

RESEARCH

Open Access

# Joint distribution of rank statistics considering the location and scale parameters and its power study

Wan-Chen Lee

Correspondence:  
umlee223@cc.umanitoba.ca  
Department of Statistics, University  
of Manitoba, Winnipeg, Canada

## Abstract

The ranking method used for testing the equivalence of two distributions has been studied for decades and is widely adopted for its simplicity. However, due to the complexity of calculations, the power of the test is either estimated by a normal approximation or found when an appropriate alternative is given. Here, via the Finite Markov chain imbedding technique, we are able to establish the marginal and joint distributions of the rank statistics considering the shift and scale parameters, respectively and simultaneously, under two different continuous distribution functions. Furthermore, the procedures of distribution equivalence tests and their power functions are discussed. Numerical results of a joint distribution of rank statistics under the standard normal distribution and the powers for a sequence of alternative normal distributions with means from -20 to 20 and standard deviations from 1 to 9 and their reciprocal are presented. In addition, we discuss the powers of the rank statistics under the Lehmann alternatives.

**2010 Mathematics Subject Classification:** Primary 62G07; Secondary 62G10

**Keywords:** FMCI; Lehmann alternative; Rank statistic; Rank-sum test; Power

## 1 Introduction

Suppose that on the basis of observations  $X_1, \dots, X_m; Y_1, \dots, Y_n$  from the cumulative distribution functions  $F$  and  $G$ , two major topics in the hypothesis testing are to test the equivalence of either the center or the dispersion of the two populations of interest. The hypotheses are stated, for some  $\theta \neq 0$ ,

$$H_0 : F(x) = G(x) \quad \text{versus} \quad H_a : F(x) = G(x - \theta), \quad \text{for all } x,$$

which is known as the shift alternative and, for some  $\sigma \neq 1$ ,

$$H_0 : F(x) = G(x) \quad \text{versus} \quad H_a : F(x) = G(x\sigma^{-1}), \quad \text{for all } x.$$

Wilcoxon (1945) proposed the ranking method for testing the significance of the difference of the two populations means, also known as the Wilcoxon rank-sum test, and defined a statistic  $W_Y$ , as the sum of the ranks of the  $y$ 's in the combined and ordered sequence of  $x$ 's and  $y$ 's, equivalent to

$$\sum_{j=1}^n \{ \# \text{ of } x'_i < y_j \} + \frac{n(n+1)}{2}.$$

Mann and Whitney (1947) introduced an elaboration of the ranking test, proposed the statistic  $U_X = mn - W_Y + \frac{n(n+1)}{2}$ , and proved that the limiting distribution of the test statistic  $U_X$  is

$$\frac{U_X - E(U_X)}{\sqrt{Var(U_X)}} \xrightarrow{L} N(0, 1)$$

as  $m$  and  $n$  go to infinity in any arbitrary manner where

$$E(U_X) = mnp_1$$

and

$$Var(U_X) = mnp_1(1 - p_1) + mn(n - 1)(p_2 - p_1^2) + mn(m - 1)(p_3 - p_1^2),$$

with

$$\begin{aligned} p_1 &= P(X > Y), \\ p_2 &= P(X > Y \text{ and } X > Y'), \\ p_3 &= P(X > Y \text{ and } X' > Y), \end{aligned} \tag{1}$$

where  $X, X'$  and  $Y, Y'$  are independently distributed,  $X, X'$  with the distribution  $F$ , and  $Y, Y'$  with the distribution  $G$ . Intuitively, the power for the right-sided test can be found as

$$P\left(\frac{U_X - E(U_X)}{\sqrt{Var(U_X)}} > \frac{c - E(U_X)}{\sqrt{Var(U_X)}} \mid H_a\right), \tag{2}$$

where  $c$  is the value such that

$$\Phi\left(\frac{c - \frac{1}{2}mn}{\sqrt{\frac{1}{12}mn(m+n+1)}} \mid H_o\right) \geq 1 - \alpha.$$

Over the years, there have been studies on finding the exact or approximate power for the rank-sum test. By choosing an appropriate alternative distribution function, Shieh et al. (2006) derived the exact power for the uniform, normal, double exponential and exponential shift models. Rosner and Glynn (2009) discussed power against the family of alternatives of the form

$$\Phi^{-1}(F_Y(y)) = \Phi^{-1}(F_X(y)) + \mu \text{ for some } \mu \neq 0,$$

where the underlying distributions  $F_X$  and  $F_Y$  are normal. Collings and Hamilton (1988) presented a bootstrap method to find the empirical distribution functions in order to approximate the power against the shift alternative. Lehmann (1953) derived the power function as

$$P(S_1 = s_1, S_2 = s_2, \dots, S_n = s_n) = \frac{k^n}{\binom{m+n}{m}} \prod_{j=1}^n \frac{\Gamma(s_j + jk - j)}{\Gamma(s_j)} \frac{\Gamma(s_{j+1})}{\Gamma(s_{j+1} + jk - j)},$$

where  $s_j$  is the rank of  $y_j$  in the combined samples for the alternative hypothesis of

$$G_Y(x) = F_X(x)^k, \text{ for all } x,$$

where  $k$  is a positive integer. However, Lehmann (1998) pointed out that the power function of the rank-sum test, Equation (2), was only qualitative. Since the numerical values for

assessing the probabilities in Equation (1) are considerably complicated in computation when  $F$  and  $G$  are continuous distributions with  $F \neq G$ .

As the rank-sum test is widely adopted for testing the center differences of two distributions, it is natural to study the efficiency of a rank-sum test for variability (Ansari and Bradley 1960). For decades, studies have focused on proposing new definitions of the rank statistic and using the methods of Chernoff and Savage to show the relative efficiency of the proposed statistic to the F-test, see for example Mood (1954), Siegel and Tukey (1960), Ansari and Bradley (1960), and Klotz (1962). Ansari and Bradley (1960) mentioned that if the means of the  $X$  and  $Y$  samples cannot be considered equal, differences in location have a severe impact on all the tests of dispersion. Klotz (1962) showed the power of a rank test can be found by integrating the joint density of  $X$  and  $Y$  samples over that part of the  $m + n$  dimensional space defined by the alternative orderings which lie in the critical region of the test, for which conditions are very strict.

Our approach aims at releasing some of the conditions for finding the distribution of the proposed rank statistic. We systematically imbed the random vector  $\mathbf{U}_n$  into a Markov chain to induce the marginal and joint distributions of the rank statistics considering the shift and scale parameter, respectively, under any form of two distribution functions. A joint distribution of rank statistics, to the best of our knowledge, has not been studied in the literature. The main strength of using the finite Markov chain imbedding approach (FMCI) is to derive the distribution of the rank statistic without giving any conditions. Therefore, under the null hypothesis of  $F = G$ , we are able to identify a proper critical region and, under the alternative assumption, the power of the test can be determined naturally. The distribution of the random vector  $\mathbf{U}_n$ , independent of the form of the distribution function  $F$ , is also demonstrated under the null hypothesis of the distribution equivalence.

The main contributions of this paper are as follows. In Section 2.1, we introduce the procedures of deriving the distribution of the rank statistic considering the shift parameter and its power function by using FMCI. The procedures are general and can be applied to either two identical distribution functions of interest or two different continuous density functions. In Section 2.2, we address the steps for finding the distribution of the rank statistic considering the scale parameter and its power function. In Section 2.3, we retrieve the joint distribution of the rank statistics considering the location and scale parameters simultaneously as well as its power function. Numerical results of a joint distribution and some powers of the rank statistics against shift parameter and scale parameter, individually and simultaneously, are presented in Section 3. We also discuss the powers of the rank statistics under the Lehmann alternatives. We end this paper with a short conclusion in Section 4.

## 2 Methods

### 2.1 Distributions of the rank statistic in the shift case

Let  $\{X_1, \dots, X_m\}$  and  $\{Y_1, \dots, Y_n\}$  be two independent samples from the continuous cumulative density distributions  $F(x)$  and  $G(x - \theta)$ , respectively. Given  $\mathbf{x} = \{x_1, \dots, x_m\}$  and  $x_{[i]}$  is the  $i^{th}$  smallest number in the sample, we have

$$p_i = P(x_{[i-1]} < Y < x_{[i]}) = \int_{x_{[i-1]}}^{x_{[i]}} g(y) dy = G(x_{[i]}) - G(x_{[i-1]}),$$

for  $i = 1, 2, \dots, m + 1$  where  $x_{[0]} = -\infty$  and  $x_{[m+1]} = \infty$ . Therefore, we define the sampling distribution of  $Y$  in the  $(m + 1)$  intervals as

$$\begin{aligned} \mathbf{p} &= (G(x_{[1]}) - G(x_{[0]}), \dots, G(x_{[m+1]}) - G(x_{[m]})) \\ &= (p_1, p_2, \dots, p_{m+1}). \end{aligned} \quad (3)$$

Given  $m$ , for  $t = 1, 2, \dots, n$ , let

$$\Omega_t = \left\{ \mathbf{u}_t = (u_1(t), \dots, u_{m+1}(t)) : \sum_{i=1}^{m+1} u_i(t) = t \text{ and } u_i(t) \geq 0, \quad i = 1, \dots, m + 1 \right\},$$

where  $u_i(t)$  is the number of  $y$ 's in the interval  $[x_{[i-1]}, x_{[i]}]$  among  $y_1, \dots, y_t$ . For each  $\mathbf{u}_n = (u_1(n), \dots, u_{m+1}(n))$ , we have a corresponding rank-sum of  $y$ 's in the combined sample

$$R_l(\mathbf{U}_n = \mathbf{u}_n | \mathbf{X}) = \frac{\sum_{i=1}^{m+1} u_i^2(n) + \sum_{i=1}^{m+1} u_i(n)}{2} + \sum_{i=1}^m (u_i(n) + 1) \left( \sum_{j=i+1}^{m+1} u_j(n) \right). \quad (4)$$

**Theorem 1.** *The statistic  $R_l$  is equivalent to the statistic  $W_Y$ , which is addressed by Wilcoxon in 1945.*

*Proof.* Let

$$I(x_i, y_j) = \begin{cases} 1 & \text{if } x_i < y_j \\ 0 & \text{otherwise.} \end{cases}$$

The rank statistic  $W_Y$ , sum of the ranks of  $y$ 's observations, can be determined by

$$\begin{aligned} \sum_{j=1}^n \left( \sum_{i=1}^m I(x_i, y_j) + j \right) &= \sum_{j=1}^n \sum_{i=1}^m I(x_i, y_j) + \sum_{j=1}^n j \\ &= \sum_{i=1}^m \sum_{j=1}^n I(x_i, y_j) + \frac{n(n+1)}{2}. \end{aligned} \quad (5)$$

The first summation of the first term in Equation (5) can be interpreted as the number of  $y$  observations larger than  $x_{[i]}$  which is  $\sum_{j=i+1}^{m+1} u_j(n)$  in our expression. It is not difficult to see that  $\sum_{i=1}^{m+1} u_i(n)$  equals  $n$ , the size of  $y$  sample. Therefore, the equation can be rewritten as

$$\sum_{i=1}^m \left( \sum_{j=i+1}^{m+1} u_j(n) \right) + \frac{\sum_{i=1}^{m+1} u_i(n)^2 + 2 \sum_{i=1}^m u_i(n) \left( \sum_{j=i+1}^{m+1} u_j(n) \right) + \sum_{i=1}^{m+1} u_i(n)}{2}.$$

It is then easy to see that

$$\sum_{i=1}^m (u_i(n) + 1) \left( \sum_{j=i+1}^{m+1} u_j(n) \right) + \frac{\sum_{i=1}^{m+1} u_i(n)^2 + \sum_{i=1}^{m+1} u_i(n)}{2} = R_l.$$

□

Next, we demonstrate that for two random samples from the same population, the distribution of the random vector  $\mathbf{U}_n$  is independent of the form of the distribution function.

**Theorem 2.** *Distribution-free property of  $\mathbf{U}_n$ .*

$$P(\mathbf{U}_n = \mathbf{u}_n | H_0) = \frac{1}{\text{Card}(\Omega_n)} = \frac{1}{\binom{m+n}{n}}. \quad (6)$$

*Proof.* We know the joint PDF of the ordered sample of  $x$ 's is given by

$$f(x_{[1]}, \dots, x_{[m]}) = m! \prod_{i=1}^m f(x_i)$$

and, when  $F = G$ , the conditional probability of the random vector  $\mathbf{U}_n$  given  $\mathbf{X} = (x_1, x_2, \dots, x_m)$  is

$$P(\mathbf{U}_n = \mathbf{u}_n | x_1, x_2, \dots, x_m) = \frac{n!}{\prod_{i=1}^{m+1} u_i(n)!} \prod_{i=1}^{m+1} \left( \int_{x_{[i-1]}}^{x_{[i]}} f(y) dy \right)^{u_i(n)}, \quad (7)$$

where  $x_{[0]} = -\infty$  and  $x_{[m+1]} = \infty$ . By taking the expected value of the conditional probability, we have

$$\begin{aligned} & P(\mathbf{U}_n = \mathbf{u}_n | H_0) \\ &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} P(\mathbf{u}_n | x_1, \dots, x_m) f(x_{[1]}, \dots, x_{[m]}) dx_{[1]} \cdots dx_{[m]} \\ &= \int_{-\infty}^{\infty} \int_{x_{[1]}}^{\infty} \cdots \int_{x_{[m-1]}}^{\infty} \frac{n!}{\prod_{i=1}^{m+1} u_i(n)!} (F(x_{[1]}))^{u_1(n)} (F(x_{[2]}) - F(x_{[1]}))^{u_2(n)} \\ &\quad \cdots (1 - F(x_{[m]}))^{u_{m+1}(n)} m! dF(x_{[1]}) \cdots dF(x_{[m]}). \end{aligned} \quad (8)$$

Using variable transformation, it is clear to see that the random variables  $F(x_{[1]}), \dots, F(x_{[m]})$  have a Dirichlet distribution with parameters  $u_1(n) + 1, u_2(n) + 1, \dots, u_{m+1}(n) + 1$ . Therefore, we have

$$P(\mathbf{U}_n = \mathbf{u}_n | H_0) = \frac{n! m!}{(n+m)!} = \frac{1}{\text{Card}(\Omega_n)}$$

which is independent of the distribution function.  $\square$

This is the reason that the distribution of the rank statistic  $\mathbf{U}_n$  is distribution-free under the null hypothesis. However, the distribution of the random vector  $\mathbf{U}_n$  is discrete uniform with the mass function one over the number of possible outcomes of the random vector  $\mathbf{U}_n$  only when assuming  $F = G$ . In other words, the distribution of the random vector  $\mathbf{U}_n$  can be found by the traditional combinatorial analysis when  $F = G$ . Unfortunately, when  $F \neq G$ , we will not be able to establish the distribution of  $\mathbf{U}_n$  through Equation (7) as solving the multiple integral in Equation (8) is either tedious given some appropriate alternative distribution function or difficult. Our understanding is that finding the power of the test has not been solved in most cases. To overcome this situation, we bring in the finite Markov chain imbedding approach.

Let  $\Omega_t, t = 0, 1, \dots, n$ , be the state space which has

$$\binom{m+t}{t}$$

possible states,  $\Gamma_n = \{0, 1, \dots, n\}$  be an index set, and  $\{Z_t : t \in \Gamma_n\}$  be a non-homogeneous Markov chain on the state space  $\Omega_t$ . As a transition probability matrix  $M_t$  for this chain,  $t = 1, \dots, n$ , consider

$$M_t = \Omega_{t-1} \begin{bmatrix} & \Omega_t \\ p_{u_{t-1}, u_t} & \end{bmatrix}_{\binom{m+t-1}{t-1} \times \binom{m+t}{t}},$$

where

$$\begin{aligned} p_{u_{t-1}, u_t} &= P(Z_t = u_t | Z_{t-1} = u_{t-1}) \\ &= \begin{cases} p_i & \text{if } u_i(t-1) + 1 = u_i(t) \text{ and } u_j(t-1) = u_j(t) \forall j \neq i \\ 0 & \text{otherwise} \end{cases}, \end{aligned}$$

and  $p_i$  is defined in Equation (3).

**Theorem 3.**  $R_l(\mathbf{U}_n | X)$  is finite Markov chain imbeddable, and

$$P(R_l(\mathbf{U}_n) = r | X) = \xi \left( \prod_{t=1}^n M_t \right) \mathbf{B}'(C_r),$$

where  $\mathbf{B}(C_r) = \sum_{k: R_l(\mathbf{U}_n) = r} e_k$ ,  $e_k$  is a  $1 \times \binom{m+n}{n}$  unit row vector corresponding to state  $\mathbf{u}_n$ ,  $\xi (= P(Z_0 = 1) = 1)$  is the initial probability and  $M_t$ ,  $t = 1, \dots, n$ , are the transition probability matrices of the imbedded Markov chain defined on the state space  $\Omega_t$ .

*Proof.* For each  $\mathbf{u}_n = (u_1(n), \dots, u_{m+1}(n))$  in the state space  $\Omega_n$ , we have a corresponding rank  $R_l$  as shown in Equation (4). Intuitively, the minimum rank  $r_{ls}$  is  $n(n+1)/2$  and the maximum rank  $r_{lb}$  is  $n(2m+n+1)/2$ . In accordance with the possible values of the rank  $R_l$ , we define a finite partition  $\{C_r : r = r_{ls}, \dots, r_{lb}\}$  such that

$$P(Z_n \in C_r | \mathbf{p}) = \xi \left( \prod_{t=1}^n M_t \right) \mathbf{B}'(C_r) \quad (9)$$

where  $\mathbf{B}(C_r) = \sum_{k: R_l(\mathbf{U}_n) = r} e_k$ ,  $e_k$  is a  $1 \times \binom{m+n}{n}$  unit row vector corresponding to state  $\mathbf{U}_n$ , we then obtain the conditional probability of the rank  $R_l$ .  $\square$

Then, the Law of Large Numbers is used to determine the probability of  $\mathbf{U}_n$  for any continuous  $F$  and  $G$

$$\frac{1}{N} \sum_{i=1}^N P(\mathbf{U}_n = \mathbf{u}_n | \mathbf{X}_i) \xrightarrow{p} P(\mathbf{U}_n = \mathbf{u}_n)$$

where  $\mathbf{X}_i$  is the  $i^{th}$  sample of size  $m$  from the distribution function  $F$ . It is easy to see that

$$P(R_l(\mathbf{U}_n) = r) = \sum_{\mathbf{u}_n: R_l(\mathbf{u}_n) = r} P(\mathbf{U}_n = \mathbf{u}_n). \quad (10)$$

To test

$$H_0 : F(x) = G(x) \text{ versus } H_a : F(x) = G(x - \theta),$$

for some  $\theta \neq 0$ , the power function is approximated by

$$\begin{aligned}
& P(R_l(\mathbf{U}_n) \leq r_{1\alpha} | H_a) + P(R_l(\mathbf{U}_n) \geq r_{2\alpha} | H_a) \\
&= \sum_{r=r_{ls}}^{r_{1\alpha}} P(R_l(\mathbf{U}_n) = r | H_a) + \sum_{r=r_{2\alpha}}^{r_{lb}} P(R_l(\mathbf{U}_n) = r | H_a) \\
&= \sum_{r=r_{ls}}^{r_{1\alpha}} \sum_{\mathbf{u}_n: R(\mathbf{u}_n)=r} P(\mathbf{U}_n = \mathbf{u}_n | H_a) + \sum_{r=r_{2\alpha}}^{r_{lb}} \sum_{\mathbf{u}_n: R(\mathbf{u}_n)=r} P(\mathbf{U}_n = \mathbf{u}_n | H_a) \\
&\approx \sum_{r=r_{ls}}^{r_{1\alpha}} \sum_{\mathbf{u}_n: R(\mathbf{u}_n)=r} \frac{1}{N} \sum_{i=1}^N P(\mathbf{U}_n | H_a; \mathbf{X}_i) + \sum_{r=r_{2\alpha}}^{r_{lb}} \sum_{\mathbf{u}_n: R(\mathbf{u}_n)=r} \frac{1}{N} \sum_{i=1}^N P(\mathbf{U}_n | H_a; \mathbf{X}_i) \\
&= \frac{1}{N} \left( \sum_{r=r_{ls}}^{r_{1\alpha}} \sum_{i=1}^N \sum_{\mathbf{u}_n: R(\mathbf{u}_n)=r} P(\mathbf{U}_n | H_a; \mathbf{X}_i) + \sum_{r=r_{2\alpha}}^{r_{lb}} \sum_{i=1}^N \sum_{\mathbf{u}_n: R(\mathbf{u}_n)=r} P(\mathbf{U}_n | H_a; \mathbf{X}_i) \right) \\
&= \frac{1}{N} \sum_{i=1}^N \left( \sum_{r=r_{ls}}^{r_{1\alpha}} P(R_l(\mathbf{U}_n) = r | H_a; \mathbf{X}_i) + \sum_{r=r_{2\alpha}}^{r_{lb}} P(R_l(\mathbf{U}_n) = r | H_a; \mathbf{X}_i) \right),
\end{aligned}$$

where

$$P(R_l(\mathbf{U}_n) \leq r_{1\alpha} | H_o) + P(R_l(\mathbf{U}_n) \geq r_{2\alpha} | H_o) \leq \alpha.$$

Note that the alternative hypothesis is subject to the purpose of the test. This simply needs to be slightly modified if a one-sided test is adopted.

## 2.2 Distributions of the rank statistic in the scale case

We studied the distribution and the power function of the rank statistic  $R_l$  considering a shift in location. Now, the distribution and the power function of the rank statistic considering the scale parameter will be addressed. For this purpose, we consider  $F(x) = G(x\sigma^{-1})$  and state the null and alternative hypotheses as

$$H_0 : \sigma = 1 \text{ versus } H_a : \sigma \neq 1.$$

To do so, we begin with the procedure of finding the distribution of the rank statistic, denoted  $R_s$ , considering the scale parameter through the random vector  $\mathbf{U}_n$ . The array of ranks are given by

$$(m+n)/2, \dots, 3, 2, 1, \quad 1, 2, 3, \dots, (m+n)/2;$$

if  $m+n$  is even, and

$$(m+n-1)/2, \dots, 3, 2, 1, \quad 0 \quad 1, 2, 3, \dots, (m+n-1)/2$$

if  $m+n$  is odd. We first introduce how to determine the rank-sum of  $y$ 's observations in the combined samples,  $R_s$ , with respect to

$$\Omega_n = \left\{ \mathbf{u}_n = (u_1(n), \dots, u_{m+1}(n)) : \sum_{i=1}^{m+1} u_i(n) = n \right\}$$

where  $u_i(n)$  means the number of  $y$  observations belonging to  $[x_{[i-1]}, x_{[i]})$ . Let  $med(x, y)$  be the median among  $x$ 's and  $y$ 's and belongs to  $[x_{[i]}, x_{[i+1]})$  which will

then break  $\mathbf{U}_n$  into two parts  $\mathbf{U}_n^-$  and  $\mathbf{U}_n^+$ . If  $m + n$  is odd and  $med(x, y) = x_{[i]}$ , then

$$\mathbf{U}_n^- = (u_1^- = u_i(n), u_2^- = u_{i-1}(n), \dots, u_i^- = u_1(n))$$

is a  $1 \times i$  vector and

$$\mathbf{U}_n^+ = (u_1^+ = u_{i+1}(n), u_2^+ = u_{i+2}(n), \dots, u_{m+1-i}^+ = u_{m+1}(n))$$

is a  $1 \times (m + 1 - i)$  vector. The second possible case is, if  $m + n$  is odd and  $med(x, y) = y_{[\sum_{k=1}^i u_k(n)+j]}$ , then  $\mathbf{U}_n^-$ , a row vector with length  $i + 1$ , has the form

$$(u_1^- = j - 1, u_2^- = u_i(n), \dots, u_{i+1}^- = u_1(n))$$

and  $\mathbf{U}_n^+$ , a row vector with length  $m + 1 - i$ , is given by

$$(u_1^+ = u_{i+1}(n) - j, u_2^+ = u_{i+2}(n), \dots, u_{m+1-i}^+ = u_{m+1}(n)).$$

The third possible case is, if  $m + n$  is even and  $x_{[i]}$  is the smallest number larger than  $med(x, y)$ , the vectors are now defined as

$$\mathbf{U}_n^- = (u_1^- = u_i(n), u_2^- = u_{i-1}(n), \dots, u_i^- = u_1(n))$$

and

$$\mathbf{U}_n^+ = (u_1^+ = 0, u_2^+ = u_{i+1}(n), \dots, u_{m+1-i}^+ = u_{m+1}(n)).$$

The last possibility is, if  $m + n$  is even,  $y_{[\sum_{k=1}^i u_k(n)+j]}$  is the smallest number larger than  $med(x, y)$ . The vectors are now defined as

$$\mathbf{U}_n^- = (u_1^- = j - 1, u_2^- = u_i(n), \dots, u_{i+1}^- = u_1(n))$$

and

$$\mathbf{U}_n^+ = (u_1^+ = u_{i+1}(n) - j + 1, u_2^+ = u_{i+2}(n), \dots, u_{m+1-i}^+ = u_{m+1}(n)).$$

Let  $n^-$  be the length of the vector  $\mathbf{U}_n^-$  and  $n^+$  be the length of the vector  $\mathbf{U}_n^+$ .

**Theorem 4.**  $R_s(\mathbf{U}_n | X)$  is finite Markov chain imbeddable, and

$$P(R_s(\mathbf{U}_n) = r | X) = \xi \left( \prod_{t=1}^n M_t \right) \mathbf{B}'(C_r),$$

where  $\mathbf{B}(C_r) = \sum_{k: R_s(\mathbf{U}_n) = r} e_k$ ,  $e_k$  is a  $1 \times \binom{m+n}{n}$  unit row vector corresponding to state  $\mathbf{U}_n$ ,  $\xi (= P(Z_0 = 1) = 1)$  is the initial probability and  $M_t$ ,  $t = 1, \dots, n$  are the transition probability matrices of the imbedded Markov chain defined on the state space  $\Omega_t$ .

*Proof.* For each  $\mathbf{U}_n$  in the state space  $\Omega_n$ , we have a corresponding

$$\begin{aligned} R_s(\mathbf{U}_n | X) &= R_s(\mathbf{U}_n^- | X) + R_s(\mathbf{U}_n^+ | X) \\ &= \frac{\sum_{k=1}^{n^-} (u_k^-)^2 + \sum_{k=1}^{n^-} u_k^-}{2} + \sum_{k=1}^{n^- - 1} (u_k^- + 1) \left( \sum_{j=k+1}^{n^-} u_j^- \right) \\ &\quad + \frac{\sum_{k=1}^{n^+} (u_k^+)^2 + \sum_{k=1}^{n^+} u_k^+}{2} + \sum_{k=1}^{n^+ - 1} (u_k^+ + 1) \left( \sum_{j=k+1}^{n^+} u_j^+ \right). \end{aligned} \quad (11)$$

The smallest possible value of  $R_s(\mathbf{U}_n)$  is

$$r_{ss} = \begin{cases} \frac{n(n+2)}{4} & \text{if } m+n \text{ is even and } n \text{ is even} \\ \frac{(n+1)(n+3)}{4} & \text{if } m+n \text{ is even and } n \text{ is odd} \\ \frac{n^2}{4} & \text{if } m+n \text{ is odd and } n \text{ is even} \\ \frac{(n+1)(n-1)}{4} & \text{if } m+n \text{ is odd and } n \text{ is odd} \end{cases} \quad (12)$$

and the largest possible value is

$$r_{sb} = \begin{cases} \frac{n(2m+n+2)}{4} & \text{if } m+n \text{ is even and } n \text{ is even} \\ \frac{n(2m+n+2)-1}{4} & \text{if } m+n \text{ is even and } n \text{ is odd} \\ \frac{n(2m+n-1)}{4} & \text{if } m+n \text{ is odd and } n \text{ is even} \\ \frac{n(2m+n)-1}{4} & \text{if } m+n \text{ is odd and } n \text{ is odd} \end{cases} \quad (13)$$

In accordance with Equation (11), we use the possible value of  $R_s$  as a rule of the partition. The rest of the proof follows along the same line as that of Theorem 3, and here, is omitted.

□

Similarly, we apply the LLN to conclude that

$$\frac{1}{N} \sum_{i=1}^N P(R_s | \mathbf{X}_i) \xrightarrow{P} P(R_s)$$

which establishes the distribution of  $R_s$ .

Through FMCI we, again, successfully retrieved the distribution of  $R_s$  under selected alternative distributions, for which the procedures are similar to those in the previous section. In addition, it is quite intuitive to approximate the power function by

$$\frac{1}{N} \sum_{i=1}^N \left( \sum_{s=r_{ss}}^{s_{1\alpha}} P(R_s(\mathbf{U}_n) = s | \mathbf{X}_i) + \sum_{s=s_{2\alpha}}^{r_{sb}} P(R_s(\mathbf{U}_n) = s | \mathbf{X}_i) \right),$$

where

$$P(R_s(\mathbf{U}_n) \leq s_{1\alpha} | H_0) + P(R_s(\mathbf{U}_n) \geq s_{2\alpha} | H_0) \leq \alpha.$$

### 2.3 Joint distributions of the rank statistics in the shift and scale case

We have derived the marginal distributions of  $R_l$  and  $R_s$  in terms of  $\mathbf{U}_n$ , respectively, which yield the following theorem.

**Theorem 5.**  $(R_l(\mathbf{U}_n | \mathbf{X}), R_s(\mathbf{U}_n | \mathbf{X}))$  is finite Markov chain imbeddable, and

$$P(R_l(\mathbf{U}_n) = r_1; R_s(\mathbf{U}_n) = r_2 | \mathbf{X}) = \xi \left( \prod_{t=1}^n M_t \right) \mathbf{B}'(C_r)$$

where  $\mathbf{B}(C_r) = \sum_{k: R_l(\mathbf{U}_n) = r_1 \& R_s(\mathbf{U}_n) = r_2} e_k$ ,  $e_k$  is a  $1 \times \binom{m+n}{n}$  unit row vector corresponding to state  $\mathbf{u}_n$ ,  $\xi (= P(Z_0 = 1) = 1)$  is the initial probability and  $M_t$ ,  $t = 1, \dots, n$  are the transition probability matrices of the imbedded Markov chain defined on the state space  $\Omega_t$ .

*Proof.* By Equations (4) and (11), we know each  $\mathbf{u}_n$  in the state space  $\Omega_n$  has corresponding values of  $R_l$  and  $R_s$ . The combinations of the values  $R_l$  and  $R_s$  are used to be

the standard of the partition. The rest of the proof follows along the same line as that of Theorem 3.  $\square$

The joint distribution of the ranks considering both the location and scale parameters which can be determined through our algorithm is yet to be studied in the literature. Our result allows us to test the homogeneity of the distribution functions  $F(x) = G((x - \theta)\sigma^{-1})$ . We state the hypotheses as follows

$$H_0 : \theta = 0 \text{ and } \sigma = 1 \text{ v.s. } H_a : \theta \neq 0 \text{ or } \sigma \neq 1. \quad (14)$$

Also we are able to identify a proper critical region under the null hypothesis and discuss its power when  $F \neq G$ . For example, a rectangular critical region can be

$$C_\alpha = \{R_l \leq r_{1l}, R_l \geq r_{2l}, R_s \leq r_{1s} \text{ or } R_s \geq r_{2s}\}$$

where  $r_{1l}$ ,  $r_{2l}$ ,  $r_{1s}$  and  $r_{2s}$  are the critical values such that

$$\begin{aligned} P(R_l \leq r_{1l}|H_0) + P(R_l \geq r_{2l}|H_0) + P(r_{1l} < R_l < r_{2l}, R_s \leq r_{1s}|H_0) \\ + P(r_{1l} < R_l < r_{2l}, R_s \geq r_{2s}|H_0) \leq \alpha \end{aligned}$$

or an elliptic critical region

$$C'_\alpha = \left\{ \frac{R_l^2}{a} + \frac{R_s^2}{b} > C \right\}$$

for some positive constants  $a$  and  $b$  such that

$$P\left(\frac{R_l^2}{a} + \frac{R_s^2}{b} > C|H_0\right) \leq \alpha.$$

According to the above defined rejection region, the power of the test can be found as

$$\begin{aligned} P(R_l \leq r_{1l}|H_a) + P(R_l \geq r_{2l}|H_a) + P(r_{1l} < R_l < r_{2l}, R_s \leq r_{1s}|H_a) \\ + P(r_{1l} < R_l < r_{2l}, R_s \geq r_{2s}|H_a) \end{aligned} \quad (15)$$

or

$$P\left(\frac{R_l^2}{a} + \frac{R_s^2}{b} > C|H_a\right). \quad (16)$$

Note that unless having a conjecture about the values of  $\theta$  and  $\sigma$ , we tend to use a two-sided test. However, with the knowledge of the center and shape of the distribution of interest, deciding a sectorial critical region is a better choice, for which an example is demonstrated in the numerical studies.

### 3 Numerical results and discussion

#### 3.1 A joint distribution of $R_l$ and $R_s$

Let  $\{X_1, \dots, X_5\} \sim N(0, 1)$  and  $\{Y_1, \dots, Y_7\} \sim N(\theta, \sigma)$ . Figure 1 gives the joint distribution of the random variables  $R_l$  and  $R_s$  under the null hypothesis of  $\theta = 0$  and  $\sigma = 1$ . The marginal distributions of  $R_l$  and  $R_s$  can be easily established from their joint distribution. Figure 1 also shows that the two random variables  $R_l$  and  $R_s$  are dependent. We construct two critical regions as shown in Figure 2, according to their joint distribution. Outside the yellow area in Figure 2 is the selected rectangular critical region  $C_{0.1738}$  and outside the red shadow is the elliptic one  $C'_{0.1738}$ .

#### 3.2 Powers for a joint test using $R_l$ and $R_s$

The alternative of interest is stated in the preceding section (see Equation (14)). The power functions of the test statistics  $R_l$  and  $R_s$  for a sequence of normally distributed populations with  $\theta$  from -20 to 20 with an increment of 0.5 and  $\sigma$  from 1 to 10 with an increment of 1, and its reciprocal under two types of critical regions are provided in Figures 3 and 4. We adopt a two-sided test because of the selected values of the parameters. It should be slightly modified the critical region in the previous step in order to calculate the powers if a one-sided test is adopted. Both critical regions roughly perform equally well as shown in Figures 3 and 4. Figure 5 presents the performance of the two critical regions for given various parameter settings. Figures 5(a) and (b) show that given a standard deviation of 1 or a mean of 0, the powers of the two critical regions, rectangular and elliptic, are high and similar. However, when the variation of the alternative population reduces ( $\sigma = 1/10$ ) or increases ( $\sigma = 10$ ), the elliptic critical region performs better than the rectangular one as shown in Figures 5(c) and (d). Therefore, we suggest that when conducting a test for the equivalence of two distributions, an elliptic rejection area should be used.

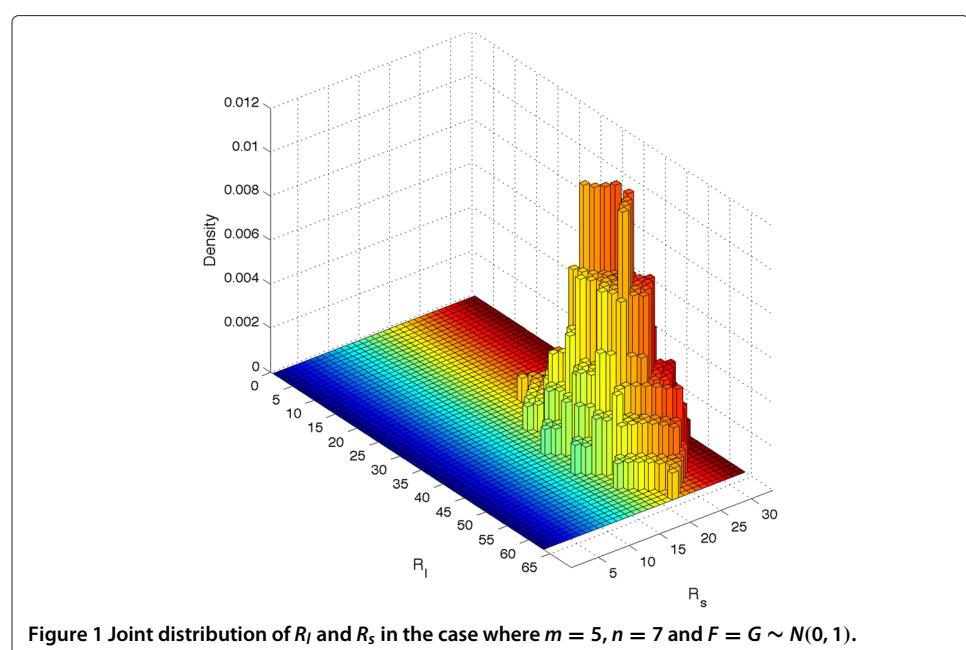


Figure 1 Joint distribution of  $R_l$  and  $R_s$  in the case where  $m = 5$ ,  $n = 7$  and  $F = G \sim N(0, 1)$ .

Figure 2 Critical Regions at size 17.38% for  $R_I$  and  $R_S$  for  $m = 5$  and  $n = 7$ .

Next, we consider the problem of determining an optimum rank test. To conduct a test of distributions equivalency, we can use either  $R_l$  or  $R_s$  as the test statistic. As mentioned earlier, the marginal distribution  $R_l$  or  $R_s$  can be easily established from their joint distribution. Figures 6 and 7 provide the power functions for the test statistics  $R_l$  and  $R_s$  at the level of significance 17.38%, respectively. Figure 7 shows that the rank test against scale parameter is badly effected by the centre of the alternative population. This was seen before by Ansari and Bradley (1960). By comparing Figures 6 and 7 with Figure 4, it seems that the joint test would be much more reliable than either  $R_l$  or  $R_s$  alone for distributions equivalence tests. A joint test for distributions equivalency would like a better option under most circumstances.

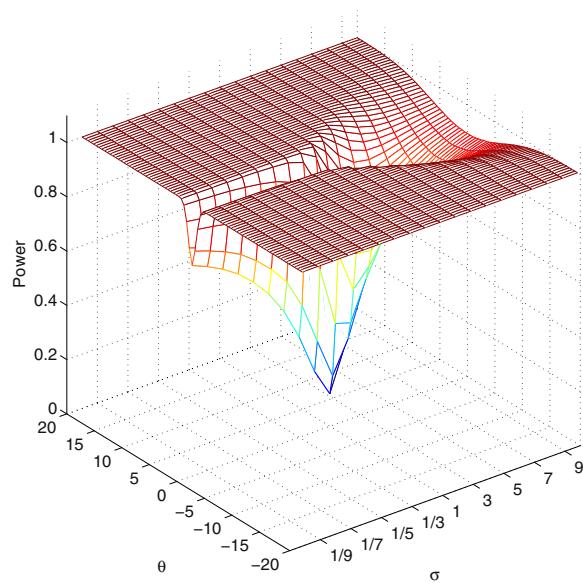


Figure 3 Power functions of  $R_I$  and  $R_S$  for  $m = 5$  and  $n = 7$  under  $C_\alpha$ .

### 3.3 Lehmann alternatives

Consider the one-sided alternative  $F(x; \theta, \sigma) > G(x; \theta, \sigma)$ , Lehmann (1953) proposed a test of  $H_0 : F(x; \theta, \sigma) = G(x; \theta, \sigma)$  against  $H_a : F(x; \theta, \sigma)^k = G(x; \theta, \sigma)$  which is known as the family of Lehmann alternative. Note  $F(x; \theta, \sigma)^k$  is the cumulative distribution of  $\max_{1 \leq i \leq k}(x_i)$  when  $X_i \sim F$  and, under the alternative hypothesis,  $G(x; \theta, \sigma)$  is stochastically larger than  $F(x; \theta, \sigma)$ . First of all, we know

$$\begin{aligned} E_k(X) &= \int_{-\infty}^0 -G(x)dx + \int_0^\infty 1 - G(x)dx \\ &> \int_{-\infty}^0 -F(x)dx + \int_0^\infty 1 - F(x)dx = E(X). \end{aligned} \quad (17)$$

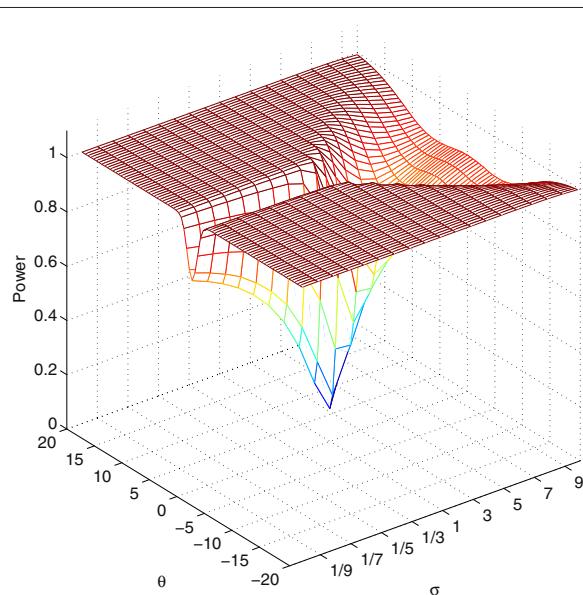
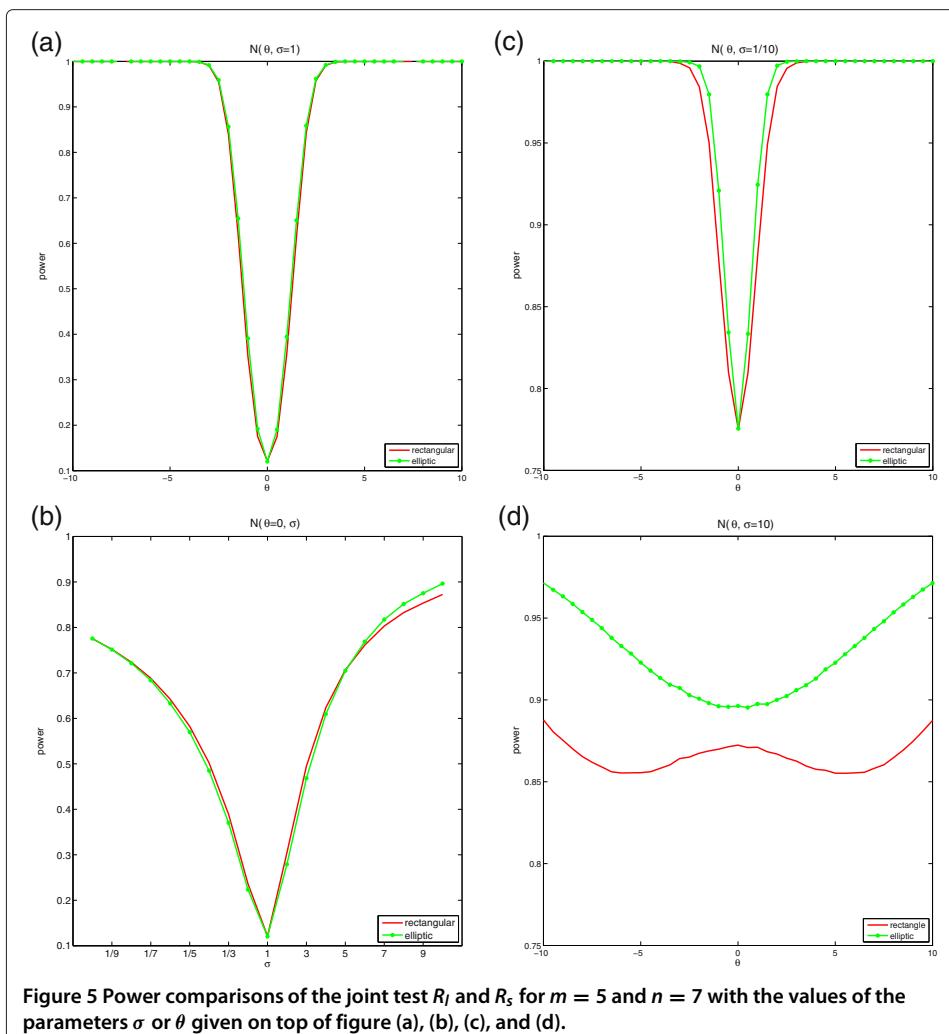
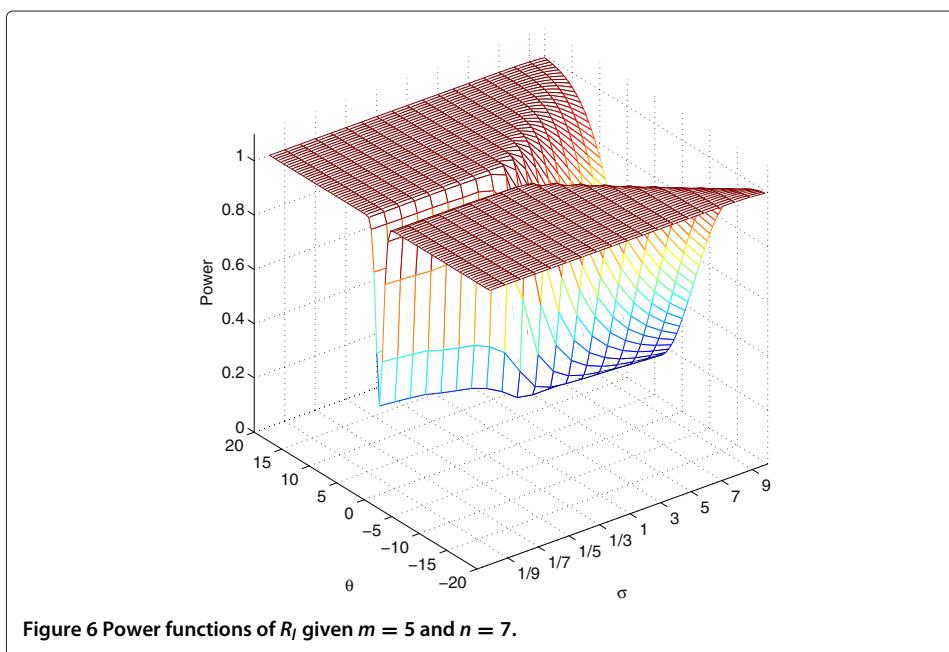


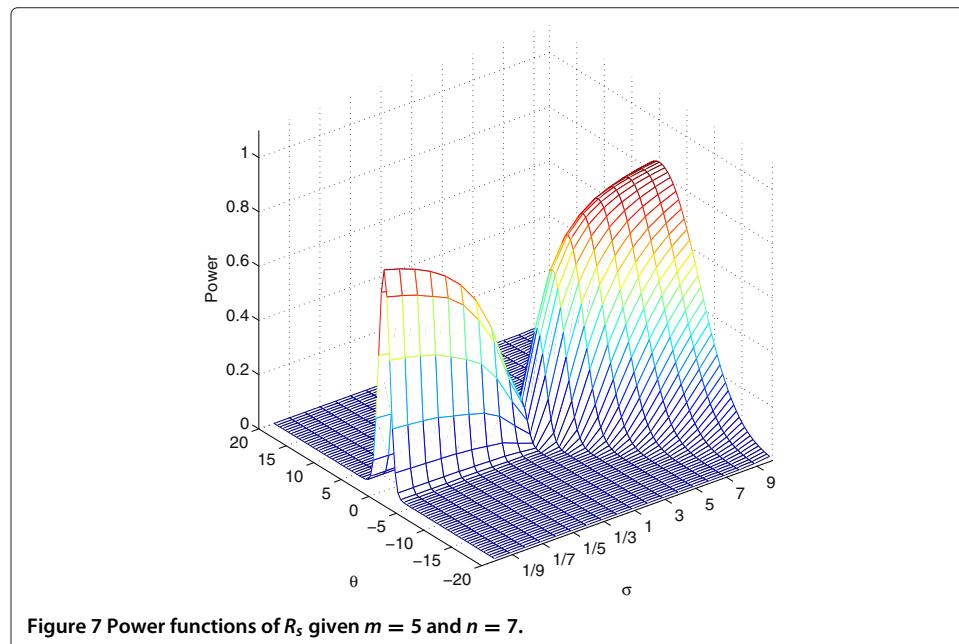
Figure 4 Power functions of  $R_I$  and  $R_S$  for  $m = 5$  and  $n = 7$  under  $C'_\alpha$ .



**Figure 5** Power comparisons of the joint test  $R_l$  and  $R_s$  for  $m = 5$  and  $n = 7$  with the values of the parameters  $\sigma$  or  $\theta$  given on top of figure (a), (b), (c), and (d).



**Figure 6** Power functions of  $R_l$  given  $m = 5$  and  $n = 7$ .



Therefore, the larger the  $R_l$  is, the stronger the evidence against the null hypothesis will be. For the variation of the distribution per se, the codomain of the density function is compressed to larger numbers; therefore, in most cases, we have  $\text{Var}(X_k) < \text{Var}(X)$ . We then propose to reject the null hypothesis when  $R_s$  is large. For example, given  $F \sim U(0, 1)$  and  $G = F^k$ , it is easy to see

$$\frac{E_{k+1}(X)}{E_k(X)} = \frac{(k+1)^2}{k(k+2)} > 1 \quad (18)$$

and

$$\frac{\text{Var}_{k+1}(X)}{\text{Var}_k(X)} = \frac{(k+1)^3}{k(k+2)(k+3)} < 1 \quad (19)$$

for all  $k$ . We first find the marginal and joint distributions of the ranks  $R_l$  and  $R_s$  in order to define critical regions for  $R_l$  and  $R_s$  individually and simultaneously. Due to the properties

**Table 1 Power comparisons for a one-sided rank test  $H_0 : F(x; \theta_o, \sigma_o) = G(x; \theta_a, \sigma_a)$  v.s.  $H_a : F^k(x; \theta_o, \sigma_o) = G(x; \theta_a, \sigma_a)$**

F	Test	$m = 6 \ n = 10$				$m = 10 \ n = 10$				$m = 10 \ n = 20$			
		$\beta(F)$	$\beta(F^2)$	$\beta(F^3)$	$\beta(F^6)$	$\beta(F)$	$\beta(F^2)$	$\beta(F^3)$	$\beta(F^6)$	$\beta(F)$	$\beta(F^2)$	$\beta(F^3)$	$\beta(F^6)$
$U(0, 1)$	$R_l$	.090	.411	.647	.900	.096	.496	.761	.967	.099	.591	.845	.984
	$R_s$	.080	.152	.193	.218	.076	.137	.149	.123	.100	.236	.370	.638
	$R_l \& R_s$	.100	.452	.699	.934	.100	.531	.799	.981	.100	.622	.878	.992
$t(3)$	$R_l$	0.090	.412	.639	.897	0.096	.493	.756	.965	0.099	.574	.841	.987
	$R_s$	0.080	.150	.197	.217	0.076	.137	.152	.121	0.100	.234	.367	.634
	$R_l \& R_s$	0.100	.453	.696	.932	0.100	.528	.798	.980	0.100	.606	.874	.993
$Exp(1)$	$R_l$	0.090	.411	.650	.899	0.096	.490	.764	.967	0.099	.579	.841	.987
	$R_s$	0.080	.149	.195	.217	0.076	.140	.152	.122	0.100	.232	.376	.641
	$R_l \& R_s$	0.100	.451	.702	.933	0.100	.525	.805	.982	0.100	.607	.875	.993

Note: A sectorial critical region is chosen for a simultaneous testing.

of the mean and variance of the alternative distribution, as shown in Equations (17), (18) and (19), we are cautious to define the critical regions. Table 1 provides powers for the tests as we choose uniform, standard Normal, student-t with 3 degrees of freedom, exponential distributions for the hypothesized distribution, a couple of different settings for sample sizes  $m$  and  $n$ , and 2, 3, 6 for  $k$ . Clearly, a joint test considering both  $R_l$  and  $R_s$  for the equality of distributions is best suited in comparison with tests considering only one of the rank statistics.

#### 4 Conclusion

Our proposed algorithm provides a solution for finding the power of distribution equivalence tests considering the shift and scale parameters, respectively and simultaneously. Numerical studies show that a joint test should be adopted for the test homogeneity of distributions as well as under Lehmann alternatives. Also an elliptic critical region is a better choice rather than a rectangular one for a joint test. In practice, it is reasonable to have neither the normality assumption nor equal mean/variance of the interested distributions. However, our algorithm highly depends on the technology equipments as the possible states in  $\Omega_n$  grow rapidly when the sample sizes increase. Therefore, we can, so far, only target small sample sizes in our work.

#### Competing interests

The author declares that she has no competing interests.

#### Acknowledgments

The author would like to thank James C. Fu and anonymous referee whose comments led to significant improvements of this manuscript.

Received: 7 August 2013 Accepted: 10 February 2014

Published: 11 June 2014

#### References

- Ansari, AR, Bradley, RA: Rank-Sum Tests for Dispersions. *Ann. Math. Stat.* **31**, 1174–1189 (1960)
- Collings, BJ, Hamilton, MA: Estimating the power of the two-sample Wilcoxon Test for location shift. *Biometrics* **44**, 847–860 (1988)
- Klotz, J: Nonparametric test for scale. *Ann. Math. Stat.* **33**, 498–512 (1962)
- Lehmann, EL: The power for rank tests. *Ann. Math. Stat.* **24**, 23–43 (1953)
- Lehmann, EL: Nonparametrics: Statistical Methods Based on Ranks. Revised ed. Prentice-Hall, New Jersey (1998)
- Mann, HB, Whitney, DR: On a test of whether one of two random variables is stochastically larger than the other. *Ann. Math. Stat.* **18**, 50–60 (1947)
- Mood, AM: On the asymptotic efficiency of certain nonparametric two-sample tests. *Ann. Math. Stat.* **25**, 514–522 (1954)
- Rosner, B, Glynn, RJ: Power and sample size estimation for the Wilcoxon rank sum test with application to comparisons of C statistics from alternative prediction models. *Biometrics* **65**, 188–197 (2009)
- Sheeh, G, Jan, SL, Randles, RH: On power and sample size determinations for the Wilcoxon-Mann-Whitney test. *Nonparametric Stat.* **18**, 33–43 (2006)
- Siegel, S, Tukey, JW: A nonparametric sum of ranks procedure for relative spread in unpaired samples. *J. Am. Stat. Assoc.* **55**, 429–445 (1960)
- Wilcoxon, F: Individual comparisons by ranking methods. *Biometrics* **1**, 80–83 (1945)

doi:10.1186/2195-5832-1-6

Cite this article as: Lee: Joint distribution of rank statistics considering the location and scale parameters and its power study. *Journal of Statistical Distributions and Applications* 2014 **1**:6.