightarrow +\infty} P \left(\max_{a \leq t \leq b} \left| E \left(\exp \left(it(T_v - a_v(Z_v)) / s_v(Z_v) \right) \middle| Z_v \right) - e^{-t^2/2} \right| < \varepsilon \right) = 1;$$

(B) $a_v(Z_v)$ is asymptotically normal (μ, b^2);

(C) $s_v^2(Z_v) \xrightarrow{P} \sigma^2, 0 < \sigma^2 < +\infty$.

Then T_v is asymptotically normal ($\mu, \sigma^2 + b^2$).

Proof of Lemma 1. Define $R_\nu(t, Z_\nu)$ by the relation

$$E\left(\exp(it(T_\nu - a_\nu)/s_\nu) | Z_\nu\right) = e^{-t^2/2} + R_\nu(t, Z_\nu).$$

Writing here ts_ν instead of t , then multiplying both sides by e^{ita_ν} and taking expectations, we get

$$E(e^{itT_\nu}) = E(e^{ita_\nu - \frac{1}{2}t^2s_\nu^2}) + E(e^{ita_\nu} R_\nu(ts_\nu, Z_\nu)).$$

The first term on the right-hand side tends to $e^{it\mu - \frac{1}{2}t^2(b^2 + \sigma^2)}$, which is easily seen from the inequality

$$\begin{aligned} |E(e^{ita_\nu - \frac{1}{2}t^2s_\nu^2}) - e^{it\mu - \frac{1}{2}t^2(\sigma^2 + b^2)}| &\leq e^{-\frac{1}{2}t^2\sigma^2} |E(e^{ita_\nu}) - e^{it\mu + \frac{1}{2}t^2b^2}| + \\ &\quad + E(|e^{-\frac{1}{2}t^2s_\nu^2} - e^{-\frac{1}{2}t^2\sigma^2}|) \end{aligned}$$

and from (B) and (C). The second term is less or equal to $E(|R_\nu(ts_\nu, Z_\nu)|)$ in absolute value. To prove its convergence to zero, it suffices — because of $|R_\nu| \leq 2$ — to show that $R_\nu(ts_\nu, Z_\nu) \xrightarrow{P} 0$. But this follows from the inequality $P(|R_\nu(ts_\nu, Z_\nu)| > \varepsilon) \leq P(\max_{t\sigma - \delta \leq u \leq t\sigma + \delta} |R(u, Z_\nu)| > \varepsilon, |ts_\nu - t\sigma| \leq \delta) + P(|ts_\nu - t\sigma| > \delta)$ and from (A) and (C).

Proof of Theorem 5. We confine ourselves to the case $\Delta \rightarrow 1-$, for analogous reasons as in the proof of Theorem 4.

(i) Suppose that in the interval $\langle \Delta, 1 \rangle$ there lie exactly m' values of the first sample and n' values of the second one. Denote them $X'_1, \dots, X'_{m'}$, $X'_{m'+1}, \dots, X'_{m'+n'}$; further denote $R'_i = \sum_{j=1}^{m'+n'} u(X'_i - X'_j)$, where

$$u(x) = \begin{cases} 1 & \text{for } x \geq 0, \\ 0 & \text{for } x < 0, \end{cases}$$

and put $S' = \sum_{i=1}^{m'} R'_i$. Now,

$$\begin{aligned} (8) \quad S &= S' + \frac{1}{2}(m - m')(m - m' + 1) + m'(m - m') = \\ &= S' + \frac{1}{2}(m - m')(m + m' + 1). \end{aligned}$$

Here m' and n' are independent binomial variables with parameters $(m, 1 - \Delta)$ and $(n, 1 - \Delta)$, respectively. Conditioned by $m' = k$, $n' = l$, the distribution of S' is that of the Wilcoxon statistic from samples of sizes k, l , under the null hypothesis; it is asymptotically normal $(\frac{1}{2}k(k + l + 1), \frac{1}{12}kl(k + l + 1))$ for $k, l \rightarrow +\infty$. This fact can be written as

(9)

$$\lim_{k, l \rightarrow +\infty} \max_{a \leq t \leq b} \left| E\left(\exp\left(it \frac{S' - \frac{1}{2}m'(m' + n' + 1)}{(\frac{1}{12}m'n'(m' + n' + 1))^{1/2}}\right) \middle| m' = k, n' = l\right) - e^{-t^2/2} \right| = 0,$$

for each real $a < b$.

Let us define

$$\begin{aligned} d^2 &= \frac{1}{3}mnN(1-\Delta)^3, \\ T &= (S - \frac{1}{2}\{m(m+1) + mn(1-\Delta)^2\})/d, \\ a(k, l) &= \frac{1}{2}(kl - mn(1-\Delta)^2)/d, \\ s^2(k, l) &= \frac{1}{12}kl(k+l+1)/d^2, \end{aligned}$$

(all these quantities being dependent on ν). We shall verify that T, a, s satisfy the conditions of Lemma 1.

(A) According to (8),

$$\frac{T - a(m', n')}{s(m', n')} = \frac{S' - \frac{1}{2}m'(m' + n' + 1)}{(\frac{1}{12}m'n'(m' + n' + 1))^{1/2}},$$

so that we can rewrite the assertion (9) as follows: To each $\varepsilon > 0$ and $a < b$, there exist K, L (independent of ν) such that

$$\max_{a \leq t \leq b} \left| E \left(\exp \left(it \frac{T - a(m', n')}{s(m', n')} \right) \middle| m' = k, n' = l \right) - e^{-t^2/2} \right| < \varepsilon$$

holds for all $k \geq K, l \geq L$.

Now, from $m(1-\Delta) \rightarrow +\infty$ (and from $m < n$) it easily follows that $m' \xrightarrow{P} +\infty, n' \xrightarrow{P} +\infty$, which means that $P(m' \geq K, n' \geq L) \rightarrow 1$. Consequently,

$$\lim_{\nu \rightarrow +\infty} P \left(\max_{a \leq t \leq b} \left| E \left(\exp \left(it \frac{T - a(m', n')}{s(m', n')} \right) \middle| m', n' \right) - e^{-t^2/2} \right| < \varepsilon \right) = 1.$$

(B) We have

$$\begin{aligned} m'n' - mn(1-\Delta)^2 &= (m' - m(1-\Delta))n(1-\Delta) + (n' - n(1-\Delta))m(1-\Delta) + \\ &+ (m' - m(1-\Delta))(n' - n(1-\Delta)) = U + V + W, \end{aligned}$$

say. The assumptions $\Delta \rightarrow 1$ and $m(1-\Delta) \rightarrow +\infty$, together with the central limit theorem, entail that $U+V$ is asymptotically normal $(0, mnN(1-\Delta)^3)$; further,

$$EW^2 = mn\Delta^2(1-\Delta)^2.$$

Hence,

$$a(m', n') = (U + V + W)/2d$$

is asymptotically normal $(0, \frac{3}{4})$, considering that $E(W^2/4d^2) \rightarrow 0$.

(C) We have

$$s^2(m', n') = \frac{m'n'(m' + n' + 1)}{4mnN(1-\Delta)^3} \xrightarrow{P} \frac{1}{4},$$

since the coefficients of variation of m' and n' tend to 0 and thus each of the sequences

$$\frac{m'}{m(1-\Delta)}, \quad \frac{n'}{n(1-\Delta)}, \quad \frac{m'+n'+1}{N(1-\Delta)}$$

tends to 1 in probability.

The random variables T, a, s satisfy thus the conditions of Lemma 1 with constants $\mu = 0, b^2 = \frac{3}{4}, \sigma^2 = \frac{1}{4}$. Consequently, T is asymptotically normal $(0, 1)$, which is equivalent to the assertion (i) [for $\Delta \rightarrow 1$].

(ii) Put $\lambda = m(1-\Delta)$ and denote by C its upper bound (which exists by assumption). Then apply Theorem 2 with $A = \Delta$; we get

$$\begin{aligned} P(X_i < A < X_j, 1 \leq i \leq m, m+1 \leq j \leq N) = \\ = \Delta^m = \left(1 - \frac{\lambda}{m}\right)^m > (1-\varepsilon)e^{-C} \end{aligned}$$

for each $\varepsilon > 0$ and large m . This completes the proof.

REMARK ADDED IN PROOF. The proof of Theorem 1 is based on the original version of [2]. From Theorem 2.1 of the published revised version, the theorem follows immediately.

References

- [1] H. Chernoff and I. R. Savage, *Asymptotic normality and efficiency of certain nonparametric test statistics*, Ann. Math. Statist. 29 (1958), pp. 972-994.
 [2] J. Hájek, *Asymptotic normality of simple linear rank statistics under alternatives*, Ann. Math. Statist. 39 (1968), pp. 325-346.

CHARLES UNIVERSITY, PRAGUE

Received on 29. 5. 1967