

Testing Multivariate Normality Based on t -Representative Points

Jiajuan Liang ^{1,2,*}, Ping He ^{1,2} and Jun Yang ^{1,†}

¹ Department of Statistics and Data Science, BNU-HKBU United International College, Zhuhai 519087, China

² Guangdong Provincial Key Laboratory of Interdisciplinary Research and Application for Data Science, BNU-HKBU United International College, Zhuhai 519087, China

* Correspondence: jjliang@uic.edu.cn

† J.L. and P.H. are deeply indebted to our Ph.D. student J.Y. who finished all simulation work in this paper but unfortunately passed away on 1 January 2022 due to illness.

Abstract: Testing multivariate normality is an ever-lasting interest in the goodness-of-fit area since the classical Pearson's chi-squared test. Among the numerous approaches in the construction of tests for multivariate normality, normal characterization is one of the common approaches, which can be divided into the necessary and sufficient characterization and necessary-only characterization. We construct a test for multivariate normality by combining the necessary-only characterization and the idea of statistical representative points in this paper. The main idea is to transform a high-dimensional sample into a one-dimensional one through the necessary normal characterization and then employ the representative-point-based Pearson's chi-squared test. A limited Monte Carlo study shows a considerable power improvement of the representative-point-based chi-square test over the traditional one. An illustrative example is given to show the supplemental function of the new test when used together with existing ones in the literature.

Keywords: chi-squared test; multivariate normality; representative points; spherical distribution; Student's t -distribution

MSC: 62H15; 62E10



Citation: Liang, J.; He, P.; Yang, J. Testing Multivariate Normality Based on t -Representative Points. *Axioms* **2022**, *11*, 587. <https://doi.org/10.3390/axioms11110587>

Academic Editor: Hans J. Haubold

Received: 4 September 2022

Accepted: 21 October 2022

Published: 24 October 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

The multivariate normal distribution can be characterized many different ways. For example, (1) one of the most well-known characterizations is that a p -dimensional random vector $x = (X_1, \dots, X_p)'$ ($p \times 1$) has a p -dimensional normal distribution if and only if all of its linear combinations $a'x$ ($a \in R^p$, the p -dimensional Euclidean space) has a univariate normal distribution; (2) Anderson [1] gave the multivariate normal characterizations by its mean vector and its covariance matrix, respectively; (3) Shao and Zhou [2] gave a characterization for a random vector without a multivariate normal distribution; to name a few. Most of the existing tests for MVN (multivariate normality) in the literature are more or less based on some kind of characterizations for MVN ([3,4]); some are based on necessary and sufficient characterizations for MVN; some others are only based on necessary characterizations for MVN ([5–17]). An MVN test based on a necessary characterization is called a necessary test in the literature. By reviewing a wide class of MVN tests, Ebner and Henze [18] discuss many of the most updated tests for MVN, which consists of mostly necessary ones. Necessary tests for MVN are usually easier to construct and their exact null distributions or asymptotic null distributions are easier to obtain with simple analytical expressions. However, a common drawback of all necessary tests for MVN is that if the null hypothesis of MVN is not rejected, there is no guarantee of MVN for the original data. Therefore, a characterization-based MVN test may have better power performance than many non-characterization-based MVN-tests. This is our motivation for developing a new characterization-based MVN test.

Yang, Fang, and Liang [19] gave a characterization for MVN by using the the sampling distributions of a series of specially constructed random vectors from an i.i.d. (independently identically distributed) normal sample from a zero-mean p -dimensional normal distribution $N_p(\mathbf{0}, \Sigma)$ with an unknown covariance matrix $\Sigma > \mathbf{0}$ (positive definite). This characterization can be employed to construct Q-Q (quantile-quantile) plots for detecting non-multinormality ([20]). In this paper, we will develop further application of the MVN characterization in [19] through employing the idea of statistical representative points (RP for short, [21]), the properties of spherical distributions ([22]), and the classical Pearson–Fisher statistic ([23,24]). This paper is organized as follows. A brief review on the MVN characterization in [19] is given in Section 2. The development of the new test is given in Section 3. Section 4 is devoted to a Monte Carlo study on the empirical performance of the new test and its application. Some concluding remarks are given in the last section.

2. A Brief Review on the MVN Characterization

Let continuous random vectors $(d \times 1)$ x_1, \dots, x_n be independently identically distributed (i.i.d.) according to a probability density function (p.d.f.) $f(x)$, $x \in R^d$ (the d -dimensional Euclidean space). For simplicity, we write x_1, \dots, x_n i.i.d. $\sim f(x)$. It is assumed that the p.d.f. $f(x)$ has a zero mean and finite second-order moments. Let:

$$S_k = \sum_{i=1}^k x_i x_i', \quad y_k = S_k^{-\frac{1}{2}} x_k, \quad z_k = \frac{y_k}{\sqrt{1 - y_k' y_k}}, \tag{1}$$

for $k = d + 1, \dots, n$, where $S_k^{-\frac{1}{2}} = (S_k^{\frac{1}{2}})^{-1}$ and $S_k^{\frac{1}{2}}$ stands for the positive definite square root of S_k .

Theorem 1 ([19]). *Let x_1, \dots, x_n be i.i.d. $\sim N_d(\mathbf{0}, \Sigma)$. Define the random vectors y_k and z_k by (1). Then,*

1. y_{d+1}, \dots, y_n are mutually independent and y_k ($k \geq d + 1$) has a symmetric multivariate Pearson Type II distribution with a p.d.f.,

$$f_k(y) = \frac{\Gamma(\frac{k}{2})}{\pi^{\frac{d}{2}} \Gamma(\frac{k-d}{2})} (1 - y' y)^{\frac{k-d-2}{2}}, \quad y \in R^d, \quad y' y < 1. \tag{2}$$

2. z_{d+1}, \dots, z_n are mutually independent and z_k ($k \geq d + 1$) has a symmetric multivariate Pearson Type VII distribution with a p.d.f.,

$$h_k(z) = \frac{\Gamma(\frac{k}{2})}{\pi^{\frac{d}{2}} \Gamma(\frac{k-d}{2})} (1 + z' z)^{-\frac{k}{2}}, \quad z \in R^d. \tag{3}$$

3. Let x_1, \dots, x_n be i.i.d. in R^d with a p.d.f. $v(x)$, which is continuous in $x \in R^d$ and $v(\mathbf{0}) > 0$. Define the random vectors z_k by (1). If z_k and z_{k-1} ($k \geq d + 2$) have p.d.f.'s $h_k(z)$ and $h_{k-1}(z)$ defined by (3), respectively, then x_i ($i = 1, \dots, n$) has a multivariate normal distribution.

From the above theorem, we can derive the following corollary.

Corollary 1. *Assume that x_1, \dots, x_n are i.i.d. $\sim N_d(\mu, \Sigma)$. Define the random vectors:*

$$u_i = \frac{x_1 + \dots + x_i - i x_{i+1}}{\sqrt{i(i+1)}}, \quad i = 1, \dots, n - 1. \tag{4}$$

and,

$$S_k = \sum_{i=1}^k u_i u_i', \quad y_k = S_k^{-\frac{1}{2}} u_k, \quad k = d + 1, \dots, n - 1. \tag{5}$$

Let $\mathbf{y}_k = (y_{k1}, \dots, y_{kd})'$: $d \times 1$ ($k = d + 1, \dots, n - 1$) and:

$$\bar{y}_k = \frac{1}{d} \sum_{i=1}^d y_{ki}, \quad t_k = \frac{\sqrt{d}\bar{y}_k}{v_k}, \quad v_k^2 = \frac{1}{d-1} \sum_{i=1}^d (y_{ki} - \bar{y}_k)^2; \tag{6}$$

Then $\{t_k : k = d + 1, \dots, n - 1\}$ are i.i.d. $\sim t(d - 1)$, the Student's t -distribution with d.f. (degrees of freedom) $d - 1$.

Proof. By the definition of the \mathbf{u}_i in (4), it is easy to verify that $\{\mathbf{u}_i : i = 1, \dots, n - 1\}$ are i.i.d. $\sim N_d(\mathbf{0}, \Sigma)$. By the above theorem, $\{\mathbf{y}_k : k = d + 1, \dots, n - 1\}$ are mutually independent and \mathbf{y}_k has a Pearson Type II distribution with a p.d.f. given by (2), which is a spherical distribution. For each fixed k , the random variable t_k only depends on \mathbf{y}_k . We write it as:

$$t_k = t_k(\mathbf{y}_k), \quad k = d + 1, \dots, n - 1.$$

The independence of $\{\mathbf{y}_k : k = d + 1, \dots, n - 1\}$ results in the independence of $\{t_k(\mathbf{y}_k) : k = d + 1, \dots, n - 1\}$. Note that for each fixed k , $t_k(\mathbf{y}_k)$ is scale invariant. That is, $t_k(a\mathbf{y}_k) = t_k(\mathbf{y}_k)$ holds for any constant $a > 0$. According to Theorem 2.22 of [22] (p. 51), we have:

$$t_k(\mathbf{y}_k) \stackrel{d}{=} t_k(\mathbf{z}_0),$$

where $\mathbf{z}_0 \sim N_d(\mathbf{0}, \mathbf{I}_d)$ (the d -dimensional standard normal), and the sign " $\stackrel{d}{=}$ " means that both sides of the equality have the same distribution. By the definitions of the t -distribution it is obvious that:

$$t_k(\mathbf{y}_k) \sim t(d - 1).$$

This completes the proof. \square

3. The RP-Based Chi-Square Test

The above corollary provides a way to construct a necessary test for MVN of the original i.i.d. sample $\{x_1, \dots, x_n\}$. Suppose that we want to test the hypothesis:

$$H_0 : \{x_1, \dots, x_n\} \text{ is a sample from } N_d(\boldsymbol{\mu}, \Sigma) \tag{7}$$

against the alternative hypothesis that $\{x_1, \dots, x_n\}$ is not a normal sample. This hypothesis can be transferred to testing:

$$H_0 : \{t_k : k = d + 1, \dots, n - 1\} \text{ is a sample from } t(d - 1) \tag{8}$$

versus the alternative that H_0 in (8) is not true. It is obvious that a test for (8) is a necessary one for (7); that is, if hypothesis (8) is rejected, hypothesis (7) is also rejected. However, if hypothesis (8) is not rejected, there is no guarantee for the truth of hypothesis (7).

To test hypothesis (8), we can employ the Pearson–Fisher test (simply called PF-test) for assessing if the set of i.i.d. t -type variates $\{t_k : k = d + 1, \dots, n - 1\}$ in (6) is from $t(d - 1)$. The traditional PF-test is facing with numerous choices of cells for grouping an i.i.d. sample. The choice of equiprobable cells is recommended in [25], which means that each cell is assigned an equal probability. If the number of cells, m , is pre-assigned, the probability for each cell is $1/m$. Under the null hypothesis, the transformed sample $\{t_k : k = d + 1, \dots, n - 1\}$ in (6) is from $t(d - 1)$. Then the endpoints $\{a_1 = -\infty, a_2, \dots, a_{m-1}, a_m = +\infty\}$ of all cells can be computed as follows.

$$\int_{a_i}^{a_{i+1}} f_t(x; d - 1) dx = \frac{1}{m}, \quad i = 1, \dots, m - 1,$$

where $f_t(x; d - 1)$ stands for the density function of the Student's t -distribution $t(d - 1)$.

Because the representative points (RP) ([18]; or called principal points in [26]) of a probability distribution have the property of minimizing some kind of quadratic loss function, we propose to employ the RP of the Student’s t -distribution $t(d - 1)$ as the endpoints of intervals for the Pearson–Fisher statistic. The t -RP are a set of points $\{R_1, \dots, R_m\}$ (for a selected number of points m) that minimize the quadratic loss function:

$$\phi(u_1, \dots, u_m) = \int_{-\infty}^{+\infty} \min_{1 \leq i \leq m} \{ (u_i - u)^2 \} f_t(u; d - 1) du, \tag{9}$$

$$\phi(R_1, \dots, R_m) = \min_{1 \leq i \leq m} \{ \phi(u_1, \dots, u_m) : -\infty < u_1 < \dots < u_m < +\infty \}.$$

Zhou and Wang [27] gave an algorithm for computing the t -RP $\{R_1, \dots, R_m\}$. Define the following intervals:

$$I_1 = \left(-\infty, \frac{R_1 + R_2}{2} \right), I_2 = \left[\frac{R_1 + R_2}{2}, \frac{R_2 + R_3}{2} \right), \dots, \tag{10}$$

$$I_{m-1} = \left[\frac{R_{m-2} + R_{m-1}}{2}, \frac{R_{m-1} + R_m}{2} \right), I_m = \left[\frac{R_{m-1} + R_m}{2}, +\infty \right)$$

and the probabilities:

$$p_i = \int_{I_i} f_t(x; d - 1) dx, \quad i = 1, \dots, m. \tag{11}$$

According to [21], $\{p_1, \dots, p_m\}$ can be considered as a set of “representative probabilities” for the Student’s t density function $f_t(\cdot; d - 1)$. The χ^2 -statistic for testing hypothesis (8) is computed by:

$$\chi_R^2 = \sum_{i=1}^m \frac{(n_i - np_i)^2}{np_i}, \tag{12}$$

where n_i is the frequency of the transformed sample points $\{t_k : k = d + 1, \dots, n - 1\}$ given by (6) that are located in the interval I_i in (10). It is known that $\chi_R^2 \rightarrow \chi^2(m - 1)$ ($n \rightarrow \infty$) in the distribution under some regular conditions. The p -value for testing (8) is computed by:

$$P(\chi_R^2, \nu) = K \int_{\chi_R^2}^{\infty} z^{\frac{\nu}{2}-1} \exp\left(-\frac{z}{2}\right) dz, \text{ with } \nu = m - 1, K = \left[2^{\frac{\nu}{2}} \Gamma\left(\frac{\nu}{2}\right)\right]^{-1}. \tag{13}$$

4. A Monte Carlo Study and an Illustrative Example

In order to compare the χ_R^2 -test (12) under the “representative probabilities” $\{p_1, \dots, p_m\}$ in (11) with the traditional chi-squared test, we choose the equiprobable cells for computing the traditional chi-square test. For a selected number of representative points m , define the interval endpoints:

$$\begin{aligned}
 a_1 & \text{ satisfies } P(\chi^2(m - 1) < a_1) = \frac{1}{m}; \\
 a_2 & \text{ satisfies } P(a_1 < \chi^2(m - 1) < a_2) = \frac{1}{m}; \\
 & \vdots \\
 a_{m-1} & \text{ satisfies } P(a_{m-2} < \chi^2(m - 1) < a_{m-1}) = \frac{1}{m}; \\
 a_m & \text{ satisfies } P(\chi^2(m - 1) > a_m) = \frac{1}{m}.
 \end{aligned} \tag{14}$$

Denote the traditional chi-square test based on the interval endpoints (14) by χ_T^2 :

$$\chi_T^2 = \sum_{i=1}^m \frac{(N_i - n/m)^2}{n/m}, \tag{15}$$

which is also an approximate $\chi^2(m - 1)$ for a large sample size, where N_i stands for the frequency of the observed sample points that are located in the intervals defined by the endpoints (14).

4.1. A Comparison between Empirical Type I Error Rates

Because the chi-square test based on the transformed sample points $\{t_k : k = d + 1, \dots, n - 1\}$ given by (6) is affine invariant under any nonsingular linear transformation of the original i.i.d. sample $\{x_1, \dots, x_n\}$, we only need to generate samples from a d -dimensional standard normal $N_d(\mathbf{0}, \mathbf{I}_d)$ (\mathbf{I}_d stands for the $d \times d$ identity matrix). The simulation results under 2000 replications for each case are summarized in Tables 1–3 for significance levels $\alpha = 0.01, 0.05,$ and $0.10,$ respectively. It can be roughly concluded that both the traditional chi-square statistic χ_T^2 and the RP chi-square statistic χ_R^2 show reasonable control of type I error rates under different choices of the number of cells and relatively large sample sizes. Because the Pearson–Fisher chi-square test is an approximate test for goodness-of-fit under its asymptotic null distribution, we do not show its small-sample empirical performance.

Table 1. Empirical type I error rates ($\alpha = 0.01$).

n	m	χ^2	$d = 5$	$d = 10$	$d = 15$	$d = 20$	
$n = 50$	$m = 5$	χ_R^2	0.0155	0.0110	0.0075	0.0070	
		χ_T^2	0.0075	0.0070	0.0065	0.0075	
	$m = 10$	χ_R^2	0.0320	0.0145	0.0120	0.0125	
		χ_T^2	0.0090	0.0105	0.0075	0.0085	
	$m = 15$	χ_R^2	0.0525	0.0285	0.0210	0.0235	
		χ_T^2	0.0080	0.0105	0.0115	0.0090	
		$m = 20$	χ_R^2	0.0835	0.0465	0.0255	0.0340
			χ_T^2	0.0115	0.0110	0.0090	0.0130
$n = 100$	$m = 5$	χ_R^2	0.0130	0.0100	0.0120	0.0105	
		χ_T^2	0.0070	0.0125	0.0105	0.0150	
	$m = 10$	χ_R^2	0.0350	0.0170	0.0130	0.0125	
		χ_T^2	0.0115	0.0105	0.0110	0.0075	
	$m = 15$	χ_R^2	0.0810	0.0230	0.0180	0.0185	
		χ_T^2	0.0105	0.0125	0.0120	0.0090	
		$m = 20$	χ_R^2	0.0370	0.0225	0.0245	0.0190
			χ_T^2	0.0110	0.0135	0.0110	0.0100

Table 1. Cont.

n	m	χ^2	$d = 5$	$d = 10$	$d = 15$	$d = 20$	
200	5	χ_R^2	0.0075	0.0045	0.0060	0.0065	
		χ_T^2	0.0105	0.0090	0.0075	0.0085	
	10	χ_R^2	0.0240	0.0150	0.0130	0.0130	
		χ_T^2	0.0115	0.0135	0.0110	0.0105	
	15	χ_R^2	0.0430	0.0125	0.0155	0.0155	
		χ_T^2	0.0120	0.0080	0.0125	0.0110	
	20	χ_R^2	0.0455	0.0145	0.0150	0.0170	
		χ_T^2	0.0075	0.0120	0.0095	0.0120	
	400	5	χ_R^2	0.0110	0.0155	0.0140	0.0085
			χ_T^2	0.0105	0.0080	0.0115	0.0090
10		χ_R^2	0.0185	0.0105	0.0140	0.0105	
		χ_T^2	0.0070	0.0090	0.0150	0.0095	
15		χ_R^2	0.0230	0.0140	0.0155	0.0145	
		χ_T^2	0.0140	0.0105	0.0115	0.0125	
20		χ_R^2	0.0570	0.0140	0.0120	0.0175	
		χ_T^2	0.0085	0.0150	0.0075	0.0130	

Table 2. Empirical type I error rates ($\alpha = 0.05$).

n	m	χ^2	$d = 5$	$d = 10$	$d = 15$	$d = 20$	
50	5	χ_R^2	0.0430	0.0550	0.0460	0.0450	
		χ_T^2	0.0400	0.0555	0.0390	0.0460	
	10	χ_R^2	0.0730	0.0505	0.0565	0.0450	
		χ_T^2	0.0350	0.0440	0.0480	0.0340	
	15	χ_R^2	0.0660	0.0750	0.0725	0.0675	
		χ_T^2	0.0445	0.0575	0.0455	0.0545	
	20	χ_R^2	0.0930	0.1115	0.0790	0.0710	
		χ_T^2	0.0490	0.0570	0.0495	0.0365	
	100	5	χ_R^2	0.0555	0.0435	0.0465	0.0410
			χ_T^2	0.0470	0.0485	0.0540	0.0505
10		χ_R^2	0.0750	0.0530	0.0500	0.0560	
		χ_T^2	0.0595	0.0510	0.0515	0.0480	
15		χ_R^2	0.0970	0.0685	0.0610	0.0560	
		χ_T^2	0.0530	0.0550	0.0540	0.0525	
20		χ_R^2	0.0755	0.0735	0.0695	0.0665	
		χ_T^2	0.0540	0.0465	0.0545	0.0420	

Table 2. Cont.

n	m	χ^2	$d = 5$	$d = 10$	$d = 15$	$d = 20$	
$n = 200$	$m = 5$	χ_R^2	0.0460	0.0485	0.0450	0.0465	
		χ_T^2	0.0565	0.0495	0.0465	0.0505	
	$m = 10$	χ_R^2	0.0580	0.0530	0.0400	0.0495	
		χ_T^2	0.0490	0.0520	0.0425	0.0530	
	$m = 15$	χ_R^2	0.1135	0.0625	0.0550	0.0565	
		χ_T^2	0.0530	0.0480	0.0485	0.0470	
	$m = 20$	χ_R^2	0.0715	0.0635	0.0560	0.0595	
		χ_T^2	0.0485	0.0450	0.0600	0.0505	
	$n = 400$	$m = 5$	χ_R^2	0.0550	0.0470	0.0495	0.0475
			χ_T^2	0.0485	0.0520	0.0450	0.0375
$m = 10$		χ_R^2	0.0590	0.0525	0.0515	0.0475	
		χ_T^2	0.0545	0.0565	0.0510	0.0460	
$m = 15$		χ_R^2	0.0740	0.0460	0.0475	0.0495	
		χ_T^2	0.0465	0.0535	0.0460	0.0520	
$m = 20$		χ_R^2	0.0880	0.0670	0.0580	0.0505	
		χ_T^2	0.0515	0.0470	0.0495	0.0475	

Table 3. Empirical type I error rates ($\alpha = 0.10$).

n	m	χ^2	$d = 5$	$d = 10$	$d = 15$	$d = 20$	
$n = 50$	$m = 5$	χ_R^2	0.0835	0.0885	0.0915	0.0880	
		χ_T^2	0.0855	0.1005	0.0845	0.0865	
	$m = 10$	χ_R^2	0.1325	0.1045	0.0970	0.0985	
		χ_T^2	0.0895	0.0885	0.0930	0.0960	
	$m = 15$	χ_R^2	0.0870	0.1065	0.1055	0.1130	
		χ_T^2	0.0940	0.0890	0.0985	0.1080	
	$m = 20$	χ_R^2	0.1155	0.1375	0.1330	0.1290	
		χ_T^2	0.0875	0.0815	0.0735	0.0835	
	$n = 100$	$m = 5$	χ_R^2	0.0980	0.0965	0.0965	0.1010
			χ_T^2	0.1155	0.1085	0.1035	0.0885
$m = 10$		χ_R^2	0.1065	0.1010	0.0950	0.0940	
		χ_T^2	0.0985	0.0955	0.0905	0.0905	
$m = 15$		χ_R^2	0.1130	0.0895	0.1100	0.1070	
		χ_T^2	0.0985	0.0930	0.1035	0.1010	
$m = 20$		χ_R^2	0.1135	0.1000	0.1140	0.1095	
		χ_T^2	0.1015	0.0980	0.1005	0.0885	

Table 3. Cont.

<i>n</i>	<i>m</i>	χ^2	<i>d</i> = 5	<i>d</i> = 10	<i>d</i> = 15	<i>d</i> = 20	
<i>n</i> = 200	<i>m</i> = 5	χ_R^2	0.0965	0.0980	0.0965	0.1050	
		χ_T^2	0.0895	0.0880	0.1055	0.0940	
	<i>m</i> = 10	χ_R^2	0.1035	0.1040	0.0925	0.0965	
		χ_T^2	0.0980	0.0970	0.1010	0.1040	
	<i>m</i> = 15	χ_R^2	0.1575	0.1040	0.0865	0.1015	
		χ_T^2	0.1040	0.0980	0.0940	0.0930	
	<i>m</i> = 20	χ_R^2	0.1005	0.1050	0.0965	0.1035	
		χ_T^2	0.0940	0.0965	0.1110	0.0990	
	<i>n</i> = 400	<i>m</i> = 5	χ_R^2	0.0865	0.0995	0.0940	0.0930
			χ_T^2	0.1010	0.0945	0.0975	0.0955
<i>m</i> = 10		χ_R^2	0.1020	0.1055	0.0975	0.1095	
		χ_T^2	0.0980	0.1000	0.1080	0.0975	
<i>m</i> = 15		χ_R^2	0.1045	0.0950	0.0970	0.0930	
		χ_T^2	0.0900	0.1025	0.1050	0.1005	
<i>m</i> = 20		χ_R^2	0.1210	0.0965	0.1085	0.1025	
		χ_T^2	0.1035	0.0910	0.1075	0.1055	

4.2. A Simple Power Comparison

To show the benefit of employing the RP-idea for grouping cells for the Pearson–Fisher chi-square test, we carry out a simple Monte Carlo study by selecting the following alternative distributions, which consist of three types of distributions (symmetric about the origin; skewed distributions; distributions with normal marginals):

- (1) [symmetric] The multivariate Cauchy distribution ([22]) has a density function of the form:

$$f_c(\|x\|) = C_1 \left(1 + \frac{\|x\|^2}{m} \right)^{-\frac{d+1}{2}},$$

where “ $\| \cdot \|$ ” stands for the Euclidean norm of a vector, C_1 is a normalizing constant depending on the dimension d .

- (2) [symmetric] The β -generalized normal distribution $N_d(0, I_d, 1/2)$ with $\beta = 1/4$ has a density function of the form by ([28]):

$$f(x_1, \dots, x_d) = \frac{\beta^d r^{d/\beta}}{2^d \Gamma^d(1/\beta)} \cdot \exp \left\{ -r \sum_{i=1}^d |x_i|^\beta \right\}, \quad (x_1, \dots, x_d)' \in R^d,$$

where $r > 0$ is a parameter. Let $r = 1/2$ in the simulation and denote it by β -normal.

- (3) [symmetric] Multivariate double Weibull distribution consisting of i.i.d. univariate double Weibull distributions ([29]), its density function is given by:

$$f_d(x) = \left(\frac{\alpha}{2\beta} \right)^d \prod_{i=1}^d \left(\frac{|x_i|}{\beta} \right)^{\alpha-1} \exp \left\{ -\sum_{i=1}^d \left(\frac{|x_i|}{\beta} \right)^\alpha \right\}, \quad x = (x_1, \dots, x_d)' \in R^d,$$

where α and β are the shape parameter and scale parameter, respectively. Let $\alpha = 1/2$ and $\beta = 1$ in the simulation.

- (4) [skewed] The shifted i.i.d. $\chi^2(1)$ with i.i.d. marginals, each marginal has the same distribution as that of the random variable $Y = X - E(X)$, where $X \sim \chi^2(1)$, the univariate chi-square distribution with 1 degree of freedom and $E(X) = 1$.

- (5) [skewed] The shifted i.i.d. $\exp(1)$ with i.i.d. marginals, each marginal has the same distribution as that of the random variable $Y = X - E(X)$, where $X \sim \exp(1)$, the univariate exponential distribution.
- (6) [skewed] The shifted i.i.d. F -distribution with i.i.d. marginals $F(4,3)$ with i.i.d. marginals, $Y = X - E(X)$, where $X \sim F(4,3)$, $E(X) = E[F(4,3)] = 3$.
- (7) [A distribution with normal marginals] The distribution $N(0,1) + \chi^2(2)$ consists of i.i.d. $[d/2]$ normal $N(0,1)$ marginals and $d - [d/2]$ i.i.d. $\chi^2(2) - 2$ marginals, where $[d/2]$ stands for the integer part of $d/2$.
- (8) [A distribution with normal marginals] The distribution $N(0,1) + \exp(1)$ consists of i.i.d. $[d/2]$ normal $N(0,1)$ marginals and $d - [d/2]$ i.i.d. $\exp(1) - 1$ marginals.
- (9) [A distribution with normal marginals] The distribution $N(0,1) + F(4,3)$ consists of i.i.d. $[d/2]$ normal $N(0,1)$ marginals and $d - [d/2]$ i.i.d. $F(4,3) - 3$ marginals.

For each of these alternative distributions, choose the sample size $n = 50, 70, \dots, 400$. Plot the power values versus the sample size n for both statistics χ^2_R in (12) and χ^2_T in (15).

A visual observation on the following Figures 1–9 immediately leads to the following two empirical conclusions:

- (1) The RP chi-square test χ^2_R is comparable to (or slightly better than) the traditional test χ^2_T for symmetric alternative distributions;
- (2) The RP chi-square test χ^2_R is able to improve the traditional test χ^2_T significantly for both skewed and normal+skewed alternative distributions.

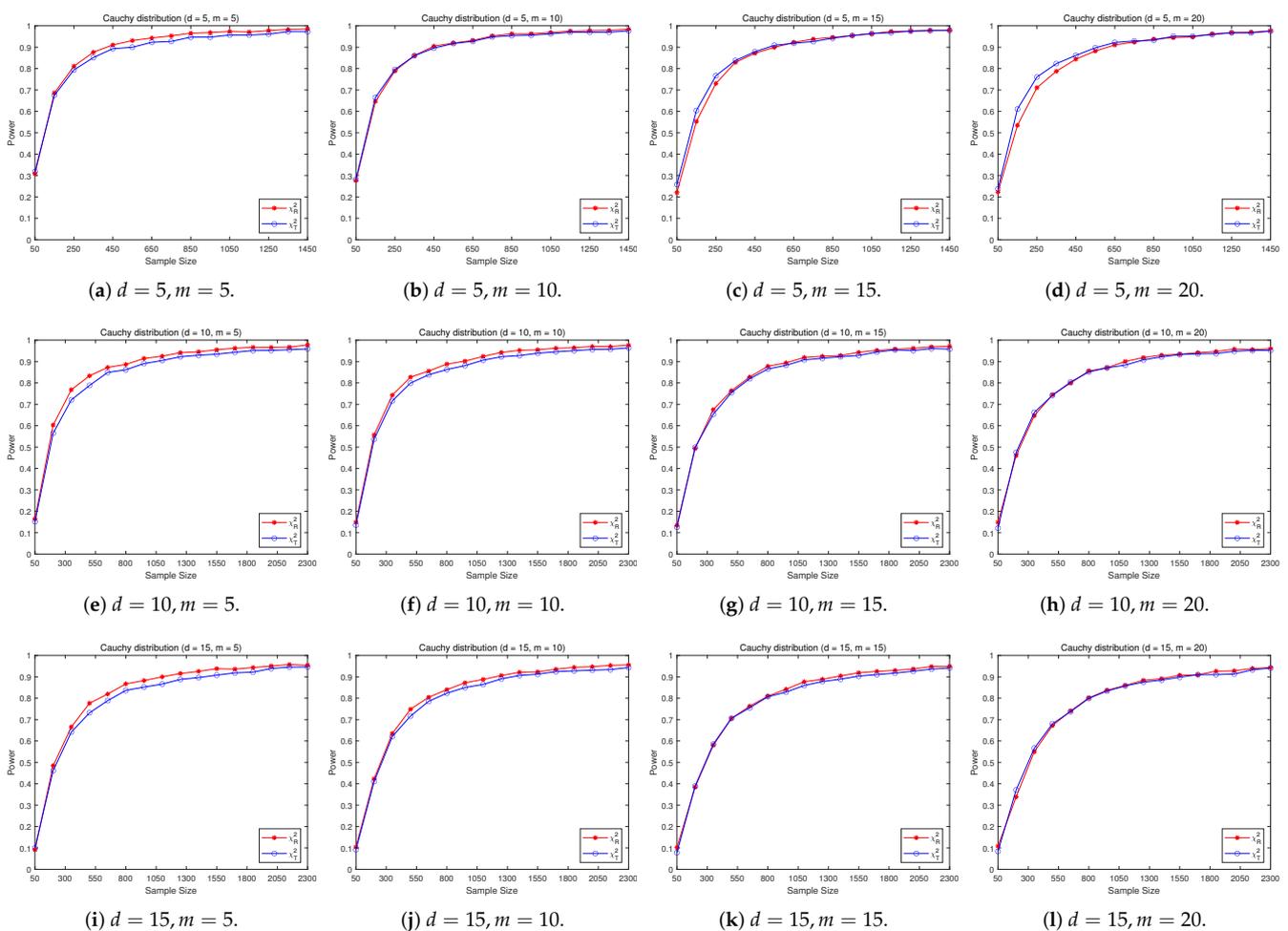


Figure 1. Cont.

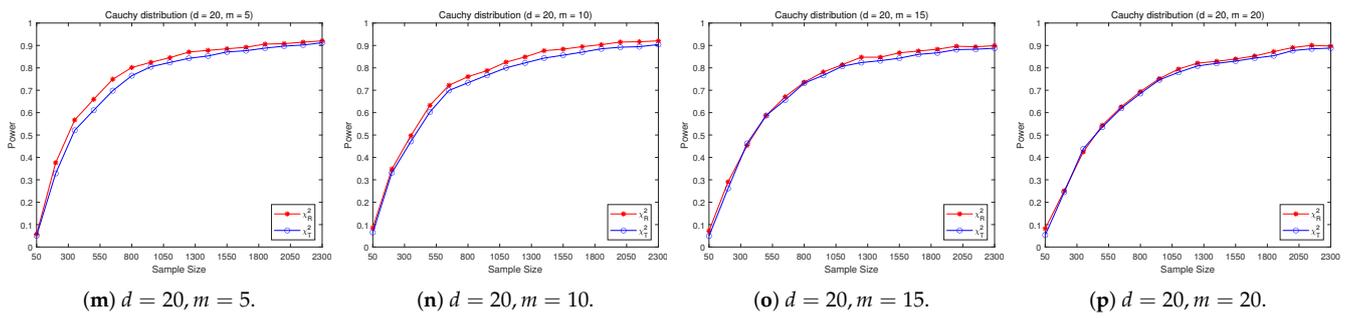


Figure 1. The Cauchy distribution. Red line for RP-chi-square, blue line for traditional chi-square.

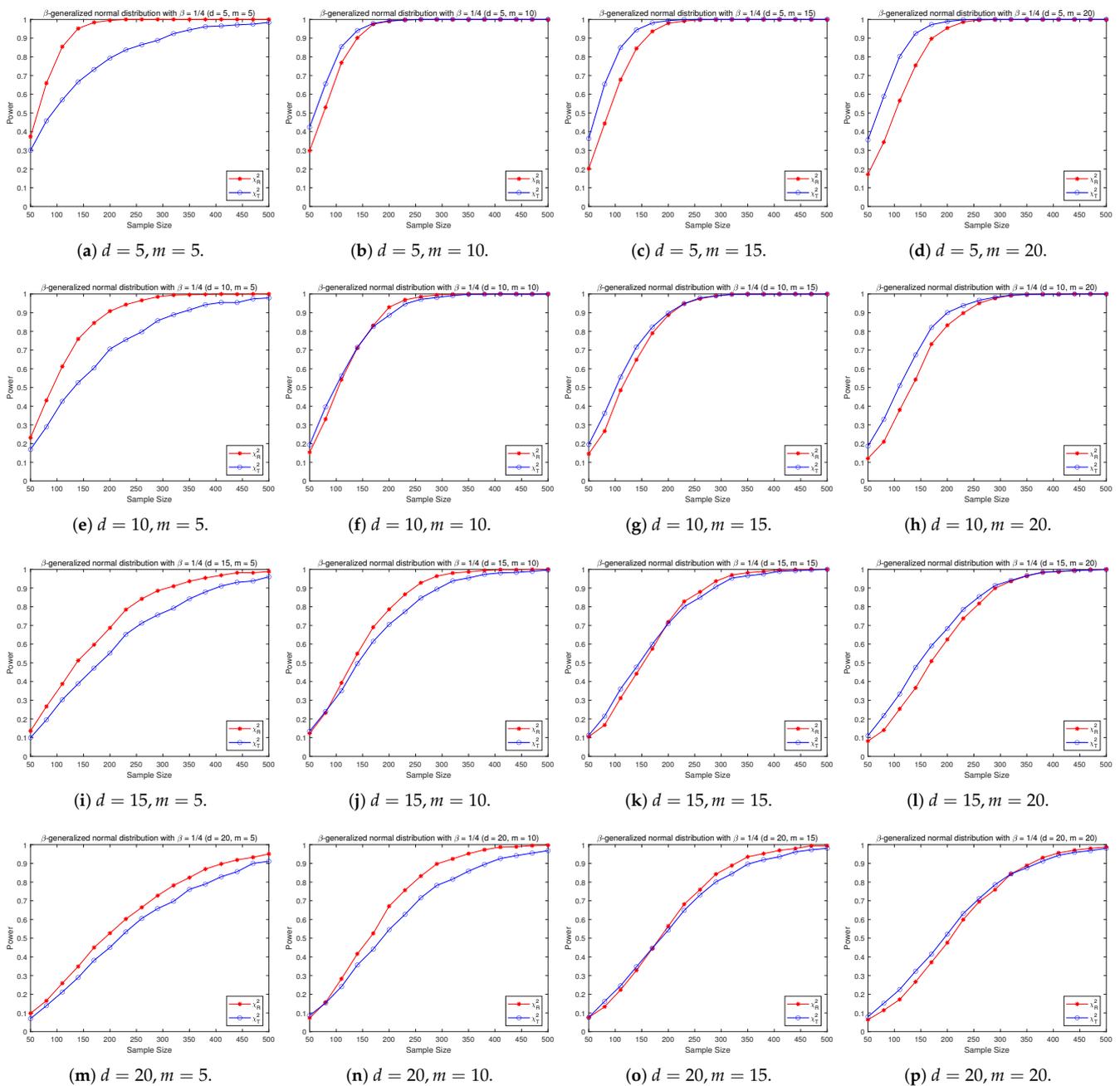


Figure 2. The β -generalized normal distribution with $\beta = 1/4$. Red line for RP-chi-square, blue line for traditional chi-square.

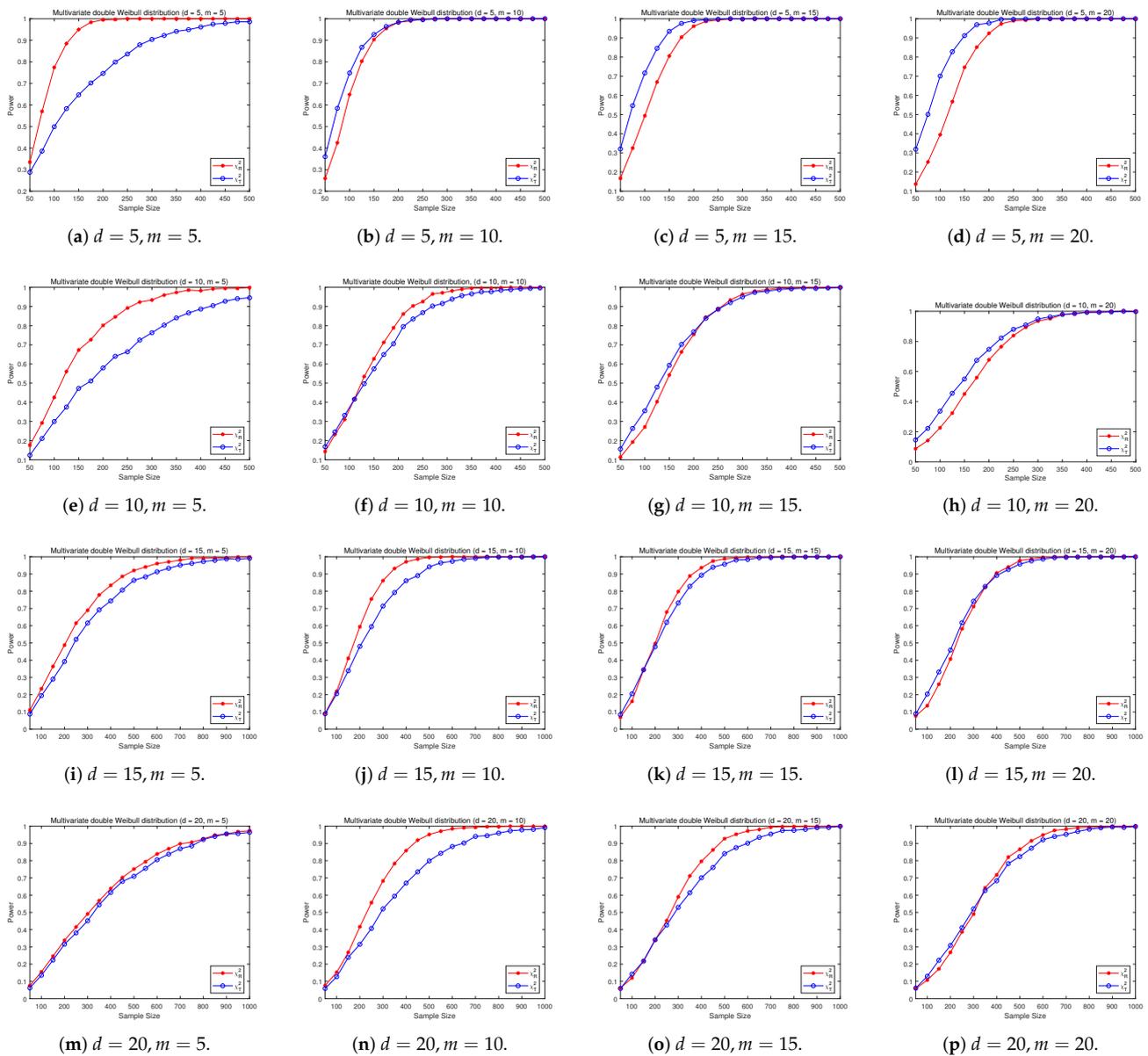


Figure 3. The double Weibull distribution. Red line for RP-chi-square, blue line for traditional chi-square.

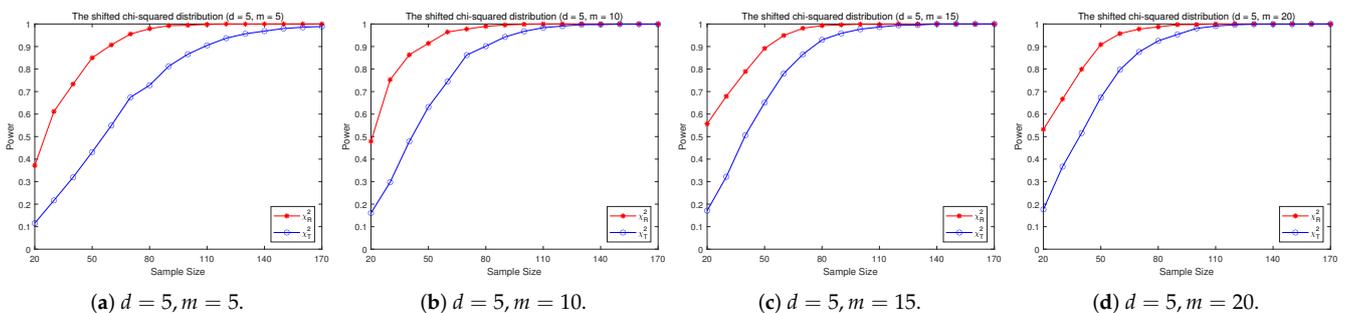


Figure 4. Cont.

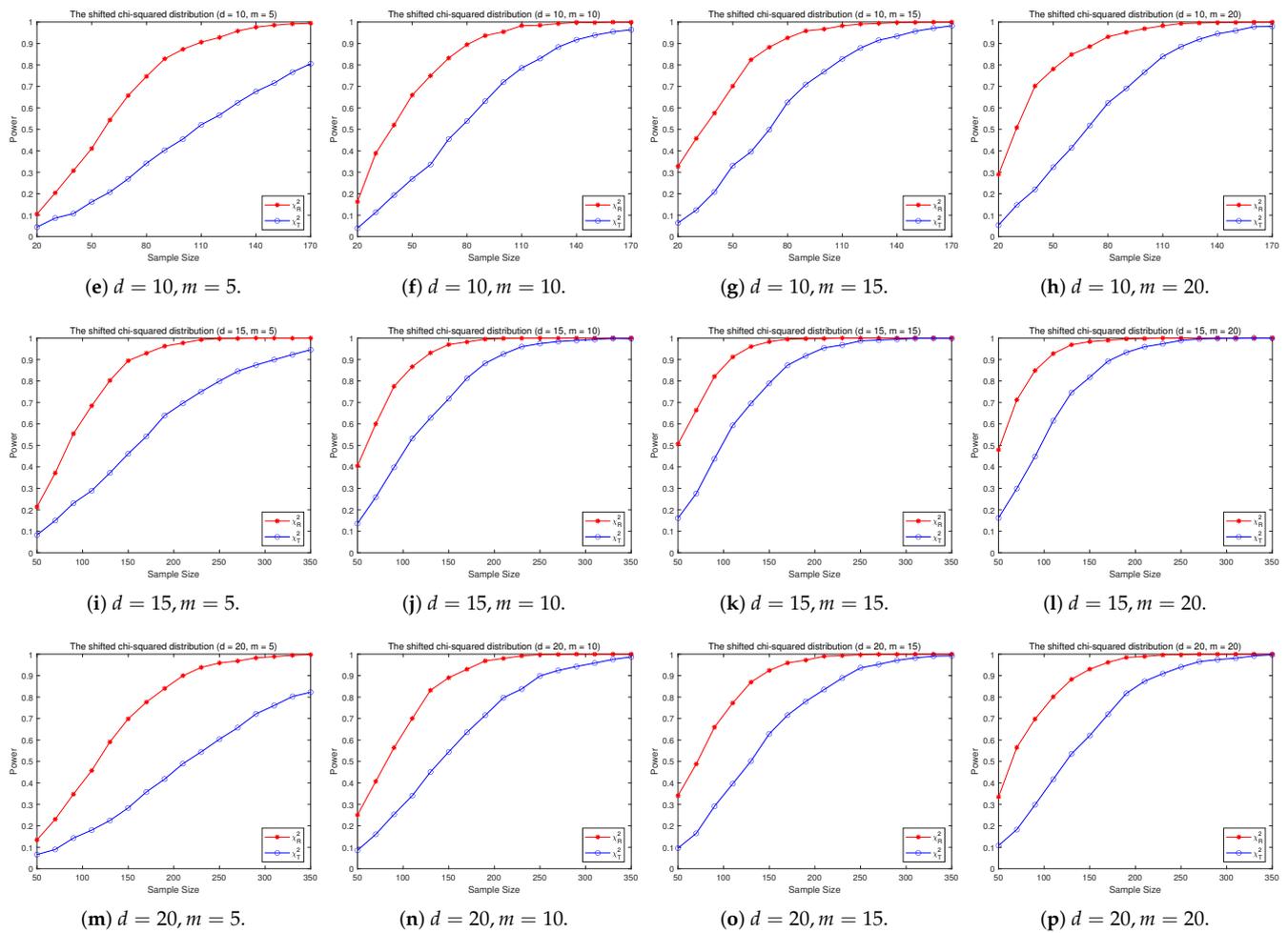


Figure 4. The shifted chi-square distribution with degree of freedom 1. Red line for RP-chi-square, blue line for traditional chi-square.

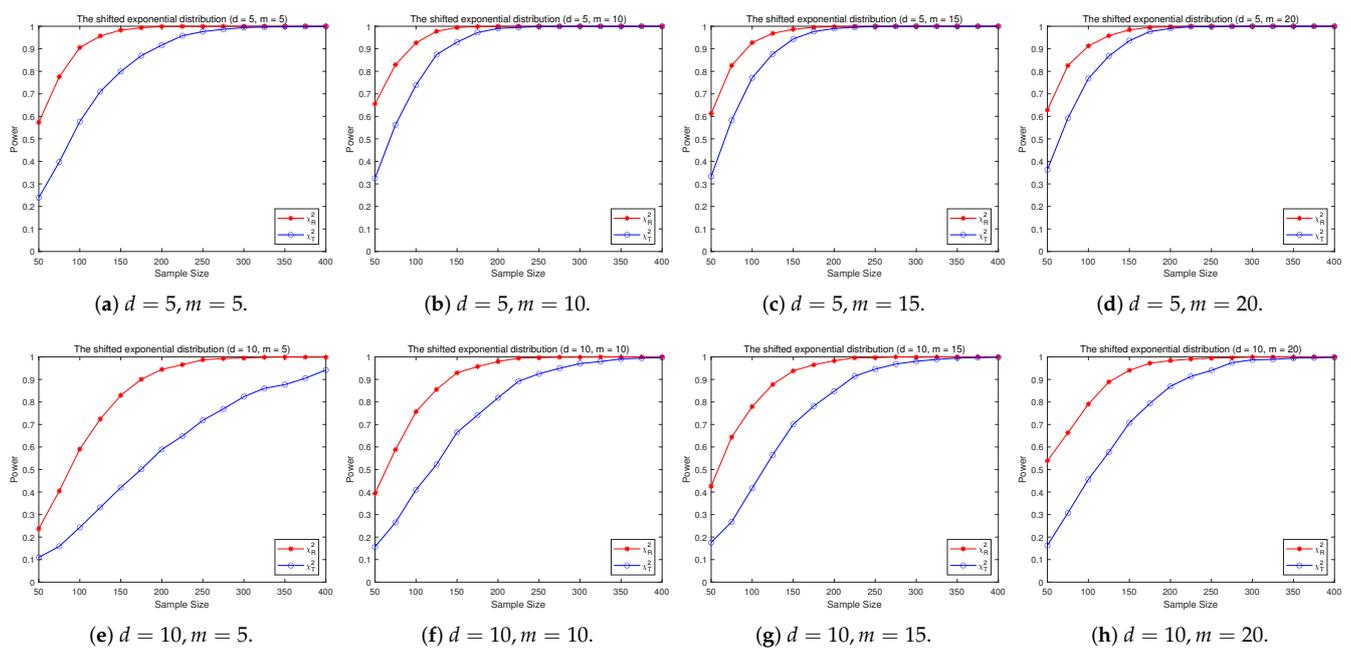


Figure 5. Cont.

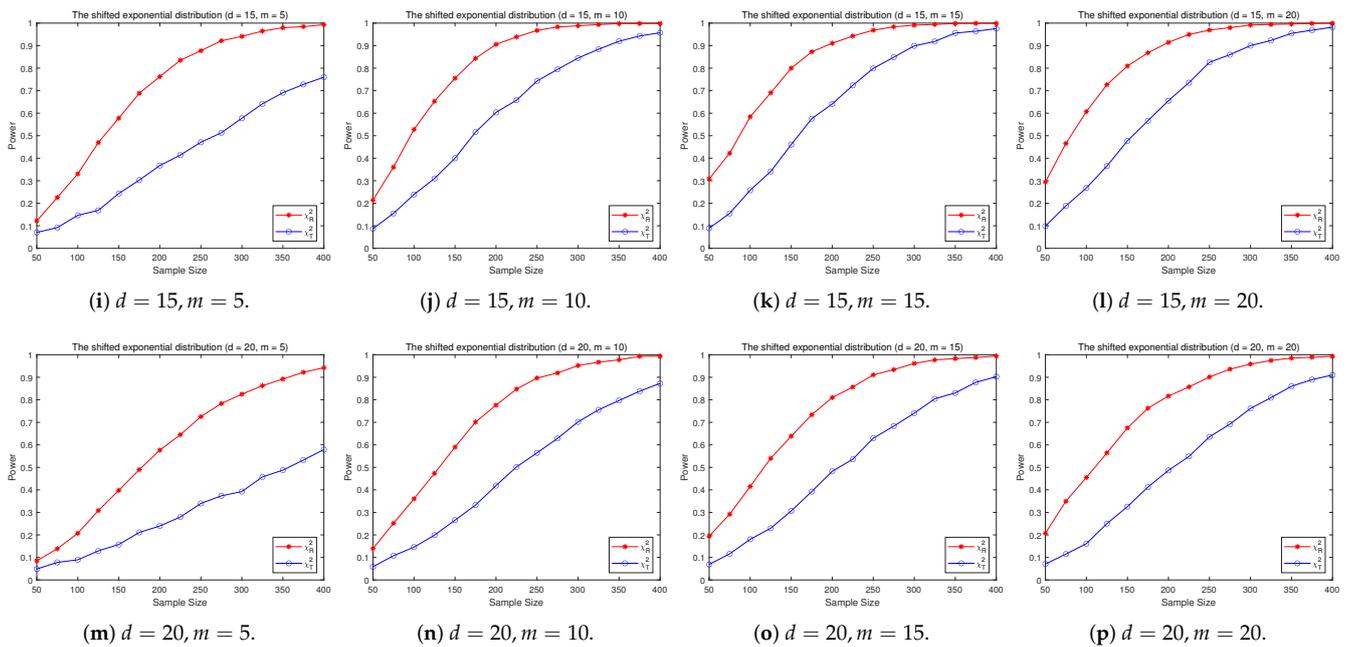


Figure 5. The shifted exponential distribution $\exp(1)$. Red line for RP-chi-square, blue line for traditional chi-square.

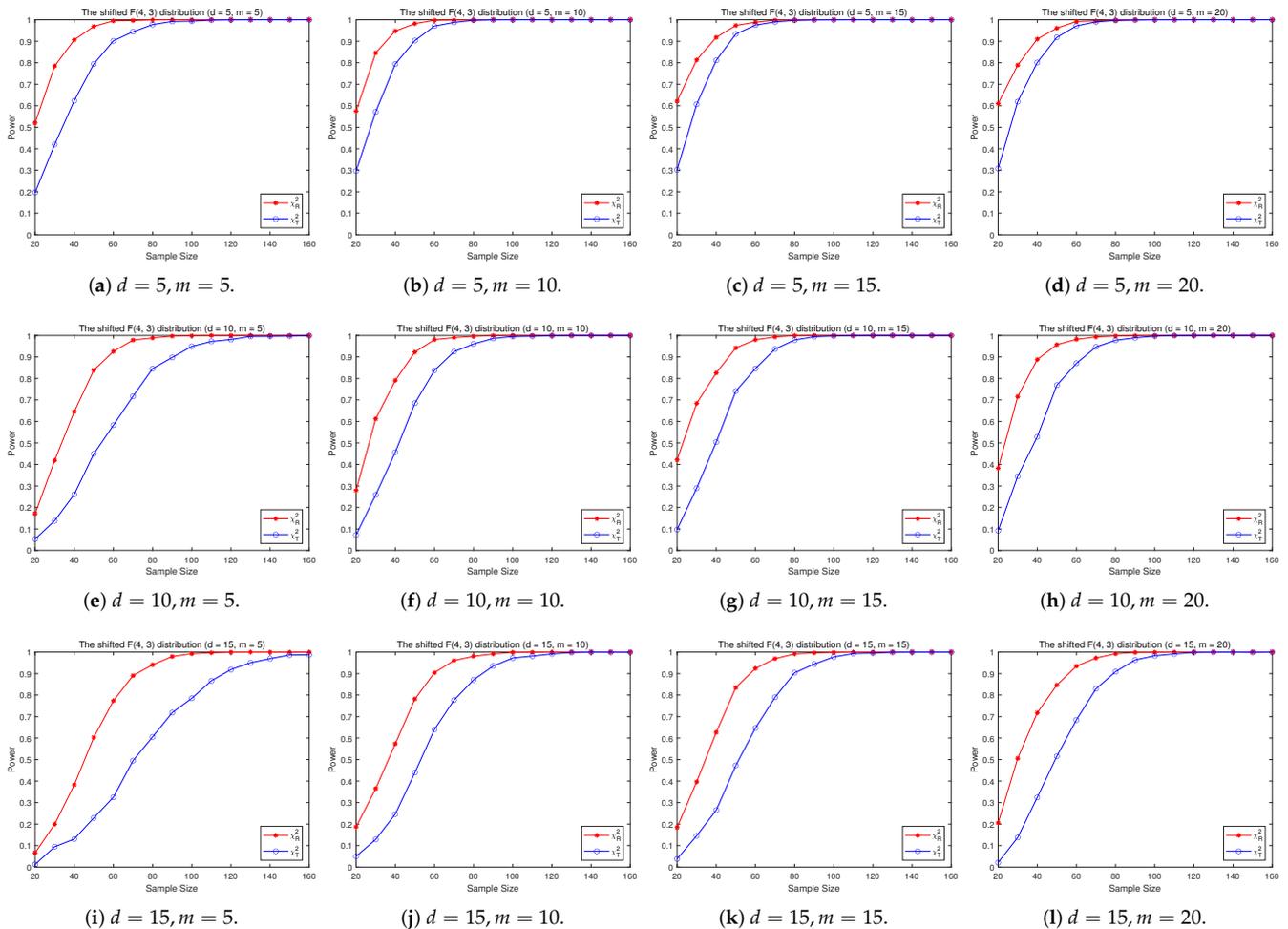


Figure 6. Cont.

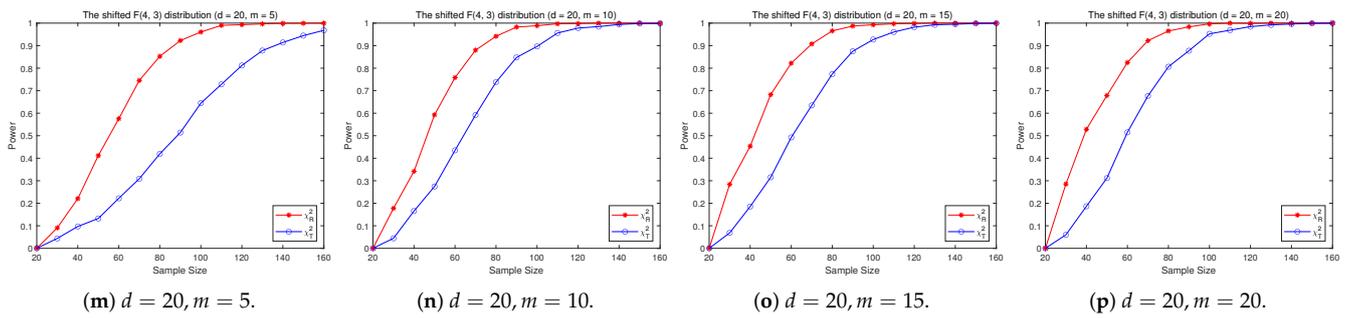


Figure 6. The shifted F distribution $F(4,3)$. Red line for RP-chi-square, blue line for traditional chi-square.

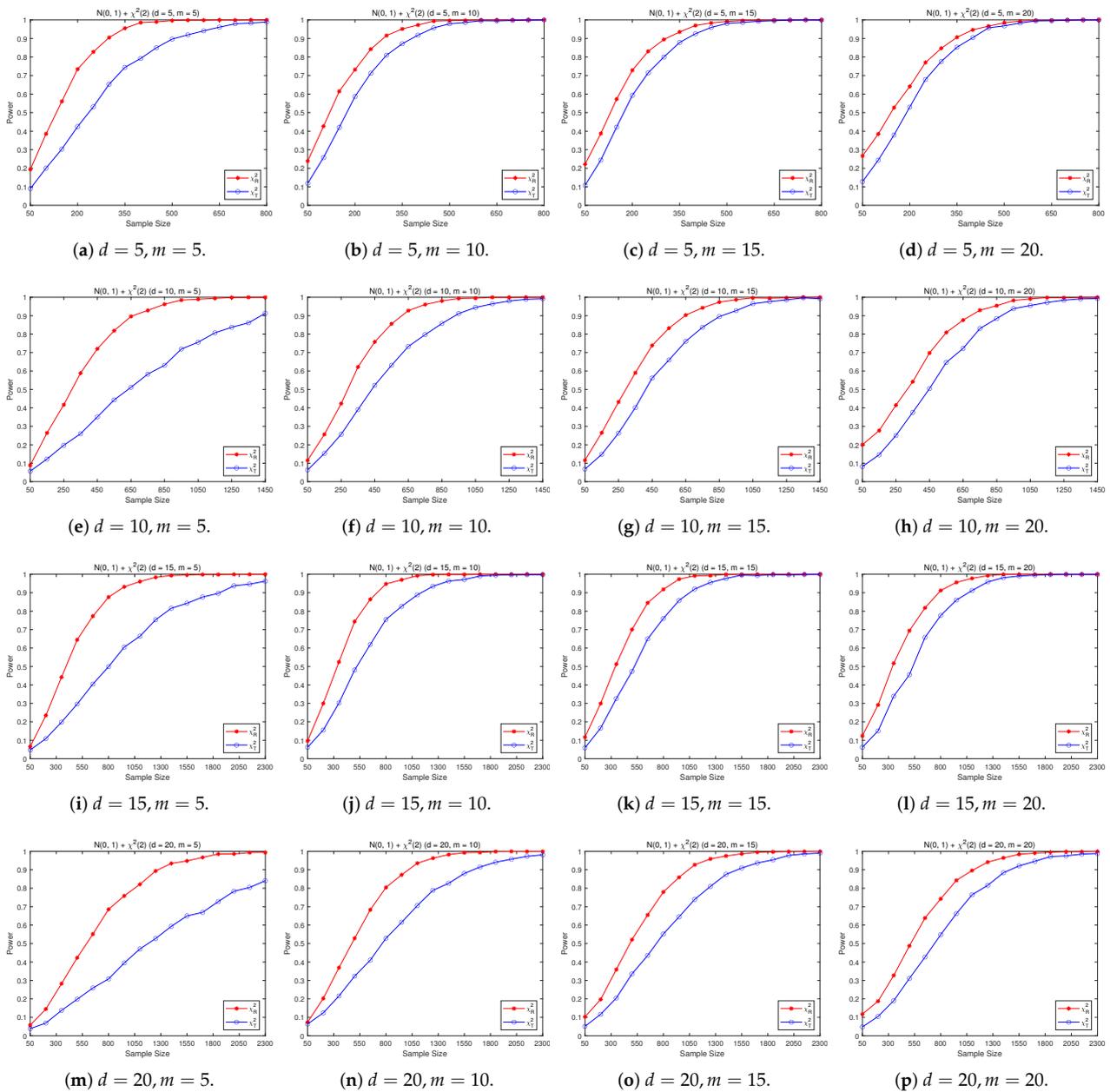


Figure 7. The distribution $N(0, 1) + \chi^2(2)$. Red line for RP-chi-square, blue line for traditional chi-square.

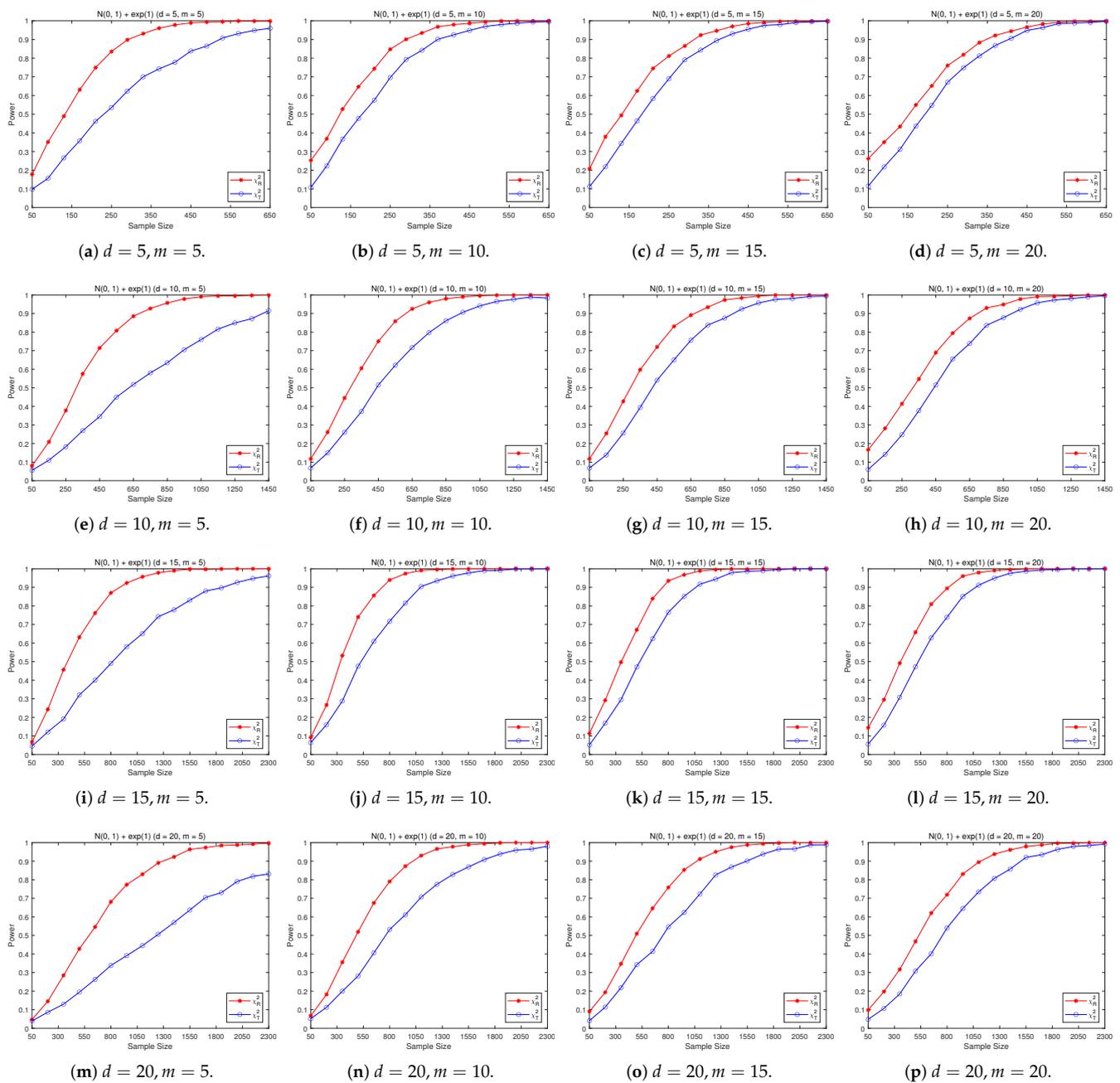


Figure 8. The distribution $N(0, 1) + \exp(1)$. Red line for RP-chi-square, blue line for traditional chi-square.

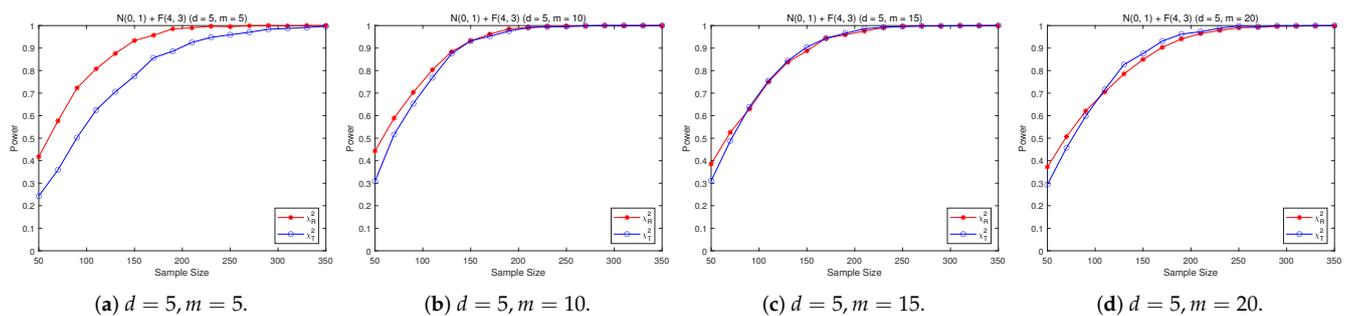


Figure 9. Cont.

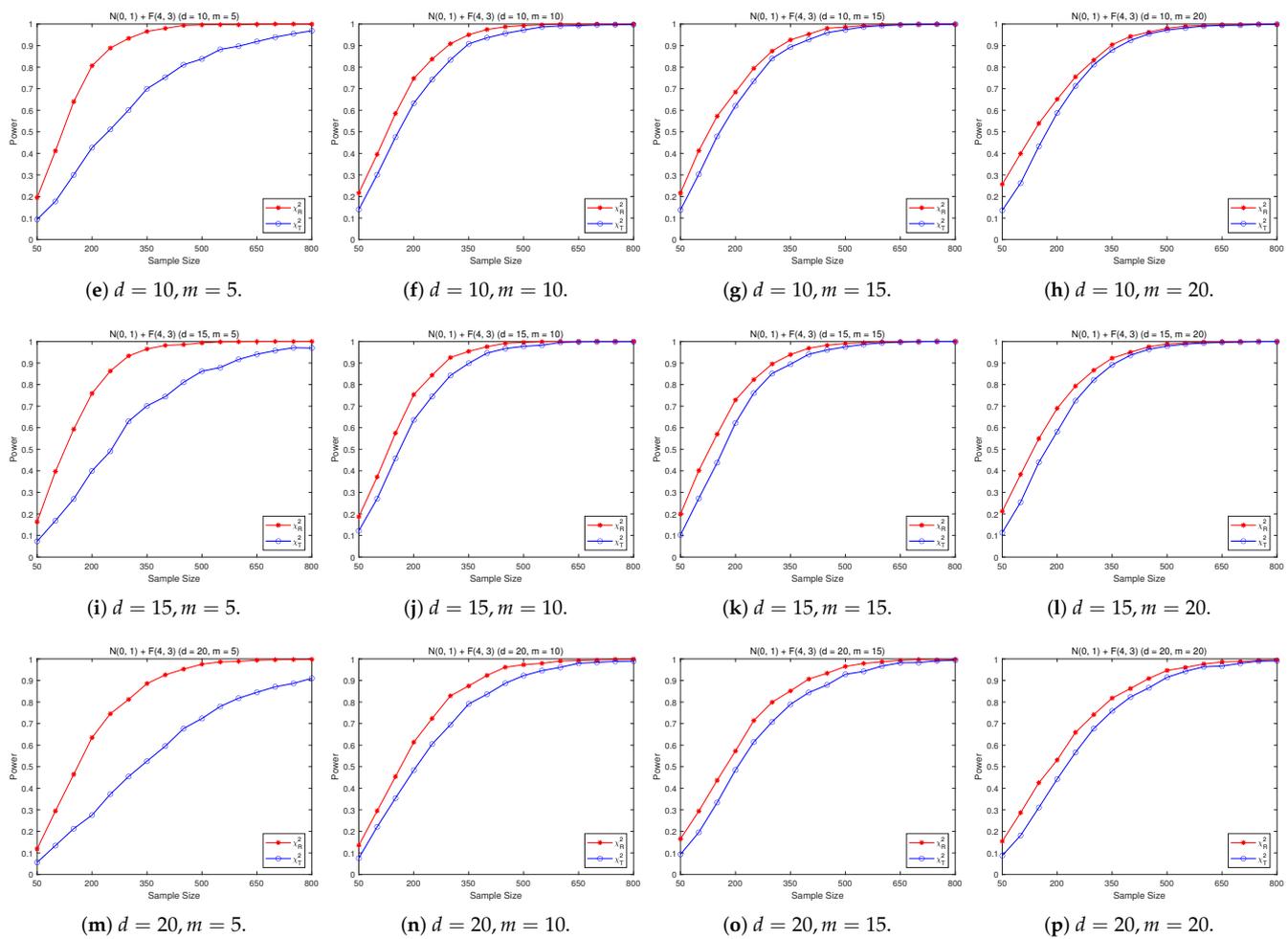


Figure 9. The distribution $N(0, 1) + F(4, 3)$. Red line for RP-chi-square, blue line for traditional chi-square.

4.3. An Illustrative Example

The data were collected from the measurements of various body circumferences of 252 men and were supplied by Dr. A. Garth Fisher. A complete list of the data can be found on the website: <http://lib.stat.cmu.edu/datasets/> accessed on 23 October 2022 (bodyfat). In our analysis of the data, the variables are listed as follows:

- X_1 : density determined from underwater weighing
- X_2 : percent body fat
- X_3 : age (years)
- X_4 : weight (lbs)
- X_5 : height (inches)
- X_6 : neck circumference (cm)
- X_7 : chest circumference (cm)
- X_8 : abdomen 2 circumference (cm)
- X_9 : hip circumference (cm)
- X_{10} : thigh circumference (cm)
- X_{11} : knee circumference (cm)
- X_{12} : ankle circumference (cm)
- X_{13} : biceps (extended) circumference (cm)
- X_{14} : forearm circumference (cm)
- X_{15} : wrist circumference (cm)

The observations from the plotting method for detecting non-multinormality in [30] are summarized as follows.

- (1) The 5-dimensional random vector $(X_4, \dots, X_8)'$ can be approximately considered as 5-dimensional normal;
- (2) The 5-dimensional random vector $(X_4, X_5, X_8, X_9, X_{10})'$ shows evidence of non-MVN;
- (3) The 5-dimensional random vector $(X_4, X_5, X_9, X_{10}, X_{11})'$ shows evidence of non-MVN;

- (4) The 5-dimensional random vector $(X_8, \dots, X_{12})'$ shows evidence of non-MVN;
- (5) The 10-dimensional random vector $(X_4, \dots, X_{13})'$ shows evidence of non-MVN;
- (6) The 10-dimensional random vector $(X_5, \dots, X_{14})'$ can be approximately considered as 10-dimensional normal;
- (7) The 10-dimensional random vector $(X_6, \dots, X_{15})'$ shows evidence of non-MVN;
- (8) The 10-dimensional random vector $(X_4, X_6, \dots, X_{11}, X_{13}, X_{14}, X_{15})'$ can be approximately considered as 10-dimensional normal.

The chi-square analysis by the two statistics χ_R^2 and χ_T^2 with three different choices of the number of cells ($m = 5, 10$, and 15) is presented in the following Table 4. We obtain mostly consistent results with those from [30] under the significance level 0.05:

- (1) The 5-dimensional random vector $(X_4, \dots, X_8)'$ can be approximately considered as 5-dimensional normal by χ_R^2 for $m = 10$ and $m = 15$, and by χ_T^2 for $m = 15$;
- (2) The 5-dimensional random vector $(X_4, X_5, X_8, X_9, X_{10})'$ shows evidence of non-MVN by χ_R^2 for $m = 15$. χ_T^2 fails to detect the non-MVN for all three choices of m ;
- (3) The 5-dimensional random vector $(X_4, X_5, X_9, X_{10}, X_{11})'$ shows evidence of non-MVN by both χ_R^2 and χ_T^2 for all three choices of m ;
- (4) The 5-dimensional random vector $(X_8, \dots, X_{12})'$ shows evidence of non-MVN by χ_R^2 for $m = 5$ and $m = 10$, and by χ_T^2 for $m = 10$ and $m = 15$;
- (5) The 10-dimensional random vector $(X_4, \dots, X_{13})'$ shows evidence of non-MVN by χ_R^2 for $m = 15$. χ_T^2 fails to detect the non-MVN for all three choices of m ;
- (6) The 10-dimensional random vector $(X_5, \dots, X_{14})'$ shows evidence of non-MVN by χ_R^2 for $m = 10$ and $m = 15$. χ_T^2 fails to detect the non-MVN for all three choices of m ;
- (7) The 10-dimensional random vector $(X_6, \dots, X_{15})'$ can be approximately considered as 10-dimensional normal by both χ_R^2 and χ_T^2 for all three choices of m ;
- (8) The 10-dimensional random vector $(X_4, X_6, \dots, X_{11}, X_{13}, X_{14}, X_{15})'$ can be approximately considered as 10-dimensional normal by both χ_R^2 and χ_T^2 .

The inconsistency from two different tests in the above conclusions (5) and (6) may come from the general drawback for the chi-square test in optimal cell selection. Some good properties were discussed for equiprobable cell selection compared to some random cell selection ([31]) under some conditions. RP-cell selection is based on the idea of minimizing some kind of expected quadratic loss when quantizing a continuous probability distribution ([32]). This kind of quantization gives some good properties in approximating a continuous probability distribution by a set of discrete points. Both methods for cell selection in the chi-square statistic construction are valid in capturing lack of fit between a set of observed frequencies and the set of expected frequencies. A captured lack of fit by any chi-square test always indicates some kind of discrepancy between the null hypothesis and the underlying distribution of sample data. Therefore, the sample data in the above observations (5) and (6) indicate evidence of non-MVN. It is also observed that different choices of the number of cells also result in different p -values. It is possible that a different number of cells could give completely different conclusions in applying the chi-square test. However, a lack of fit from any grouping of data by the chi-square test always indicates some kind of discrepancy between the null hypothesis and the underlying distribution of sample data. More discussion on selecting the number of cells can refer to some early studies on the application of Pearson's chi-square test ([33–35]). A more complete illustration on applying various MVN tests to real data analysis by the R language can be found in [36].

Table 4. *p*-values from the two chi-square tests (data: bodyfat).

Subsets	χ^2 -Test	<i>m</i> = 5	<i>m</i> = 10	<i>m</i> = 15
(X_4, \dots, X_8)	χ^2_R	0.0141	0.0982	0.1707
	χ^2_T	0.0258	0.0205	0.0507
$(X_4, X_5, X_8, X_9, X_{10})$	χ^2_R	0.1414	0.0641	0.0012
	χ^2_T	0.1822	0.2599	0.5342
$(X_4, X_5, X_9, X_{10}, X_{11})$	χ^2_R	4.0736×10^{-6}	9.0387×10^{-13}	5.4589×10^{-11}
	χ^2_T	0.0783	0.0067	6.7792×10^{-5}
(X_8, \dots, X_{12})	χ^2_R	9.8107×10^{-4}	8.6962×10^{-4}	0.0775
	χ^2_T	0.2734	0.0051	0.0022
(X_4, \dots, X_{13})	χ^2_R	0.0980	0.1271	1.1012×10^{-4}
	χ^2_T	0.7114	0.4925	0.2635
(X_5, \dots, X_{14})	χ^2_R	0.3617	0.0258	0.0274
	χ^2_T	0.5782	0.3183	0.3534
(X_6, \dots, X_{15})	χ^2_R	0.3435	0.9409	0.3285
	χ^2_T	0.2159	0.2542	0.5657
$(X_4, X_6, \dots, X_{11}, X_{13}, X_{14}, X_{15})$	χ^2_R	0.2029	0.0362	0.3270
	χ^2_T	0.1998	0.1173	0.2191

5. Concluding Remarks

The RP-based chi-square test in this paper was developed for the purpose of demonstrating the application of statistical representative points (or principal points) in goodness-of-fit problems. It shows a competitive benefit compared to the same goodness-of-fit methods without employing the RP idea. It can be considered as a successful improvement from the point of view of the significant power increase in the Monte Carlo study. Because the RP-based chi-squared test is a necessary one for testing MVN, it cannot avoid the common weakness of all necessary tests for MVN in the literature. The real-data analysis in the illustrative example shows that the RP-based chi-square test can be a good supplemental test when used together with some existing tests in the literature as reviewed in the introduction section. While many statistics for testing MVN against some general alternative distributions reviewed by [18] perform very well, none of them are perfect. As summarized in [18], MVN test statistics can be classified into two major types: univariate and multivariate approaches. Univariate approaches are based on transformed data from multivariate observations. Therefore, univariate approaches are all necessary ones, implying that no rejection of the null hypothesis cannot conclude MVN. The RP-based chi-square test in this paper belongs to this family. One of the benefits of univariate approaches may be the simple asymptotic null distributions of test statistics with relative fast convergence. This can be found from the simulation of the type I error rates (see Tables 1–3 in this paper) based on the critical values of the asymptotic null distributions of the test statistics. However, loss of the original data information seems to be another common drawback of univariate approaches. As a result, univariate approaches usually lead to more power loss than do many multivariate approaches. The most representative multivariate approaches may be the Marida’s [7] multivariate skewness and kurtosis statistics. Many subsequent multivariate approaches were developed after [7], which are more or less related to the sample covariance matrix ([14,18]). One of the drawbacks of the sample-covariance matrix-related statistics is that convergence of the sample covariance matrix to the true population covariance matrix is very slow, and it becomes slower with the increase in data dimensions. As a result, almost all multivariate approaches require

very large sample size to control type I error rates and their asymptotic null distributions do not help very much in real applications with finite sample sizes ([9]). The RP-based chi-square test for MVN in this paper can be considered as one of the miscellaneous results for testing MVN reviewed in [18]. It was developed through a necessary and sufficient characterization of MVN and data transformation. This unique characterization-based data transformation guarantees that non-normal multivariate data will not result in the same set of transformed data with the Student's t -distribution. This is the motivation of proposing the RP-based chi-square test for MVN in this paper. The choice of the number of RPs for constructing the RP chi-square test is not unique. A large number of RPs may result in a zero frequency of transformed data points in some cells of the RP chi-square statistic, causing the numerical computation of the chi-square statistic broken under a small size. This is a common drawback of Pearson–Fisher's chi-square statistics. After more MNV tests were developed since Pearson's [23] initial chi-square test, and with the extensive Monte Carlo studies available in the literature, it is arguable that the Pearson chi-square test has become out-of-date. The RP-idea seems to inject new energy into the old Pearson chi-square test. A challenging application of the Pearson chi-square test is the situation that the null distribution contains unknown parameter(s), which has (have) to be estimated before implementing the chi-square test. As a result of the parameter estimation, the asymptotic null distribution of Pearson's [23] classical chi-square test is no longer an exact chi-square distribution, but a linear combination of independent chi-squares [37]. Our future research direction is to find a way to employ the RP idea to improve other types of chi-square tests, as studied in [38]. Although it is not the purpose of this paper to develop a superior MVN test to any existing MVN test in the literature, which requires an extensive Monte Carlo study, the RP-idea to improve the oldest goodness-of-fit test in statistical history may shed some additional light to the nonparametric statistical inference.

Author Contributions: J.L. developed the theory for the RP chi-squared test and wrote the initial draft. P.H. supervised J.Y. to finish the simulation and the real data analysis and edited all figures. J.Y. finished all simulation and real data analysis. All authors have read and agreed to the published version of the manuscript.

Funding: This work was partially supported by a UIC New Faculty Start-up Research Fund R72021106, and in part by the Guangdong Provincial Key Laboratory of Interdisciplinary Research and Application for Data Science, BNU-HKBU United International College (UIC), project code 2022B1212010006.

Data Availability Statement: Not applicable.

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations

The following abbreviations are used in this manuscript:

i.i.d.	Independent identically distributed
MVN	Multivariate normality
p.d.f.	Probability density function
RP	Representative points

References

1. Anderson, M.R. A characterization of the multivariate normal distribution. *Ann. Math. Stat.* **1971**, *42*, 824–827. [[CrossRef](#)]
2. Shao, Y.; Zhou, M. A characterization of multivariate normality through univariate projections. *J. Multivar. Anal.* **2010**, *101*, 2637–2640. [[CrossRef](#)] [[PubMed](#)]
3. Malkovich, J.F.; Afifi, A.A. On tests for multivariate normality. *J. Am. Stat. Assoc.* **1973**, *68*, 176–179. [[CrossRef](#)]
4. Cox, D.R.; Small, N.J.H. Testing multivariate normality. *Biometrika* **1978**, *65*, 263–272. [[CrossRef](#)]
5. Andrews, D.F.; Gnanadesikan, R.; Warner, J.L. Methods for assessing multivariate normality. In Proceedings of the Third International Symposium on Multivariate Analysis, Dayton, OH, USA, 19–24 June 1972; Volume 3, pp. 95–116.
6. Gnanadesikan, R. *Methods for Statistical Data Analysis of Multivariate Observations*; Wiley: New York, NY, USA, 1977.
7. Mardia, K.V. Measures of multivariate skewness and kurtosis with applications. *Biometrika* **1970**, *57*, 519–530. [[CrossRef](#)]

8. Mardia, K.V. Tests of univariate and multivariate normality. In *Handbook of Statistics*; Krishnaiah, P.R., Ed.; North-Holland Publishing Company: Amsterdam, The Netherlands, 1980; Volume 1, pp. 279–320.
9. Romeu, J.L.; Ozturk, A. A comparative study of goodness-of-fit tests for multivariate normality. *J. Multivar. Anal.* **1993**, *46*, 309–334. [[CrossRef](#)]
10. Horswell, R.L.; Looney, S.W. A comparison of tests for multivariate normality that are based on measures of multivariate skewness and kurtosis. *J. Stat. Comput. Simul.* **1992**, *42*, 21–38. [[CrossRef](#)]
11. Looney, S.W. How to use tests for univariate normality to assess multivariate normality. *Am. Stat.* **1995**, *39*, 75–79.
12. Liang, J.; Li, R.; Fang, H.; Fang, K.T. Testing multinormality based on low-dimensional projection. *J. Stat. Plann. Inference.* **2000**, *86*, 129–141. [[CrossRef](#)]
13. Srivastava, D.K.; Mudholkar, G.S. Goodness-of-fit tests for univariate and multivariate normal models. *Handb. Stat.* **2003**, *22*, 869–906.
14. Mecklin, C.J.; Mundfrom, D.J. An appraisal and bibliography of tests for multivariate normality. *Int. Stat. Rev.* **2004**, *72*, 123–138. [[CrossRef](#)]
15. Batsidis, A.; Martin, N.; Pardo, L.; Zografos, K. A Necessary power divergence type family tests of multivariate normality. *Commun. Stat. Simul. Comput.* **2013**, *42*, 2253–2271. [[CrossRef](#)]
16. Al-Labadi, L.; Fazeli Asl, F.; Saberi, Z. A necessary Bayesian nonparametric test for assessing multivariate normality. *Math. Methods Stat.* **2021**, *30*, 64–81. [[CrossRef](#)]
17. Doornik, J.A.; Hansen, H. An omnibus test for univariate and multivariate normality. *Oxf. Bull. Econ. Stat.* **2008**, *70*, 927–939.
18. Ebner, B.; Henze, N. Tests for multivariate normality? a critical review with emphasis on weighted L^2 -statistics. *Test* **2020**, *29*, 845–892. [[CrossRef](#)]
19. Yang, Z.H.; Fang, K.T.; Liang, J. A characterization of multivariate normal distribution and its application. *Stat. Prob. Lett.* **1996**, *30*, 347–352. [[CrossRef](#)]
20. Liang, J.; Pan, W.; Yang, Z.H. Characterization-based Q-Q plots for testing multinormality. *Stat. Prob. Lett.* **2004**, *70*, 183–190. [[CrossRef](#)]
21. Fang, K.T.; He, S.D. *The Problem of Selecting a Given Number of Representative Points in a Normal Distribution and a Generalized Mill's Ratio*; Technical Report; Department of Statistics, Stanford University: Stanford, CA, USA, 1982.
22. Fang, K.T.; Kotz, S.; Ng, K.W. *Symmetric Multivariate and Related Distributions*; Chapman and Hall: London, UK; New York, NY, USA, 1990.
23. Pearson, K. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *Philos. Mag.* **1900**, *50*, 157–175. [[CrossRef](#)]
24. Fisher, R.A. The condition under which χ^2 measures the discrepancy between observation and hypothesis. *J. R. Stat. Soc.* **1924**, *87*, 442–450.
25. Voinov, V.; Pya, N.; Alloyarova, R. A comparative study of some modified chi-squared tests. *Commun. Stat. Simul. Comput.* **2009**, *38*, 355–367. [[CrossRef](#)]
26. Flury, B. Principal points. *Biometrika* **1990**, *77*, 33–41. [[CrossRef](#)]
27. Zhou, M.; Wang, W. Representative points of the Student's t_n distribution and their applications in statistical simulation. *Acta Math. Appl. Sin.* **2016**, *39*, 620–640. (In Chinese)
28. Goodman, I.R.; Kotz, S. Multivariate θ -generalized normal distribution. *J. Multivar. Anal.* **1973**, *3*, 204–219. [[CrossRef](#)]
29. Perveen, Z.; Munir, M.; Ahmad, M. Double Weibull distribution: Properties and application. *Pak. J. Sci.* **2017**, *69*, 95–100.
30. Liang, J.; Bentler, P.M. A t -distribution plot to detect non-multinormality. *Comput. Stat. Data Anal.* **1999**, *30*, 31–44. [[CrossRef](#)]
31. Koehler, K.; Gann, F. Chi-squared goodness-of-fit tests: Cell Selection and Power. *J. Commun. Stat. Simul.* **1990**, *19*, 1265–1278. [[CrossRef](#)]
32. Graf, S.; Luschgy, H. *Foundations of Quantization for Probability Distributions*; Springer: Berlin, Germany, 2000.
33. Mann, H.; Wald, A. On the choice of the number of class intervals in the application of the chi-square test. *Ann. Math. Stat.* **1942**, *13*, 306–317. [[CrossRef](#)]
34. Dahiya, R.C.; Gurland, J. How many classes in the Pearson chi-square test? *J. Am. Stat. Assoc.* **1973**, *68*, 707–712.
35. Kallenberg, W.; Oosterhoff, J.; Schriever, B. The number of classes in chi-squared goodness-of-fit tests. *J. Am. Stat. Assoc.* **1985**, *80*, 959–968. [[CrossRef](#)]
36. Korkmaz, S.; Goksuluk, D.; Zararsiz, G. MVN: An R package for assessing multivariate normality. *R J.* **2014**, *6*, 151–162. [[CrossRef](#)]
37. Chernoff, H.; Lehmann, E.L. The use of maximum likelihood estimates in tests for goodness of fit. *Ann. Math. Stat.* **1954**, *25*, 579–589. [[CrossRef](#)]
38. Voinov, V.; Nikulin, M.; Balakrishnan, N. *Chi-Squared Goodness of Fit Tests with Applications*; Academic Press: New York, NY, USA, 2013.