

Revisiting the two-sample runs test

Ludwig Baringhaus¹ · Norbert Henze²

Received: 8 January 2015 / Accepted: 5 October 2015 / Published online: 13 October 2015
© Sociedad de Estadística e Investigación Operativa 2015

Abstract We give new representations for the two-sample runs test statistic, derive explicit expressions of its mean and variance also in the general non-null case, and present an elementary proof of the consistency of the runs test.

Keywords Two-sample tests · Runs · Consistency · Fixed alternative

Mathematics Subject Classification 62G10 · 62G20

1 Introduction

The runs test is one of the most popular two-sample rank tests. Proposed by [Wald and Wolfowitz \(1940\)](#), it is still presented and discussed in modern textbooks, see, e.g., [Gibbons and Chakraborti \(2011\)](#), [Govindarajulu \(2007\)](#). Irrespective of its efficiency deficits detected by [Mood \(1954\)](#), it has its merits as being a procedure based on a very simple test statistic, the number of runs in the series of zeros and ones associated with the variables of the two samples. Like its common competitors, the Kolmogorov–Smirnov test and the Cramér–von Mises test, the runs test is an omnibus test which is consistent against a broad class of alternative distributions. But unlike the former that

✉ Ludwig Baringhaus
lbaring@stochastik.uni-hannover.de

Norbert Henze
norbert.henze@kit.edu

¹ Institut für Mathematische Stochastik, Leibniz Universität Hannover, Welfengarten 1, 30167 Hannover, Germany

² Karlsruher Institut für Technologie (KIT), Institut für Stochastik, Englerstraße 2, 76133 Karlsruhe, Germany

come along with complicated asymptotic null distributions of the test statistics, the null distribution of the runs test statistic is asymptotically normal. Wald and Wolfowitz sketch the proof which is along the lines of that of the De Moivre–Laplace central limit theorem, see also [Morgenstern \(1962\)](#) for a simplification of the proof. The major part of the paper of Wald and Wolfowitz is concerned with the proof of the consistency of the runs test, which is quite long and technically involved. A shorter proof is presented in [Blum and Weiss \(1957\)](#). However, these authors use a result of [Weiss \(1955\)](#) on the almost sure uniform convergence of certain proportions of spacings associated with a sequence of independent and identically distributed real random variables, the proof of which is not less challenging. Using the Efron–Stein inequality, [Henze and Voigt \(1992\)](#) give valuable bounds of the test statistic and obtain, as a consequence, the (even strong) consistency of the runs test. [Blumenthal \(1963\)](#) using results of [Weiss \(1955\)](#) too, proves the limiting normality of the test statistic for a broad class of alternative distributions. In what follows, on the basis of new representations of the runs statistic involving the number of ‘interior’ runs of zeros and the number of ‘interior’ runs of ones, we give explicit expressions for the expectation and variance of the runs statistic in the general non-null case. From these, we derive the desired limit expressions to establish the consistency using Pratt’s extended version of the dominated convergence theorem.

Runs tests have also been studied in other contexts, such as testing for symmetry [[Cohen and Menjoge \(1988\)](#), [McWilliams \(1990\)](#) and [Henze \(1993a\)](#)] or testing for the proportional hazards model of random censorship [[Henze \(1993b\)](#)]. A multivariate generalization of the runs test for the two-sample problem has been studied by [Friedman and Rafsky \(1979\)](#), [Henze and Penrose \(1999\)](#), [Paindaveine \(2009\)](#), building on runs proposal by [Marden \(1999\)](#), [Biswas et al. \(2014\)](#), and [Dyckerhoff et al. \(2015\)](#).

2 Expectation and variance of the runs statistic

Let X_1, \dots, X_m and Y_1, \dots, Y_n be independent random samples of independent real-valued random variables, where the X_j have an unknown continuous distribution (function) F and the Y_k have an unknown continuous distribution (function) G . Put $Z_j = X_j$, $j = 1, \dots, m$, and $Z_{m+k} = Y_k$, $k = 1, \dots, n$. Due to the continuity of F and G , the Z_j are pairwise distinct with probability one, so that we will assume without loss of generality that $Z_{(1)} < \dots < Z_{(n+m)}$ where $Z_{(j)}$ is the j -th order statistic of Z_1, \dots, Z_{m+n} . Denoting by D_1, \dots, D_{m+n} the antiranks of Z_1, \dots, Z_{m+n} , i.e., the random permutation of the indices $1, \dots, m+n$ such that $Z_{D_1} < \dots < Z_{D_{m+n}}$, and by $A_j = \mathbf{1}\{1 \leq D_j \leq m\}$ the indicator of the event that the j -th order statistic is an element of the first sample X_1, \dots, X_m , $j = 1, \dots, m+n$, the Wald–Wolfowitz runs statistic $W_{m,n}$ is defined as the number of runs in the random sequence A_1, \dots, A_{m+n} of zeros and ones, i.e., we have

$$W_{m,n} = 1 + \sum_{j=2}^{m+n} \mathbf{1}\{A_{j-1} \neq A_j\}, \quad (2.1)$$

see [Gibbons and Chakraborti \(2011\)](#), page 81. Treating the testing problem

$$H : F = G \text{ versus } K : F \neq G$$

[Wald and Wolfowitz \(1940\)](#) suggest to reject the hypothesis H if $W_{m,n} < c_{m,n;\alpha}$, where for given $\alpha \in (0, 1)$ the constant $c_{m,n;\alpha}$ is the α -quantile of $W_{m,n}$ under H . Being a rank statistic, $W_{m,n}$ is distribution-free under H . Denoting by

$$\Delta_{m,n}(F, G) = \mathbb{E}(W_{m,n}) \text{ and } \sigma_{m,n}^2(F, G) = \mathbb{V}(W_{m,n})$$

the expectation and the variance of $W_{m,n}$, it is easily verified that, under H ,

$$\Delta_{m,n} = \Delta_{m,n}(F, F) = 1 + 2 \frac{mn}{m+n}$$

and

$$\sigma_{m,n}^2 = \sigma_{m,n}^2(F, F) = \frac{2mn(2mn - (m+n))}{(m+n)^2(m+n-1)},$$

see, e.g., [Gibbons and Chakraborti \(2011\)](#), p 82. In the sequel, we derive $\mathbb{E}(W_{m,n})$ and $\mathbb{V}(W_{m,n})$ in the general setting stated at the beginning of this section. The key idea is to introduce so-called ‘interior’ runs of zeros and ones. To be precise, let $X_{(1)} < \dots < X_{(m)}$ and $Y_{(1)} < \dots < Y_{(n)}$ be the order statistics of X_1, \dots, X_m and Y_1, \dots, Y_n , respectively. For a (possibly random) subset B of \mathbb{R} , let

$$N_X(B) = \sum_{i=1}^m \mathbf{1}\{X_i \in B\}, \quad N_Y(B) = \sum_{j=1}^n \mathbf{1}\{Y_j \in B\}$$

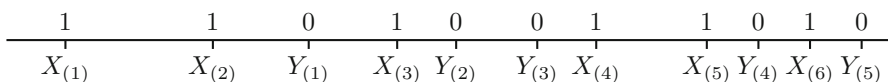
be the number of X_1, \dots, X_m and the number of Y_1, \dots, Y_n falling into B , respectively. If B is an open interval, $B = (a, b)$, say, we write $N_X(a, b)$ and $N_Y(a, b)$ instead of $N_X((a, b))$ and $N_Y((a, b))$, respectively. Then,

$$W_{m,n} = 2 + W_{m,n}(0) + W_{m,n}(1),$$

where

$$W_{m,n}(0) = \sum_{i=2}^m \mathbf{1}\{N_Y(X_{(i-1)}, X_{(i)}) \geq 1\}, \quad W_{m,n}(1) = \sum_{i=2}^n \mathbf{1}\{N_X(Y_{(i-1)}, Y_{(i)}) \geq 1\}$$

are the so-called numbers of interior runs of zeros and interior runs of ones, respectively. These concepts are illustrated below. Here, $m = 6$, $n = 5$, $W_{m,n} = 8$, $W_{m,n}(0) = W_{m,n}(1) = 3$.



In what follows, we deal with the statistic

$$V_{m,n}(0) = m - 1 - W_{m,n}(0) = \sum_{i=2}^m \mathbf{1}\{N_Y(X_{(i-1)}, X_{(i)}) = 0\}$$

that counts the number of intervals $(X_{(i-1)}, X_{(i)})$, $i = 2, \dots, m$, which do not contain any of the Y_j 's. Likewise,

$$V_{m,n}(1) = n - 1 - W_{m,n}(1) = \sum_{i=2}^n \mathbf{1}\{N_X(Y_{(i-1)}, Y_{(i)}) = 0\}$$

is the number of intervals $(Y_{(i-1)}, Y_{(i)})$ that do not contain any of the X_j 's. Notice that $V_{m,n}^2(i) = V_{m,n}(i) + 2V_{m,n}(i, 1) + 2V_{m,n}(i, 2)$, $i \in \{0, 1\}$, where

$$\begin{aligned} V_{m,n}(0, 1) &= \sum_{i=2}^{m-1} \mathbf{1}\{N_Y(X_{(i-1)}, X_{(i+1)}) = 0\}, \\ V_{m,n}(0, 2) &= \sum_{i=2}^{m-2} \sum_{j=i+2}^m \mathbf{1}\{N_Y((X_{(i-1)}, X_{(i)}) \cup (X_{(j-1)}, X_{(j)})) = 0\}, \\ V_{m,n}(1, 1) &= \sum_{i=2}^{n-1} \mathbf{1}\{N_X(Y_{(i-1)}, Y_{(i+1)}) = 0\}, \\ V_{m,n}(1, 2) &= \sum_{i=2}^{n-2} \sum_{j=i+2}^n \mathbf{1}\{N_X((Y_{(i-1)}, Y_{(i)}) \cup (Y_{(j-1)}, Y_{(j)})) = 0\}. \end{aligned} \quad (2.2)$$

We define $V_{m,n} = V_{m,n}(0) + V_{m,n}(1)$ and note that

$$W_{m,n} = m + n - V_{m,n}.$$

We first give the expectations of the random variables $V_{m,n}(i)$, $V_{m,n}(i, 1)$, $V_{m,n}(i, 2)$, $i = 0, 1$, and $V_{m,n}(0)V_{m,n}(1)$. For a real number x and an integer k write

$$x_{(k)} = x(x-1) \cdots (x-k+1)$$

for the k -th descending factorial of x .

Proposition 1 *If $m \geq 2$, then $\mathbb{E}V_{m,n}(0) = m_{(2)}\mu_{m,n}(0)$, where*

$$\mu_{m,n}(0) = \iint_{u < v} (1 - (G(v) - G(u)))^n (1 - (F(v) - F(u)))^{m-2} dF(u) dF(v).$$

If $n \geq 2$, we have $\mathbb{E}V_{m,n}(1) = n_{(2)}\mu_{m,n}(1)$, where

$$\mu_{m,n}(1) = \iint_{u < v} (1 - (F(v) - F(u)))^m (1 - (G(v) - G(u)))^{n-2} dG(u)dG(v).$$

If $m \geq 3$, then $\mathbb{E}V_{m,n}(0, 1) = m_{(3)}\mu_{m,n}(0, 1)$, where

$$\begin{aligned} \mu_{m,n}(0, 1) &= \iint_{u < v} (1 - (G(v) - G(u)))^n (1 - (F(v) - F(u)))^{m-3} \\ &\quad \times (F(v) - F(u)) dF(u)dF(v). \end{aligned}$$

If $n \geq 3$, we have $\mathbb{E}V_{m,n}(1, 1) = n_{(3)}\mu_{m,n}(1, 1)$, where

$$\begin{aligned} \mu_{m,n}(1, 1) &= \iint_{u < v} (1 - (F(v) - F(u)))^m (1 - (G(v) - G(u)))^{n-3} \\ &\quad \times (G(v) - G(u)) dG(u)dG(v). \end{aligned}$$

If $m \geq 4$, then $\mathbb{E}V_{m,n}(0, 2) = m_{(4)}\mu_{m,n}(0, 2)$, where

$$\begin{aligned} \mu_{m,n}(0, 2) &= \iiint_{u < v < w < z} (1 - ((G(z) - G(w)) + (G(v) - G(u))))^n \\ &\quad \cdot (1 - ((F(z) - F(w)) + (F(v) - F(u))))^{m-4} dF(u)dF(v)dF(w)dF(z). \end{aligned}$$

If $n \geq 4$, we have $\mathbb{E}V_{m,n}(1, 2) = n_{(4)}\mu_{m,n}(1, 2)$, where

$$\begin{aligned} \mu_{m,n}(1, 2) &= \iiint_{u < v < w < z} (1 - ((F(z) - F(w)) + (F(v) - F(u))))^m \\ &\quad \cdot (1 - ((G(z) - G(w)) + (G(v) - G(u))))^{n-4} dG(u)dG(v)dG(w)dG(z). \end{aligned}$$

If $m \geq 2$ and $n \geq 2$, then $\mathbb{E}[V_{m,n}(0)V_{m,n}(1)] = m_{(2)}n_{(2)}(\rho_{m,n}(0) + \rho_{m,n}(1))$, where

$$\begin{aligned} \rho_{m,n}(0) &= \iiint_{u < v < w < z} (1 - ((F(z) - F(w)) + (F(v) - F(u))))^{m-2} \\ &\quad \cdot (1 - ((G(z) - G(w)) + (G(v) - G(u))))^{n-2} dF(u)dF(v)dG(w)dG(z) \end{aligned} \quad (2.3)$$

and

$$\begin{aligned} \rho_{m,n}(1) = & \int \int \int \int_{u < v < w < z} (1 - ((F(z) - F(w)) + (F(v) - F(u))))^{m-2} \\ & \cdot (1 - ((G(z) - G(w)) + (G(v) - G(u))))^{n-2} dG(u) dG(v) dF(w) dF(z). \end{aligned} \quad (2.4)$$

If one of the conditions on the sample sizes given above is not true, the corresponding expectation is 0.

Proof In what follows, we write $x \wedge y := \min(x, y)$, $x \vee y := \max(x, y)$. It is convenient to use alternative representations of the statistics. Let $m \geq 2$. Putting

$$I_{ij} = \mathbf{1}\{N_X(X_i \wedge X_j, X_i \vee X_j) = 0, N_Y(X_i \wedge X_j, X_i \vee X_j) = 0\}$$

for $1 \leq i < j \leq m$ (notice that $X_i, X_j \notin (X_i \wedge X_j, X_i \vee X_j)$), we have

$$V_{m,n}(0) = \sum_{1 \leq i < j \leq m} I_{ij}.$$

This equation is readily established since $V_{m,n}(0)$ counts the intervals $(X_{(i-1)}, X_{(i)})$, $i = 2, \dots, m$, that do not contain any of the Y_j 's. Notice that, although the sum on the right-hand side has $\binom{m}{2}$ summands, only those pairs i, j contribute to the sum for which the points X_i, X_j are adjacent within the ordered sample $X_{(1)} < \dots < X_{(m)}$. From $\mathbb{E}I_{ij} = 2\mu_{m,n}(0)$ for each $1 \leq i < j \leq m$, we deduce that

$$\mathbb{E}V_{m,n}(0) = \binom{m}{2} 2\mu_{m,n}(0) = m_{(2)}\mu_{m,n}(0).$$

Let $m \geq 3$, and put

$$\begin{aligned} I_{ij;k} = & \mathbf{1}\{X_i \wedge X_j < X_k < X_i \vee X_j, N_Y(X_i \wedge X_j, X_i \vee X_j) = 0, \\ & X_\mu \notin (X_i \wedge X_j, X_i \vee X_j) \text{ for each } \mu \in \{1, \dots, m\} \setminus \{i, j, k\}\} \end{aligned}$$

for $1 \leq i < j \leq m$, $k \in \{1, \dots, m\} \setminus \{i, j\}$. Since, by (2.2), $V_{m,n}(0, 1)$ is the number of intervals $(X_{(i-1)}, X_{(i+1)})$, $i = 2, \dots, m-1$, that do not contain any of the Y_j 's, it follows that

$$V_{m,n}(0, 1) = \sum_{1 \leq i < j \leq m} \sum_{k \in \{1, \dots, m\} \setminus \{i, j\}} I_{ij;k} \quad (2.5)$$

(notice that only those pairs i, j contribute to the sum for which there is an $l \in \{2, \dots, m-1\}$ such that $X_i \wedge X_j = X_{(l-1)}$ and $X_i \vee X_j = X_{(l+1)}$).

We have $\mathbb{E}I_{ij;k} = 2\mu_{m,n}(0, 1)$ for each $1 \leq i < j \leq m$, $k \in \{1, \dots, m\} \setminus \{i, j\}$, whence

$$\mathbb{E}V_{m,n}(0, 1) = 2\binom{m}{2}(m-2)\mu_{m,n}(0, 1) = m_{(3)}\mu_{m,n}(0, 1).$$

Let $m \geq 4$. Putting

$$\begin{aligned} I_{ij;k\ell} = & \mathbf{1}\{X_i < X_j < X_k < X_\ell, N_Y((X_i, X_j) \cup (X_k, X_\ell)) = 0, \\ & X_\mu \notin ((X_i, X_j) \cup (X_k, X_\ell)) \text{ for each } \mu \in \{1, \dots, m\} \setminus \{i, j, k, \ell\} \\ & + \mathbf{1}\{X_i < X_j < X_\ell < X_k, N_Y((X_i, X_j) \cup (X_\ell, X_k)) = 0, \\ & X_\mu \notin ((X_i, X_j) \cup (X_\ell, X_k)) \text{ for each } \mu \in \{1, \dots, m\} \setminus \{i, j, k, \ell\} \\ & + \mathbf{1}\{X_j < X_i < X_k < X_\ell, N_Y((X_j, X_i) \cup (X_k, X_\ell)) = 0, \\ & X_\mu \notin ((X_j, X_i) \cup (X_k, X_\ell)) \text{ for each } \mu \in \{1, \dots, m\} \setminus \{i, j, k, \ell\} \\ & + \mathbf{1}\{X_j < X_i < X_\ell < X_k, N_Y((X_j, X_i) \cup (X_\ell, X_k)) = 0, \\ & X_\mu \notin ((X_j, X_i) \cup (X_\ell, X_k)) \text{ for each } \mu \in \{1, \dots, m\} \setminus \{i, j, k, \ell\} \end{aligned}$$

for $1 \leq i < j \leq m$, $1 \leq k < \ell \leq m$, $k, \ell \notin \{i, j\}$, a similar reasoning as that leading to (2.5) yields

$$V_{m,n}(0, 2) = \sum_{1 \leq i < j \leq m} \sum_{\substack{1 \leq k < \ell \leq m \\ k, \ell \notin \{i, j\}}} I_{ij;k\ell}.$$

Since $\mathbb{E}I_{ij;k\ell} = 4\mu_{m,n}(0, 2)$ for each $1 \leq i < j \leq m$, $1 \leq k < \ell \leq m$, $k, \ell \notin \{i, j\}$, it follows that

$$\mathbb{E}V_{m,n}(0, 2) = 4\binom{m}{2}\binom{m-2}{2}\mu_{m,n}(0, 2) = m_{(4)}\mu_{m,n}(0, 2).$$

The expectations of $V_{m,n}(1)$, $V_{m,n}(1, 1)$ and $V_{m,n}(1, 2)$ can be obtained in the same way. Finally, if $m \geq 2$, $n \geq 2$, then

$$V_{m,n}(0)V_{m,n}(1) = \sum_{1 \leq i < j \leq m} \sum_{1 \leq k < \ell \leq n} (J_{ij;k\ell} + K_{ij;k\ell}),$$

where

$$\begin{aligned} J_{ij;k\ell} = & \mathbf{1}\{X_i < X_j < Y_k < Y_\ell, X_\mu \notin ((X_i, X_j) \cup (Y_k, Y_\ell)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\ & Y_\nu \notin ((X_i, X_j) \cup (Y_k, Y_\ell)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\} \\ & + \mathbf{1}\{X_i < X_j < Y_\ell < Y_k, X_\mu \notin ((X_i, X_j) \cup (Y_\ell, Y_k)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\ & Y_\nu \notin ((X_i, X_j) \cup (Y_\ell, Y_k)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\} \\ & + \mathbf{1}\{X_j < X_i < Y_k < Y_\ell, X_\mu \notin ((X_j, X_i) \cup (Y_k, Y_\ell)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\ & Y_\nu \notin ((X_j, X_i) \cup (Y_k, Y_\ell)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\} \\ & + \mathbf{1}\{X_j < X_i < Y_\ell < Y_k, X_\mu \notin ((X_j, X_i) \cup (Y_\ell, Y_k)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\ & Y_\nu \notin ((X_j, X_i) \cup (Y_\ell, Y_k)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\} \end{aligned}$$

and

$$\begin{aligned}
 K_{ij;k\ell} = & \mathbf{1}\{Y_k < Y_\ell < X_i < X_j, X_\mu \notin ((X_i, X_j) \cup (Y_k, Y_\ell)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\
 & Y_\nu \notin ((X_i, X_j) \cup (Y_k, Y_\ell)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\}\} \\
 & + \mathbf{1}\{Y_k < Y_\ell < X_j < X_i, X_\mu \notin ((X_j, X_i) \cup (Y_k, Y_\ell)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\
 & Y_\nu \notin ((X_j, X_i) \cup (Y_k, Y_\ell)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\}\} \\
 & + \mathbf{1}\{Y_\ell < Y_k < X_i < X_j, X_\mu \notin ((X_i, X_j) \cup (Y_\ell, Y_k)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\
 & Y_\nu \notin ((X_i, X_j) \cup (Y_\ell, Y_k)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\}\} \\
 & + \mathbf{1}\{Y_\ell < Y_k < X_j < X_i, X_\mu \notin ((X_j, X_i) \cup (Y_\ell, Y_k)) \text{ for } \mu \in \{1, \dots, m\} \setminus \{i, j\}, \\
 & Y_\nu \notin ((X_j, X_i) \cup (Y_\ell, Y_k)) \text{ for } \nu \in \{1, \dots, n\} \setminus \{k, \ell\}\}.
 \end{aligned}$$

Obviously, for each $1 \leq i < j \leq m$, $1 \leq k < \ell \leq n$,

$$\mathbb{E}J_{ij;k\ell} = 4\rho_{m,n}(0) \text{ and } \mathbb{E}K_{ij;k\ell} = 4\rho_{m,n}(1)$$

[recall $\rho_{m,n}(0)$ and $\rho_{m,n}(1)$ from (2.3) and (2.4), respectively]. Thus,

$$\begin{aligned}
 \mathbb{E}[V_{m,n}(0)V_{m,n}(1)] &= 4 \binom{m}{2} \binom{n}{2} (\rho_{m,n}(0) + \rho_{m,n}(1)) \\
 &= m_{(2)}n_{(2)}(\rho_{m,n}(0) + \rho_{m,n}(1)).
 \end{aligned}$$

□

The following theorem gives explicit expressions of the expectation and the variance of the number of runs of zeros and ones, a result which does not seem to be stated elsewhere.

Theorem 1 (a) *We have $\mathbb{E}W_{m,n} = m + n - \mathbb{E}V_{m,n}$, where*

$$\begin{aligned}
 \mathbb{E}V_{m,n} = & \binom{m}{2} \mathbb{E} \left((1 - |G(X_1) - G(X_2)|)^n (1 - |F(X_1) - F(X_2)|)^{m-2} \right) \\
 & + \binom{n}{2} \mathbb{E} \left((1 - |F(Y_1) - F(Y_2)|)^m (1 - |G(Y_1) - G(Y_2)|)^{n-2} \right).
 \end{aligned}$$

(b) *We have $\mathbb{V}(W_{m,n}) = \mathbb{V}(V_{m,n}) = \mathbb{E}V_{m,n}^2 - (\mathbb{E}V_{m,n})^2$, where*

$$\begin{aligned}
 \mathbb{E}V_{m,n}^2 = & \mathbb{E}V_{m,n} \\
 & + m_{(3)} \mathbb{E} \left((1 - |G(X_1) - G(X_2)|)^n (1 - |F(X_1) - F(X_2)|)^{m-3} |F(X_1) - F(X_2)| \right) \\
 & + n_{(3)} \mathbb{E} \left((1 - |F(Y_1) - F(Y_2)|)^m (1 - |G(Y_1) - G(Y_2)|)^{n-3} |G(Y_1) - G(Y_2)| \right) \\
 & + \frac{m_{(4)}}{2} \mathbb{E} \left((1 - (|G(X_1) - G(X_2)| + |G(X_3) - G(X_4)|))^n \right. \\
 & \cdot (1 - (|F(X_1) - F(X_2)| + |F(X_3) - F(X_4)|))^{m-4} \mathbf{1}\{X_1 \vee X_2 < X_3 \wedge X_4\} \Big)
 \end{aligned}$$

$$\begin{aligned}
& + \frac{n_{(4)}}{2} \mathbb{E} \left((1 - (|F(Y_1) - F(Y_2)| + |F(Y_3) - F(Y_4)|))^m \right. \\
& \cdot (1 - (|G(Y_1) - G(Y_2)| + |G(Y_3) - G(Y_4)|))^{n-4} \mathbf{1}\{Y_1 \vee Y_2 < Y_3 \wedge Y_4\} \Big) \\
& + \frac{1}{2} m_{(2)} n_{(2)} \mathbb{E} \left((1 - (|F(Y_1) - F(Y_2)| + |F(X_1) - F(X_2)|))^m \right. \\
& \cdot (1 - (|G(Y_1) - G(Y_2)| + |G(X_1) - G(X_2)|))^{n-2} \mathbf{1}\{X_1 \vee X_2 < Y_1 \wedge Y_2\} \Big) \\
& + \frac{1}{2} m_{(2)} n_{(2)} \mathbb{E} \left((1 - (|F(X_1) - F(X_2)| + |F(Y_1) - F(Y_2)|))^m \right. \\
& \cdot (1 - (|G(X_1) - G(X_2)| + |G(Y_1) - G(Y_2)|))^{n-2} \mathbf{1}\{Y_1 \vee Y_2 < X_1 \wedge X_2\} \Big).
\end{aligned}$$

Proof From

$$\begin{aligned}
\mathbb{E}W_{m,n}(0) &= m - 1 - \mathbb{E}V_{m,n}(0), \\
\mathbb{E}W_{m,n}(1) &= n - 1 - \mathbb{E}V_{m,n}(1),
\end{aligned}$$

and Proposition 1, it follows that

$$\mathbb{E}W_{m,n} = m + n - (m_{(2)}\mu_{m,n}(0) + n_{(2)}\mu_{m,n}(1)),$$

and we have

$$\begin{aligned}
\mu_{m,n}(0) &= \frac{1}{2} \mathbb{E} \left((1 - |G(X_1) - G(X_2)|)^n (1 - |F(X_1) - F(X_2)|)^{m-2} \right), \\
\mu_{m,n}(1) &= \frac{1}{2} \mathbb{E} \left((1 - |F(Y_1) - F(Y_2)|)^m (1 - |G(Y_1) - G(Y_2)|)^{n-2} \right),
\end{aligned}$$

which yields assertion (a). To prove (b), notice that

$$\begin{aligned}
\mathbb{E}V_{m,n}^2 &= m_{(2)}\mu_{m,n}(0) + n_{(2)}\mu_{m,n}(1) + 2(m_{(3)}\mu_{m,n}(0, 1) + n_{(3)}\mu_{m,n}(1, 1)) \\
&+ 2(m_{(4)}\mu_{m,n}(0, 2) + n_{(4)}\mu_{m,n}(1, 2)) + 2m_{(2)}n_{(2)}(\rho_{m,n}(0) + \rho_{m,n}(1)).
\end{aligned}$$

Putting the corresponding pieces together, b) follows. \square

In the hypothesis case, the result of Theorem 1 specializes to the formulas for the expectation and the variance of the runs statistic $W_{m,n}$ stated in Gibbons and Chakraborti (2011). In fact, in this case, the various integrals appearing in Proposition 1 simplify and can be carried out easily; see, for example, the identities (3.2) and (3.3) given below.

3 Limits

Suppose that X_1 and Y_1 have absolutely continuous distribution functions F and G having densities f and g , respectively. In what follows, each of the asymptotic results

given is derived under the usual limiting regime for two-sample problems, namely $m, n \rightarrow \infty$ in such a way that, for some $\tau \in (0, 1)$, the limit

$$\lim_{m, n \rightarrow \infty} \frac{m}{m+n} = \tau \quad (3.1)$$

exists. We will make use of the identities

$$\iint_{u_1 < u_2} (1 - (F(u_2) - F(u_1)))^{m-2} dF(u_1) dF(u_2) = \frac{1}{m}, \quad m \geq 2, \quad (3.2)$$

and

$$\begin{aligned} & \iiint_{u_1 < u_2 < u_3 < u_4} (1 - ((F(u_4) - F(u_3)) + (F(u_2) - F(u_1))))^{m-4} \\ & \quad \times dF(u_1) dF(u_2) dF(u_3) dF(u_4) \\ & = \frac{1}{2m(m-1)}, \quad m \geq 4. \end{aligned} \quad (3.3)$$

These follow readily by introducing the independent random variables $U_j = F(X_j)$, each of which has the uniform distribution on the unit interval $[0, 1]$. For,

$$\begin{aligned} & \iint_{u_1 < u_2} (1 - (F(u_2) - F(u_1)))^{m-2} dF(u_1) dF(u_2) \\ & = \mathbb{E} \left((1 - (F(X_2) - F(X_1)))^{m-2} \mathbf{1}\{X_1 < X_2\} \right) \\ & = \mathbb{E} \left((1 - (U_2 - U_1))^{m-2} \mathbf{1}\{U_1 < U_2\} \right) \\ & = \iint_{u_1 < u_2} (1 - (u_2 - u_1))^{m-2} du_1 du_2 \\ & = \frac{1}{m}. \end{aligned}$$

Analogously,

$$\begin{aligned} & \iiint_{u_1 < u_2 < u_3 < u_4} (1 - ((F(u_4) - F(u_3)) + (F(u_2) - F(u_1))))^{m-4} \\ & \quad \times dF(u_1) dF(u_2) dF(u_3) dF(u_4) \\ & = \iiint_{u_1 < u_2 < u_3 < u_4} (1 - ((u_4 - u_3) + (u_2 - u_1)))^{m-4} du_1 du_2 du_3 du_4 \\ & = \frac{1}{2m(m-1)}. \end{aligned}$$

Proposition 2 *If the densities f and g are continuous almost everywhere then, under the limiting regime (3.1),*

$$\lim_{m,n \rightarrow \infty} \frac{\mathbb{E}V_{m,n}(0)}{m+n} = \int_{-\infty}^{+\infty} \frac{\tau^2 f^2(u)}{\tau f(u) + (1-\tau)g(u)} du, \quad (3.4)$$

$$\lim_{m,n \rightarrow \infty} \frac{\mathbb{E}V_{m,n}(1)}{m+n} = \int_{-\infty}^{+\infty} \frac{(1-\tau)^2 g^2(u)}{\tau f(u) + (1-\tau)g(u)} du, \quad (3.5)$$

and

$$\lim_{m,n \rightarrow \infty} \frac{\mathbb{V}(V_{m,n}(i))}{(m+n)^2} = 0, \quad i \in \{0, 1\}. \quad (3.6)$$

Proof To derive (3.4), we define

$$\begin{aligned} h_{m,n}(u, t) &= \frac{m-1}{m} \left(1 - \left(G\left(u + \frac{t}{m}\right) - G(u) \right) \right)^n \\ &\quad \cdot \left(1 - \left(F\left(u + \frac{t}{m}\right) - F(u) \right) \right)^{m-2} f(u) f\left(u + \frac{t}{m}\right), \end{aligned}$$

and

$$b_{m,n}(u, t) = \frac{m-1}{m} \left(1 - \left(F\left(u + \frac{t}{m}\right) - F(u) \right) \right)^{m-2} f(u) f\left(u + \frac{t}{m}\right)$$

for $-\infty < u < +\infty$, $t \geq 0$, and 0 else. Then,

$$\frac{\mathbb{E}V_{m,n}(0)}{m+n} = \frac{m}{m+n} \int_{-\infty}^{+\infty} \int_0^\infty h_{m,n}(u, t) dt du.$$

Moreover, $0 \leq h_{m,n} \leq b_{m,n}$. If u is a continuity point of f then, for each $t \geq 0$,

$$\begin{aligned} \lim_{m,n \rightarrow \infty} h_{m,n}(u, t) &= \exp\left(-\left(\frac{1-\tau}{\tau}g(u) + f(u)\right)t\right) f^2(u), \\ \lim_{m,n \rightarrow \infty} b_{m,n}(u, t) &= \exp(-f(u)t) f^2(u). \end{aligned}$$

By (3.2), it follows that

$$\begin{aligned} &\lim_{m,n \rightarrow \infty} \int_{-\infty}^{+\infty} \int_0^\infty b_{m,n}(u, t) dt du \\ &= \lim_{m,n \rightarrow \infty} (m-1) \iint_{u < v} (1 - (F(v) - F(u)))^{m-2} f(u) f(v) du dv \\ &= 1 = \int_{-\infty}^{+\infty} \int_0^\infty \exp(-f(u)t) f^2(u) dt du. \end{aligned}$$

Invoking Pratt's extended version of the dominated convergence theorem, see [Pratt \(1960\)](#), we obtain

$$\begin{aligned}\lim_{m,n \rightarrow \infty} \frac{\mathbb{E}V_{m,n}(0)}{m+n} &= \tau \int_{-\infty}^{+\infty} \int_0^{\infty} \exp\left(-\left(\frac{1-\tau}{\tau}g(u) + f(u)\right)t\right) f^2(u) dt du \\ &= \tau \int_{-\infty}^{+\infty} \frac{1}{\frac{1-\tau}{\tau}g(u) + f(u)} f^2(u) du \\ &= \int_{-\infty}^{+\infty} \frac{\tau^2 f^2(u)}{\tau f(u) + (1-\tau)g(u)} du.\end{aligned}$$

In the same way, (3.5) follows. To verify (3.6) for $i = 0$, we start from

$$\mathbb{E}V_{m,n}^2(0) = \mathbb{E}V_{m,n}(0) + 2\mathbb{E}V_{m,n}(0, 1) + 2\mathbb{E}V_{m,n}(0, 2) \quad (3.7)$$

and notice that due (3.4),

$$\lim_{m,n \rightarrow \infty} \frac{\mathbb{E}V_{m,n}(0)}{(m+n)^2} = 0.$$

For $-\infty < u < +\infty$ and $t \geq 0$, put

$$\begin{aligned}k_{m,n}(u, t) &= \frac{m-2}{m} \left(1 - \left(G\left(u + \frac{t}{m}\right) - G(u)\right)\right)^n \left(1 - \left(F\left(u + \frac{t}{m}\right) - F(u)\right)\right)^{m-3} \\ &\quad \cdot \left(F\left(u + \frac{t}{m}\right) - F(u)\right) f(u) f\left(u + \frac{t}{m}\right).\end{aligned}$$

Then,

$$\frac{\mathbb{E}V_{m,n}(0, 1)}{(m+n)^2} = \frac{m_{(2)}}{(m+n)^2} \int_{-\infty}^{+\infty} \int_0^{\infty} k_{m,n}(u, t) dt du$$

and $\lim_{m,n \rightarrow \infty} k_{m,n}(u, t) = 0$ for each continuity point u of f and each $t \geq 0$. Arguing as for (3.4), we obtain

$$\lim_{m,n \rightarrow \infty} \frac{\mathbb{E}V_{m,n}(0, 1)}{(m+n)^2} = 0. \quad (3.8)$$

Putting

$$\begin{aligned}\ell_{m,n}(u, s, w, t) &= \frac{m_{(4)}}{(m+n)^2 m^2} \left(1 - \left(\left(G\left(w + \frac{t}{m}\right) - G(w)\right) + \left(G\left(u + \frac{s}{m}\right) - G(u)\right)\right)\right)^n\end{aligned}$$

$$\cdot \left(1 - \left(\left(F \left(w + \frac{t}{m} \right) - F(w) \right) + \left(F \left(u + \frac{s}{m} \right) - F(u) \right) \right) \right)^{m-4} \\ \cdot f(u) f \left(u + \frac{s}{m} \right) f(w) f \left(w + \frac{t}{m} \right)$$

and

$$d_{m,n}(u, s, w, t) \\ = \frac{m_{(4)}}{(m+n)^2 m^2} \left(1 - \left(\left(F \left(w + \frac{t}{m} \right) - F(w) \right) + \left(F \left(u + \frac{s}{m} \right) - F(u) \right) \right) \right)^{m-4} \\ \cdot f(u) f \left(u + \frac{s}{m} \right) f(w) f \left(w + \frac{t}{m} \right)$$

for $-\infty < u < w < +\infty$, $0 < s < m(w-u)$, $t > 0$, and 0 else, we have

$$\frac{2}{(m+n)^2} \mathbb{E} V_{m,n}(0, 2) = 2 \iiint \ell_{m,n}(u, s, w, t) \, ds \, du \, dt \, dw.$$

Define

$$\ell(u, s, w, t) \\ = \tau^2 \exp \left(- \left(\frac{1-\tau}{\tau} g(u) + f(u) \right) s - \left(\frac{1-\tau}{\tau} g(w) + f(w) \right) t \right) f^2(u) f^2(w)$$

and

$$d(u, s, w, t) = \tau^2 \exp(-f(u)s - f(w)t) f^2(u) f^2(w)$$

for $-\infty < u < w < +\infty$, $s, t > 0$, and 0 else. Then, $0 \leq \ell_{m,n} \leq d_{m,n}$ and

$$\lim_{m,n \rightarrow \infty} \ell_{m,n}(u, s, w, t) = \ell(u, s, v, t), \quad \lim_{m,n \rightarrow \infty} d_{m,n}(u, s, w, t) = d(u, s, v, t)$$

for continuity points $u < w$ of both f and g , and $s, t \geq 0$. By (3.3),

$$\lim_{m,n \rightarrow \infty} \int d_{m,n}(u, s, w, t) \, du \, ds \, dw \, dt = \frac{1}{2} \tau^2 = \iiint \ell(u, s, w, t) \, ds \, du \, dt \, dw.$$

Making again use of the extended dominated convergence theorem, we get

$$\lim_{m,n \rightarrow \infty} \frac{2}{(m+n)^2} \mathbb{E} V_{m,n}(0, 2) \\ = 2 \iiint \ell(u, s, w, t) \, ds \, du \, dt \, dw = \left(\int_{-\infty}^{+\infty} \frac{\tau^2 f^2(u)}{\tau f(u) + (1-\tau)g(u)} \, du \right)^2. \quad (3.9)$$

By (3.7), combining (3.4), (3.8) and (3.9), we obtain

$$\lim_{m,n \rightarrow \infty} \frac{1}{(m+n)^2} \mathbb{V}(V_{m,n}(0)) = 0. \quad (3.10)$$

The corresponding assertion for $\mathbb{V}(V_{m,n}(1))$ is proved in the same way. \square

Define

$$\Delta(F, G, \tau) = 2\tau(1-\tau) \int_{-\infty}^{+\infty} \frac{f(u)g(u)}{\tau f(u) + (1-\tau)g(u)} du. \quad (3.11)$$

In the literature on image processing, $\Delta(F, G, \tau)$ is known as the Henze–Penrose affinity between the distributions F and G with the weights τ and $1-\tau$; see, e.g., van Gemert et al. (2006).

Theorem 2 *Under the conditions of Proposition 2, we have*

$$\lim_{m,n \rightarrow \infty} \frac{1}{m+n} \Delta_{m,n}(F, G) = \Delta(F, G, \tau), \quad (3.12)$$

where

$$\Delta(F, G, \tau) \leq 2\tau(1-\tau). \quad (3.13)$$

Equality holds in (3.13) if, and only if, $F = G$. Moreover,

$$\lim_{m,n \rightarrow \infty} \frac{\sigma_{m,n}^2(F, G)}{(m+n)^2} = 0. \quad (3.14)$$

Proof Remember that $\Delta_{m,n}(F, G) = \mathbb{E}(W_{m,n})$, $\sigma_{m,n}^2(F, G) = \mathbb{V}(W_{m,n})$. From the identity $\mathbb{E}W_{m,n} = m+n - \mathbb{E}V_{m,n}$ stated in part a) of Theorem 1, we deduce that

$$\frac{1}{m+n} \Delta_{m,n}(F, G) = 1 - \frac{1}{m+n} \mathbb{E}V_{m,n}.$$

By (3.4) and (3.5),

$$\lim_{m,n \rightarrow \infty} \mathbb{E}V_{m,n} = \int_{-\infty}^{+\infty} \frac{\tau^2 f(u)^2 + (1-\tau)^2 g(u)^2}{\tau f(u) + (1-\tau)g(u)} du.$$

On writing $1 = \int_{-\infty}^{+\infty} (\tau f(u) + (1-\tau)g(u)) du$, the limit of $\frac{1}{m+n} \Delta_{m,n}(F, G)$ is seen to be $\Delta(F, G, \tau)$, which proves (3.12).

Due to

$$\mathbb{V}(W_{m,n}) = \mathbb{V}(V_{m,n}) = \mathbb{V}(V_{m,n}(0)) + \mathbb{V}(V_{m,n}(1)) + 2\text{Cov}(V_{m,n}(0), V_{m,n}(1))$$

and

$$2|\text{Cov}(V_{m,n}(0), V_{m,n}(1))| \leq \mathbb{V}(V_{m,n}(0)) + \mathbb{V}(V_{m+n}(1)),$$

(3.14) follows from (3.6). The assertion (3.13) is well known, see e.g., [Henze \(1984\)](#). The simple proof is obtained by observing that for non-negative x, y the inequality

$$\frac{xy}{\tau x + (1 - \tau)y} \leq \tau y + (1 - \tau)x$$

(defining $0/0 := 0$) is equivalent to $(x - y)^2 \geq 0$, and by noting that

$$\int_{-\infty}^{\infty} ((1 - \tau)f(u) + \tau g(u)) \, du = 1.$$

□

The consistency of the test is a simple consequence of Theorem 2. By Chebyshev's inequality,

$$\mathbb{P}\left(W_{m,n} \leq \Delta_{m,n} - \sqrt{\frac{1}{\alpha - \varepsilon} \sigma_{m,n}^2}\right) \leq \alpha - \varepsilon$$

for each $0 < \varepsilon < \alpha$. Thus,

$$\Delta_{m,n} - \sqrt{\frac{1}{\alpha - \varepsilon} \sigma_{m,n}^2} < c_{m,n;\alpha}$$

for each $0 < \varepsilon < \alpha$ and therefore

$$\Delta_{m,n} - \sqrt{\frac{1}{\alpha} \sigma_{m,n}^2} \leq c_{m,n;\alpha}.$$

If $F \neq G$ we have $\varepsilon = 2\tau(1 - \tau) - \Delta(F, G, \tau) > 0$, whence for large m, n

$$\mathbb{P}(W_{m,n} \geq c_{m,n;\alpha}) \leq \mathbb{P}\left(\frac{1}{m+n} W_{m,n} - \Delta(F, G, \tau) \geq \frac{\varepsilon}{2}\right).$$

Again applying Chebyshev's inequality, the assertion

$$\lim_{m,n \rightarrow \infty} \mathbb{P}(W_{m,n} \geq c_{m,n;\alpha}) = 0$$

follows.

4 Final remarks

1. In the null hypothesis case, the test statistic $W_{m,n}$ has a limiting normal distribution see, e.g., [Gibbons and Chakraborti \(2011\)](#). A corresponding result for fixed alternative distributions also exists, see [Blumenthal \(1963\)](#). It is of interest to obtain the latter using our approach. Although presently we are unable to do so, we believe that our method may be refined to indeed yield the appropriate weak convergence result. We hope to address this issue in our future work.
2. The main results of this paper are explicit expressions for the expectation and the variance of the celebrated runs statistic under a fixed alternative. The expectation may be regarded as a ‘measure of non-centrality’, and the power of the test should increase with increasing values of $\Delta_{m,n}(F, F) - \Delta_{m,n}(F, G)$. Although the limit $\Delta(F, G, \tau)$ figuring in (3.14) is smaller than $\Delta(F, F, \tau) = 2\tau(1-\tau)$ if $F \neq G$ [cf. (3.13)], one may conjecture that $\Delta(F, G) = \mathbb{E}W_{m,n} \leq \Delta(F, F) = 1 + 2mn/(m + m)$ for any choice of m and n . We leave this question as an open problem.

Acknowledgments The authors thank the referees for constructive comments and helpful suggestions.

References

- Biswas M, Mukhopadhyay M, Ghosh AK (2014) A distribution-free two-sample run test applicable to high-dimensional data. *Biometrika* 101:913–926
- Blum J, Weiss L (1957) Consistency of certain two-sample tests. *Ann Math Stat* 28:242–246
- Blumenthal S (1963) The asymptotic normality of two test statistics associated with the two-sample problem. *Ann Math Stat* 34:1513–1523
- Cohen J, Menjoge S (1988) One-sample-run tests of symmetry. *J Stat Plan Inference* 18:93–100
- Dyckerhoff R, Ley C, Paidaveine D (2015) Depth-based runs tests for bivariate central symmetry. *Ann Inst Statist Math* 67:917–941
- Friedman J, Rafsky L (1979) Multivariate generalizations of the Wolfowitz and Smirnov two sample tests. *Ann Stat* 7:697–717
- Gibbons J, Chakraborti S (2011) Nonparametric statistical inference. CRC Press, Boca Raton
- Govindarajulu Z (2007) Nonparametric inference. World Scientific Publishing, Singapore
- Henze N (1984) On the number of random points with nearest neighbors of the same type and a multivariate two-sample test (in German). *Metrika* 31:259–273
- Henze N (1993a) On the consistency of a test for symmetry based on a runs statistic. *J Nonparametr Stat* 3:195–199
- Henze N (1993b) A quick omnibus test for the proportional hazards model of random censorship. *Statistics* 24:253–263
- Henze N, Penrose M (1999) On the multivariate runs test. *Ann Stat* 27:290–298
- Henze N, Voigt B (1992) Almost sure convergence of certain slowly changing symmetric one- and multi-sample statistics. *Ann Probab* 20:1086–1098
- Marden JI (1999) Multivariate rank tests. In: Multivariate analysis, design of experiments, and survey sampling, volume 159 of *Statist. Textbooks Monogr.* New York: Dekker, pp 401–432
- McWilliams TP (1990) A distribution-free test for symmetry based on a runs statistic. *J Am Stat Assoc* 85:1130–1133
- Mood A (1954) On the asymptotic efficiency of certain nonparametric two-sample tests. *Ann Math Stat* 25:514–522
- Morgenstern D (1962) Zur asymptotik des run-testes. *Metrika* 5:150–153
- Paidaveine D (2009) On multivariate runs tests for randomness. *J Am Stat Assoc* 104:1525–1538
- Pratt J (1960) On interchanging limits and integrals. *Ann Math Stat* 31:74–77
- van Gemert J, Burghouts G, Seinstra F, Geusebroek J-M (2006) Color invariant object recognition using entropic graphs. *Int J Imaging Syst Technol* 16:146–153

- Wald A, Wolfowitz J (1940) On a test whether two samples are from the same population. *Ann Math Stat* 11:147–162
- Weiss L (1955) The stochastic convergence of a function of sample successive differences. *Ann Math Stat* 26:532–536