

TESTING MULTIVARIATE UNIFORMITY AND ITS APPLICATIONS

JIA-JUAN LIANG, KAI-TAI FANG, FRED J. HICKERNELL, AND RUNZE LI

ABSTRACT. Some new statistics are proposed to test the uniformity of random samples in the multidimensional unit cube $[0, 1]^d$ ($d \geq 2$). These statistics are derived from number-theoretic or quasi-Monte Carlo methods for measuring the discrepancy of points in $[0, 1]^d$. Under the null hypothesis that the samples are independent and identically distributed with a uniform distribution in $[0, 1]^d$, we obtain some asymptotic properties of the new statistics. By Monte Carlo simulation, it is found that the finite-sample distributions of the new statistics are well approximated by the standard normal distribution, $N(0, 1)$, or the chi-squared distribution, $\chi^2(2)$. A power study is performed, and possible applications of the new statistics to testing general multivariate goodness-of-fit problems are discussed.

1. INTRODUCTION

Testing uniformity in the unit interval $[0, 1]$ has been studied by many authors. Some early work in this area is [Ney37], [Pea39] and [AD54]. Quesenberry and Miller [QM77, MQ79] made a thorough Monte Carlo simulation to compare a number of existing statistics for testing uniformity in $[0, 1]$ and recommended Watson's U^2 -test [Wat62] and Neyman's smooth test [Ney37] as general choices for testing uniformity in $[0, 1]$. D'Agostino and Stephens [DS86, Chapter 6] gave a comprehensive review on tests for uniformity in $[0, 1]$.

Testing uniformity of random samples in the multidimensional unit cube ($d \geq 2$),

$$\bar{C}^d = [0, 1]^d = \{\mathbf{x} = (x_1, \dots, x_d)' \in R^d : 0 \leq x_i \leq 1, i = 1, \dots, d\},$$

seems to have received less attention in the literature. Two well-known quantities are the Kolmogorov-Smirnov type statistic,

$$(1.1) \quad KS_n = \sup_{\mathbf{x} \in R^d} |F_n(\mathbf{x}) - F(\mathbf{x})|,$$

and the Cramér-von Mises type statistic,

$$(1.2) \quad CM_n = n \int_{R^d} [F_n(\mathbf{x}) - F(\mathbf{x})]^2 \psi(\mathbf{x}) d\mathbf{x}.$$

Here, $F(\mathbf{x})$ is the null distribution function (d.f.), $F_n(\mathbf{x})$ denotes the empirical distribution function (e.d.f.) based on n independent and identically distributed

Received by the editor August 14, 1998 and, in revised form, February 11, 1999.

2000 *Mathematics Subject Classification.* Primary 65C05, 62H10, 65D30.

Key words and phrases. Goodness-of-fit, discrepancy, quasi-Monte Carlo methods, testing uniformity.

This work was partially supported by a Hong Kong Research Grants Council grant RGC/97-98/47.

(i.i.d.) samples, and $\psi(\mathbf{x}) \geq 0$ in (1.2) is a suitable weight function. Unfortunately, the Kolmogorov-Smirnov type statistic is difficult to compute for large d .

In the literature of number-theoretic methods or quasi-Monte Carlo methods [Nie92, FW94, SJ94], there are a number of criteria for measuring whether a set of points is uniformly scattered in the unit cube \bar{C}^d . These criteria are called *discrepancies*, and they arise in the error analysis of quasi-Monte Carlo methods for evaluating multiple integrals. Given an integral over the unit cube,

$$(1.3) \quad I(f) = \int_{\bar{C}^d} f(\mathbf{x}) \, d\mathbf{x},$$

quasi-Monte Carlo methods approximate this integral by the sample mean,

$$(1.4) \quad Q(f) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{z}_i),$$

over a set of uniformly scattered sample points, $\mathcal{P} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \bar{C}^d$. Examples of good sets for quasi-Monte Carlo integration are given in [HW81, Nie92, SJ94, Tez95] and the references therein. The worst-case quadrature error of a quasi-Monte Carlo method is bounded by a generalized Koksma-Hlawka inequality,

$$(1.5) \quad |I(f) - Q(f)| \leq D(\mathcal{P})V(f),$$

where $V(f)$ is a measure of the variation or fluctuation of the integrand, and the discrepancy, $D(\mathcal{P})$, is a measure of the quality of the quadrature rule, or equivalently, of the set of points, \mathcal{P} . A smaller discrepancy implies a better set of points.

The precise definitions of the discrepancy and the variation depend on the space of integrands. In the original Koksma-Hlawka inequality [Nie92, Chap. 2], $V(f)$ is the variation of the integrand in the sense of Hardy and Krause, and the discrepancy is the star discrepancy, defined as follows:

$$(1.6) \quad D^*(\mathcal{P}) = \sup_{\mathbf{x} \in \bar{C}^d} \left| \frac{|\mathcal{P} \cap [\mathbf{0}, \mathbf{x}]|}{n} - \text{Vol}([\mathbf{0}, \mathbf{x}]) \right|,$$

where $|\cdot|$ denotes the number of points in a set, and $\text{Vol}([\mathbf{0}, \mathbf{x}])$ denotes the volume of the hypercube $[\mathbf{0}, \mathbf{x}]$ ($\mathbf{x} \in R^d$). Taking the null distribution in (1.1) to be the uniform distribution, $F(\mathbf{x}) = \text{Vol}([\mathbf{0}, \mathbf{x}])$, and noting that the e.d.f. is $F_n(\mathbf{x}) = |\mathcal{P} \cap [\mathbf{0}, \mathbf{x}]|/n$, the star discrepancy, (1.6), is a special case of the Kolmogorov-Smirnov type statistic, (1.1).

This relationship between discrepancies arising in quasi-Monte Carlo quadrature error bounds and goodness-of-fit statistics is rather general and has been discussed in [Hic98b, Hic99]. If the space of integrands is a reproducing kernel Hilbert space, then one may obtain a computationally simple form for the discrepancy. A special case is the \mathcal{L}^2 -star discrepancy, which corresponds to the Cramér-von Mises statistic, (1.2), for the uniform distribution with weight function $\psi(\mathbf{x}) = 1$.

Hickernell [Hic98a] proposed a generalized discrepancy based on reproducing kernels. This discrepancy has a computationally simple formula. Three special cases, the *symmetric discrepancy*, the *centered discrepancy* and the *star discrepancy*, have interesting geometrical interpretations. Thus, all three of them may give useful statistics for testing multivariate uniformity of a set of points.

However, to perform a statistical test, one must know the probability distribution of the test statistic under the assumption of i.i.d. uniform random points. This problem is the focus of this article. Section 2 defines the new test statistics and

describes their asymptotic behavior. For statistical reasons the discrepancy itself is not the best goodness-of-fit statistic, but useful statistics are derived from the discrepancy. Section 3 presents simulation results on how well the distributions of the statistics are approximated by their asymptotic limits and on the power performance of the new statistics. Some applications are also discussed.

2. THE NEW STATISTICS AND THEIR ASYMPTOTIC PROPERTIES

2.1. **Some \mathcal{L}^2 -type discrepancies.** The generalized \mathcal{L}^2 -type discrepancy proposed in [Hic98a] is as follows:

$$(2.1) \quad [D(\mathcal{P})]^2 = M^d - \frac{2}{n} \sum_{k=1}^n \prod_{j=1}^d [M + \beta^2 \mu(z_{kj})] + \frac{1}{n^2} \sum_{k,l=1}^n \prod_{j=1}^d \left(M + \beta^2 \left[\mu(z_{kj}) + \mu(z_{lj}) + \frac{1}{2} B_2(\{z_{kj} - z_{lj}\}) + B_1(z_{kj}) B_1(z_{lj}) \right] \right),$$

where $\{ \}$ denotes the fractional part of a real number or vector, β is an arbitrarily given positive constant, and $\mu(\cdot)$ is an arbitrary function satisfying

$$\mu \in \left\{ f : \frac{df}{dx} \in \mathcal{L}^\infty([0, 1]) \text{ and } \int_0^1 f(x) dx = 0 \right\}.$$

The constant M is determined in terms of β and μ as follows:

$$(2.2) \quad M = 1 + \beta^2 \int_0^1 \left(\frac{d\mu}{dx} \right)^2 dx.$$

The $B_1(\cdot)$ and $B_2(\cdot)$ in (2.1) are the first and the second degree Bernoulli polynomials, respectively:

$$B_1(x) = x - \frac{1}{2} \quad \text{and} \quad B_2(x) = x^2 - x + \frac{1}{6}.$$

For any $z_1, z_2 \in [0, 1]$, it is true that

$$B_2(\{z_1 - z_2\}) = B_2(|z_1 - z_2|) = |z_1 - z_2|^2 - |z_1 - z_2| + \frac{1}{6}.$$

The three special cases of $[D(\mathcal{P})]^2$ (denoted by $D_s(\mathcal{P})^2$, $D_c(\mathcal{P})^2$ and $D_*(\mathcal{P})^2$, respectively) given in [Hic98a] are derived by taking three different choices of the function $\mu(\cdot)$ and the constant β in (2.2):

1) the symmetric discrepancy:

$$\mu(x) = -\frac{1}{2} B_2(x) = -\frac{1}{2} (x^2 - x + \frac{1}{6}), \quad \beta^{-1} = \frac{1}{2}, \quad M = \frac{4}{3},$$

$$(2.3) \quad D_s(\mathcal{P})^2 = \left(\frac{4}{3} \right)^d - \frac{2}{n} \sum_{k=1}^n \prod_{j=1}^d (1 + 2z_{kj} - 2z_{kj}^2) + \frac{2^d}{n^2} \sum_{k,l=1}^n \prod_{j=1}^d (1 - |z_{kj} - z_{lj}|);$$

2) the centered discrepancy:

$$\mu(x) = -\frac{1}{2}B_2\left(\left\{x - \frac{1}{2}\right\}\right) = -\frac{1}{2}\left(\left|x - \frac{1}{2}\right|^2 - \left|x - \frac{1}{2}\right| + \frac{1}{6}\right), \quad \beta^{-1} = 1, \quad M = \frac{13}{12},$$

$$(2.4) \quad D_c(\mathcal{P})^2 = \left(\frac{13}{12}\right)^d - \frac{2}{n} \sum_{k=1}^n \prod_{j=1}^d \left(1 + \frac{1}{2} \left|z_{kj} - \frac{1}{2}\right| - \frac{1}{2} \left|z_{kj} - \frac{1}{2}\right|^2\right) + \frac{1}{n^2} \sum_{k,l=1}^n \prod_{j=1}^d \left[1 + \frac{1}{2} \left|z_{kj} - \frac{1}{2}\right| + \frac{1}{2} \left|z_{lj} - \frac{1}{2}\right| - \frac{1}{2} |z_{kj} - z_{lj}|\right];$$

3) the star discrepancy:

$$\mu(x) = \frac{1}{6} - \frac{x^2}{2}, \quad \beta^{-1} = 1, \quad M = \frac{4}{3},$$

$$(2.5) \quad D_*(\mathcal{P})^2 = \left(\frac{4}{3}\right)^d - \frac{2}{n} \sum_{k=1}^n \prod_{j=1}^d \left(\frac{3 - z_{kj}^2}{2}\right) + \frac{1}{n^2} \sum_{k,l=1}^n \prod_{j=1}^d [2 - \max(z_{kj}, z_{lj})].$$

2.2. Asymptotic properties of the discrepancies. The null hypothesis for testing the uniformity of random samples $\mathcal{P} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \bar{C}^d$ can be stated as

$$(2.6) \quad H_0 : \mathbf{z}_1, \dots, \mathbf{z}_n \text{ are uniformly distributed in } \bar{C}^d.$$

The alternative hypothesis H_1 implies rejection of H_0 in (2.6). A test for (2.6), that is, a *test of multivariate uniformity*, can be performed by determining whether the value of a test statistic is unlikely under the null hypothesis. If one wishes to use one of the discrepancies described above as a test statistic, then its probability distribution under the null hypothesis must be calculated. Although, this distribution is too complicated to describe for finite sample size, it can be characterized rather simply in the limit of infinite sample size.

The main results on the asymptotic properties of the discrepancies are contained in Theorems 2.1 and 2.3 below. Their proofs rely on the theory of U -type statistics [Ser80, Chapter 5].

Theorem 2.1. *Under the null hypothesis (2.6), the statistic $[D(\mathcal{P})]^2$ given by (2.1) has the asymptotic property*

$$(2.7) \quad [D(\mathcal{P})]^2 \xrightarrow{a.s.} 0 \quad (n \rightarrow \infty),$$

for an arbitrary function $\mu(\cdot)$ and an arbitrary constant β , where “ $\xrightarrow{a.s.}$ ” means “converges almost surely”.

Proof. For $\mathcal{P} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$ with $\mathbf{z}_k = (z_{k1}, \dots, z_{kd})'$ ($k = 1, \dots, n$), under the null hypothesis (2.6), the random variables z_{kj} ($k = 1, \dots, n, j = 1, \dots, d$) have a uniform distribution $U[0, 1]$. Let

$$(2.8) \quad g_1(\mathbf{z}_k) = \prod_{j=1}^d [M + \beta^2 \mu(z_{kj})],$$

and

$$(2.9) \quad h(\mathbf{z}_k, \mathbf{z}_l) = \prod_{j=1}^d \left(M + \beta^2 \left[\mu(z_{kj}) + \mu(z_{lj}) + \frac{1}{2} B_2(|z_{kj} - z_{lj}|) + B_1(z_{kj}) B_1(z_{lj}) \right] \right)$$

for $k, l = 1, \dots, n$. Then the square discrepancy given by (2.1) can be written as

$$(2.10) \quad \begin{aligned} [D(\mathcal{P})]^2 &= M^d - \frac{2}{n} \sum_{k=1}^n g_1(\mathbf{z}_k) + \frac{1}{n^2} \sum_{k,l=1}^n h(\mathbf{z}_k, \mathbf{z}_l) \\ &= M^d - \frac{2}{n} \sum_{k=1}^n g_1(\mathbf{z}_k) + \frac{1}{n^2} \left[2 \sum_{k<l}^n h(\mathbf{z}_k, \mathbf{z}_l) + \sum_{k=1}^n h(\mathbf{z}_k, \mathbf{z}_k) \right] \\ &= M^d - \frac{2}{n} \sum_{k=1}^n g_1(\mathbf{z}_k) + \frac{n-1}{n} \cdot \frac{2}{n(n-1)} \sum_{k<l}^n h(\mathbf{z}_k, \mathbf{z}_l) + \frac{1}{n^2} \sum_{k=1}^n g_2(\mathbf{z}_k), \end{aligned}$$

where

$$(2.11) \quad g_2(\mathbf{z}_k) = \prod_{j=1}^d \left(M + \beta^2 \left[\frac{1}{12} + 2\mu(z_{kj}) + B_1(z_{kj})^2 \right] \right).$$

Note that both $\{g_1(\mathbf{z}_k)\}_{k=1}^n$ and $\{g_2(\mathbf{z}_k)\}_{k=1}^n$ are sequences of i.i.d. random variables. By the strong law of large numbers,

$$(2.12) \quad U_1 = \frac{1}{n} \sum_{k=1}^n g_1(\mathbf{z}_k) \xrightarrow{\text{a.s.}} E[g_1(\mathbf{z}_1)] = M^d,$$

and $\frac{1}{n} \sum_{k=1}^n g_2(\mathbf{z}_k) \xrightarrow{\text{a.s.}} E[g_2(\mathbf{z}_1)] < \infty$. It follows that

$$(2.13) \quad \frac{1}{n^2} \sum_{k=1}^n g_2(\mathbf{z}_k) \xrightarrow{\text{a.s.}} 0.$$

By the theory of the U -type statistics [Ser80, Chapter 5],

$$(2.14) \quad U_2 = \frac{2}{n(n-1)} \sum_{k<l}^n h(\mathbf{z}_k, \mathbf{z}_l)$$

is a second-order U -statistic. By the strong law of large numbers for general U -statistics,

$$U_2 \xrightarrow{\text{a.s.}} E[h(\mathbf{z}_1, \mathbf{z}_2)] = M^d \quad (n \rightarrow \infty).$$

Therefore, by (2.10), as $n \rightarrow \infty$, $[D(\mathcal{P})]^2 \xrightarrow{\text{a.s.}} M^d - 2M^d + M^d = 0$. □

Corollary 2.2. *Under the null hypothesis (2.6), it is true that*

$$D_s(\mathcal{P})^2 \xrightarrow{\text{a.s.}} 0, \quad D_c(\mathcal{P})^2 \xrightarrow{\text{a.s.}} 0, \quad D_*(\mathcal{P})^2 \xrightarrow{\text{a.s.}} 0,$$

as $n \rightarrow \infty$, where $D_s(\mathcal{P})^2$, $D_c(\mathcal{P})^2$ and $D_*(\mathcal{P})^2$ are given by (2.3)-(2.5), respectively.

The following theorem and corollaries consider pieces of the discrepancy defined in (2.1). These pieces are then re-combined to give new statistics for testing multivariate normality.

Theorem 2.3. *Let U_1 and U_2 be given by (2.12) and (2.14), respectively. Then, under the null hypothesis (2.6),*

$$\sqrt{n} \begin{pmatrix} U_1 - M^d \\ U_2 - M^d \end{pmatrix} \xrightarrow{\mathcal{D}} N_2(\mathbf{0}, \Sigma) \quad (n \rightarrow \infty),$$

where “ $\xrightarrow{\mathcal{D}}$ ” means “converges in distribution”, and Σ is a singular covariance matrix:

$$(2.15) \quad \Sigma = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} \zeta_1,$$

where $\zeta_1 = (M^2 + \beta^4 c^2)^d - M^{2d}$ and $c^2 = \int_0^1 \mu(x)^2 dx$.

Proof. Since the random variable U_1 given by (2.12) is a first-order U -statistic and the random variable U_2 given by (2.14) is a second-order U -statistic, by the central limit theorem for U -statistics, we have

$$(2.16) \quad \sqrt{n}(U_1 - EU_1) / \sqrt{\text{var}(\sqrt{n}U_1)} \xrightarrow{\mathcal{D}} N(0, 1),$$

where $EU_1 = M^d$ by the proof of Theorem 2.1. It is easy to calculate

$$\text{var}(\sqrt{n}U_1) = (M^2 + \beta^4 c^2)^d - M^{2d}.$$

By Lemma A of [Ser80, p. 183], we obtain the variance of U_2 :

$$(2.17) \quad \text{var}(U_2) = \frac{4(n-2)}{n(n-1)} \zeta_1 + \frac{2}{n(n-1)} \zeta_2,$$

where ζ_1 is given by (2.15) and $\zeta_2 = [M^2 + \beta^4(2c^2 + 1/90)]^d - M^{2d}$. By the theorem in [Ser80, p. 189], U_2 can be written as

$$(2.18) \quad U_2(n) = \hat{U}_2(n) + R_n,$$

where $\hat{U}_2(n)$ is a random variable that can be written as a sum of i.i.d. random variables as follows:

$$(2.19) \quad \hat{U}_2(n) - EU_2(n) = \frac{2}{n} \sum_{i=1}^n h_1(\mathbf{z}_i),$$

for some function $h_1(\cdot)$ [Ser80, p. 188, equation (2)] and $EU_2(n) = M^d$. Formula (2.19) can be written as

$$(2.20) \quad \hat{U}_2(n) = \frac{1}{n} \sum_{i=1}^n h_2(\mathbf{z}_i),$$

where $h_2(\cdot) = 2h_1(\cdot) + M^d$. In (2.18), R_n is a residual term, $R_n = o_p(n^{-1}(\log n)^\delta)$ ($n \rightarrow \infty$), which implies that R_n tends to zeros in probability, where $\delta > 1/v$ ($v > 0$) if $E\{h(\mathbf{z}_1, \mathbf{z}_2)\}^v < \infty$. $h(\cdot, \cdot)$ is defined by (2.9). Under the null hypothesis (2.6), $E\{h(\mathbf{z}_1, \mathbf{z}_2)\}^v < \infty$ for any $v > 0$. Combining (2.18) and (2.20), we can write

$$U_2 = \frac{1}{n} \sum_{i=1}^n h_2(\mathbf{z}_i) + R_n.$$

Then

$$(2.21) \quad \begin{aligned} \sqrt{n} \begin{pmatrix} U_1 \\ U_2 \end{pmatrix} &= \sqrt{n} \begin{pmatrix} \frac{1}{n} \sum_{k=1}^n g_1(\mathbf{z}_k) \\ \frac{1}{n} \sum_{k=1}^n h_2(\mathbf{z}_k) + R_n \end{pmatrix} \\ &= \frac{1}{\sqrt{n}} \sum_{k=1}^n \begin{pmatrix} g_1(\mathbf{z}_k) \\ h_2(\mathbf{z}_k) \end{pmatrix} + \begin{pmatrix} 0 \\ \sqrt{n}R_n \end{pmatrix}. \end{aligned}$$

Under the null hypothesis (2.6), by the multivariate central limit theorem, we have

$$\sqrt{n} \begin{pmatrix} U_1 \\ U_2 \end{pmatrix} \xrightarrow{\mathcal{D}} \text{2-dimensional normal distribution}$$

because of the independence of the \mathbf{z}_k ($k = 1, \dots, n$) and the fact that

$$\begin{pmatrix} 0 \\ \sqrt{n}R_n \end{pmatrix} \xrightarrow{\mathcal{P}} \mathbf{0} \quad (n \rightarrow \infty),$$

where “ $\xrightarrow{\mathcal{P}}$ ” means “converges in probability”. It is easy to obtain the covariance between $\sqrt{n}U_1$ and $\sqrt{n}U_2$:

$$\text{cov}(\sqrt{n}U_1, \sqrt{n}U_2) = 2[(M^2 + \beta^4 c^2)^d - M^{2d}] = 2\zeta_1,$$

where ζ_1 is the same as in (2.15). Then we have

$$\sqrt{n} \begin{pmatrix} U_1 - M^d \\ U_2 - M^d \end{pmatrix} \xrightarrow{\mathcal{D}} N_2(\mathbf{0}, \Sigma),$$

with Σ given by (2.15). This completes the proof. □

Theorem 2.3 implies that the asymptotic distribution of $(\sqrt{n}U_1, \sqrt{n}U_2)$ is a singular (degenerate) normal distribution. This property results in the following corollary, which shows why the discrepancy itself is not necessarily a suitable goodness-of-fit statistic.

Corollary 2.4. *Under the null hypothesis (2.6), the generalized \mathcal{L}^2 -type discrepancy $[D(\mathcal{P})]^2$ given by (2.1) has a further asymptotic property:*

$$\sqrt{n}[D(\mathcal{P})]^2 \xrightarrow{\mathcal{P}} 0 \quad (n \rightarrow \infty).$$

Proof. By the proof of Theorem 2.1, $[D(\mathcal{P})]^2$ can be written as

$$[D(\mathcal{P})]^2 = M^d - 2U_1 + \frac{n-1}{n}U_2 + \frac{1}{n^2} \sum_{k=1}^n g_2(\mathbf{z}_k),$$

where $g_2(\cdot)$ is given by (2.11). By (2.13) and (2.14), we can write

$$\frac{1}{n^2} \sum_{k=1}^n g_2(\mathbf{z}_k) = O_p\left(\frac{1}{n}\right),$$

where the notation “ $f(n) = O_p(\frac{1}{n})$ ” means $nf(n) \xrightarrow{\mathcal{P}}$ a constant. Then we have

$$(2.22) \quad \sqrt{n}[D(\mathcal{P})]^2 = -2\sqrt{n}(U_1 - M^d) + \sqrt{n}(U_2 - M^d) - \frac{1}{\sqrt{n}}U_2 + O_p\left(\frac{1}{\sqrt{n}}\right).$$

Since $U_2 \xrightarrow{\text{a.s.}} M^d < \infty$ ($n \rightarrow \infty$) by Theorem 2.1, we have $(1/\sqrt{n})U_2 = O_p(1/\sqrt{n})$. Then we can write (2.22) as

$$(2.23) \quad \begin{aligned} \sqrt{n}[D(\mathcal{P})]^2 &= -2\sqrt{n}(U_1 - M^d) + \sqrt{n}(U_2 - M^d) + O_p\left(\frac{1}{\sqrt{n}}\right) \\ &= \sqrt{n}(-2, 1) \begin{pmatrix} U_1 - M^d \\ U_2 - M^d \end{pmatrix} + O_p\left(\frac{1}{\sqrt{n}}\right). \end{aligned}$$

By Theorem 2.3, we obtain

$$\sqrt{n}[D(\mathcal{P})]^2 \xrightarrow{\mathcal{D}} N(0, \mathbf{a}'\Sigma\mathbf{a}),$$

where $\mathbf{a}' = (-2, 1)$ and Σ is given by (2.15). It turns out that $\mathbf{a}'\Sigma\mathbf{a} = 0$. Therefore, $\sqrt{n}[D(\mathcal{P})]^2 \xrightarrow{\mathcal{D}} 0$. This implies $\sqrt{n}[D(\mathcal{P})]^2 \xrightarrow{\mathcal{P}} 0$. □

A better goodness-of-fit statistic than the discrepancy can be obtained by a linear combination of U_1 and U_2 that is not a degenerate normal distribution. This is the idea behind the statistic A_n defined below.

Corollary 2.5. *Under the null hypothesis (2.6), the statistic*

$$(2.24) \quad A_n = \sqrt{n}[(U_1 - M^d) + 2(U_2 - M^d)]/(5\sqrt{\zeta_1}) \xrightarrow{\mathcal{D}} N(0, 1) \quad (n \rightarrow \infty),$$

where U_1 and U_2 are defined by (2.12) and (2.14), respectively, and ζ_1 is given by (2.15).

Proof. It was noted in the proof of Corollary 2.4 that $\mathbf{a}' = (-2, 1)$ is an eigenvector associated with the zero-eigenvalue of Σ given by (2.15). On the other hand, $\mathbf{b} = (1, 2)'$ ($\mathbf{a}'\mathbf{b} = 0$) is an eigenvector associated with the eigenvalue $5\zeta_1$ of Σ . By Theorem 2.3, $E(A_n) = 0$ and the variance

$$\text{var}(A_n) = \mathbf{b}'\Sigma_n\mathbf{b}/(25\zeta_1) \rightarrow \mathbf{b}'\Sigma\mathbf{b}/(25\zeta_1) = 1 \quad (n \rightarrow \infty)$$

under the null hypothesis (2.6), where Σ_n is the covariance matrix of $(\sqrt{n}U_1, \sqrt{n}U_2)$, which turns out to be

$$(2.25) \quad \Sigma_n = \begin{pmatrix} \zeta_1 & 2\zeta_1 \\ 2\zeta_1 & \frac{4(n-2)}{n-1}\zeta_1 + \frac{2}{n-1}\zeta_2 \end{pmatrix},$$

and ζ_2 is the same as in (2.17). Assertion (2.24) holds as a result of Theorem 2.3. □

The statistic A_n can be employed for testing hypothesis (2.6). Larger values of $|A_n|$ imply rejection for the null hypothesis (2.6). It can be verified that Σ_n given by (2.25) tends to singularity very slowly. For example, under the symmetric discrepancy, when $n = 1000$, the condition number of $\Sigma_n = 1467$ ($d = 2$), 1196 ($d = 5$) and 841 ($d = 10$). Therefore, for finite sample size n (e.g., $n \leq 1000$), we recommend using the normal distribution $N_2(\mathbf{0}, \Sigma_n)$ as the approximate joint distribution of $(\sqrt{n}U_1, \sqrt{n}U_2)$. Based on this idea, we propose the following χ^2 -type statistic for testing the null hypothesis (2.6):

$$(2.26) \quad T_n = n[(U_1 - M^d), (U_2 - M^d)]\Sigma_n^{-1}[(U_1 - M^d), (U_2 - M^d)]'.$$

The approximate null distribution of T_n can be taken as the chi-squared distribution $\chi^2(2)$. Larger values of T_n imply rejection for the null hypothesis (2.6). The

convergence rate of $\sqrt{n}[D(\mathcal{P})]^2 \xrightarrow{\mathcal{P}} 0$ in Corollary 2.4 is also very slow. For example, for the symmetric discrepancy, we can write

$$(2.27) \quad \sqrt{n}D_s(\mathcal{P})^2 = \left(-2, \frac{2^d(n-1)}{n}\right) [\sqrt{n}(U_1 - EU_1), \sqrt{n}(U_2 - EU_2)]' + R_n.$$

The residual term $R_n = [2^d - (\frac{4}{3})^d]/\sqrt{n}$ in (2.27) tends to zero very slowly. When $n = 10,000$, $R_n = 0.0222$ ($d = 2$), 0.2779 ($d = 5$) and 10.0624 ($d = 10$).

For the three special cases of the generalized \mathcal{L}^2 -type discrepancy, we can easily obtain the parameters needed to define the statistics A_n and T_n in (2.24) and (2.26):

- 1) the symmetric discrepancy:

$$U_1 = \frac{1}{n} \sum_{k=1}^n \prod_{j=1}^d (1 + 2z_{kj} - 2z_{kj}^2),$$

$$U_2 = \frac{2^{d+1}}{n(n-1)} \sum_{k < l}^n \prod_{j=1}^d (1 - |z_{kj} - z_{lj}|),$$

$$M = 4/3, \zeta_1 = (9/5)^d - (6/9)^d \text{ and } \zeta_2 = 2^d - (16/9)^d;$$

- 2) the centered discrepancy:

$$U_1 = \frac{1}{n} \sum_{k=1}^n \prod_{j=1}^d \left(1 + \frac{1}{2}|z_{kj} - \frac{1}{2}| - \frac{1}{2}|z_{kj} - \frac{1}{2}|^2\right),$$

$$U_2 = \frac{2}{n(n-1)} \sum_{k < l}^n \prod_{j=1}^d \left[1 + \frac{1}{2}|z_{kj} - \frac{1}{2}| + \frac{1}{2}|z_{lj} - \frac{1}{2}| - \frac{1}{2}|z_{kj} - z_{lj}|\right],$$

$$M = 13/12, \zeta_1 = (47/40)^d - (13/12)^{2d} \text{ and } \zeta_2 = (57/48)^d - (13/12)^{2d};$$

- 3) the star discrepancy:

$$U_1 = \frac{1}{n} \sum_{k=1}^n \prod_{j=1}^d \left(\frac{3 - z_{kj}}{2}\right),$$

$$U_2 = \frac{2}{n(n-1)} \sum_{k < l}^n \prod_{j=1}^d [2 - \max(z_{kj}, z_{lj})],$$

$$M = 4/3, \zeta_1 = (9/5)^d - (16/9)^d \text{ and } \zeta_2 = (11/6)^d - (16/9)^d.$$

3. MONTE CARLO STUDY AND APPLICATIONS

The exact finite-sample distributions of the statistics A_n (given by (2.24)) and T_n (given by (2.26)) under the null hypothesis (2.6) are not readily obtained. However, the effectiveness of the approximation of their finite-sample distributions by their asymptotic distributions can be studied by Monte Carlo simulation. The approximation of the distribution of T_n by $\chi^2(2)$ is influenced not only by the convergence of U_1 and U_2 but by the convergence of Σ_n as well, while the convergence of A_n to normal $N(0, 1)$ is influenced only by the convergence of U_1 and U_2 . Therefore, it is expected that the approximation of the distribution of A_n by $N(0, 1)$ is better than that of T_n by $\chi^2(2)$.

3.1. Numerical comparisons between the finite-sample distributions of A_n and $N(0, 1)$, and T_n and $\chi^2(2)$. In the simulation, for each sample size n ($n = 25, 50, 100, 200$), we generate 10,000 uniform samples $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)'$. The elements of \mathbf{Z} are i.i.d. $U(0, 1)$ -variates. Then we obtain 10,000 values of A_n and T_n under the three (symmetric, centered and star) discrepancies, respectively. These values are sorted and the ordered samples of A_n and T_n are obtained. The empirical quantiles of A_n and T_n , respectively, are calculated from 10,000 order statistics of A_n and T_n . Tables 1, 2, and 3 list the numerical comparisons between some selected $100(1 - \alpha)$ -percentiles ($\alpha = 1\%, 5\%$ and 10%) of A_n and T_n and the $100(1 - \alpha)$ -percentiles of $N(0, 1)$ and $\chi^2(2)$. Since a test using A_n is two-sided, and a test using T_n is one-sided, we list the upper (U) and lower (L) percentiles of both A_n and $N(0, 1)$, and only the upper percentiles of both T_n and $\chi^2(2)$.

The closer the percentiles of A_n and the percentiles of $N(0, 1)$ are, the better approximation we obtain by using $N(0, 1)$ as the approximate finite-sample distribution of A_n . The same is true for the numerical comparison between the statistic T_n and $\chi^2(2)$. Several empirical conclusions can be summarized from the numerical results in Tables 1–3:

- a) the standard normal $N(0, 1)$ approximates the finite-sample distribution of A_n better than the $\chi^2(2)$ approximates the finite-sample distribution of T_n ;
- b) the approximation of the finite-sample distribution of A_n by $N(0, 1)$, and the approximation of the finite-sample distribution of T_n by $\chi^2(2)$, appear to be the best for the symmetric discrepancy;
- c) the approximation of the percentiles of the finite-sample distribution of A_n and the approximation of the percentiles of the finite-sample distribution of T_n for $\alpha = 5\%$ and $\alpha = 10\%$ are much better than for $\alpha = 1\%$; and
- d) the approximation for the case $n = 25$ is almost as good as those cases for $n = 200$.

3.2. Type I error rates. Based on the numerical comparisons shown in Tables 1–3, we perform a simulation of the empirical type I error rates of the two statistics A_n and T_n under the three discrepancies. For convenience, we choose the null distribution of the random vectors \mathbf{z}_i to be composed of i.i.d. $U(0, 1)$ marginals. In the simulation, we generate 2,000 sets of $\{\mathbf{z}_1, \dots, \mathbf{z}_n\}$ for each n with the components of \mathbf{z}_i consisting of i.i.d. $U(0, 1)$ variates. Tables 4, 5, and 6 summarize the simulation results on the type I error rates of A_n and T_n under the three discrepancies, where the percentiles of A_n are chosen as those of $N(0, 1)$ and the percentiles of T_n as those of $\chi^2(2)$. It shows that for $\alpha = 5\%$ and $\alpha = 10\%$, the type I error rates are better controlled for the statistics A_n and T_n , while the type I error rates for $\alpha = 1\%$ tend to be large by using the percentiles of $N(0, 1)$ for A_n and the percentiles of $\chi^2(2)$ for T_n in most cases.

3.3. Power study. Now we turn to study the power of A_n and T_n in testing hypothesis (2.6). The alternative distributions are chosen to be *meta-type uniform distributions*. The theoretical background of some meta-type multivariate distributions is given in [KS91] and [FFK97]. The idea for constructing the alternative distributions is as follows. Let the random vector $\mathbf{x} = (X_1, \dots, X_d)$ have a d.f. $F(\mathbf{x})$ ($\mathbf{x} \in R^d$). Denote by $F_i(x_i)$ the marginal d.f. of X_i ($i = 1, \dots, d$), which is assumed to be continuous. Define the random vector $\mathbf{u} = (U_1, \dots, U_d)$ by

$$(3.1) \quad U_i = F_i(X_i), \quad i = 1, \dots, d.$$

TABLE 1. Comparisons between the empirical percentiles of A_n , T_n and the percentiles of $N(0, 1)$ and $\chi^2(2)$ for the symmetric discrepancy (U=upper, L=lower)

| | | A_n | | | | | T_n | | |
|---------------|--|---------|---------|---------|-----------------|--|---------|--------|--------|
| α | | 1% | 5% | 10% | α | | 1% | 5% | 10% |
| $N(0, 1)$ (U) | | 2.5758 | 1.9600 | 1.6449 | $\chi^2(2)$ (U) | | 9.2103 | 5.9915 | 4.6052 |
| (L) | | -2.5758 | -1.9600 | -1.6449 | | | | | |
| $d = 2$ | | | | | | | | | |
| $n = 25$ (U) | | 2.9759 | 2.1836 | 1.7969 | $n = 25$ | | 13.7011 | 6.7329 | 4.3474 |
| (L) | | -2.6414 | -1.9566 | -1.6775 | | | | | |
| 50(U) | | 2.6455 | 2.0585 | 1.7237 | 50 | | 13.5989 | 6.3139 | 4.3748 |
| (L) | | -2.4507 | -1.9969 | -1.6847 | | | | | |
| 100(U) | | 2.6058 | 2.0078 | 1.6935 | 100 | | 15.2375 | 6.2400 | 4.3013 |
| (L) | | -2.4772 | -1.9102 | -1.6327 | | | | | |
| 200(U) | | 2.5923 | 1.9897 | 1.6601 | 200 | | 15.2146 | 6.7055 | 4.3251 |
| (L) | | -2.6379 | -1.9279 | -1.6301 | | | | | |
| $d = 5$ | | | | | | | | | |
| $n = 25$ (U) | | 2.9520 | 2.1698 | 1.8390 | $n = 25$ | | 10.6401 | 6.0031 | 4.2624 |
| (L) | | -2.4989 | -1.9750 | -1.7051 | | | | | |
| 50(U) | | 2.7856 | 2.0814 | 1.7396 | 50 | | 11.4226 | 6.1015 | 4.4108 |
| (L) | | -2.3915 | -1.9143 | -1.6484 | | | | | |
| 100(U) | | 2.6219 | 2.0271 | 1.6849 | 100 | | 11.4069 | 6.0066 | 4.2354 |
| (L) | | -2.5679 | -1.9372 | -1.6497 | | | | | |
| 200(U) | | 2.6534 | 2.0671 | 1.7026 | 200 | | 12.7450 | 6.0279 | 4.4207 |
| (L) | | -2.5295 | -1.9199 | -1.6241 | | | | | |
| $d = 10$ | | | | | | | | | |
| $n = 25$ (U) | | 3.2566 | 2.3549 | 1.9631 | $n = 25$ | | 11.1983 | 6.3805 | 4.6083 |
| (L) | | -2.4658 | -1.9266 | -1.6559 | | | | | |
| 50(U) | | 2.9148 | 2.0965 | 1.7282 | 50 | | 10.6401 | 5.7312 | 4.2321 |
| (L) | | -2.5183 | -1.8729 | -1.5886 | | | | | |
| 100(U) | | 2.7635 | 2.0861 | 1.7492 | 100 | | 10.0934 | 5.7531 | 4.3382 |
| (L) | | -2.4521 | -1.8981 | -1.6210 | | | | | |
| 200(U) | | 2.8986 | 2.0855 | 1.6790 | 200 | | 10.7926 | 6.3940 | 4.6259 |
| (L) | | -2.5858 | -1.9743 | -1.6473 | | | | | |

TABLE 2. Comparisons between the empirical percentiles of A_n , T_n and the percentiles of $N(0, 1)$ and $\chi^2(2)$ for the centered discrepancy (U=upper, L=lower)

| | | A_n | | | | | T_n | | |
|---------------|--|---------|---------|---------|-----------------|--|---------|--------|--------|
| α | | 1% | 5% | 10% | α | | 1% | 5% | 10% |
| $N(0, 1)$ (U) | | 2.5758 | 1.9600 | 1.6449 | $\chi^2(2)$ (U) | | 9.2103 | 5.9915 | 4.6052 |
| (L) | | -2.5758 | -1.9600 | -1.6449 | | | | | |
| $d = 2$ | | | | | | | | | |
| $n = 25$ (U) | | 3.0050 | 2.1787 | 1.8135 | $n = 25$ | | 14.7350 | 6.6225 | 4.4123 |
| (L) | | -2.5991 | -1.9844 | -1.6711 | | | | | |
| 50(U) | | 2.8674 | 2.1098 | 1.8002 | 50 | | 13.2973 | 6.4543 | 4.4933 |
| (L) | | -2.4697 | -1.9843 | -1.6570 | | | | | |
| 100(U) | | 2.7243 | 2.0026 | 1.6594 | 100 | | 15.1737 | 6.4344 | 4.2733 |
| (L) | | -2.4646 | -1.9056 | -1.6336 | | | | | |
| 200(U) | | 2.6659 | 2.0101 | 1.6705 | 200 | | 15.9685 | 6.5391 | 4.4612 |
| (L) | | -2.4212 | -1.9303 | -1.6483 | | | | | |
| $d = 5$ | | | | | | | | | |
| $n = 25$ (U) | | 2.9690 | 2.1944 | 1.8036 | $n = 25$ | | 12.0965 | 6.2195 | 4.2602 |
| (L) | | -2.6154 | -1.9494 | -1.6214 | | | | | |
| 50(U) | | 2.7896 | 2.0864 | 1.7549 | 50 | | 12.3615 | 6.3261 | 4.3271 |
| (L) | | -2.5491 | -1.9160 | -1.6482 | | | | | |
| 100(U) | | 2.7273 | 2.0579 | 1.6885 | 100 | | 12.0372 | 6.1285 | 4.3838 |
| (L) | | -2.4673 | -1.9303 | -1.6325 | | | | | |
| 200(U) | | 2.6119 | 1.9727 | 1.6686 | 200 | | 14.1589 | 6.3164 | 4.3577 |
| (L) | | -2.4947 | -1.9312 | -1.6596 | | | | | |
| $d = 10$ | | | | | | | | | |
| $n = 25$ (U) | | 3.0643 | 2.1765 | 1.7999 | $n = 25$ | | 11.8336 | 6.2523 | 4.4722 |
| (L) | | -2.5317 | -1.9699 | -1.7057 | | | | | |
| 50(U) | | 2.7659 | 2.0547 | 1.7327 | 50 | | 10.3546 | 5.8925 | 4.4789 |
| (L) | | -2.4328 | -1.9181 | -1.6338 | | | | | |
| 100(U) | | 2.7047 | 1.9789 | 1.6812 | 100 | | 10.6855 | 5.9538 | 4.3766 |
| (L) | | -2.5490 | -1.9115 | -1.6321 | | | | | |
| 200(U) | | 2.6627 | 2.0203 | 1.7219 | 200 | | 11.4535 | 6.3443 | 4.4515 |
| (L) | | -2.6466 | -1.9017 | -1.5883 | | | | | |

TABLE 3. Comparisons between the empirical percentiles of A_n , T_n and the percentiles of $N(0, 1)$ and $\chi^2(2)$ for the star discrepancy (U=upper, L=lower)

| | | A_n | | | | | T_n | | |
|---------------|--|---------|---------|---------|-----------------|---------|--------|--------|--|
| α | | 1% | 5% | 10% | α | 1% | 5% | 10% | |
| $N(0, 1)$ (U) | | 2.5758 | 1.9600 | 1.6449 | $\chi^2(2)$ (U) | 9.2103 | 5.9915 | 4.6052 | |
| (L) | | -2.5758 | -1.9600 | -1.6449 | | | | | |
| $d = 2$ | | | | | | | | | |
| $n = 25$ (U) | | 2.7798 | 2.1064 | 1.7502 | $n = 25$ | 18.1738 | 7.4153 | 4.2674 | |
| (L) | | -2.4210 | -1.8260 | -1.5894 | | | | | |
| 50 (U) | | 2.6203 | 2.0494 | 1.7041 | 50 | 17.7953 | 7.3728 | 4.3048 | |
| (L) | | -2.4262 | -1.8957 | -1.6033 | | | | | |
| 100 (U) | | 2.7023 | 2.0041 | 1.6674 | 100 | 20.2587 | 7.0995 | 4.2731 | |
| (L) | | -2.5098 | -1.9783 | -1.6809 | | | | | |
| 200 (U) | | 2.7373 | 2.0500 | 1.7079 | 200 | 21.0088 | 7.6134 | 4.3452 | |
| (L) | | -2.5060 | -1.9245 | -1.6407 | | | | | |
| $d = 5$ | | | | | | | | | |
| $n = 25$ (U) | | 2.9385 | 2.0791 | 1.7058 | $n = 25$ | 14.3399 | 6.3745 | 4.0116 | |
| (L) | | -2.2055 | -1.7930 | -1.5380 | | | | | |
| 50 (U) | | 2.8831 | 2.0920 | 1.7123 | 50 | 16.6810 | 6.7896 | 4.1579 | |
| (L) | | -2.4423 | -1.8550 | -1.5617 | | | | | |
| 100 (U) | | 2.7979 | 2.1148 | 1.7166 | 100 | 14.6865 | 6.7067 | 4.1141 | |
| (L) | | -2.3267 | -1.8339 | -1.6227 | | | | | |
| 200 (U) | | 2.8597 | 2.0473 | 1.6776 | 200 | 18.5595 | 6.6619 | 4.1685 | |
| (L) | | -2.4238 | -1.8852 | -1.5846 | | | | | |
| $d = 10$ | | | | | | | | | |
| $n = 25$ (U) | | 2.9127 | 2.1215 | 1.7109 | $n = 25$ | 14.8453 | 6.4506 | 4.1141 | |
| (L) | | -2.1692 | -1.8174 | -1.5444 | | | | | |
| 50 (U) | | 2.9363 | 2.1483 | 1.7379 | 50 | 15.4874 | 6.7275 | 4.1528 | |
| (L) | | -2.3410 | -1.8569 | -1.5453 | | | | | |
| 100 (U) | | 2.7731 | 2.0901 | 1.7077 | 100 | 15.2233 | 6.4890 | 4.0312 | |
| (L) | | -2.4133 | -1.8551 | -1.5742 | | | | | |
| 200 (U) | | 2.8309 | 2.0886 | 1.6720 | 200 | 16.2075 | 6.8058 | 4.1076 | |
| (L) | | -2.4539 | -1.8862 | -1.5914 | | | | | |

TABLE 4. Empirical type I error rates of A_n and T_n under the symmetric discrepancy

| α | A_n | | | T_n | | |
|----------|--------|--------|--------|--------|--------|--------|
| | 1% | 5% | 10% | 1% | 5% | 10% |
| $d = 2$ | | | | | | |
| $n = 25$ | 0.0165 | 0.0660 | 0.1265 | 0.0245 | 0.0550 | 0.0905 |
| 50 | 0.0130 | 0.0575 | 0.1095 | 0.0265 | 0.0615 | 0.0915 |
| 100 | 0.0150 | 0.0615 | 0.1120 | 0.0255 | 0.0595 | 0.0860 |
| 200 | 0.0095 | 0.0585 | 0.1010 | 0.0270 | 0.0580 | 0.0920 |
| $d = 5$ | | | | | | |
| $n = 25$ | 0.0160 | 0.0695 | 0.1315 | 0.0195 | 0.0525 | 0.1035 |
| 50 | 0.0140 | 0.0605 | 0.1095 | 0.0205 | 0.0545 | 0.0970 |
| 100 | 0.0150 | 0.0525 | 0.1060 | 0.0190 | 0.0515 | 0.0815 |
| 200 | 0.0125 | 0.0540 | 0.0925 | 0.0175 | 0.0535 | 0.0930 |
| $d = 10$ | | | | | | |
| $n = 25$ | 0.0235 | 0.0650 | 0.1245 | 0.0240 | 0.0620 | 0.1005 |
| 50 | 0.0105 | 0.0550 | 0.1095 | 0.0105 | 0.0435 | 0.0860 |
| 100 | 0.0145 | 0.0595 | 0.1145 | 0.0155 | 0.0550 | 0.1010 |
| 200 | 0.0075 | 0.0410 | 0.0835 | 0.0135 | 0.0420 | 0.0825 |

TABLE 5. Empirical type I error rates of A_n and T_n under the centered discrepancy

| α | A_n | | | T_n | | |
|----------|--------|--------|--------|--------|--------|--------|
| | 1% | 5% | 10% | 1% | 5% | 10% |
| $d = 2$ | | | | | | |
| $n = 25$ | 0.0180 | 0.0655 | 0.1240 | 0.0275 | 0.0610 | 0.0950 |
| 50 | 0.0120 | 0.0565 | 0.1020 | 0.0270 | 0.0640 | 0.0940 |
| 100 | 0.0170 | 0.0640 | 0.1210 | 0.0260 | 0.0595 | 0.0950 |
| 200 | 0.0090 | 0.0495 | 0.0955 | 0.0280 | 0.0525 | 0.0890 |
| $d = 5$ | | | | | | |
| $n = 25$ | 0.0135 | 0.0665 | 0.1195 | 0.0200 | 0.0575 | 0.0915 |
| 50 | 0.0135 | 0.0625 | 0.1235 | 0.0210 | 0.0570 | 0.0960 |
| 100 | 0.0095 | 0.0530 | 0.1070 | 0.0185 | 0.0515 | 0.0895 |
| 200 | 0.0130 | 0.0475 | 0.1050 | 0.0145 | 0.0575 | 0.0910 |
| $d = 10$ | | | | | | |
| $n = 25$ | 0.0150 | 0.0690 | 0.1325 | 0.0230 | 0.0650 | 0.1095 |
| 50 | 0.0100 | 0.0460 | 0.1045 | 0.0125 | 0.0410 | 0.0835 |
| 100 | 0.0105 | 0.0495 | 0.1125 | 0.0145 | 0.0555 | 0.1045 |
| 200 | 0.0100 | 0.0455 | 0.0970 | 0.0170 | 0.0500 | 0.0850 |

It is obvious that each U_i in (3.1) has a uniform distribution $U(0, 1)$, but the joint distribution of $\mathbf{u} = (U_1, \dots, U_d)$ may be quite different from the uniform distribution in \bar{C}^d . If the joint d.f. $F(\mathbf{x})$ of $\mathbf{x} = (X_1, \dots, X_d)$ possesses a density function $f(\mathbf{x}) = f(x_1, \dots, x_d)$, where $f_i(x_i)$ denotes the marginal density function of X_i , then the joint density function of $\mathbf{u} = (U_1, \dots, U_d)$ can be obtained by a

TABLE 6. Empirical type I error rates of A_n and T_n under the star discrepancy

| α | A_n | | | T_n | | |
|----------|--------|--------|--------|--------|--------|--------|
| | 1% | 5% | 10% | 1% | 5% | 10% |
| $d = 2$ | | | | | | |
| $n = 25$ | 0.0105 | 0.0470 | 0.0955 | 0.0325 | 0.0630 | 0.0845 |
| 50 | 0.0105 | 0.0450 | 0.0895 | 0.0385 | 0.0610 | 0.0860 |
| 100 | 0.0075 | 0.0465 | 0.0930 | 0.0300 | 0.0550 | 0.0800 |
| 200 | 0.0100 | 0.0510 | 0.1030 | 0.0360 | 0.0630 | 0.0920 |
| $d = 5$ | | | | | | |
| $n = 25$ | 0.0100 | 0.0560 | 0.1165 | 0.0260 | 0.0665 | 0.0975 |
| 50 | 0.0050 | 0.0530 | 0.0980 | 0.0255 | 0.0665 | 0.0920 |
| 100 | 0.0105 | 0.0575 | 0.1090 | 0.0325 | 0.0675 | 0.0940 |
| 200 | 0.0120 | 0.0555 | 0.0985 | 0.0290 | 0.0585 | 0.0835 |
| $d = 10$ | | | | | | |
| $n = 25$ | 0.0115 | 0.0460 | 0.0940 | 0.0245 | 0.0630 | 0.0890 |
| 50 | 0.0200 | 0.0525 | 0.1060 | 0.0275 | 0.0605 | 0.0890 |
| 100 | 0.0120 | 0.0525 | 0.0895 | 0.0265 | 0.0550 | 0.0810 |
| 200 | 0.0110 | 0.0470 | 0.0945 | 0.0300 | 0.0545 | 0.0780 |

direct calculation:

$$(3.2) \quad p(u_1, \dots, u_d) = f(F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)) / \prod_{i=1}^d f_i(F_i^{-1}(u_i)),$$

where $(u_1, \dots, u_d) \in \bar{C}^d$ and $F_i^{-1}(\cdot)$ denotes the inverse function of $F_i(\cdot)$. It is clear that the complexity of (3.2) is determined by the joint distribution of $\mathbf{x} = (X_1, \dots, X_d)$ in (3.1). In particular, if the random variables X_i in (3.1) are independent, then the U_i 's given by (3.1) will be i.i.d. $U(0, 1)$ variates. In this case, the random vector $\mathbf{u} = (U_1, \dots, U_d)$ is uniformly distributed in \bar{C}^d . This can be seen from (3.2).

In the simulation, we choose the random vector $\mathbf{x} = (X_1, \dots, X_d)$ in (3.1) to have a joint distribution belonging to the subclasses of elliptical distributions [FKN90, Chapter 3], where the parameters μ and Σ are chosen as $\mu = \mathbf{0}$ and $\Sigma = (\sigma_{ij})$, where $\sigma_{ii} = 1$ and $\sigma_{ij} = \sigma_{ji} = \rho = 0.5$ for $1 \leq i \neq j \leq d$. Except the normal distribution $N_d(\mu, \Sigma)$, we give general expressions for the density functions of the selected subclasses of elliptical distributions below:

- 1) the multivariate t -distribution $\mathbf{x} \sim Mt_d(m, \mu, \Sigma)$, with density function given by

$$C|\Sigma|^{-1/2} (1 + m^{-1}(\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu))^{-(d+m)/2}, \quad m > 0, \mathbf{x} \in R^d,$$

where C is a normalizing constant (the following C 's have a similar meaning, but their values may be different);

- 2) the Kotz type distribution, with density function given by

$$C|\Sigma|^{-1/2} [(\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu)]^{N-1} \exp \{-r[(\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu)]^s\},$$

$r, s > 0, 2N + d > 2, \mathbf{x} \in R^d;$

3) the Pearson type VII distribution, with density function given by

$$C|\Sigma|^{-1/2} (1 + m^{-1}(\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu))^{-N}, \quad N > d/2, m > 0, \mathbf{x} \in R^d;$$

4) the Pearson type II distribution, with density function given by

$$C|\Sigma|^{-1/2} (1 - (\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu))^m, \quad m > -1, (\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu) < 1;$$

5) the multivariate Cauchy distribution $MC_d(\mu, \Sigma)$, with density function given by

$$\frac{\Gamma((d+1)/2)}{\pi^{(d+1)/2}} (1 + (\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu))^{-(d+1)/2}, \quad \mathbf{x} \in R^d.$$

When the random vector $\mathbf{x} = (X_1, \dots, X_d)'$ is generated by one of the elliptical distributions above, then the random vector $\mathbf{u} = (U_1, \dots, U_d)'$ given by (3.1) is considered to have a meta-type uniform distribution denoted as follows:

- 0) $\mathbf{u} \sim MNU$ when $\mathbf{x} \sim N_d(\mathbf{0}, \Sigma)$;
- 1) $\mathbf{u} \sim MTU$ when \mathbf{x} has a multivariate t -distribution with $m = 5$;
- 2) $\mathbf{u} \sim MKU$ when \mathbf{x} has a Kotz type distribution with $N = 2$, $r = 1$ and $s = 0.5$;
- 3) $\mathbf{u} \sim MPVIIU$ when \mathbf{x} has a Pearson type VII distribution with $N = 10$ and $m = 2$;
- 4) $\mathbf{u} \sim MPIIU$ when \mathbf{x} has a Pearson type II distribution with $m = 3/2$;
- 5) $\mathbf{u} \sim MCU$ when \mathbf{x} has a Cauchy distribution.

The power of the multivariate test for uniformity is the probability that the statistical test correctly identifies a sample coming from one of the above distributions as being non-uniform. Table 7 summarizes the simulation results on the power of A_n and T_n , where the simulation is done with 2,000 replications, the critical points of both A_n and T_n are chosen as those of $N(0, 1)$ and $\chi^2(2)$, respectively, and the empirical samples from the elliptical distributions are generated by the *TFWW algorithm* [Tas77] and [FW94, pp. 160-170]. It shows that the two statistics A_n and T_n under the three discrepancies are powerful for testing uniformity in \bar{C}^d in most cases. The χ^2 -type statistic T_n seems to be more powerful than the normal-type statistic A_n . For both A_n and T_n , they seem to be more powerful in the higher dimensional case ($d = 10$) than in the lower dimensional case ($d = 5$). It is also noticed that the two statistics A_n and T_n under the symmetric discrepancy seem to be more powerful than under the centered discrepancy and the star discrepancy in most cases.

A power comparison between the two statistics A_n, T_n and some existing statistics, such as the Kolmogorov-Smirnov type statistic (1.1), seems to be infeasible in high dimensions because of computational difficulties. Thus, we have not performed such comparisons.

3.4. Applications. By using the Rosenblatt transformation [Ros52], we can transfer a test for the simple hypothesis

$$(3.3) \quad H_0 : \mathbf{x}_1, \dots, \mathbf{x}_n \text{ have a known d.f. } F(\mathbf{x}),$$

to a test for the uniformity of random points in \bar{C}^d , and then apply the two statistics A_n and T_n to test for uniformity of the transformed variates in \bar{C}^d . Denote by

$$(3.4) \quad \begin{aligned} F_1(x_1) &= \text{the marginal d.f. of } X_1, \\ F_{2|1}(x_2|x_1) &= \text{the conditional d.f. of } X_2 \text{ given } X_1 = x_1, \\ &\vdots \quad \quad \quad \vdots \\ F_{k+1|(1,\dots,k)}(x_{k+1}|x_1, \dots, x_k) &= \text{the conditional d.f. of } X_{k+1} \\ &\quad \quad \quad \text{given } (X_1, \dots, X_k) = (x_1, \dots, x_k), \end{aligned}$$

where $k = 1, \dots, d-1$. Perform the following series of Rosenblatt transformations on each observation $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$ ($i = 1, \dots, n$):

$$(3.5) \quad \begin{aligned} U_{i1} &= F_1(x_{i1}) \\ U_{i2} &= F_{2|1}(x_{i2}|x_{i1}) \\ &\vdots \quad \quad \quad \vdots \\ U_{i,k+1} &= F_{k+1|(1,\dots,k)}(x_{i,k+1}|(x_{i1}, \dots, x_{ik})), \end{aligned}$$

where $i = 1, \dots, n$ and $k = 1, \dots, d-1$. If the null hypothesis (3.3) is true, the random vectors $\mathbf{u}_i = (U_{i1}, \dots, U_{id})$ ($i = 1, \dots, n$) are i.i.d. and the components of \mathbf{u}_i have a uniform distribution $U[0, 1]$.

The more frequent cases in applications are those hypotheses which involve unknown parameters in the null distributions. Justel, Peña and Zamar [JPZ97] proposed a multivariate version of the Kolmogorov-Smirnov type statistic (1.1) for testing the simple hypothesis (3.3). Their statistics are difficult to compute for large dimensions ($d \geq 3$) and always require estimating unknown parameters in the null distribution. The two statistics A_n and T_n developed in this paper are easy to compute in arbitrary dimension, and estimation of unknown parameters can be avoided when the null distribution belongs to the class of the multivariate normal distributions, the spherically symmetric distributions, the l_1 -norm symmetric distributions [FKN90, Chapter 5], the l_p -norm symmetric distributions [YM95], or the L_p -norm spherical distributions [GS97].

REFERENCES

- [AD54] T. W. Anderson and D. A. Darling, *A test of goodness-of-fit*, J. Amer. Statist. Assoc. **49** (1954), 765–769. MR **16**:1039h
- [DS86] R. B. D'Agostino and M. A. Stephens (eds.), *Goodness-of-fit techniques*, Marcel Dekker, Inc., New York and Basel, 1986. MR **88c**:62075
- [FFK97] H. B. Fang, K. T. Fang, and S. Kotz, *The meta-elliptical distributions with given marginals*, Tech. Report MATH-165, Hong Kong Baptist University, 1997.
- [FKN90] K. T. Fang, S. Kotz, and K. W. Ng, *Symmetric multivariate and related distributions*, Chapman and Hall, London and New York, 1990. MR **91i**:62070
- [FW94] K. T. Fang and Y. Wang, *Number-theoretic methods in statistics*, Chapman and Hall, London, 1994. MR **95g**:65189
- [GS97] A. K. Gupta and D. Song, *L_p -norm spherical distributions*, J. Statist. Plann. Inference **60** (1997), 241–260. MR **98h**:62092
- [Hic98a] F. J. Hickernell, *A generalized discrepancy and quadrature error bound*, Math. Comp. **67** (1998), 299–322. MR **98c**:65032
- [Hic98b] F. J. Hickernell, *Lattice rules: How well do they measure up?*, Random and Quasi-Random Point Sets (P. Hellekalek and G. Larcher, eds.), Lecture Notes in Statistics, vol. 138, Springer-Verlag, New York, 1998, pp. 109–166. CMP 99:06

- [Hic99] F. J. Hickernell, *Goodness-of-fit statistics, discrepancies and robust designs*, Statist. Probab. Lett. **44** (1999), 73–78. CMP 99:17
- [HW81] L. G. Hua and Y. Wang, *Applications of number theory to numerical analysis*, Springer-Verlag, Berlin, and Science Press Beijing, 1981. MR **83g**:10034
- [JPZ97] A. Justel, D. Peña, and R. Zamar, *A multivariate Kolmogorov-Smirnov test of goodness of fit*, Statist. Probab. Lett. **35** (1997), 251–259. MR **98k**:62087
- [KS91] S. Kotz and J. P. Seeger, *A new approach to dependence in multivariate distributions*, Advances in Probability Distributions with Given Marginals (S. Kotz G. Dall’Aglio and G. Salinetti, eds.), Kluwer, Dordrecht, Netherlands, 1991, pp. 113–127. MR **95k**:62162
- [MQ79] F. L. Miller Jr. and C. P. Quesenberry, *Power studies of tests for uniformity II*, Comm. Statist. Simulation Comput. **B8(3)** (1979), 271–290.
- [Ney37] J. Neyman, “Smooth” test for goodness of fit, J. Amer. Statist. Assoc. **20** (1937), 149–199.
- [Nie92] H. Niederreiter, *Random number generation and quasi-Monte Carlo methods*, CBMS-NSF Regional Conf. Ser. Appl. Math., vol. 63, SIAM, Philadelphia, 1992. MR **93h**:65008
- [Pea39] E. S. Pearson, *The probability transformation for testing goodness of fit and combining independent tests of significance*, Biometrika **30** (1939), 134–148.
- [QM77] C. P. Quesenberry and F. L. Miller Jr., *Power studies of some tests for uniformity*, J. Statist. Comput. Simulation **5** (1977), 169–192.
- [Ros52] M. Rosenblatt, *Remarks on a multivariate transformation*, Ann. Math. Statist. **23** (1952), 470–472. MR **14**:189j
- [Ser80] R. J. Serfling, *Approximation theorems of mathematical statistics*, John Wiley & Sons inc., New York, 1980. MR **82a**:62003
- [SJ94] I. H. Sloan and S. Joe, *Lattice methods for multiple integration*, Oxford University Press, Oxford, 1994. MR **98a**:65026
- [Tas77] D. Tashiro, *On methods for generating uniform points on the surface of a sphere*, Ann. Inst. Statist. Math. **29** (1977), 295–300. MR **58**:8144
- [Tez95] S. Tezuka, *Uniform random numbers: Theory and practice*, Kluwer Academic Publishers, Boston, 1995.
- [Wat62] G. S. Watson, *Goodness-of-fit tests on a circle. II*, Biometrika **49** (1962), 57–63. MR **25**:1626
- [YM95] X. Yue and C. Ma, *Multivariate l_p -norm symmetric distributions*, Statist. Probab. Lett. **24** (1995), 281–288. MR **96i**:62057

DEPARTMENT OF MATHEMATICS, HONG KONG BAPTIST UNIVERSITY, KOWLOON TONG, HONG KONG SAR, CHINA, AND INSTITUTE OF APPLIED MATHEMATICS, CHINESE ACADEMY OF SCIENCES, BEIJING, CHINA

E-mail address: jjliang@hkbu.edu.hk

DEPARTMENT OF MATHEMATICS, HONG KONG BAPTIST UNIVERSITY, KOWLOON TONG, HONG KONG SAR, CHINA, AND INSTITUTE OF APPLIED MATHEMATICS, CHINESE ACADEMY OF SCIENCES, BEIJING, CHINA

E-mail address: ktfang@hkbu.edu.hk

DEPARTMENT OF MATHEMATICS, HONG KONG BAPTIST UNIVERSITY, KOWLOON TONG, HONG KONG SAR, CHINA

E-mail address: fred@hkbu.edu.hk

URL: <http://www.math.hkbu.edu.hk/~fred>

DEPARTMENT OF STATISTICS, UNIVERSITY OF NORTH CAROLINA, CHAPEL HILL, NC, 27599-3260, UNITED STATES OF AMERICA

E-mail address: lirz@email.unc.edu