# Asymptotic Properties of the Wilcoxon-Mann-Whitney Statistics

## Shomurodov Nozmbek * To'rayev Alimardon

### Abstract

Random variables seen in many practical problems of statistical physics, quantum field theory, and reliability theory are associated connected random variables. This article focuses on nonparametric estimates for statistics constructed by associated random variables. It proves a theorem for a sequence of stationary associated random variables with two identical marginal distributions.

*Keywords:* associated; statistics; stationary; sequence.

## 1. INTRODUCTION

It is well-known that independent random variables have been extensively studied in science. However, in nature and technology, random variables are often dependent. Therefore, the study of dependent random variables, specifically associated random variables, under certain conditions, and demonstrating their applications in practical problems has been the focus of many prominent experts. In this field, renowned mathematicians such as Newman, Prakasa Rao, Harris, Fortuin, Lebowitz, Hoeffding, Wilcoxon, Mann, Whitney, and their students have achieved fundamental results. Currently, with the development of several directions in mathematical statistics, the importance of the theory of associated random variables has significantly increased, which is well-known among specialists. The topic of this master's thesis is dedicated to gathering future-relevant results on associated random variables, which have been relatively less studied compared to dependent variables, and to studying nonparametric estimators for statistics constructed based on associated variables.

*1.1. Preliminaries*

**Definition 1.1.** Let $(X, Y)$ be a random vector with $E[X^2] < \infty$ and $E[Y^2] < \infty$. Define

$$H(x, y) = P(X \le x, Y \le y) - P(X \le x)P(Y \le y). \tag{1.1}$$

Recall the Hoeffding identity [8]:

$$\text{cov}(X, Y) = \int_{\mathbb{R}^2} H(x, y) \, dx \, dy. \tag{1.2}$$

This identity was extended to the multivariate case by Block and Fang (1988) using the concept of cumulants for random vectors. Yu (1993) generalized Newman's (1984) earlier work by extending the covariance identity to

absolutely continuous functions of the components of the random vector $X$. Cuesta-Molina (1992) generalized the Hoeffding identity to semi-monotonic functions $K(\cdot)$ in the following form:

$$K(x', y') - K(x, y') - K(x', y) + K(x, y) \geq 0 \quad [9]$$

for all $x \leq x'$ and $y \leq y'$. This was proven as:

$$E[K(X, Y)] - E[K(X^*, Y^*)] = \int_{\mathbb{R}^2} H(x, y) K(dx, dy),$$

where $X^*$ and $Y^*$ are independent random variables with the same marginal distributions as $X$ and $Y$, respectively. These results were further generalized by Yu (1993), Cuesta-Molina (1992), and Prakasa Rao (1998) to the multivariate case. Cuadras (2002) showed that if $\alpha(x)$ and $\beta(y)$ are functions with finite variation, then:

$$\text{cov}(\alpha(X), \beta(Y)) = \int_{\mathbb{R}^2} H(x, y) \alpha(dx) \beta(dy).$$

This result is a special case of (1.2.). From this, we can see that $\text{cov}(X_1, X_n) \to 0$ as $n \to \infty$. In particular, we have:

$$\sup_n |\text{cov}(X_1, X_n)| < \infty.$$

Using the association property of $X_1, \ldots, X_n$, we observe that $\text{cov}(X_1, X_n) > 0$ and obtain:

$$0 \leq \text{cov}(X_1, X_j) = [\text{cov}(X_1, X_j)]^{2/3} [\text{cov}(X_1, X_j)]^{1/3} \leq [\sup \text{cov}(X_1, X_n)]^{2/3} [\text{cov}(X_1, X_j)]^{1/3}.$$

Therefore,

$$\sum_{j=2}^{n} \text{cov}(X_1, X_j) \leq [\sup \text{cov}(X_1, X_n)]^{2/3} \sum_{j=2}^{n} [\text{cov}(X_1, X_j)]^{1/3} < \infty. \tag{1.3}$$

Let $R_1, R_2, \ldots, R_n$ be the ranks of $X_1, X_2, \ldots, X_n$. The Wilcoxon signed-rank statistic is defined as...

## 2. Main part

Let $\{X_n, n \geq 1\}$ be a sequence of stationary random variables. We can express $T$ as a linear combination of two U-statistics (Hettmansperger (1984)):

$$T = nU_n^{(1)} + \binom{n}{2} U_n^{(2)}, \tag{2.1}$$

where

$$nU_n^{(1)} = \sum_{i=1}^{n} \phi(X_i),$$

$$\binom{n}{2} U_n^{(2)} = \sum_{1 \leq i < j \leq n} \psi(X_i, X_j), \tag{2.2}$$

and

$$\psi(x, y) = I(x + y > 0). \tag{2.3}$$

For a stationary sequence $\{X_n, n \geq 1\}$, we have:

$$E(U_n^{(2)}) = \frac{1}{\binom{n}{2}} \sum_{1 \le i < j \le n} P_{ij} = \frac{1}{\binom{n}{2}} \sum_{j=2}^{n} (n - j + 1) p_{1,j},$$

where $p_{ij} = P[X_i + X_j > 0]$. Define:

$$\theta = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \psi(x, y) dF(x) dF(y),$$

$$\theta = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \psi(x, y) dF(x) dF(y) = 1 - \int_{-\infty}^{\infty} F(-x) dF(x),$$

$$\psi_1(x_1) = E(\psi(x_1, x_2)) = \int_{-\infty}^{\infty} \psi(x_1, x_2) dF(x_2) = 1 - F(-x_1).$$

Then,

$$h^{(1)}(x_1) = \psi_1(x_1) - \theta, \tag{2.4}$$

and

$$h^{(2)}(x_1, x_2) = \psi(x_1, x_2) - \psi_1(x_1) - \psi_1(x_2) + \theta = \psi(x_1, x_2) + F(-x_1) + F(-x_2) - 2 + \theta. \tag{2.5}$$

The Hoeffding decomposition (H-decomposition) for $U_n^{(2)}$ is given by (Lee (1990)):

$$U_n^{(2)} = \theta + 2H_n^{(1)} + H_n^{(2)}, \tag{2.6}$$

where $H_n^{(j)}$ is the $j$-th degree U-statistic based on the kernel $h^{(j)}$, $j = 1, 2$:

$$H_n^{(j)} = \frac{1}{\binom{n}{j}} \sum h^{(j)}(X_{i_1}, \dots, X_{i_j}), \tag{2.7}$$

with the sum taken over all subsets $1 \le i_1 < \dots < i_j \le n$ of $\{1, \dots, n\}$.

## 2.1. Variance Decomposition

Here, the sum is taken over all subsets $\{1, \dots, n\}$ for $1 \le i_1 < \dots < i_j \le n$. Considering the H-decomposition, we obtain the following:

$$\text{Var}\left(U_n^{(2)}\right) = 4\text{Var}\left(H_n^{(1)}\right) + \text{Var}\left(H_n^{(2)}\right) + 4\text{Cov}\left(H_n^{(1)}, H_n^{(2)}\right). \tag{2.8}$$

Now, consider the following equality (Dewan and Prakasa Rao (2001)):

$$\text{Var}\left(H_n^{(1)}\right) = \frac{1}{n}\left(\sigma_1^2 + 2\sum_{j=2}^{\infty} \sigma_{1j}^2\right) + o\left(\frac{1}{n}\right), \tag{2.9}$$

where

$$\sigma_1^2 = \text{Var}\left(F(-X_1)\right),$$

$$\sigma_{1j}^2 = \text{Cov}\left(F(-X_1), F(-X_{1+j})\right). \tag{2.10}$$

Using Newman's inequality and (1.3), we can write:

$$\sum_{j=2}^{\infty} \sigma_{1j}^2 = \sum_{j=2}^{\infty} \text{Cov}\left(F(-X_1), F(-X_{1+j})\right) < \infty.$$

Additionally,

$$\text{Var}\left(H_n^{(2)}\right) = \binom{n}{2}^{-2} \sum_{1 \le i < j \le n} \sum_{1 \le k < l \le n} \text{Cov}\left\{h^{(2)}(X_i, X_j), h^{(2)}(X_k, X_l)\right\},$$

where

$$\text{Cov}\left\{h^{(2)}(X_i, X_j), h^{(2)}(X_k, X_l)\right\} = \text{Cov}\left\{\psi(X_i, X_j), \psi(X_k, X_l)\right\} +$$

$$+\text{Cov}\left\{\psi(X_i, X_j), F(-X_k)\right\} + \text{Cov}\left\{\psi(X_i, X_j), F(-X_l)\right\} +$$
$$+\text{Cov}\left\{\psi(X_k, X_l), F(-X_k)\right\} + \text{Cov}\left\{\psi(X_k, X_l), F(-X_l)\right\} +$$

$$+\text{Cov}\left\{F(-X_k), F(-X_l)\right\} + \text{Cov}\left\{F(-X_k), F(-X_l)\right\}.$$

Using Newman's (1980) inequality, we obtain: Using Newman's (1980) inequality, we obtain:

$$|\text{Cov}\left(F(-X_k), F(-X_l)\right)| \le \sup_x \left(f(x)\right)^2 \text{Cov}\left(X_k, X_l\right). \tag{2.11}$$

Due to the boundedness of the density function, the following result from Bagai and Prakasa Rao (1991) holds:

$$\left|\text{Cov}\left(\psi(X_i, X_j), \psi(X_k, X_l)\right)\right| =$$

$$= \left|P\left[X_i + X_j > 0, X_k + X_l > 0\right] - P\left[X_i + X_j > 0\right]P\left[X_k + X_l > 0\right]\right| \le$$

$$\le C\left[\text{Cov}\left(X_i + X_j, X_k + X_l\right)\right]^{1/3} =$$

$$= C\left[\text{Cov}\left(X_i, X_k\right) + \text{Cov}\left(X_j, X_k\right) + \text{Cov}\left(X_i, X_l\right) + \text{Cov}\left(X_j, X_l\right)\right]^{1/3}. \tag{2.12}$$

Let $Z = X_i + X_j$. Then, $\psi(X_i, X_j) = I(X_i + X_j > 0) = I(Z > 0)$. Note that this function has a jump at $z = 0$. From equation (2.12), we can conclude that:

$$\left|\text{Cov}\left(\psi(X_i, X_j), F(X_k)\right)\right| =$$

$$= \left|\int_{-\infty}^{\infty} \left(P\left[X_i + X_j \le 0, X_k \le x\right] - P\left[X_i + X_j \le 0\right]P\left[X_k \le x\right]\right)dF(x)\right| \le$$

$$\le \int_{-\infty}^{\infty} |P[X_i + X_j \le 0, X_k \le x] - P[X_i + X_j \le 0]P[X_k \le x]|dF(x) \le$$

$$\le C\int_{-\infty}^{\infty} [\text{cov}(X_i + X_j, X_k)]^{1/3}dF(x) = C[\text{cov}(X_i + X_j, X_k)]^{1/3} =$$

$$= C[\text{cov}(X_i, X_k) + \text{cov}(X_j, X_k)]^{1/3}.$$

Using equations (2.11), (2.12), and (2.13) in equation (2.10), we obtain the following:

$$\left| \text{cov} \left\{ h^{(2)}(X_i, X_j), h^{(2)}(X_k, X_l) \right\} \right| \leq$$

$$\leq C \left[ \text{cov}(X_i, X_k) + \text{cov}(X_j, X_k) + \text{cov}(X_i, X_l) + \text{cov}(X_j, X_l) \right]^{\frac{1}{3}} +$$

$$+ \left[ \text{cov}(X_i, X_k) + \text{cov}(X_j, X_k) \right]^{\frac{1}{3}} + \left[ \text{cov}(X_i, X_l) + \text{cov}(X_j, X_l) \right]^{\frac{1}{3}} +$$

$$+ \left[ \text{cov}(X_k, X_i) + \text{cov}(X_l, X_i) \right]^{\frac{1}{3}} + \left[ \text{cov}(X_k, X_j) + \text{cov}(X_l, X_j) \right]^{\frac{1}{3}} +$$

$$+ \text{cov}(X_i, X_k) + \text{cov}(X_j, X_k) + \text{cov}(X_i, X_l) + \text{cov}(X_j, X_l) \leq$$

$$\leq C \left[ \text{cov}(X_i, X_k) + \text{cov}(X_j, X_k) + \text{cov}(X_i, X_l) + \text{cov}(X_j, X_l) \right]^{\frac{1}{3}} +$$

$$+ \left[ \text{cov}(X_k, X_i) + \text{cov}(X_l, X_i) \right]^{\frac{1}{3}} + \left[ \text{cov}(X_k, X_j) + \text{cov}(X_l, X_j) \right]^{\frac{1}{3}} +$$

$$= C \left[ \text{cov}(X_i, X_k) + \text{cov}(X_j, X_k) + \text{cov}(X_i, X_l) + \text{cov}(X_j, X_l) \right]^{\frac{1}{3}} +$$

$$= \left[ \text{cov}(X_k, X_i) + \text{cov}(X_l, X_i) \right]^{\frac{1}{3}} + \left[ \text{cov}(X_k, X_j) + \text{cov}(X_l, X_j) \right]^{\frac{1}{3}}$$

$$r = (|i - k|) + r(|j - k|) + r(|i - l|) + r[|j - l|]$$
$$\sum r(k) < \infty.$$

Therefore, from Serfling's (1968) theorem, we obtain as $n \to \infty$

$$operatornamevar\left( H_n^{(2)} \right) = o\left( \frac{1}{n} \right). \tag{2.13}$$

Using the Cauchy-Schwarz inequality, the following follows

$$\text{cov}\left( H_n^{(1)}, H_n^{(2)} \right) = o\left( \frac{1}{n} \right).$$

Using equations (2.8), (2.9), (2.16), and (2.17), we can write

$$\text{var}\left( U_n^{(2)} \right) = 4 \left[ \sigma_1^2 + 2 \sum_{j=1}^{\infty} \sigma_{1j}^2 \right] + o\left( \frac{1}{n} \right).$$

To obtain the limit distribution of the U-statistic, we introduce the following theorem.

**Theorem 2.1.** *Let $\{X_n, n \geq 1\}$ be a sequence of associated random variables. Suppose*

$$\sum_{k=1}^{\infty} r(k) < \infty \text{ holds. Then, as } n \to \infty,$$

$$n^{1/2}\left( U_n^{(2)} - \theta \right) \xrightarrow{e} N(0, 1)$$

*where $\sigma_U^2 = \sigma_1^2 + 2 \sum_{j=1}^{\infty} \sigma_{1j}^2$.*

Proof. Here, we also use relations (2.2)–(2.12), and make appropriate modifications in the remaining relations: Using Newman's (1980) inequality, we obtain

$$|\text{cov}\left(F(-X_i), F(-X_k)\right)| \leq \sup_x (f(x))^2 \text{cov}(X_i, X_k). \tag{2.14}$$

Due to the boundedness of the density function, the following result follows from Bagai and Prakasa Rao's (1991) theorem:

$$\left|\text{cov}\left(\psi(X_i, X_j), \psi(X_k, X_l)\right)\right| =$$

$$= \left|P\left[X_i + X_j > 0, X_i + X_k > 0\right] - P\left[X_i + X_j > 0\right] P\left[X_i + X_k > 0\right]\right| \leq$$

$$\leq Cr\left(|i - l|\right).$$

Let $Z = X_i + X_j$. Note that the function $\psi(X_i, X_j) = I(X_i + X_j > 0) = I(z > 0)$ has a discontinuity at $z = 0$. Now, from equation (2.12), it follows that

$$\left|\text{cov}\left(\psi(X_i, X_j), F(X_k)\right)\right| =$$

$$= \left|\int_{-\infty}^{\infty} \left(P\left[X_i + X_j \leq 0, X_k \leq x\right] - P\left[X_i + X_j \leq 0\right] P\left[X_k \leq x\right]\right) dF(x)\right| \leq$$

$$\leq \int_{-\infty}^{\infty} \left|P\left[X_i + X_j \leq 0, X_k \leq x\right] - P\left[X_i + X_j \leq 0\right] P\left[X_k \leq x\right]\right| dF(x) \leq$$

$$= C\left[r\left(|i - k|\right) + r\left(|j - k|\right)\right]. \tag{2.15}$$

Using equations (2.14) and (2.15) in equation (2.12)we obtain

$$\left|\text{cov}\left(h^{(2)}(X_i, X_j), h^{(2)}(X_k, X_l)\right)\right| \leq$$

$$= r\left(|i - k|\right) + r\left(|j - k|\right) + r\left(|i - l|\right) + r\left(|j - l|\right).$$

$$\sum_{k=1}^{\infty} r(k) < \infty.$$

From this and Serfling's (1968) theorem, we obtain as $n \to \infty$

$$\text{var}\left(H_n^{(2)}\right) = o\left(\frac{1}{n}\right). \tag{2.16}$$

Using the Cauchy-Schwarz inequality, the following follows

$$\text{cov}\left(H_n^{(1)}, H_n^{(2)}\right) = o\left(\frac{1}{n}\right). \tag{2.17}$$

Using equations (2.8), (2.9), (2.16), and (2.17), we can write

$$\text{var}\left(U_n^{(2)}\right) = 4\left[\sigma_1^2 + 2\sum_{j=1}^{\infty} \sigma_{1j}^2\right] + o\left(\frac{1}{n}\right).$$

## 3. Conclusion

It is well-known that independent random variables have been sufficiently studied in science. However, in nature and technology, random variables are often dependent. Therefore, the study of dependent variables, specifically associated random variables, under certain conditions, and demonstrating their applications in practical problems has been the focus of many prominent experts. The topic of this master's thesis is dedicated to gathering future-relevant results on associated random variables, which have been relatively less studied compared to dependent variables, and to studying nonparametric estimators for statistics constructed based on associated variables.

## References

[1] Bagai, I. and Prakasa Rao, B.L.S. Estimation of the survival function for stationary associated processes, *Statist. Probab. Lett.*, 12, 385-391.(1991)

[2] Baek, Jong Li., Park, Sung Tae., Chung, Sung Mo. and Seo, Hye Young. On the almost sure convergence of weighted sums of negatively associated random variables, *Commun. Korean Math. Soc.*, 20, 539-546. (2005).

[3] Bagai, I. and Prakasa Rao, B.L.S. Kernel-type density and failure rate estimation for associated sequences, *Ann. Inst. Statist. Math.*, 47, 253-266.(1995)

[4] Birkel, T. On the convergence rate in the central limit theorem for associated processes, *Ann. Probab.*, 16, 1689-1698.(1988b)

[5] Bulinski, A.V. Rates of convergence in the central limit theorem for fields of associated random variables, *Theor. Probab. Appl.*, 40, 136-144.(1995)

[6] Matula, P. On almost sure limit theorems for positively dependent random variables, *Statist. Probab. Lett.*, 74, 59-66.(2005)

[7] Newman, C.M. and Wright, A.L. Associated random variables and martingale inequalities, *Z. Wahrsch. Theorie und Verw. Gebiete*, 59, 361-371. (1982)

[8] Hoeffding, W. Masstabinvariante Korrelations-theorie, Schr. Math. Inst.*University Berlin*, 5, 181-233. (1940)

[9] Queseda- Molina, J.J. A generalization of an identity of Hoeffding and some applications, *J. Ital. Statist. Soc.*, 3, 405-411.(1992)

## Affiliations

SHOMURODOV NOZIMBEK

**ADDRESS:** Tashkent State Transport University, Tashkent.

**E-MAIL:** nozimshoh007@gmail.com

**ORCID ID:** https://orcid.org/0009-0005-2748-4484

TURAEV ALIMARDON

**ADDRESS:** Tashkent State Transport University, Tashkent.

**E-MAIL:** alimardontoxirovich0413@gmail.com

**ORCID ID:** https://orcid.org/0000-0002-0584-9738