

## Incomplete interdirections and lift-interdirections

Šárka Hudecová, Jana Klicnarová & Miroslav Šiman

To cite this article: Šárka Hudecová, Jana Klicnarová & Miroslav Šiman (2019): Incomplete interdirections and lift-interdirections, Journal of Nonparametric Statistics, DOI: [10.1080/10485252.2019.1700255](https://doi.org/10.1080/10485252.2019.1700255)

To link to this article: <https://doi.org/10.1080/10485252.2019.1700255>



Published online: 06 Dec 2019.



Submit your article to this journal [↗](#)



Article views: 23



View related articles [↗](#)



View Crossmark data [↗](#)



## Incomplete interdirections and lift-interdirections

Šárka Hudecová<sup>a</sup>, Jana Klicnarová<sup>b,c</sup> and Miroslav Šiman<sup>b</sup>

<sup>a</sup>Department of Probability and Statistics, Faculty of Mathematics and Physics, Charles University, Praha 8, Czech Republic; <sup>b</sup>The Czech Academy of Sciences, Institute of Information Theory and Automation, Praha 8, Czech Republic; <sup>c</sup>Faculty of Economics, University of South Bohemia in České Budějovice, České Budějovice, Czech Republic

### ABSTRACT

The article presents, discusses, and explores incomplete variants of interdirections, lift-interdirections, and symmetrised lift-interdirections (with a few incomplete designs). Although they are easier to compute in high-dimensional spaces than the originals, they can still replace them in many optimal statistical procedures based on signs and ranks without significantly changing their properties. This is proved theoretically and confirmed empirically in a small simulation study dealing with the canonical examples of multivariate sign and signed-rank one-sample tests applied to high-dimensional data sets.

### ARTICLE HISTORY

Received 16 May 2019  
Accepted 28 November 2019

### KEYWORDS

Interdirection; multivariate rank; multivariate sign; signed-rank test; robustness

## 1. Introduction

Simple, robust, and powerful univariate sign and rank statistical procedures with weak assumptions have already inspired many multivariate generalisations. Current literature on multivariate rank statistics is largely dominated by spatial ranks and signs (Möttönen and Oja 1995; Oja 2010), followed by component-wise rank and sign vectors (Puri and Sen 1971), Oja multivariate ranks and signs (Oja 1999), the ranks of pseudo-Mahalanobis distances (Hallin and Paindaveine 2002b), and various depth concepts (Zuo and Serfling 2000); see also Chaudhuri and Sengupta (1993) and Chernozhukov, Galichon, Hallin, and Henry (2017). Unfortunately, the signs and ranks based on the purely geometric hyperplane-based concepts of interdirections (Randles 1989) and lift-interdirections (Hettmansperger, Möttönen, and Oja 1999; Oja and Paindaveine 2005) seem to fall out of fashion due to their high computational demands, even though they avoid any estimation of the shape matrix.

This article can be viewed as an attempt to revive these two appealing concepts and to bring them back to the forefront. In particular, it presents incomplete interdirections, incomplete lift-interdirections and incomplete symmetrised lift-interdirections based on incomplete  $U$ -statistics. They allow for various random and deterministic designs, they are relatively quick to compute in high-dimensional spaces, and they can replace their

complete counterparts in many sign and signed-rank test statistics without changing their asymptotic behaviour.

In other words, this article addresses the problem of too many hyperplanes passing through certain number of observations by working only with some of them. Such an approach is hardly surprising, and it was already suggested in the concluding section of Oja and Paindaveine (2005). This article shows that there is a clever way to do the selection, analyses it, and illustrates it with multidimensional data.

The signs and ranks based on incomplete interdirections and incomplete (symmetrised) lift-interdirections might still be used in the large number of statistical tests mentioned in Oja and Paindaveine (2005), including the test of randomness against VARMA dependence, the two-sample and multi-sample tests, the Durbin-Watson test and the test of the order of a VARMA model. In particular, the incomplete interdirections might be used in any sign test where the original interdirections appear; see, e.g. Randles (1989), Randles and Peters (1990), Jan and Randles (1996), Gieser and Randles (1997), Um and Randles (1998), Ghosh and Sengupta (2001), Hallin and Paindaveine (2002a, 2002b, 2004), Taskinen, Oja, and Randles (2005), and Paindaveine (2009). Here, it is proved rigorously only for the one-sample tests of Randles (1989) and Oja and Paindaveine (2005). Note also that all the resulting tests can be expected to be robust with respect to both radial and angular outliers like other hyperplane-based testing procedures (Oja and Paindaveine 2005).

Multivariate rank tests often need some symmetry to work. The sign tests discussed here work for all distributions with elliptical directions (Randles 1989) while the signed-rank tests require elliptical symmetry as those in Oja and Paindaveine (2005); see Serfling (2006) for a review of the most important concepts of multivariate symmetry. These assumptions permit even some distributions that are not unimodal. Such distributions may arise easily in the context of mixtures even in the univariate case; see, e.g. Došlá (2009).

Next Section 2 introduces incomplete interdirections, lift-interdirections, and their symmetric variant. Section 3 then describes their properties and justifies their use in the one-sample tests of Randles (1989) and Oja and Paindaveine (2005), and Section 4 employs them in a brief comparative simulation study using also large sample sizes and dimensions. Section 5 discusses the results and achievements. Final Appendix collects the proofs.

## 2. Definitions and notation

Let  $\mathcal{X}_n$  be a random sample consisting of  $n$   $p$ -dimensional observations  $\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathbb{R}^p$ ,  $p \geq 2$ . Any ordered  $k$ -tuple  $\mathbf{q}(n, k) = (q_1, \dots, q_k)$  of distinct integer indices  $1 \leq q_1 < q_2 < \dots < q_k \leq n$  then defines subsample  $\mathcal{X}^{\mathbf{q}(n, k)} = (\mathbf{X}_{q_1}, \dots, \mathbf{X}_{q_k})$ . The set of all possible  $k$ -tuples  $\mathbf{q}(n, k)$  will be denoted by  $\mathcal{Q}(n, k)$ . It has exactly  $\binom{n}{k}$  elements.

If  $k = p$ , then  $H^{\mathbf{q}(n, p)} \subset \mathbb{R}^p$  is defined as the hyperplane that contains all the observations from  $\mathcal{X}^{\mathbf{q}(n, p)}$ . If  $k = p-1$ , then  $H^{\mathbf{q}(n, p-1)} \subset \mathbb{R}^p$  is defined as the hyperplane containing all the observations from  $\mathcal{X}^{\mathbf{q}(n, p-1)}$  and the origin.

Any hyperplane  $H^{\mathbf{q}(n, p-1)}$  can be expressed by the equation

$$\det(\mathbb{M}^{\mathbf{q}(n, p-1)}(\mathbf{y})) = \mathbf{d}^{\mathbf{q}(n, p-1)'} \mathbf{y} = 0$$

where  $\mathbf{d}^{\mathbf{q}(n,p-1)} = (d_1^{\mathbf{q}(n,p-1)}, \dots, d_p^{\mathbf{q}(n,p-1)})'$  and  $d_j^{\mathbf{q}(n,p-1)}$ ,  $j = 1, \dots, p$ , is the cofactor of the  $j$ th element in the last column of matrix

$$\mathbb{M}^{\mathbf{q}(n,p-1)}(\mathbf{y}) = (\mathbf{X}_{q_1}, \mathbf{X}_{q_2}, \dots, \mathbf{X}_{q_{p-1}}, \mathbf{y}).$$

Similarly, any hyperplane  $H^{\mathbf{q}(n,p)}$  can be expressed by the equation

$$\det(\mathbb{M}^{\mathbf{q}(n,p)}(\mathbf{y})) = d_0^{\mathbf{q}(n,p)} + \mathbf{d}^{\mathbf{q}(n,p)'} \mathbf{y} = 0$$

where  $\mathbf{d}^{\mathbf{q}(n,p)} = (d_1^{\mathbf{q}(n,p)}, \dots, d_p^{\mathbf{q}(n,p)})'$  and  $d_j^{\mathbf{q}(n,p)}$ ,  $j = 0, \dots, p$ , is the cofactor of the  $(j + 1)$ th element in the last column of matrix

$$\mathbb{M}^{\mathbf{q}(n,p)}(\mathbf{y}) = \begin{pmatrix} 1 & 1 & \dots & 1 & 1 \\ \mathbf{X}_{q_1} & \mathbf{X}_{q_2} & \dots & \mathbf{X}_{q_p} & \mathbf{y} \end{pmatrix}.$$

The signs

$$S^{\mathbf{q}(n,p-1)}(\mathbf{y}) = \text{sign}(\mathbf{d}^{\mathbf{q}(n,p-1)'} \mathbf{y}) \quad \text{and} \quad S^{\mathbf{q}(n,p)}(\mathbf{y}) = \text{sign}(d_0^{\mathbf{q}(n,p)} + \mathbf{d}^{\mathbf{q}(n,p)'} \mathbf{y})$$

then reveal the position of  $\mathbf{y} \in \mathbb{R}^p$  with respect to  $H^{\mathbf{q}(n,p-1)}$  and  $H^{\mathbf{q}(n,p)}$ , respectively.

If  $\mathcal{X}_n$  comes from the spherically symmetric distribution centred around the origin, then the angular distance

$$\alpha(\mathbf{y}_1, \mathbf{y}_2) = \arccos\left(\frac{\mathbf{y}'_1 \mathbf{y}_2}{\|\mathbf{y}_1\| \|\mathbf{y}_2\|}\right)$$

between any two points  $\mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^p$  can be measured by means of (affine invariant) interdirections

$$C_{\mathbf{y}_1, \mathbf{y}_2}(\mathcal{X}_n) = \sum_{\mathbf{q} \in \mathcal{Q}(n,p-1)} (1 - S^{\mathbf{q}}(\mathbf{y}_1) S^{\mathbf{q}}(\mathbf{y}_2)) / 2;$$

see Randles (1989). (They are defined by means of all hyperplanes passing through the origin and  $p-1$  observations, but only those hyperplanes separating  $\mathbf{y}_1$  from  $\mathbf{y}_2$  are actually counted.) That is to say that

$$a_{\mathbf{y}_1, \mathbf{y}_2} = \pi C_{\mathbf{y}_1, \mathbf{y}_2}(\mathcal{X}_n) / \binom{n}{p-1}$$

is then an affine invariant consistent estimator of  $\alpha(\mathbf{y}_1, \mathbf{y}_2)$ . Needless to say that the angles and their estimators play a crucial role in various multivariate sign and signed-rank tests.

Oja and Paindaveine (2005) suggest to measure the empirical distance between two points  $\mathbf{y}_1$  and  $\mathbf{y}_2$  by means of (affine invariant) lift-interdirection

$$L_{\mathbf{y}_1, \mathbf{y}_2}(\mathcal{X}_n) = \sum_{\mathbf{q} \in \mathcal{Q}(n,p)} (1 - S^{\mathbf{q}}(\mathbf{y}_1) S^{\mathbf{q}}(\mathbf{y}_2)) / 2$$

that counts all hyperplanes passing through  $p$  observations and separating  $\mathbf{y}_1$  from  $\mathbf{y}_2$ . They consider especially the particular case  $L_{\mathbf{y}, -\mathbf{y}}(\mathcal{X}_n) = L_{\mathbf{y}, -\mathbf{y}}(\mathcal{X}_n)$ ,  $\mathbf{y} \in \mathbb{R}^p$ , and its symmetrised

version

$$\underline{L}_y(\mathcal{X}_n) = \sum_{\substack{q \in \mathcal{Q}(n,p) \\ s \in \{-1,1\}^p}} \frac{1 - \text{sign}(d_{0s}^{q(n,p)} + \mathbf{d}_s^{q(n,p)'} \mathbf{y}) \text{sign}(d_{0s}^{q(n,p)} - \mathbf{d}_s^{q(n,p)'} \mathbf{y})}{2}$$

where  $(d_{0s}^{q(n,p)}, \mathbf{d}_s^{q(n,p)'})'$  is the vector of cofactors of the last column of matrix

$$\mathbb{M}_s^{q(n,p)}(\mathbf{y}) = \begin{pmatrix} 1 & 1 & \cdots & 1 & 1 \\ s_1 \mathbf{X}_{q_1} & s_2 \mathbf{X}_{q_2} & \cdots & s_p \mathbf{X}_{q_p} & \mathbf{y} \end{pmatrix},$$

and  $\{-1, 1\}^p$  is the set of all  $2^p$   $p$ -dimensional vectors  $\mathbf{s} = (s_1, \dots, s_p)'$  with individual coordinates equal to either 1 or  $-1$ . That is to say that the symmetrised lift-interdirections  $\underline{L}_y(\mathcal{X}_n)$  are invariant with respect to affine data transformations, permutations of the observations, and reflections of the observations with respect to the origin; see Oja and Paindaveine (2005). The ranks of  $\underline{L}_y(\mathcal{X}_n)$  can thus replace the ranks of pseudo-Mahalanobis distances in many tests based on them; see ibidem.

The cardinality of  $\mathcal{Q}(n, p-1)$  may quickly become impractical for the computation of  $a_{y_1, y_2}$  with growing  $n$  and  $p$ , which is why it seems advantageous to define *incomplete interdirections* as incomplete  $U$ -statistics (see, e.g. Lee 1990)

$$\tilde{C}_{y_1, y_2}(\mathcal{X}_n) = \sum_{q \in \mathcal{Q}_S(n, p-1)} (1 - S^q(\mathbf{y}_1) S^q(\mathbf{y}_2)) / 2$$

considering only some hyperplanes determined by some (design)  $\mathcal{Q}_S(n, p-1) \subset \mathcal{Q}(n, p-1)$  with  $m_{n, p-1}$  elements.

One can also analogously define *incomplete lift-interdirections*  $\tilde{L}_{y_1, y_2}(\mathcal{X}_n)$

$$\tilde{L}_{y_1, y_2}(\mathcal{X}_n) = \sum_{q \in \mathcal{Q}_S(n, p)} (1 - S^q(\mathbf{y}_1) S^q(\mathbf{y}_2)) / 2$$

for some (design)  $\mathcal{Q}_S(n, p) \subset \mathcal{Q}(n, p)$  with  $m_{n, p}$  elements, and *incomplete symmetrised lift-interdirections*:

$$\tilde{\underline{L}}_y(\mathcal{X}_n) = \sum_{\substack{q \in \mathcal{Q}_S^s(n, p) \\ s \in \{-1, 1\}^p}} \left( 1 - \text{sign}(d_{0s}^{q(n,p)} + \mathbf{d}_s^{q(n,p)'} \mathbf{y}) \text{sign}(d_{0s}^{q(n,p)} - \mathbf{d}_s^{q(n,p)'} \mathbf{y}) \right) / 2$$

for some (design)  $\mathcal{Q}_S^s(n, p)$  of elements from  $\mathcal{Q}(n, p)$  with  $m_{n, p}^s$  elements.

To be more specific, the designs  $\mathcal{Q}_S(n, p-1)$ ,  $\mathcal{Q}_S(n, p)$ , and  $\mathcal{Q}_S^s(n, p)$  are arbitrary collections of predefined deterministic sizes  $m_{n, p-1}$ ,  $m_{n, p}$ , and  $m_{n, p}^s$  containing (not necessarily distinct) elements from  $\mathcal{Q}(n, p-1)$ ,  $\mathcal{Q}(n, p)$ , and  $\mathcal{Q}(n, p)$  in this order. They are called sets or denoted as subsets by slight but useful abuse of terminology. The subscript  $S$  indicates the (design) subsets and the index  $s$  indicates a relation to the symmetrised lift-interdirections.

These considerations lead to a few meaningful rank concepts in the multivariate space because one can easily write  $R_i$ ,  $\tilde{R}_i$ , and  $\tilde{\underline{R}}_i$  for the ranks of  $\mathbf{X}_i' \boldsymbol{\Sigma}^{-1} \mathbf{X}_i$ ,  $\tilde{L}_{X_i}(\mathcal{X}_n)$ , and  $\tilde{\underline{L}}_{X_i}(\mathcal{X}_n)$  in the corresponding samples of the same quantities for  $i = 1, \dots, n$ . Of course, these ranks would be meaningful only for centred observations.

### 3. Theoretical considerations

Fortunately, both  $\tilde{C}_{y_1, y_2}(\mathcal{X}_n)/m_{n, p-1}$  and  $\tilde{L}_{y_1, y_2}(\mathcal{X}_n)/m_{n, p}$ ,  $y_1 \neq y_2$ , are non-degenerate incomplete  $U$ -statistics (with kernel  $h_{y_1, y_2}$ ) in the form

$$\tilde{U}_n^{y_1, y_2} = \frac{1}{m_{n, k}} \sum_{\mathbf{q}(n, k) \in \mathcal{Q}_S(n, k)} h_{y_1, y_2}(\mathbf{q}(n, k)), \quad m_{n, k} = |\mathcal{Q}_S(n, k)|,$$

(Note also that even  $\tilde{L}_{y_1, y_2}(\mathcal{X}_n)$  may be viewed as an incomplete  $U$ -statistic if the design set is really in the form specified above.)

The properties of  $\tilde{U}_n^{y_1, y_2}$  and relations to the corresponding complete  $U$ -statistics

$$U_n^{y_1, y_2} = \frac{1}{\binom{n}{k}} \sum_{\mathbf{q}(n, k) \in \mathcal{Q}(n, k)} h_{y_1, y_2}(\mathbf{q}(n, k))$$

have already been thoroughly investigated in the literature; see Blom (1976) and Janson (1984) for most of the results presented here. All the references below nevertheless point to the book Lee (1990) for simplicity. Note also that both  $\tilde{U}_n^{y_1, y_2}$  and  $U_n^{y_1, y_2}$  are unbiased.

Here, the kernel is

$$h_{y_1, y_2}(\mathbf{q}(n, k)) = \left(1 - S^{\mathbf{q}(n, k)}(y_1)S^{\mathbf{q}(n, k)}(y_2)\right) / 2$$

and  $k$  stands for the kernel order equal to the dimension of  $\mathbf{q}(n, k)$ :  $k = p - 1$  for  $\tilde{C}_{y_1, y_2}(\mathcal{X}_n)$ , and  $k = p$  for  $\tilde{L}_{y_1, y_2}(\mathcal{X}_n)$ . The kernel, despite the simplified notation, actually depends on  $k$  independent and identically distributed random variables from the random sample  $\mathcal{X}_n$ . Similarly, all the arguments and indices may be omitted when no confusion is possible. In particular,  $\mathcal{Q}_S(n, k)$  may be shortened to  $\mathcal{Q}_S$  and  $m_{n, k}$  to  $m_n$ .

The first question concerns the choice of the design set  $\mathcal{Q}_S$  and its influence on the asymptotic relative efficiency (ARE) of  $\tilde{U}_n$  with respect to  $U_n$ :

$$ARE(\tilde{U}_n, U_n) = \lim_{n \rightarrow \infty} \frac{\text{var} U_n}{\text{var} \tilde{U}_n}.$$

The following discussion is limited to the three most important cases when the design elements of  $\mathcal{Q}_S$  are set deterministically or chosen randomly with or without replacement. Then the theory of incomplete  $U$ -statistics answers the problem satisfactorily (Lee 1990, Section 4.3).

Define

$$h_c(\mathbf{x}_1, \dots, \mathbf{x}_c) := E h(\mathbf{x}_1, \dots, \mathbf{x}_c, \mathbf{X}_{c+1}, \dots, \mathbf{X}_k),$$

$$p_c := P(h_c(\mathbf{X}_1, \dots, \mathbf{X}_c) = 1), \quad \text{and}$$

$$\sigma_c^2 := \text{var}(h_c(\mathbf{X}_1, \dots, \mathbf{X}_c)),$$

$c = 1, \dots, k$ . Generally,  $\sigma_d^2 \geq (d/c)\sigma_c^2$  for any  $k \geq d \geq c \geq 1$  (Lee 1990, Theorem 4 on p. 15). In the special context considered here, each  $h_c$  follows the Bernoulli distribution  $B(1, p_c)$  and, therefore,  $\sigma_c^2 = p_c(1 - p_c)$ , which also depends on  $y_1$  and  $y_2$ . These quantities influence the asymptotic relative efficiencies stated below.

Recall that  $\lim_{n \rightarrow \infty} n \text{var} U_n = k^2 \sigma_1^2$  (Lee 1990, Theorem 3 on p. 12). Furthermore, Theorem 4 on p. 193 there makes it possible to obtain the variance of  $\tilde{U}_n$  for the two special random designs considered. Consequently, if  $\mathcal{Q}_S$  is chosen from  $\mathcal{Q}$  by random sampling with replacement, then

$$ARE(\tilde{U}_n, U_n) = \lim_{n \rightarrow \infty} \left( \frac{n}{k^2 m_n} \frac{\sigma_k^2}{\sigma_1^2} + 1 - \frac{1}{m_n} \right)^{-1},$$

and if  $\mathcal{Q}_S$  is chosen from  $\mathcal{Q}$  by random sampling without replacement, then

$$ARE(\tilde{U}_n, U_n) = \lim_{n \rightarrow \infty} \left( \frac{(C_{n,k} - m_n)}{(C_{n,k} - 1)} \frac{n}{k^2 m_n} \frac{\sigma_k^2}{\sigma_1^2} + \frac{C_{n,k}}{C_{n,k} - 1} \left( 1 - \frac{1}{m_n} \right) \right)^{-1},$$

where  $C_{n,k} = \binom{n}{k}$ . The difference between the two sampling schemes is thus minimal if both  $m_n$  and  $n$  are large but  $m_n \ll C_{n,k}$ .

Fixed designs are interesting only when they are simple or optimal. In the latter case, they should minimise the variance of  $\tilde{U}_n$  and, therefore, maximise its  $ARE(\tilde{U}_n, U_n)$ . In view of symmetry, it seems natural to consider only equireplicate (or, balanced) designs when each index appears in the same number, say  $r$ , of the elements of  $\mathcal{Q}_S$ . Then  $mk = nr$ .

If the simplest balanced design  $\mathcal{Q}_S = \{(1, \dots, k), (k + 1, \dots, 2k), \dots, (mk - k + 1, \dots, n)\}$  is considered for  $n = mk$ , then

$$ARE(\tilde{U}_n, U_n) = k \frac{\sigma_1^2}{\sigma_k^2} \leq 1.$$

In general,  $\mathcal{Q}_S$  is a minimum variance design, if, for each  $v = 1, 2, \dots, k - 1$ , every  $v$ -subset of  $\{1, 2, \dots, n\}$  is contained in the same number  $m \binom{k}{v} / \binom{n}{v}$  of elements of  $\mathcal{Q}_S$  (Lee 1990, Theorem 1 on p. 195). For example, if  $\mathcal{Q}_S(n, 2)$  is a balanced design of  $\tilde{U}_n$ , then it must be optimal and

$$ARE(\tilde{U}_n, U_n) = \frac{2r}{2(r - 1) + \frac{\sigma_2^2}{\sigma_1^2}} \leq 1$$

increases with  $r$  (Lee 1990, p. 196).

Write  $\theta$  for the common expectation of unbiased  $\tilde{U}_n$  and  $U_n$ :  $\theta = E\tilde{U}_n = EU_n$ .

For the two random designs considered here, the justification for the use of incomplete statistics instead of the complete ones follows from Theorem 1 on p. 200 of Lee (1990) restated here:

**Theorem 3.1:** *Assume  $\min(n, m_n) \rightarrow \infty$ ,  $m_n/n \rightarrow \infty$ , and that  $\mathcal{Q}_S$  results from the random sampling with or without replacement. Then  $\sqrt{n}(\tilde{U}_n - \theta)$  and  $\sqrt{n}(U_n - \theta)$  have the same asymptotic distribution.*

Analogous Theorem 3 on p. 211 of Lee (1990) for fixed balanced designs denotes with  $f_{cn}$  the number of pairs  $q_1, q_2 \in \mathcal{Q}_S$  that have exactly  $c$  elements in common,  $c = 0, 1, \dots, k$ :

**Theorem 3.2:** *If  $\mathcal{Q}_S$  is a balanced design,  $\min(n, m_n) \rightarrow \infty$ ,  $nf_{cn}/m_n^2 \rightarrow 1$  for  $c = 0$ , and  $nf_{cn}/m_n^2 \rightarrow 0$  for  $c = 1, \dots, k$ , then  $\sqrt{n}(\tilde{U}_n - \theta)$  and  $\sqrt{n}(U_n - \theta)$  have the same asymptotic distribution.*

Note that the assumption of Theorem 3.1 is satisfied even for  $m_n \ll \binom{n}{k}$  such as  $m_n = O(n \log(\log(\log(n))))$ , irrespective of the fixed kernel order  $k$ . This fact seems rather unappreciated in the field of multivariate statistics.

Now assume that  $X_1, \dots, X_n$  is a random sample from an elliptical distribution with median vector  $\theta \in \mathbb{R}^p$ , positive definite scatter matrix  $\Sigma \in \mathbb{R}^{p \times p}$ , and density proportionate to

$$f(\sqrt{(\mathbf{x} - \theta)' \Sigma^{-1} (\mathbf{x} - \theta)}), \quad \mathbf{x} \in \mathbb{R}^p,$$

where function  $f : [0, \infty) \rightarrow [0, \infty)$  satisfies  $\int_0^\infty z^{p-1} f(z) dz < \infty$ . Then  $\|\Sigma^{-1/2}(X_1 - \theta)\|$  has cumulative distribution function  $F_r$  with density  $f_r(z) \propto z^{p-1} f(z) I[z > 0]$ .

Without any loss of generality, consider the one-sample testing problem with the null hypothesis  $H_0 : \theta = \mathbf{0}$  and alternative  $H_1 : \theta \neq \mathbf{0}$ . Randles (1989) proposed a test for  $H_0$  using interdirections and showed that the test statistic

$$V_B := \frac{p}{n} \sum_{i=1}^n \sum_{j=1}^n \cos(a_{X_i, X_j})$$

has an asymptotic (distribution-free)  $\chi_p^2$  null distribution with  $p$  degrees of freedom in the whole class of elliptical distributions.

It appears that the test remains valid even when the interdirections are replaced with incomplete interdirections.

**Proposition 3.3:** *Assume a random sample  $X_1, \dots, X_n$  from an elliptical distribution and consider incomplete interdirections  $\tilde{C}_{X_i, X_j}(\mathcal{X}_n)$ ,  $i, j = 1, \dots, n$ , based on  $\mathcal{Q}_S(n, p - 1)$  (with cardinality  $|\mathcal{Q}_S(n, p - 1)| = m_{n, p-1}$ ) chosen randomly with or without replacement from  $\mathcal{Q}(n, p - 1)$ . If  $H_0$  holds and  $m_{n, p-1}/n \rightarrow \infty$ , then*

$$S_B := \frac{p}{n} \sum_{i=1}^n \sum_{j=1}^n \cos(\tilde{a}_{X_i, X_j}^{m_{n, p-1}}) \xrightarrow{\mathcal{D}} \chi_p^2 \quad \text{for } n \rightarrow \infty \tag{1}$$

where  $\tilde{a}_{X_i, X_j}^{m_{n, p-1}} = \pi \tilde{C}_{X_i, X_j}(\mathcal{X}_n) / m_{n, p-1}$ .

Note that Proposition 3.3 holds even for the distributions with elliptical directions (Randles 1989) because of the same arguments contained in that article.

Also the three types of ranks considered here are asymptotically equivalent.

**Proposition 3.4:** *If  $X_1, \dots, X_n$  is a random sample of size  $n$  from an elliptical distribution,  $m_{n, p}/n \rightarrow \infty$ , and  $m_{n, p}^s/n \rightarrow \infty$ , then*

$$\frac{\tilde{R}_i}{n+1} = \frac{R_i}{n+1} + o_p(1), \quad \text{and} \quad \frac{\tilde{\underline{R}}_i}{n+1} = \frac{R_i}{n+1} + o_p(1)$$

as  $n \rightarrow \infty$ .

Oja and Paindaveine (2005) proposed a test statistic for  $H_0$  that also has the asymptotic  $\chi_p^2$  null distribution for all elliptical distributions and that uses both interdirections and ranks of symmetrised lift-interdirections which can still be replaced with their incomplete variants without changing the limiting null distribution under mild assumptions.

**Assumption A:** Assume a function  $K$  that is continuous almost everywhere on  $(0,1)$  and satisfies  $(1/n) \sum_{i=1}^n |K(i/(n+1))|^{2+\delta} \rightarrow \int_0^1 |K(u)|^{2+\delta} du < \infty$  for some  $\delta > 0$ .

**Proposition 3.5:** Assume that a random sample  $X_1, \dots, X_n$  comes from an elliptical distribution and that a function  $K$  satisfies Assumption A. Consider incomplete interdirections  $\tilde{C}_{X_i, X_j}(\mathcal{X}_n)$  and incomplete symmetrised lift-interdirections  $\tilde{L}_{X_i}(\mathcal{X}_n)$ ,  $i, j = 1, \dots, n$ , that are based, respectively, on the design sets  $\mathcal{Q}_S(n, p-1)$ ,  $|\mathcal{Q}_S(n, p-1)| = m_{n,p-1}$ , and  $\mathcal{Q}_S^s(n, p)$ ,  $|\mathcal{Q}_S^s(n, p)| = m_{n,p}^s$ , chosen independently and randomly with or without replacement. If  $H_0$  holds,  $m_{n,p-1}/n \rightarrow \infty$ , and  $m_{n,p}^s/n \rightarrow \infty$ , then

$$S_A := \frac{p}{n E K^2(V)} \sum_{i=1}^n \sum_{j=1}^n K\left(\frac{\tilde{R}_i}{n+1}\right) K\left(\frac{\tilde{R}_j}{n+1}\right) \cos(\tilde{a}_{X_i, X_j}^{m_{n,p-1}}) \xrightarrow{\mathcal{D}} \chi_p^2 \quad (2)$$

as  $n \rightarrow \infty$  where  $\tilde{a}_{X_i, X_j}^{m_{n,p-1}} = \pi \tilde{C}_{X_i, X_j}(\mathcal{X}_n)/m_{n,p-1}$  and  $V$  is uniformly distributed on  $[0, 1]$ .

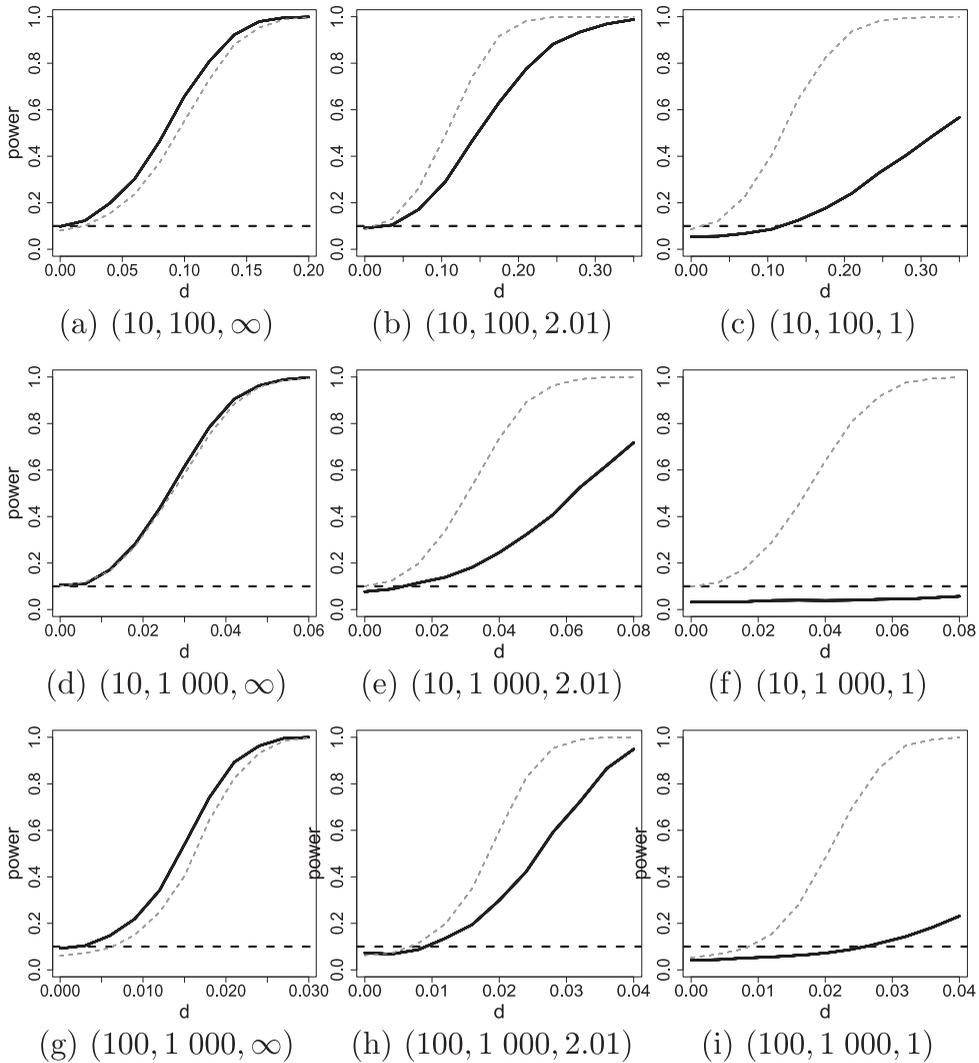
All the propositions stated above are proved in the Appendix.

#### 4. Simulation study

Consider a  $p$ -dimensional random sample  $\mathcal{X}_n$  of size  $n$ ,  $X_1, \dots, X_n$ , that comes from an elliptical distribution with location parameter  $\theta$ . Use random selection with replacement to independently construct design sets  $\mathcal{Q}_S(n, p-1)$  (of size  $N_H$ ) and  $\mathcal{Q}_S^s(n, p)$  (also of size  $N_H$ ) and compute incomplete interdirections  $\tilde{C}_{X_i, X_j}(\mathcal{X}_n)$  and incomplete symmetrised lift-interdirections  $\tilde{L}_{X_i}(\mathcal{X}_n)$ ,  $i, j = 1, \dots, n$ . Use van der Waerden's scores, i.e., set  $K = \sqrt{F^{-1}}$  where  $F^{-1}$  stands for the quantile function of the  $\chi_p^2$  distribution (with  $p$ -degrees of freedom).

Then the null hypothesis  $H_0 : \theta = \mathbf{0}$  can be checked by means of the  $\chi_p^2$  signed-rank test  $T_A$  and sign test  $T_B$  that are based on the statistics  $S_A$  (2) and  $S_B$  (1), respectively. Their asymptotic null distributions are assumed  $\chi_p^2$  in both cases (for  $p$  fixed and  $n \rightarrow \infty$ ), which is justified in the previous section and proved in the Appendix.

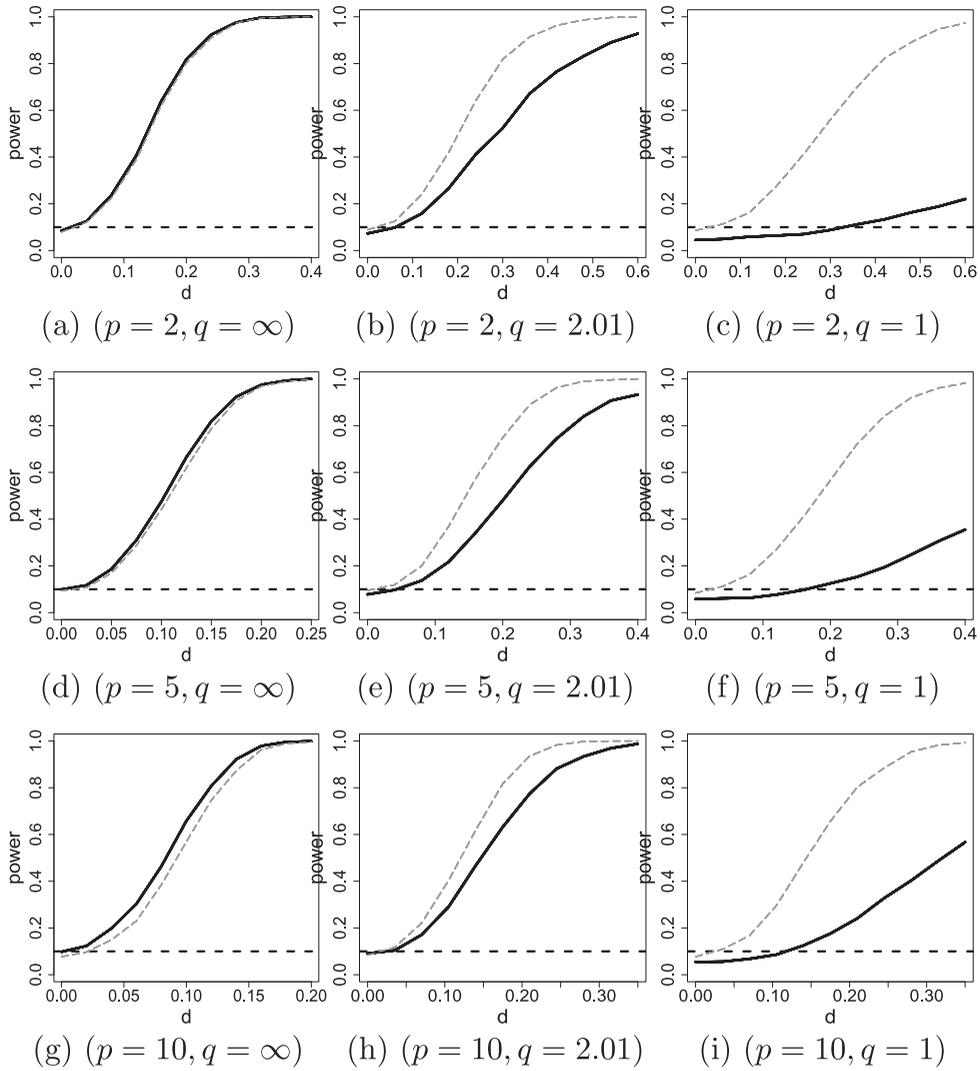
The tests  $T_A$  and  $T_B$  are compared with the benchmark Hotelling's  $T^2$  test (as implemented in the ICSNP package (Nordhausen, Sirkia, Oja, and Tyler 2015) for R (R Development Core Team 2008)) in Figures 1 and 2 for a few elliptical distributions and parameters  $n$ ,  $p$ , and  $N_H$  in terms of empirical power based on 1000 independent replications. In particular, the simulation study uses multivariate canonical Student  $t$  null distributions with  $q = 1$ ,  $q = 2.01$ , and  $q = \infty$  degrees of freedom,  $n = 100$  or  $n = 1000$   $p$ -dimensional observations with fixed  $p$  such as  $p = 10$  or  $p = 100$  always less than  $n$ , and  $N_H = 5n$ . Every simulated data set was used for evaluating the tests under both  $H_0$  and all ten shift alternatives considered.



**Figure 1.** Power comparison of two one-sample tests at significance level  $\alpha = 0.10$ . The plots display empirical powers of the sign test  $T_B$  (dotted line) and the benchmark Hotelling's test (thick solid line) with common significance level  $\alpha = 0.10$  (horizontal thick dashed line). The sign test based on  $N_H = 5n$  hyperplanes was applied to 1000  $p$ -dimensional random samples of size  $n$  from the multivariate canonical Student  $t$  distribution with  $q$  degrees of freedom shifted by  $(d, \dots, d)'$ . The parametric combination  $(p, n, q)$  characterising each experiment is stated below individual pictures. Each column of subfigures corresponds to one null distribution.

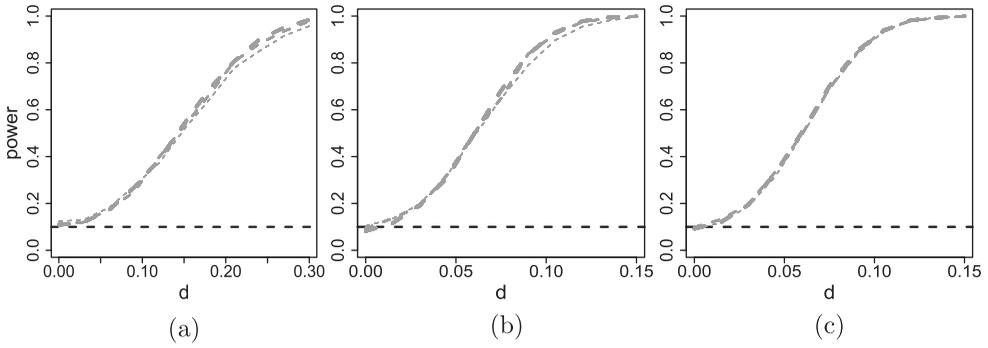
Apparently, even  $N_H$  mildly higher than  $n$  already leads to tests  $T_A$  and  $T_B$  with power comparable to the benchmark for normally distributed data and with far better performance for heavy-tailed distributions.

In fact, Figure 3 illustrates the sensitivity of test  $T_B$  to the choice of  $N_H$  and confirms that even  $N_H = 25$  may be satisfactory for common sample sizes ( $n = 100$  or  $n = 500$ ) to achieve a desirable test performance, at least in case of not too large  $p$  such as  $p = 5$  and the multivariate canonical Student  $t$  null distribution with only one degree of freedom.



**Figure 2.** Power comparison of two one-sample tests at significance level  $\alpha = 0.10$ . The plots display empirical powers of the signed-rank test  $T_A$  (dashed line) and the benchmark Hotelling's test (thick solid line) with common significance level  $\alpha = 0.10$  (horizontal thick dashed line). The signed-rank test based on  $2^p N_H$  hyperplanes,  $N_H = 5n$ , was applied to 1000  $p$ -dimensional random samples of size  $n = 100$  from the multivariate canonical Student  $t$  distribution with  $q$  degrees of freedom shifted by  $(d, \dots, d)'$ . The parametric combination  $(p, q)$  characterising each experiment is stated below individual pictures. Each column of subfigures corresponds to one null distribution.

The simulation studies do not include the classic tests (based on complete interdirections/symmetrised lift-interdirections), because their computation would be too time consuming or even impossible for most settings considered. Speed comparison is omitted for the same reason. Except for some rather special cases of low  $n$  and  $p$ , the difference would be not only in the speed of computation but also in its feasibility.



**Figure 3.** Test dependence on  $N_H$ . The plots display empirical powers of the sign test  $T_B$  (dashed line) with common significance level  $\alpha = 0.10$  (horizontal thick dashed line). The test based on  $N_H^1$  (thin),  $N_H^2$  (normal) or  $N_H^3$  (thick) hyperplanes was applied to 1000 five-dimensional ( $p = 5$ ) random samples of size (a)  $n = 100$  or (b,c)  $n = 500$  from the multivariate canonical Student  $t$  distribution with one degree of freedom shifted by  $(d, \dots, d)'$ . In particular,  $N_H^1 = cn/4$ ,  $N_H^2 = cn/2$  and  $N_H^3 = cn$  where  $c = 1/5$  in (b) and  $c = 1$  in (a) and (c), which leads to the same  $N_H^i$  in (a) and (b),  $i = 1, 2, 3$ . Apparently, there is hardly any visible difference in the test performance if the number of included hyperplanes is not too small.

To sum up the results, the incomplete interdirections and incomplete lift-interdirections make it possible to obtain powerful and computationally feasible tests of high-dimensional data.

### 5. Discussion

The classic tests based on interdirections or (symmetrised) lift-interdirections are computationally feasible only if the number  $n$  or the dimension  $p$  of the observations is small. The article uses incomplete  $U$ -statistics to define incomplete versions of those hyperplane-based concepts and shows their usefulness for nonparametric statistical inference in spaces with large but fixed dimensions.

In particular, the number of hyperplanes considered by the incomplete interdirections does not grow with the dimension of observations at all while it is equal to  $\binom{n}{p-1}$  in case of their complete counterparts. Consequently, even dimensions like  $p = 100$  can be handled easily.

As for the incomplete symmetrised lift-interdirections, the number of processed hyperplanes grows with the dimension by factor  $2^p$ , which still may make the computation challenging for  $p > 10$  or so. But the improvement achieved by them is still considerable.

The classic tests should be preferable due to their non-stochastic nature given the observations. When they are not feasible to compute, then the incomplete variants provide a way how to proceed further. The results for hyperplanes sampled at random with or without replacement are then virtually the same because the chance of choosing a hyperplane at least twice in the former case is negligible. Therefore, the first option is recommended due to its simplicity.

The brief simulation study confirms that the exact number  $N_H$  of hyperplanes considered by the incomplete variants is not too important and that even  $N_H$  as low as  $N_H = 5n$  can still produce reliable results in common settings with  $n \leq 1000$  and  $p \leq 100$  as far as

$p \leq n/10$  (ideally  $p \leq n/20$ ). In fact, even  $N_H = 25$  was shown satisfactory for small  $p$  and not too large  $n$  in a very limited comparison. Nevertheless, further simulations are needed to confirm and extend the rules to other contexts and to explore the incomplete variants in small samples.

This article also shows that incomplete  $U$ -statistics and simple geometric concepts of multivariate signs and ranks may be very useful for nonparametric inference and play important role in the future.

## Disclosure statement

No potential conflict of interest was reported by the authors.

## Funding

The research of Miroslav Šiman and Jana Klicnarová was supported by the Czech Science Foundation (Grantová agentura České republiky) project GA17-07384S. The work of Šárka Hudecová was supported by the Czech Science Foundation project GA18-01781Y.

## References

- Blom, G. (1976), 'Some Properties of Incomplete  $U$ -Statistics', *Biometrika*, 63, 573–580.
- Chaudhuri, P., and Sengupta, D. (1993), 'Sign Tests in Multidimension: Inference Based on the Geometry of the Data Cloud', *Journal of the American Statistical Association*, 88, 1363–1370.
- Chernozhukov, V., Galichon, A., Hallin, M., and Henry, M. (2017), 'Monge-Kantorovich Depth, Quantiles, Ranks, and Signs', *The Annals of Statistics*, 45, 223–256.
- Došlá, Š. (2009), 'Conditions for Bimodality and Multimodality of a Mixture of Two Unimodal Densities', *Kybernetika*, 45, 279–292.
- Ghoush, S.K., and Sengupta, D. (2001), 'Testing for Proportionality of Multivariate Dispersion Structures Using Interdirections', *Journal of Nonparametric Statistics*, 13, 331–349.
- Gieser, P.W., and Randles, R.H. (1997), 'A Nonparametric Test of Independence Between Two Vectors', *Journal of the American Statistical Association*, 92, 561–567.
- Hallin, M., and Paindaveine, D. (2002a), 'Optimal Procedures Based on Interdirections and Pseudo-Mahalanobis Ranks for Testing Multivariate Elliptic White Noise Against ARMA Dependence', *Bernoulli*, 8, 787–815.
- Hallin, M., and Paindaveine, D. (2002b), 'Optimal Tests for Multivariate Location Based on Interdirections and Pseudo-Mahalanobis Ranks', *The Annals of Statistics*, 30, 1103–1133.
- Hallin, M., and Paindaveine, D. (2004), 'Multivariate Signed-Rank Tests in Vector Autoregressive Order Identification', *Statistical Science*, 19, 697–711.
- Hettmansperger, T.P., Möttönen, J., and Oja, H. (1999), 'The Geometry of the Affine Invariant Multivariate Sign and Rank Methods', *Journal of Nonparametric Statistics*, 11, 271–285.
- Jan, S.L., and Randles, R.H. (1996), 'Interdirection Tests for Simple Repeated-Measures Designs', *Journal of the American Statistical Association*, 91, 1611–1618.
- Janson, S. (1984), 'The Asymptotic Distributions of Incomplete  $U$ -Statistics', *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 66, 495–505.
- Lee, A.J. (1990), *U-Statistics: Theory and Practice*, New York: CRC Press.
- Möttönen, J., and Oja, H. (1995), 'Multivariate Spatial Sign and Rank Methods', *Journal of Nonparametric Statistics*, 5, 201–213.
- Nordhausen, K., Sirkia, S., Oja, H., and Tyler, D.E. (2015), *ICSNP: Tools for Multivariate Nonparametrics*, R package version 1.1-0, Available at <https://CRAN.R-project.org/package=ICSNP>.
- Oja, H. (1999), 'Affine Invariant Multivariate Sign and Rank Tests and Corresponding Estimates: A Review', *Scandinavian Journal of Statistics*, 26, 319–343.
- Oja, H. (2010), *Multivariate Nonparametric Methods With R: An Approach Based on Spatial Signs and Ranks*, New York: Springer. ISBN 9781441904676.

Oja, H., and Paindaveine, D. (2005), ‘Optimal Signed-Rank Tests Based on Hyperplanes’, *Journal of Statistical Planning and Inference*, 135, 300–323.

Paindaveine, D. (2009), ‘On Multivariate Runs Tests for Randomness’, *Journal of the American Statistical Association*, 104, 1525–1538.

Puri, M.L., and Sen, P.K. (1971), *Nonparametric Methods in Multivariate Analysis*, New York: Wiley. ISBN 0471702404.

R Development Core Team (2008), *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria. Available at <http://www.R-project.org>, ISBN 3900051070.

Randles, R.H. (1989), ‘A Distribution-Free Multivariate Sign Test Based on Interdirections’, *Journal of the American Statistical Association*, 84, 1045–1050.

Randles, R.H., and Peters, D. (1990), ‘Multivariate Rank Tests for the Two-Sample Location Problem’, *Communications in Statistics – Theory and Methods*, 19, 4225–4238.

Serfling, R.J. (2006), ‘Multivariate Symmetry and Asymmetry’, in *Encyclopedia of Statistical Sciences* (Vol. 8, 2nd ed.), eds. S. Kotz, N. Balakrishnan, C.B. Read, and B. Vidakovic, pp. 5338–5345. New York: Wiley.

Taskinen, S., Oja, H., and Randles, R.H. (2005), ‘Multivariate Nonparametric Tests of Independence’, *Journal of the American Statistical Association*, 100, 916–925.

Um, Y., and Randles, R.H. (1998), ‘Nonparametric Tests for the Multivariate Multi-Sample Location Problem’, *Statistica Sinica*, 8, 801–812.

Zuo, Y., and Serfling, R. (2000), ‘General Notions of Statistical Depth Functions’, *The Annals of Statistics*, 28, 461–482.

## Appendix. Proofs

All the expectations involved in the proofs are taken with respect to the probability distribution of  $\mathbf{X}$  and with respect to the random designs involved due to the use of incomplete interdirections and symmetrised lift-interdirections. For the sake of simplicity, this is stressed out with a subscript only where confusion would be possible.

**Proof of Proposition 3.3:** Thanks to the invariance properties of incomplete interdirections, it is sufficient to consider the situation when  $\mathbf{X}_i = \mathbf{U}_i$ ,  $i = 1, \dots, n$ , where  $\mathbf{U}_1, \dots, \mathbf{U}_n$  are independent variables uniformly distributed on the unit sphere in  $\mathbb{R}^p$ . The same arguments as in the proof of Theorem A.1. in Randles (1989) imply

$$E \left[ S_B - \frac{p}{n} \sum_{i=1}^n \sum_{j=1}^n \mathbf{U}_i^\top \mathbf{U}_j \right]^2 \leq 2\pi^2 p^2 E \left[ \frac{\tilde{C}_{\mathbf{U}_1, \mathbf{U}_2}(\mathcal{X}_n)}{m_{n,p-1}} - \frac{1}{\pi} \arccos(\mathbf{U}_1^\top \mathbf{U}_2) \right]^2.$$

The right unconditional expectation converges to zero for  $n \rightarrow \infty$  because all the conditional expectations given  $\mathbf{U}_1, \mathbf{U}_2$  and  $\mathcal{Q}_S(n, p - 1)$  do so owing to Theorem 3.1 and Lemma 1 in Hallin and Paindaveine (2002b). Consequently,  $S_B$  must have the same asymptotic distribution as  $(p/n) \sum_{i=1}^n \sum_{j=1}^n \mathbf{U}_i^\top \mathbf{U}_j$ , which is known to be  $\chi_p^2$ . ■

**Proof of Proposition 3.4.:** Consider only  $\tilde{R}_i$  (as the proof for  $\tilde{R}_j$  is analogous), arbitrarily fix  $\mathbf{x} \in \mathbb{R}^p$ , and assume  $\Sigma = \mathbf{I}_p$  without any loss of generality. If  $\mathbf{X}_i = \mathbf{x}$ , then  $\tilde{R}_i \equiv \tilde{R}_i(\mathbf{x}) = \sum_{j=1}^n I[\tilde{L}_x(\mathcal{X}_n) \geq \tilde{L}_{\mathbf{X}_j}(\mathcal{X}_n)]$  and  $R_i \equiv R_i(\mathbf{x}) = \sum_{j=1}^n I[\|\mathbf{x}\| \geq \|\mathbf{X}_j\|]$ . It suffices to show that  $n^{-2} E[\tilde{R}_i(\mathbf{x}) - R_i(\mathbf{x})]^2 \rightarrow 0$  as  $n \rightarrow \infty$ , which follows, as in the proof of Proposition 2 in Oja and Paindaveine (2005), from

$$\int_{\mathbb{R}^p} g_j(\mathbf{x}, \mathbf{z}) f(\mathbf{z}) \, d\mathbf{z} = o(1), \quad j = 1, 2, \tag{A1}$$

where  $g_1(\mathbf{x}, \mathbf{z}) = E[I[\tilde{L}_x(\mathcal{X}_n) \geq \tilde{L}_z(\mathcal{X}_n), \|\mathbf{x}\| < \|\mathbf{z}\|]$  and  $g_2(\mathbf{x}, \mathbf{z}) = E[I[\tilde{L}_x(\mathcal{X}_n) < \tilde{L}_z(\mathcal{X}_n), \|\mathbf{x}\| \geq \|\mathbf{z}\|]$ . This is shown in the following.

If  $\|\mathbf{x}\| \geq \|\mathbf{z}\|$ , then  $g_1(\mathbf{x}, \mathbf{z}) = 0$ . Furthermore, Theorem 3.1 and Oja and Paindaveine (2005) imply that  $m_{n,p}^{-1}\tilde{L}_x(\mathcal{X}_n)$  converges in quadratic mean to the theoretical lift-interdirection  $l(\mathbf{x}) = Em_{n,p}^{-1}\tilde{L}_x(\mathcal{X}_n)$ , which is a strictly increasing function of  $\|\mathbf{x}\|$ . If  $\|\mathbf{x}\| < \|\mathbf{z}\|$ , then the Chebyshev inequality consequently results in

$$\begin{aligned} g_1(\mathbf{x}, \mathbf{z}) &= \mathbf{P}(m_{n,p}^{-1}[\tilde{L}_x(\mathcal{X}_n) - \tilde{L}_z(\mathcal{X}_n)] \geq 0) \\ &\leq \mathbf{P}\left(\left|m_{n,p}^{-1}[\tilde{L}_x(\mathcal{X}_n) - \tilde{L}_z(\mathcal{X}_n)] - (l(\mathbf{x}) - l(\mathbf{z}))\right| \geq l(\mathbf{z}) - l(\mathbf{x})\right) \\ &\leq \frac{\text{var}[m_{n,p}^{-1}\tilde{L}_x(\mathcal{X}_n) - m_{n,p}^{-1}\tilde{L}_z(\mathcal{X}_n)]}{(l(\mathbf{z}) - l(\mathbf{x}))^2} = O(1/n) = o(1) \end{aligned}$$

as  $n \rightarrow \infty$  for all  $\mathbf{x}, \mathbf{z} \in \mathbb{R}^p$ , which implies the result because  $g_1$  is bounded. The proof for  $g_2$  is analogous.  $\blacksquare$

**Proof of Proposition 3.5:** The proof mimics that of Proposition 3 in Hallin and Paindaveine (2002b), but some steps require different justification due to the use of incomplete versions of both interdirections and symmetrised lift-interdirections. These arguments are gathered in the concluding Lemma A.1.

Consider distances  $d_i = \|\Sigma^{-1/2}\mathbf{X}_i\|$  and independent random vectors  $\mathbf{U}_i = d_i^{-1}\Sigma^{-1/2}\mathbf{X}_i$  uniformly distributed on the unit sphere in  $\mathbb{R}^p$ ,  $i = 1, \dots, n$ . Assume  $\Sigma = \mathbf{I}_p$  without any loss of generality and define:

$$\begin{aligned} T_1^n &= \frac{1}{n} \sum_{i,j=1}^n K\left(\frac{\tilde{R}_i}{n+1}\right) K\left(\frac{\tilde{R}_j}{n+1}\right) \left[\cos(\tilde{a}_{\mathbf{X}_i, \mathbf{X}_j}^{m_{n,p}-1}) - \mathbf{U}_i^\top \mathbf{U}_j\right], \\ T_2^n &= \frac{1}{n} \sum_{i,j=1}^n \left[ K\left(\frac{\tilde{R}_i}{n+1}\right) K\left(\frac{\tilde{R}_j}{n+1}\right) - K(F_r(d_i))K(F_r(d_j)) \right] \mathbf{U}_i^\top \mathbf{U}_j, \\ \mathbf{T}_n &= \frac{1}{\sqrt{n}} \sum_{i=1}^n K(F_r(d_i))\mathbf{U}_i, \\ \mathbf{S}_n &= \frac{1}{\sqrt{n}} \sum_{i=1}^n K(R_i/(n+1))\mathbf{U}_i, \quad \text{and} \quad \hat{\mathbf{S}}_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n K(\tilde{R}_i/(n+1))\mathbf{U}_i \end{aligned}$$

where  $R_i$  is the rank of  $d_i$ ,  $i = 1, \dots, n$ .

To sum up, the proof approximates  $S_A$  with

$$S_A^0 = \frac{p}{nE K^2(V)} \sum_{i,j=1}^n K(F_r(d_i))K(F_r(d_j))\mathbf{U}_i^\top \mathbf{U}_j$$

whose limit distribution follows easily from the central limit theorem for  $\mathbf{T}_n$ . The difference  $S_A - S_A^0 \equiv (p/E[K^2(V)])(T_1^n + T_2^n)$  must be  $o_p(1)$  because both  $T_1^n \rightarrow 0$  and  $T_2^n \rightarrow 0$  in probability as  $n \rightarrow \infty$ , which is proved below.

As for  $T_2^n$ , Hallin and Paindaveine (2002b) proved that  $E\|\mathbf{T}_n - \mathbf{S}_n\|^2 = o(1)$  and  $E[K(R_i/(n+1)) - K(F_r(d_i))]^2 = o(1)$  as  $n \rightarrow \infty$ . Lemma A.1(1) further implies that

$$E\|\mathbf{S}_n - \hat{\mathbf{S}}_n\|^2 = \frac{1}{n} \sum_{i=1}^n E\left[ K\left(\frac{R_i}{n+1}\right) - K\left(\frac{\tilde{R}_i}{n+1}\right) \right]^2 \rightarrow 0$$

because the squared differences in the summands are both uniformly integrable thanks to Assumption A and  $o(1)$  in probability thanks to Proposition 3.4 and the continuity of  $K$  almost everywhere. They are thus  $o(1)$  in quadratic mean. Hence,  $E\|\mathbf{T}_n - \hat{\mathbf{S}}_n\|^2 \rightarrow 0$ , too. Therefore,  $E\|\mathbf{T}_n\|^2 < \infty$

implies  $E\|\widehat{\mathbf{S}}_n\|^2 < \infty$ . The Cauchy-Schwartz inequality then concludes that

$$E|T_2^n| = E|\widehat{\mathbf{S}}_n^\top \widehat{\mathbf{S}}_n - \mathbf{T}_n^\top \mathbf{T}_n| \leq \sqrt{E\|\widehat{\mathbf{S}}_n + \mathbf{T}_n\|^2} \sqrt{E\|\widehat{\mathbf{S}}_n - \mathbf{T}_n\|^2} \rightarrow 0$$

for  $n \rightarrow \infty$ , which proves that  $T_2^n \rightarrow 0$  both in  $L^1$  and in probability.

Now turn to  $T_1^n$  and define

$$D_i := K\left(\frac{\widetilde{R}_i}{n+1}\right) \quad \text{and} \quad C_{ij} := \cos(\widetilde{a}_{\mathbf{X}_i, \mathbf{X}_j}^{m_{n,p-1}}) - \mathbf{U}_i^\top \mathbf{U}_j.$$

Then

$$\begin{aligned} E\|T_1^n\|^2 &= \frac{1}{n^2} E\left(\sum_{i,j=1}^n D_i D_j C_{ij}\right)^2 = \frac{2}{n^2} \sum_{i \neq j} E[D_i D_j C_{ij}]^2 = \frac{2(n-1)}{n} E[D_1 D_2 C_{12}]^2 \\ &\leq \frac{2(n-1)}{n} (E|D_1 D_2|^{2+\delta})^{2/(2+\delta)} (E|C_{12}|^{2(2+\delta)/\delta})^{\delta/(2+\delta)} \end{aligned}$$

for the particular  $\delta > 0$  from Assumption A thanks to  $C_{ii} = 0$  almost surely,  $i = 1, \dots, n$ , Lemma A.1(2,3) and the Hölder inequality.

The conditional expectation of bounded  $C_{12}$  given the design,  $\mathbf{X}_1$  and  $\mathbf{X}_2$  always converges to 0 in quadratic mean for the same reasons as in Proposition 3.3. Therefore,  $E|C_{12}|^2 \rightarrow 0$  and also  $E|C_{12}|^{2(2+\delta)/\delta} \rightarrow 0$  as  $n \rightarrow \infty$ . Furthermore,  $E|D_1 D_2|^{2+\delta} < \infty$  due to Lemma A.1(2). Consequently,  $T_1^n \rightarrow 0$  in  $L^2$  (and in probability) as  $n \rightarrow \infty$ .  $\blacksquare$

**Lemma A.1:** (1) If  $i \neq j$ , then

$$E\left[K\left(\frac{R_i}{n+1}\right) - K\left(\frac{\widetilde{R}_i}{n+1}\right)\right] \left[K\left(\frac{R_j}{n+1}\right) - K\left(\frac{\widetilde{R}_j}{n+1}\right)\right] \mathbf{U}_i^\top \mathbf{U}_j = 0.$$

(2)  $E|D_1 D_2|^{2+\delta} < \infty$  and  $ED_i D_j C_{ij} = ED_1 D_2 C_{12}$  for  $i \neq j$ .

(3) If  $(i, j) \neq (k, l)$  and  $(i, j) \neq (l, k)$ , then

$$ED_i D_j C_{ij} D_k D_l C_{kl} = 0.$$

**Proof:** Recall that  $\widetilde{R}_i$  and  $\widetilde{C}_{\mathbf{X}_i, \mathbf{X}_j}(\mathcal{X}_n)$  respectively depend on independent and equally probable design sets  $Q_S^s(n, p)$  and  $Q_S(n, p-1)$  that are also independent of  $\mathcal{X}$ .

(1) Let  $i \neq j$ . The statement holds when

$$EK(R_i/(n+1))K(R_j/(n+1))\mathbf{U}_i^\top \mathbf{U}_j = 0,$$

$$EK(\widetilde{R}_i/(n+1))K(\widetilde{R}_j/(n+1))\mathbf{U}_i^\top \mathbf{U}_j = 0, \quad \text{and}$$

$$EK(\widetilde{R}_i/(n+1))K(R_j/(n+1))\mathbf{U}_i^\top \mathbf{U}_j = 0.$$

Set  $B_i := \text{sgn}(X_{i1})$ ,  $\mathbf{Y}_i := B_i \mathbf{X}_i$ , and note that  $B_i$  are (for spherically distributed  $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})'$ ) independent of  $\mathbf{Y}_i$ , mutually independent and identically distributed with  $\mathbf{P}(B_i = 1) = \mathbf{P}(B_i = -1) = 1/2$ ,  $i = 1, \dots, n$ . Recall that  $\mathbf{X}_i = d_i \mathbf{U}_i$  and  $\mathbf{U}_i = \mathbf{X}_i/d_i = \mathbf{Y}_i B_i/d_i$  where  $d_i = \|\mathbf{X}_i\| = \|\mathbf{Y}_i\|$ . Hence,  $d_i$  is a function of  $\mathbf{Y}_i$  and the ranks  $R_i$ 's are functions of the sample  $\mathcal{Y}$  consisting of  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ . Consequently,

$$\begin{aligned} &EK(R_i/(n+1))K(R_j/(n+1))\mathbf{U}_i^\top \mathbf{U}_j \\ &= E_{\mathcal{Y}} E_B [K(R_i/(n+1))K(R_j/(n+1))\mathbf{U}_i^\top \mathbf{U}_j | \mathbf{Y}_1, \dots, \mathbf{Y}_n] = \\ &= E_{\mathcal{Y}} \frac{\mathbf{Y}_i^\top \mathbf{Y}_j}{d_i d_j} K(R_i/(n+1))K(R_j/(n+1)) E_B [B_i B_j | \mathbf{Y}_1, \dots, \mathbf{Y}_n] = 0. \end{aligned}$$

Any design  $\mathcal{Q}_S^s(n, p)$  independent of  $\mathcal{X}$  leads to the symmetrised lift-interdirections  $\tilde{\mathbf{L}}_{X_i}(\mathcal{X}_n)$ ,  $i = 1, \dots, n$ , that are the same for any  $s_1 \mathbf{X}_1, \dots, s_n \mathbf{X}_n$ ,  $(s_1, \dots, s_n) \in \{-1, 1\}^p$ . Therefore, they are only functions of  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ , and the same thus holds even for their ranks.

Consequently,

$$\begin{aligned} & EK(\tilde{\mathbf{R}}_i/(n+1))K(\tilde{\mathbf{R}}_j/(n+1))\mathbf{U}_i^\top \mathbf{U}_j \\ &= E_{\mathcal{Y}}E_{\mathcal{Q}_S^s}E_B[K(\tilde{\mathbf{R}}_i/(n+1))K(\tilde{\mathbf{R}}_j/(n+1))\mathbf{U}_i^\top \mathbf{U}_j | \mathbf{Y}_1, \dots, \mathbf{Y}_n, \mathcal{Q}_S^s] \\ &= E_{\mathcal{Y}}E_{\mathcal{Q}_S^s} \frac{\mathbf{Y}_i^\top \mathbf{Y}_j}{d_i d_j} K(\tilde{\mathbf{R}}_i/(n+1))K(\tilde{\mathbf{R}}_j/(n+1))E_B[B_i B_j | \mathbf{Y}_1, \dots, \mathbf{Y}_n, \mathcal{Q}_S^s] = 0. \end{aligned}$$

The proof of the last equality is analogous.

(2) Obviously,

$$\begin{aligned} ED_i D_j C_{ij} &= \frac{1}{|\{\mathcal{Q}_S^s\}| |\{\mathcal{Q}_S\}|} \sum_{\mathcal{Q}_S^s} \sum_{\mathcal{Q}_S} E_{\mathcal{X}} K\left(\frac{\tilde{\mathbf{R}}_i}{n+1}\right) K\left(\frac{\tilde{\mathbf{R}}_j}{n+1}\right) \\ &\quad \times \left( \cos(\tilde{a}_{\mathbf{X}_i, \mathbf{X}_j}^{m_n, p-1}(\mathcal{X})) - \frac{\mathbf{X}_i^\top \mathbf{X}_j}{d_i d_j} \right). \end{aligned}$$

The expectation is the same for all pairs  $(i, j)$  of distinct indices because the  $\mathbf{X}_i$ 's are i.i.d. and the design sets are properly independent and equiprobable. The following expectation is independent of  $(i, j)$ ,  $i \neq j$ , for the same reason:

$$\begin{aligned} E|D_1 D_2|^{2+\delta} &= \frac{1}{|\{\mathcal{Q}_S^s\}|} \sum_{\mathcal{Q}_S^s} E_{\mathcal{X}} |K(\tilde{\mathbf{R}}_1/(n+1))K(\tilde{\mathbf{R}}_2/(n+1))|^{2+\delta} \\ &= \sum_{l \neq m} |K(l/(n+1))|^{2+\delta} |K(m/(n+1))|^{2+\delta} \underbrace{\sum_{\mathcal{Q}_S^s} \frac{1}{|\{\mathcal{Q}_S^s\}|} \mathbf{P}(\tilde{\mathbf{R}}_1 = l, \tilde{\mathbf{R}}_2 = m)}_{p(l, m)}. \end{aligned}$$

As  $p(l, m)$  does not depend on  $l$  and  $m$ , necessarily  $p = [n(n-1)]^{-1}$ . The finiteness of  $E|D_1 D_2|^{2+\delta}$  then follows from the arguments in the proof of Proposition 3 in Hallin and Paindaveine (2002b).

(3) Define  $C_{ij}^{\mathcal{Y}} := \cos(\pi \tilde{C}_{\mathbf{Y}_i, \mathbf{Y}_j}(\mathcal{Y}_n)/m_{n, p-1}) - \mathbf{Y}_i^\top \mathbf{Y}_j / (d_i d_j)$ , and realise that

$$\cos(\pi \tilde{C}_{\mathbf{X}_i, \mathbf{X}_j}(\mathcal{X}_n)/m_{n, p-1}) = B_i B_j \cos(\pi \tilde{C}_{\mathbf{Y}_i, \mathbf{Y}_j}(\mathcal{Y}_n)/m_{n, p-1}).$$

Then

$$\begin{aligned} & ED_i D_j C_{ij} D_k D_l C_{kl} \\ &= E_{\mathcal{Y}} E_{\mathcal{Q}_S^s} E_{\mathcal{Q}_S} C_{ij}^{\mathcal{Y}} C_{kl}^{\mathcal{Y}} E_B[D_i D_j D_k D_l B_i B_j B_k B_l | \mathbf{Y}_1, \dots, \mathbf{Y}_n, \mathcal{Q}_S^s, \mathcal{Q}_S] \\ &= E_{\mathcal{Y}} E_{\mathcal{Q}_S^s} E_{\mathcal{Q}_S} D_i D_j C_{ij}^{\mathcal{Y}} D_k D_l C_{kl}^{\mathcal{Y}} E_B[B_i B_j B_k B_l], \end{aligned}$$

which is indeed zero if  $(i, j)$  is different from both  $(k, l)$  and  $(l, k)$ . ■