

MANN-WHITNEY TEST FOR ASSOCIATED SEQUENCES

ISHA DEWAN AND B. L. S. PRAKASA RAO

*Department of Statistics and Mathematics, Indian Statistical Institute,
7 S.J.S. Sansanwal Marg, New Delhi 110016, India*

(Received May 21, 2001; revised March 26, 2002)

Abstract. Let $\{X_1, \dots, X_m\}$ and $\{Y_1, \dots, Y_n\}$ be two samples independent of each other, but the random variables within each sample are stationary associated with one dimensional marginal distribution functions F and G , respectively. We study the properties of the classical Wilcoxon-Mann-Whitney statistic for testing for stochastic dominance in the above set up.

Key words and phrases: U-statistics, Mann-Whitney statistic, central limit theorem, associated random variables.

1. Introduction

Suppose that two samples $\{X_1, \dots, X_m\}$ and $\{Y_1, \dots, Y_n\}$ are independent of each other, but the random variables within each sample are stationary associated with one dimensional marginal distribution functions F and G respectively. Assume that the density functions f and g of F and G respectively, exist. We wish to test for the equality of the two marginal distribution functions F and G . A commonly used statistic for this nonparametric testing problem is the Wilcoxon-Mann-Whitney statistic when the observations X_i , $1 \leq i \leq m$ are independent and identically distributed (i.i.d.) and Y_j , $1 \leq j \leq n$ are i.i.d. However, most often the X and the Y observations are not i.i.d. Suppose the samples are from a stationary associated stochastic process.

A finite family $\{X_1, \dots, X_n\}$ of random variables is said to be *associated* if

$$\text{Cov}(h_1(X_1, \dots, X_n), h_2(X_1, \dots, X_n)) \geq 0$$

for any coordinatewise nondecreasing functions h_1, h_2 on R^n such that the covariance exists. An infinite family of random variables is said to be *associated* if every finite subfamily is associated (cf. Esary *et al.* (1967)).

We wish to test the hypothesis that

$$(1.1) \quad H_0 : F(x) = G(x) \quad \text{for all } x,$$

against the alternative

$$(1.2) \quad H_1 : F(x) \geq G(x) \quad \text{for all } x,$$

with strict inequality for some x . We can test the above hypothesis conservatively by testing

$$(1.3) \quad H'_0 : \gamma = 0,$$

against the alternative

$$(1.4) \quad H_1' : \gamma > 0,$$

where $\gamma = 2P(Y > X) - 1 = P(Y > X) - P(Y < X)$.

Probabilistic aspects of associated random variables have been extensively studied (see, for example, Prakasa Rao and Dewan (2001) and Roussas (1999)). Here we extend the Wilcoxon-Mann-Whitney statistic to stationary sequences of associated variables. Serfling (1968) studied the Wilcoxon statistic when the samples are from stationary mixing processes. Louhichi (2000) gave an example of a sequence of random variables which is associated but not mixing. This shows that tests for samples from stationary associated random sequences need to be studied separately.

In Section 2 we state some results that are used to study the properties of Wilcoxon statistic for associated random variables. In Section 3 we discuss the asymptotic normality of the Wilcoxon statistic based on independent sequences of stationary associated variables.

2. Preliminaries

We state some theorems that are used in proving the main results in the next section.

THEOREM 2.1. (Bagai and Prakasa Rao (1991)) *Suppose X and Y are associated random variables with bounded continuous densities f_X and f_Y , respectively. Then there exists an absolute constant $C > 0$ such that*

$$(2.1) \quad \sup_{x,y} |P[X \leq x, Y \leq y] - P[X \leq x]P[Y \leq y]| \\ \leq C \left\{ \max \left(\sup_x f_X(x), \sup_x f_Y(x) \right) \right\}^{2/3} (\text{Cov}(X, Y))^{1/3}.$$

The following theorem gives the asymptotic normality of a sequence of associated variables.

THEOREM 2.2. (Newman (1980, 1984)) *Let $\{X_n, n \geq 1\}$ be a stationary associated sequence of random variables with $E[X_1^2] < \infty$ and $0 < \sigma^2 = V(X_1) + 2 \sum_{j=2}^{\infty} \text{Cov}(X_1, X_j) < \infty$. Then, $n^{-1/2}(S_n - E(S_n)) \xrightarrow{\mathcal{L}} N(0, \sigma^2)$ as $n \rightarrow \infty$.*

Assume that

$$(2.2) \quad \sup_x f(x) < c \quad \sup_x g(x) < c.$$

Further assume that

$$(2.3) \quad \sum_{j=2}^{\infty} \text{Cov}^{1/3}(X_1, X_j) < \infty,$$

and

$$(2.4) \quad \sum_{j=2}^{\infty} \text{Cov}^{1/3}(Y_1, Y_j) < \infty.$$

This would imply

$$(2.5) \quad \sum_{j=2}^{\infty} \text{Cov}(X_1, X_j) < \infty,$$

and

$$(2.6) \quad \sum_{j=2}^{\infty} \text{Cov}(Y_1, Y_j) < \infty.$$

THEOREM 2.3. (Peligard and Suresh (1995)) *Let $\{X_n, n \geq 1\}$ be a stationary associated sequence of random variables with $E(X_1) = \mu$, $E(X_1^2) < \infty$. Let $\{\ell_n, n \geq 1\}$ be a sequence of positive integers with $1 \leq \ell_n \leq n$. Let $S_j(k) = \sum_{i=j+1}^{j+k} X_i$, $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Let $\ell_n = o(n)$ as $n \rightarrow \infty$. Assume that (2.5) holds. Then, with $\ell = \ell_n$*

$$(2.7) \quad B_n = \frac{1}{n - \ell} \left(\sum_{j=0}^{n-\ell} \frac{|S_j(\ell) - \ell \bar{X}_n|}{\sqrt{\ell}} \right) \\ \rightarrow \left(\text{Var}(X_1) + 2 \sum_{i=2}^{\infty} \text{Cov}(X_1, X_i) \right) \sqrt{\frac{2}{\pi}} \quad \text{in } L_2\text{-mean as } n \rightarrow \infty.$$

In addition assume that $\ell_n = O(n/(\log n)^2)$ as $n \rightarrow \infty$, the convergence above holds in the almost sure sense.

THEOREM 2.4. (Roussas (1993)) *Let $\{X_n, n \geq 1\}$ be a stationary associated sequence of random variables with bounded one dimensional probability density function. Suppose*

$$(2.8) \quad u(n) = 2 \sum_{j=n+1}^{\infty} \text{Cov}(X_1, X_j) \\ = O(n^{-(s-2)/2}) \quad \text{for some } s > 2.$$

Let ψ_n be any positive norming factor. Then, for any bounded interval $I_M = [-M, M]$, we have

$$(2.9) \quad \sup_{x \in I_M} \psi_n |F_n(x) - F(x)| \rightarrow 0,$$

almost surely as $n \rightarrow \infty$, provided

$$(2.10) \quad \sum_{n=1}^{\infty} n^{-s/2} \psi_n^{s+2} < \infty.$$

3. Wilcoxon statistic

The Wilcoxon two-sample statistic is the U-statistic given by

$$(3.1) \quad U = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \phi(Y_j - X_i),$$

where

$$\phi(u) = \begin{cases} 1 & \text{if } u > 0, \\ 0 & \text{if } u = 0, \\ -1 & \text{if } u < 0. \end{cases}$$

Note that ϕ is a kernel of degree (1, 1) with $E\phi(Y - X) = \gamma$. We now obtain the limiting distribution of the statistic U under some conditions. Let

$$(3.2) \quad \sigma_X^2 = 4 \int_{-\infty}^{\infty} G^2(x) dF(x) - 4 \int_{-\infty}^{\infty} G(x) dF(x) + 1 + 8 \sum_{j=2}^{\infty} \text{Cov}(G(X_1), G(X_j))$$

and

$$(3.3) \quad \sigma_Y^2 = 4 \int_{-\infty}^{\infty} F^2(x) dG(x) - 4 \int_{-\infty}^{\infty} F(x) dG(x) + 1 + 8 \sum_{j=2}^{\infty} \text{Cov}(F(Y_i), F(Y_j)).$$

THEOREM 3.1. *Let $\{X_i, i \geq 1\}$ and $\{Y_j, j \geq 1\}$ be independent sequences of random variables with one dimensional distribution functions F and G , respectively, such that each sequence is stationary associated satisfying conditions (2.2) to (2.4). Then, as $m, n \rightarrow \infty$ such that $\frac{m}{n} \rightarrow c \in (0, \infty)$, we have*

$$\sqrt{m}(U - \gamma) \xrightarrow{\mathcal{L}} N(0, A^2) \quad \text{as } n \rightarrow \infty,$$

where

$$(3.4) \quad A^2 = \sigma_X^2 + c\sigma_Y^2.$$

PROOF. Following Hoeffding's decomposition (Lee (1990)), we can write U as

$$(3.5) \quad U = \gamma + H_{m,n}^{(1,0)} + H_{m,n}^{(0,1)} + H_{m,n}^{(1,1)},$$

where

$$\begin{aligned} H_{m,n}^{(1,0)} &= \frac{1}{m} \sum_{i=1}^m h^{(1,0)}(X_i), \\ h^{(1,0)}(x) &= \phi_{10}(x) - \gamma, \quad \phi_{10}(x) = 1 - 2G(x), \\ H_{m,n}^{(0,1)} &= \frac{1}{n} \sum_{j=1}^n h^{(0,1)}(Y_j), \\ h^{(0,1)}(y) &= \phi_{01}(y) - \gamma, \quad \phi_{01}(y) = 2F(y) - 1, \end{aligned}$$

and

$$H_{m,n}^{(1,1)} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n h^{(1,1)}(X_i, Y_j),$$

where

$$h^{(1,1)}(x, y) = \phi(x - y) - \phi_{10}(x) - \phi_{01}(y) + \gamma.$$

It is easy to see that

$$\begin{aligned} E(\phi_{10}(X)) &= \gamma, \\ E(\phi_{10}^2(X)) &= 4 \int_{-\infty}^{\infty} G^2(x) dF(x) - 4 \int_{-\infty}^{\infty} G(x) dF(x) + 1, \end{aligned}$$

and

$$(3.6) \quad \text{Cov}(\phi_{10}(X_i), \phi_{01}(X_j)) = 4 \text{Cov}(G(X_i), G(X_j)).$$

Since the random variables X_1, \dots, X_m are associated, so are $\phi_{10}(X_1), \dots, \phi_{10}(X_m)$ since ϕ is monotone (see, Esary *et al.* (1967)). Furthermore conditions (2.2), (2.5) and (2.6) imply that

$$\sum_{j=2}^{\infty} \text{Cov}(G(X_1), G(X_j)) < \infty,$$

and

$$\sum_{j=2}^{\infty} \text{Cov}(F(Y_1), F(Y_j)) < \infty,$$

since

$$|\text{Cov}(G(X_1), G(X_j))| < \left(\sup_x g \right) \text{Cov}(X_1, X_j),$$

and

$$|\text{Cov}(F(Y_1), F(Y_j))| < \left(\sup_x f \right) \text{Cov}(Y_1, Y_j),$$

by Newman's inequality (1980). Following Newman (1980, 1984), we get that

$$(3.7) \quad m^{-1/2} \sum_{i=1}^m (\phi_{10}(X_i) - \gamma) \xrightarrow{\mathcal{L}} N(0, \sigma_X^2) \quad \text{as } n \rightarrow \infty.$$

Similarly, we see that

$$(3.8) \quad n^{-1/2} \sum_{j=1}^n (\phi_{01}(Y_j) - \gamma) \xrightarrow{\mathcal{L}} N(0, \sigma_Y^2) \quad \text{as } n \rightarrow \infty.$$

Note that $E(H_{m,n}^{(1,1)}) = 0$. Consider

$$(3.9) \quad \begin{aligned} \text{Var}(H_{m,n}^{(1,1)}) &= E(H_{m,n}^{(1,1)})^2 \\ &= \frac{\Delta}{m^2 n^2}, \end{aligned}$$

where

$$(3.10) \quad \Delta = \sum_{i=1}^m \sum_{j=1}^n \sum_{i'=1}^m \sum_{j'=1}^n \Delta(i, j; i', j'),$$

and

$$(3.11) \quad \Delta(i, j; i', j') = \text{Cov}(h^{(1,1)}(X_i, Y_j), h^{(1,1)}(X_{i'}, Y_{j'})).$$

Following Serfling (1968),

$$(3.12) \quad \begin{aligned} \Delta(i, j; i', j') &= 4(E(F_{i,i'}(Y_j, Y_{j'}) - F(Y_j)F(Y_{j'})) \\ &\quad - \text{Cov}(G(X_i, X_{i'}))) \\ &= 4(E(G_{j,j'}(X_i, X_{i'}) - G(X_i)G(X_{i'})) \\ &\quad - \text{Cov}(F(Y_j, Y_{j'}))), \end{aligned}$$

where $F_{i,i'}$ is the joint distribution function of $(X_i, X_{i'})$ and $G_{j,j'}$ is the joint distribution function of $(Y_j, Y_{j'})$.

Then, by Theorem 2.1, there exists a constant $C > 0$ such that

$$(3.13) \quad \begin{aligned} \Delta(i, j; i', j') &\leq C[\text{Cov}^{1/3}(X_i, X_{i'}) + \text{Cov}(X_i, X_{i'})] \\ &= r_1(|i - i'|) \quad (\text{say}), \end{aligned}$$

by stationarity and

$$(3.14) \quad \begin{aligned} \Delta(i, j; i', j') &\leq C[\text{Cov}^{1/3}(Y_j, Y_{j'}) + \text{Cov}(Y_j, Y_{j'})] \\ &= r_2(|j - j'|) \quad (\text{say}), \end{aligned}$$

by stationarity. Note that

$$(3.15) \quad \sum_{k=1}^{\infty} r_1(k) < \infty, \quad \sum_{k=1}^{\infty} r_2(k) < \infty,$$

by (2.3)–(2.6). Then, following Serfling (1968), we have

$$(3.16) \quad \Delta = o(mn^2)$$

as m and $n \rightarrow \infty$ such that $\frac{m}{n}$ has a limit $c \in (0, \infty)$.

Hence, from (3.4), we have

$$(3.17) \quad \begin{aligned} \sqrt{m}(U - \gamma) &= \sqrt{m} \frac{1}{m} \sum_{i=1}^m h^{(1,0)}(X_i) + \sqrt{\frac{m}{n}} \frac{1}{\sqrt{n}} \sum_{j=1}^n h^{(0,1)}(Y_j) + \sqrt{m} H_{m,n}^{(1,1)} \\ &\xrightarrow{\mathcal{L}} N(0, A^2), \end{aligned}$$

since $E(H_{m,n}^{(1,1)}) = 0$ and $\text{Var}(\sqrt{m}H_{m,n}^{(1,1)}) \rightarrow 0$ as $m, n \rightarrow \infty$ such that $\frac{m}{n} \rightarrow c \in (0, \infty)$. This completes the proof of the theorem.

COROLLARY 3.1. *Suppose the conditions of Theorem 3.1 hold. If $F = G$, then,*

$$(3.18) \quad \begin{aligned} \sigma_X^2 &= \sigma_Y^2 \\ &= 4 \left(\frac{1}{12} + 2 \sum_{j=2}^{\infty} \text{Cov}(F(X_1), F(X_j)) \right). \end{aligned}$$

Then, as $m, n \rightarrow \infty$ such that $\frac{m}{n} \rightarrow c \in (0, \infty)$, we have

$$\sqrt{m}(U - \gamma) \xrightarrow{L} N(0, A^2) \quad \text{as } n \rightarrow \infty,$$

where

$$(3.19) \quad A^2 = 4(1 + c) \left(\frac{1}{12} + 2 \sum_{j=2}^{\infty} \text{Cov}(F(X_1), F(X_j)) \right).$$

Estimation of the limiting variance

Note that the limiting variance A^2 depends on the unknown distribution F even under the null hypothesis. We need to estimate it so that the proposed test statistic can be used for testing purposes. The unknown variance A^2 can be estimated using the estimators given by Peligard and Suresh (1995). We now give a consistent estimator of the unknown variance A^2 under some conditions.

Let $N = m + n$. Under the hypothesis $F = G$, the random variables $X_1, \dots, X_m, Y_1, \dots, Y_n$ are associated with the one-dimensional marginal distribution function F . Denote Y_1, \dots, Y_n as X_{m+1}, \dots, X_N . Then X_1, \dots, X_N are associated as independent sets of associated random variables are associated (cf. Esary *et al.* (1967)).

Let $\{\ell_N, N \geq 1\}$ be a sequence of positive integers with $1 \leq \ell_N \leq N$. Let $S_j(k) = \sum_{i=j+1}^{j+k} \phi_{10}(X_i)$, $\bar{\phi}_N = \frac{1}{N} \sum_{i=1}^N \phi_{10}(X_i)$. Define $\ell = \ell_N$ and

$$(3.20) \quad B_N = \frac{1}{N - \ell} \left[\sum_{j=0}^{N-\ell} \frac{|S_j(\ell) - \ell \bar{\phi}_N|}{\sqrt{\ell}} \right].$$

Note that B_N depends on the unknown function F . Let $\hat{\phi}_{10}(x) = 1 - 2F_N(x)$ where F_N is the empirical distribution function corresponding to F based on the associated random variables X_1, \dots, X_N . Let $\hat{S}_j(k)$, $\hat{\phi}_N$ and \hat{B}_N be expressions analogous to $S_j(k)$, $\bar{\phi}_N$ and B_N with ϕ_{10} replaced by $\hat{\phi}_{10}$. Let $Z_i = \phi_{10}(X_i) - \hat{\phi}_{10}(X_i)$. Then

$$(3.21) \quad |B_N - \hat{B}_N| = \left| \frac{1}{N - \ell} \sum_{j=0}^{N-\ell} \frac{|S_j(\ell) - \ell \bar{\phi}|}{\sqrt{\ell}} - \frac{1}{N - \ell} \sum_{j=0}^{N-\ell} \frac{|\hat{S}_j(\ell) - \ell \hat{\phi}|}{\sqrt{\ell}} \right|$$

$$\leq \frac{1}{(N - \ell)\sqrt{\ell}} \sum_{j=0}^{N-\ell} |S_j(\ell) - \hat{S}_j(\ell) - \ell(\bar{\phi} - \hat{\phi})|$$

$$= \frac{1}{(N - \ell)\sqrt{\ell}} \sum_{j=0}^{N-\ell} \left| \sum_{i=j+1}^{j+\ell} Z_i - \ell \frac{1}{N} \sum_{i=1}^N Z_i \right|$$

$$\leq \frac{1}{(N - \ell)\sqrt{\ell}} \sum_{j=0}^{N-\ell} \left\{ \sum_{i=j+1}^{j+\ell} |Z_i| + \ell \frac{1}{N} \sum_{i=1}^N |Z_i| \right\}.$$

Note that

$$|Z_i| = 2|F_N(X_i) - F(X_i)|.$$

Suppose that the density function corresponding to F has a bounded support. Then, for sufficiently large $M > 0$, with probability 1,

$$\begin{aligned}
 (3.22) \quad & \sup_{x \in R} |F_N(x) - F(x)| \\
 &= \max \left\{ \sup_{x \in [-M, M]} |F_N(x) - F(x)|, \sup_{x \in [-M, M]^c} |F_N(x) - F(x)| \right\} \\
 &= \sup_{x \in [-M, M]} |F_N(x) - F(x)|.
 \end{aligned}$$

Hence, from (3.21) and Theorem 2.4 we get

$$\begin{aligned}
 (3.23) \quad |B_N - \hat{B}_N| &\leq \frac{2}{(N - \ell)\sqrt{\ell}} (N - \ell)\ell \sup_x |F_N(x) - F(x)| \\
 &= 2\sqrt{\ell}\psi_N^{-1} \sup_x \psi_N |F_N(x) - F(x)| \\
 &\rightarrow 0 \quad \text{as } N \rightarrow \infty
 \end{aligned}$$

provided $\sqrt{\ell}\psi_N^{-1} = O(1)$ or $\ell_N = O(\psi_N^2)$. Therefore we get,

$$(3.24) \quad |B_N - \hat{B}_N| \rightarrow 0 \quad \text{a.s. as } n \rightarrow \infty.$$

Hence, from Theorem 2.3,

$$(3.25) \quad \frac{\pi}{2} \hat{B}_N^2 \rightarrow 4 \left(\frac{1}{12} + 2 \sum_{j=2}^{\infty} \text{Cov}(F(X_1), F(X_j)) \right)$$

as $n \rightarrow \infty$. Define $J_N^2 = (1 + c)\frac{\pi}{2} \hat{B}_N^2$.

Then,

$$\frac{\sqrt{N}(U - \gamma)}{J_N} \xrightarrow{\mathcal{L}} N(0, 1) \quad \text{as } m, n \rightarrow \infty \quad \text{such that } \frac{m}{n} \rightarrow c \in (0, \infty); \quad \text{as } n \rightarrow \infty.$$

Hence the statistic $\frac{\sqrt{N}(U - \gamma)}{J_N}$ can be used as a test statistic for testing $H'_0 : \gamma = 0$ against $H'_1 : \gamma > 0$.

On the other hand, by using Newman's inequality, one could obtain an upper bound on A^2 given by

$$(3.26) \quad 4(1 + c) \left(\frac{1}{12} + 2 \sum_{j=2}^{\infty} \text{Cov}(X_1, X_j) \right)$$

and we can have conservative tests and estimates of power based on (3.26).

Acknowledgements

We thank the referees for their suggestions.

REFERENCES

- Bagai, I. and Prakasa Rao, B. L. S. (1991). Estimation of the survival function for stationary associated processes, *Statist. Probab. Lett.*, **12**, 385–391.
- Esary, J., Proschan, F. and Walkup, D. (1967). Association of random variables with applications, *Ann. Math. Statist.*, **38**, 1466–1474.
- Lee, A. J. (1990). *U-Statistics*, Marcel Dekker, New York.
- Louhichi, S. (2000). Weak convergence for empirical processes of associated sequences, *Ann. Inst. H. Poincaré Probab. Statist.*, **36**, 547–567.
- Newman, C. M. (1980). Normal fluctuations and the FKG inequalities, *Comm. Math. Phys.*, **74**, 119–128.
- Newman, C. M. (1984). Asymptotic independence and limit theorems for positively and negatively dependent random variables, *Inequalities in Statistics and Probability* (ed. Y. L. Tong), 127–140, IMS, Hayward.
- Peligiard, M. and Suresh, R. (1995). Estimation of variance of partial sums of an associated sequence of random variables, *Stochastic Process. Appl.*, **56**, 307–319.
- Prakasa Rao, B. L. S. and Dewan, I. (2001). Associated sequences and related inference problems, *Handbook of Statistics, 19, Stochastic Processes: Theory and Methods* (eds. C. R. Rao and D. N. Shanbag), 693–728, North Holland, Amsterdam.
- Roussas, G. G. (1993). Curve estimation in random field of associated processes, *J. Nonparametr. Statist.*, **2**, 215–224.
- Roussas, G. G. (1999). Positive and negative dependence with some statistical applications, *Asymptotics, Nonparametrics and Time Series* (ed. S. Ghosh), 757–788, Marcel Dekker, New York.
- Serfling, R. J. (1968). The Wilcoxon two-sample statistic on strongly mixing processes, *Ann. Math. Statist.*, **39**, 1202–1209.