

Limiting distributions of likelihood ratio test for independence of components for high-dimensional normal vectors

Yongcheng Qi 1 \cdot Fang Wang 2 \cdot Lin Zhang 3

Received: 10 May 2016 / Revised: 19 February 2018 / Published online: 14 May 2018 © The Institute of Statistical Mathematics, Tokyo 2018

Abstract Consider a p-variate normal random vector. We are interested in the limiting distributions of likelihood ratio test (LRT) statistics for testing the independence of its grouped components based on a random sample of size n. In classical multivariate analysis, the dimension p is fixed or relatively small, and the limiting distribution of the LRT is a chi-square distribution. When p goes to infinity, the chi-square approximation to the classical LRT statistic may be invalid. In this paper, we prove that the LRT statistic converges to a normal distribution under quite general conditions when p goes to infinity. We propose an adjusted test statistic which has a chi-square limit in general. Our comparison study indicates that the adjusted test statistic outperforms among the three approximations in terms of sizes. We also report some numerical results to compare the performance of our approaches and other methods in the literature.

Keywords Likelihood ratio test \cdot Covariance matrix \cdot Independence \cdot Highdimensional normal vector \cdot Central limit theorem \cdot Chi-square approximation

☑ Fang Wang fang72_wang@cnu.edu.cn

> Yongcheng Qi yqi@d.umn.edu

Lin Zhang neoivylinzhang@gmail.com

- ¹ Department of Mathematics and Statistics, University of Minnesota Duluth, 1117 University Drive, Duluth, MN 55812, USA
- ² School of Mathematical Sciences, Capital Normal University, Beijing 100048, China
- ³ School of Statistics, University of Minnesota, 224 Church Street S. E., Minneapolis, MN 55455, USA

1 Introduction

In classical statistical inference, the likelihood ratio method has been widely used for testing parametric hypotheses, assessing statistical model fittings and constructing confidence intervals/regions for parameters of interest. An advantage of using the likelihood ratio method is that one does not have to estimate the variances of the test statistics. It is well known that the limiting distributions for the likelihood ratio test (LRT) statistics are chi-square distributions under certain regularity conditions when the dimension of the data or the number of the parameters of interest is fixed.

For many modern datasets, their dimensions can be proportionally large compared with the sample size. For example, financial data, consumer data, modern manufacturing data, and multimedia data all have this feature. However, the chi-square approximation does not fit the distribution of the LRT statistics very well for the high-dimensional case, especially when the dimension of the data grows with the sample size. To deal with this feature, Schott (2001), Ledoit and Wolf (2002), Schott (2005), Schott (2007), Bai et al. (2009), Chen et al. (2010), Jiang et al. (2012), Srivastava and Reid (2012), Jiang and Yang (2013), Jiang et al. (2013), Jiang and Qi (2015a), Bao et al. (2017), and Li et al. (2017) have derived different methods to study the classical LRT or alternatives to LRT when the dimension p is large.

In this paper, we consider the LRT statistics for testing the independence of subvectors from a high-dimensional normal vector. We are devoted to deriving the limiting distributions of the LRT statistics. Our motivation is from recent papers by Jiang and Yang (2013) and Jiang and Qi (2015a) where several LRTs on high-dimensional normal random vectors have been considered when the dimension goes to infinity with the sample size and corresponding LRT statistics are shown to be asymptotically normal. From both Jiang and Yang (2013) and Jiang and Qi (2015a), one can conclude that the standard LRT statistics cannot be fitted by chi-square distribution as the dimension of the data diverges.

Our problem can be addressed as the statistical model below. For a multivariate distribution $N_p(\mu, \Sigma)$, we partition a set of p variates with a joint normal distribution into k subsets and ask whether the k subsets are mutually independent, or equivalently, we want to test whether the covariance matrix Σ is block diagonal. It is worth mentioning that Srivastava and Reid (2012) and Jiang et al. (2013) consider the same testing problem for multivariate normal distributions. Both of the two papers aim at deriving new test statistics by allowing that the dimension p can be larger than the sample size *n* but assuming that p/n has a limit and the numbers of components within subsets are proportional asymptotically. Our current paper focuses on the limiting distributions of the LRT statistics which exist only when $2 \le p < n$. We will establish the central limit theorem for the LRT statistic when the dimension of the normal random vector goes to infinity. In this paper, we allow that k changes with n and the partition can be unbalanced in the sense that numbers of components within subsets are not necessarily proportional. Our results are extensions of some results in Jiang and Yang (2013) and Jiang and Qi (2015a) where the number of the partition is fixed and the partition is well balanced.

Our central limit theorem established in (3) in Sect. 2 provides normal approximation to the LRT statistic when p is large. The classical chi-square approximation in

(9) is valid when p is fixed. To assess the normal approximation and compare it to the classical chi-square approach, we present some finite-sample simulation results in Sect. 3. The study indicates that the normal approximation outperforms the chi-square approximation when p is large or relatively large, but the chi-square approximation is better when p is small. Similar phenomena have been observed in both Jiang and Yang (2013) and Jiang and Qi (2015a). However, the theoretical results and simulation study do not suggest clear ranges of p when to use the normal approximation and when to use the chi-square approximation. Such a cutoff on p should depend on the sample size n and grouping parameters q_i 's and k. For the limiting distributions of $-2 \log A_n$, where A_n is the Wilks likelihood ratio statistic defined in (2), the transition from the chi-square to the normality seems a problem in practice. This motivates us to investigate why the chi-square approach fails when p is large. As a result, we will propose an adjusted LRT statistic which has a chi-square limit in the entire range of p. Some more insights will be given in Sect. 2.

One related problem is to test the independence of the p variates of a p-dimensional normal random vectors, which is also considered in Jiang and Yang (2013) and Jiang and Qi (2015a) as a different test problem. This is indeed a special case of our setting when the random vector is partitioned into p blocks and each block contains only one component.

The rest of the paper is organized as follows. In Sect. 2, we state our main results including the central limit theorem for the LRT statistic when dimension p goes to infinity with the sample size n and the chi-square approximation for the adjusted LRT statistic for any p in the range that the LRT can be applied. In Sect. 3, we carry out some simulation studies to compare the performance of three methods including the chi-square approximation to the LRT statistic, the normal approximation to the LRT statistics, and the chi-square approximation to the adjusted LRT statistic. We also compare our test statistics with some other approaches in the literature. All proofs are given in Sect. 4.

2 Main results

Throughout, let χ_f^2 denote the chi-square random variable with f degrees of freedom and N(0, 1) the standard normal variable. For each $\alpha \in (0, 1)$, $\chi_{f,\alpha}^2$ and z_{α} are the α level critical values of χ_f^2 and N(0, 1), respectively. $N_p(\mu, \Sigma)$ denotes p-dimensional multivariate normal distribution with mean μ and covariance matrix Σ .

For $k \ge 2$, let q_1, \ldots, q_k be k positive integers. Denote $p = q_1 + \cdots + q_k$ and let

$$\mathbf{\Sigma} = (\mathbf{\Sigma}_{ij})_{1 \le i, j \le k}$$

be a positive definite matrix, where Σ_{ij} is a $q_i \times q_j$ sub-matrix for all $1 \le i, j \le k$. Assume ξ_i is a q_i -dimensional normal random (column) vector for each $1 \le i \le k$, and the *p*-dimensional random vector $(\xi'_1, \ldots, \xi'_k)'$ has the distribution $N_p(\mu, \Sigma)$. We are interested in testing the independence of *k* random vectors ξ_1, \ldots, ξ_k , or equivalently testing

$$H_0: \Sigma_{ij} = \mathbf{0} \text{ for all } 1 \le i < j \le k \quad \text{vs} \quad H_a: H_0 \text{ is not true.}$$
(1)

🖉 Springer

Assume that $\mathbf{x}_1, \ldots, \mathbf{x}_n$ are independent and identically distributed (i.i.d.) random vectors from distribution $N_p(\mu, \mathbf{\Sigma})$. Define

$$\mathbf{A} = \sum_{i=1}^{n} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})' \text{ with } \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i,$$

and partition it as follows:

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & \cdots & \mathbf{A}_{1k} \\ \mathbf{A}_{21} & \mathbf{A}_{22} & \cdots & \mathbf{A}_{2k} \\ \vdots & \ddots & \ddots & \vdots \\ \mathbf{A}_{k1} & \mathbf{A}_{k2} & \cdots & \mathbf{A}_{kk} \end{pmatrix}$$

where A_{ij} is a $q_i \times q_j$ matrix. The likelihood ratio statistic for testing (1) is given by

$$\Lambda_n = \frac{|\mathbf{A}|^{n/2}}{\prod_{i=1}^k |\mathbf{A}_{ii}|^{n/2}} := (W_n)^{n/2},$$
(2)

see Wilks (1935) or Theorem 11.2.1 from Muirhead (1982).

In the sequel, we assume that both k and p can depend on the sample size n, that is, $k = k_n$ and $p = p_n$, and both can diverge to infinity. First, we extend the results given in Jiang and Yang (2013) and Jiang and Qi (2015a).

Theorem 1 Let $p = p_n$ satisfy $2 \le p < n$ and $p_n \to \infty$ as $n \to \infty$ and $k = k_n \ge 2$ be a sequence of integers. Assume that q_1, \ldots, q_k are k positive integers such that $p = \sum_{i=1}^{k} q_i$ and $\max_{1 \le i \le k} q_i \le \delta p$, for a fixed $\delta \in (0, 1)$ and all large n. Λ_n is the Wilks likelihood ratio statistic defined in (2). Then, under the null hypothesis in (1)

$$T_0 := \frac{-2\log \Lambda_n - \bar{\mu}_n}{\bar{\sigma}_n} \xrightarrow{d} N(0, 1)$$
(3)

as $n \to \infty$, where

$$\bar{\mu}_n = \mu_n + \frac{n}{3} \left(b(n, p) - \sum_{i=1}^k b(n, q_i) \right), \quad \bar{\sigma}_n^2 = \sigma_n^2 + 2n^2 \left(b(n, p) - \sum_{i=1}^k b(n, q_i) \right)$$
(4)

with

$$\mu_n = n \sum_{i=1}^k \left(q_i - n + \frac{3}{2} \right) \log \left(1 - \frac{q_i}{n} \right) - n \left(p - n + \frac{3}{2} \right) \log \left(1 - \frac{p}{n} \right), \quad (5)$$

$$\sigma_n^2 = 2n^2 \left(-\log\left(1 - \frac{p}{n}\right) + \sum_{i=1}^k \log\left(1 - \frac{q_i}{n}\right) \right)$$
(6)

🖉 Springer

and

$$b(n,q) = \sum_{j=1}^{q} \frac{1}{(n-j)^2} \quad \text{for } 1 \le q < n.$$
⁽⁷⁾

Remark 1 The central limit theorem in (3) is still valid if $\bar{\mu}_n$ and $\bar{\sigma}_n$ are replaced by μ_n and σ_n , respectively. In fact, we will show in the proof of Theorem 1 in Sect. 4 that

$$\frac{\bar{\mu}_n - \mu_n}{\bar{\sigma}_n} \to 0 \quad \text{and} \quad \frac{\sigma_n}{\bar{\sigma}_n} \to 1 \quad \text{as } n \to \infty.$$
(8)

In general, $\bar{\mu}_n$ and $\bar{\sigma}_n^2$ are better approximations to the mean and variance of $-2 \log \Lambda_n$, but μ_n and σ_n^2 have simpler expressions. Our simulation study indicates little difference in the convergence rate to normality between two different selections if p is not too close to n. The situation is quite different when p is close to n such as p = n - 1. One can verify that in this case the convergence rates in (8) are of order $(\log n)^{-1/2}$. Therefore, with n = 100 or even n = 1000, using μ_n and σ_n in (3) can result in a serious bias in the normal approximation.

Remark 2 When the number of the partition *k* is fixed and $p = p_n$ goes to infinity as $n \to \infty$, Jiang and Yang (2013) and Jiang and Qi (2015a) prove the central limit theorem for log W_n under additional assumptions including: (A) $2 \le p < n - 1$; (B) for some $\delta \in (0, 1), \delta \le q_i/q_j \le \delta^{-1}$ for all $1 \le i, j \le n$ and all *n*. Note that $-2 \log \Lambda_n = -n \log W_n$ from (2). After taking into account the constant scale -n, one can find out that the estimates for the mean and variance of $-2 \log \Lambda_n$ in this paper are slightly different from those in Jiang and Qi (2015a). Our estimates are more accurate in the sense that our estimates catch the leading terms not only for large *p* but also for small *p*. This will be very critical in establishing our Theorem 2.

Remark 3 Our restriction $2 \le p < n$ is a natural one since the matrix **A** is not of full rank and the LRT fails when $p \ge n$. See, for example, Jiang and Yang (2013). Compared with Jiang and Yang (2013) and Jiang and Qi (2015a), our present paper has removed several constraints in the following three aspects: a). the number of the partition k can depend on n; b). the numbers of the components $\{q_j\}$ do not have to be comparable; c). the dimension p can be any integer in the range $2 \le p < n$. In particular, it follows from Theorem 1 that the central limit theorem holds for $-2 \log A_n$ when p = n - 1.

Remark 4 The classical likelihood method handles the case when both p and k are fixed integers. When q_1, q_2, \ldots, q_k are fixed as n goes to infinity, the standard LRT statistic of (1) has a chi-square limit:

$$-2\rho\log\Lambda_n \xrightarrow{d} \chi_f^2, \tag{9}$$

where

$$f = \frac{1}{2} \left(p^2 - \sum_{i=1}^k q_i^2 \right),$$
(10)

🖄 Springer

$$\rho = 1 - \frac{2\left(p^3 - \sum_{i=1}^k q_i^3\right) + 9\left(p^2 - \sum_{i=1}^k q_i^2\right)}{6n\left(p^2 - \sum_{i=1}^k q_i^2\right)},\tag{11}$$

see, for example, Theorem 11.2.5 in Muirhead (1982). Note that ρ is introduced to achieve the second-order accuracy and ρ converges to one as *n* goes to infinity.

As we have known, the chi-square approximation and the normal approximation apply in different ranges of p. But for moderate values of p, it seems difficult to tell which approximation method should be used, and therefore, some guidelines in this direction are desirable in practice. To better understand the chi-square approximation, one needs to look at the asymptotic mean and variance of the test statistic. It is well known that the ratio between the mean and the variance of a chi-square distribution is 1:2. One can verify that the ratio of the asymptotic mean $\bar{\mu}_n$ and the asymptotic variance $\bar{\sigma}_n^2$ of $-2\log \Lambda_n$ is close to 1:2 when p is fixed but this ratio can be quite different from 1:2 when p goes to infinity from Theorem 1. The same is true for $-2\rho \log \Lambda_n$. This explains why the chi-square approximation for $-2\log \Lambda_n$ works only for small p. On the other hand, a sequence of random variables with a normal limit can be also approximated in distribution by linear functions of some chi-square random variables, and this means that when p is large, the LRT statistic $-2 \log A_n$, after proper normalization, can also be approximated by some chi-square distribution. Since we aim at a unified chi-square approximation in the full range $2 \le p < n$, we should also take into account the case when p is small.

We propose an adjusted log-likelihood ratio test statistic (ALRT)

$$Z_n = (-2\log\Lambda_n)\sqrt{\frac{2f_n}{\bar{\sigma}_n^2}} + f_n - \bar{\mu}_n\sqrt{\frac{2f_n}{\bar{\sigma}_n^2}}$$
(12)

with f_n , $\bar{\mu}_n$ and $\bar{\sigma}_n^2$ being defined in (10) and (4), respectively. Note that the ALRT is essentially the LRT statistic since it is a linear combination of the log-likelihood ratio test statistic, $-2 \log \Lambda_n$. We have the following theorem regarding the chi-square approximation of Z_n statistic.

Theorem 2 Let $p = p_n$ be a sequence of integers with $2 \le p < n$ for any $n \ge 1$. Assume $k = k_n$ is also a sequence of positive integers, and q_1, \ldots, q_k are k positive integers such that $p = \sum_{i=1}^{k} q_i$. Assume there exists a constant $\delta \in (0, 1)$ such that $\max_{1 \le i \le k} q_i \le \delta p$ for all large n. Then, we have under the null hypothesis in (1)

$$\lim_{n \to \infty} \sup_{-\infty < x < \infty} |P(Z_n \le x) - P(\chi_{f_n}^2 \le x)| = 0.$$
(13)

Remark 5 The test statistic Z_n defined in (12) is a linear combination of $-2 \log A_n$, and the coefficients are selected in such a way that the asymptotic mean and variance of Z_n are close to f_n and $2f_n$, respectively, where f_n and $2f_n$ are the mean and variance of a chi-squared random variable with f_n degrees of freedom.

Remark 6 Let $\alpha \in (0, 1)$ be any given number. Based on the classical chi-square approximation in (9), the LRT rejects the null hypothesis in (1) if $-2\rho \log \Lambda_n \ge \chi^2_{f_n,\alpha}$. Based on the normal approximation in Theorem 1, the rejection region is $-2 \log \Lambda_n \ge \mu_n + \bar{\sigma}_n z_\alpha$. Based on the chi-square approximation in Theorem 2, the LRT rejects the null hypothesis in (1) if $Z_n \ge \chi^2_{f_n,\alpha}$.

Remark 7 In Theorems 1 and 2, we impose assumption that $\max_{1 \le i \le k} q_i \le \delta p$ for some $\delta \in (0, 1)$. This condition is quite mild. We notice that it is trivial when p is fixed. It rules out the extreme situation that $\max_{1 \le i \le k} q_i/p \to 1$ along the entire sequence or any subsequence. Violating this assumption may result in a non-normal or chi-square limit.

As an application, we consider the test for complete independence investigated in Jiang and Yang (2013) and Jiang and Qi (2015a). Assume that a *p*-dimensional random vector $\mathbf{x} = (x_1, \ldots, x_p)'$ has a distribution $N_p(\mu, \Sigma)$. Our interest is in testing the independence of the *p* components x_1, x_2, \ldots, x_p or equivalently testing that the covariance matrix Σ is diagonal based on a random sample $\mathbf{x}_1, \ldots, \mathbf{x}_n$ from distribution $N_p(\mu, \Sigma)$. Let $\mathbf{R} = (r_{ij})_{p \times p}$ be the correlation matrix generated from $N_p(\mu, \Sigma)$. Then, the test is equivalent to

$$H_0: \mathbf{R} = \mathbf{I}_p \text{ vs } H_a: \mathbf{R} \neq \mathbf{I}_p, \tag{14}$$

where \mathbf{I}_p denotes $p \times p$ identity matrix.

Obviously, test (14) is a special case of the test (1) with k = p and $q_1 = \ldots = q_p = 1$. In this case, W_n , defined in (2), is equal to the determinant of Pearson's correlation matrix. More specifically, let $\mathbf{x}_i = (x_{i1}, \ldots, x_{ip})'$ for $1 \le i \le n$ and $\bar{x}_j = \frac{1}{n} \sum_{m=1}^n x_{mj}$ for $1 \le j \le p$. For $1 \le i, j \le p$, define

$$\hat{r}_{ij} = \frac{\sum_{m=1}^{n} (x_{mi} - \bar{x}_i)(x_{mj} - \bar{x}_j)}{\sqrt{\sum_{m=1}^{n} (x_{mi} - \bar{x}_i)^2 \cdot \sum_{m=1}^{n} (x_{mj} - \bar{x}_j)^2}}.$$
(15)

Then, Pearson's correlation matrix is given by $\hat{\mathbf{R}}_n = (\hat{r}_{ij})_{p \times p}$.

From (2), $\log \Lambda_n = \frac{n}{2} \log |\hat{\mathbf{R}}_n|$. Under condition $2 \le p < n-4$ and $p \to \infty$, a central limit theorem is proved for $\log |\hat{\mathbf{R}}_n|$ in Jiang and Yang (2013) and Jiang and Qi (2015a), see, for example, Corollary 1 in Jiang and Qi (2015a). Since $\max_{1\le i\le p} q_i = 1 \le p/2$ for any $2 \le p < n$, both Theorems 1 and 2 are valid in this case. In our central limit theorem in Theorem 1, we have extended the range for p and allow p = n - 4, n - 3, n - 2, and n - 1 as well. From Theorems 1 and 2, we conclude the following corollary.

Corollary 1 Let $p = p_n$ be a sequence of positive integers with $2 \le p < n$. Set $\log A_n = \frac{n}{2} \log |\hat{\mathbf{R}}_n|, \bar{\sigma}_n^2 = 2n^2 \left(p \log(1 - \frac{1}{n}) - \log(1 - \frac{p}{n}) + \sum_{j=1}^p \frac{1}{(n-j)^2} - \frac{p}{(n-1)^2} \right),$

🖄 Springer

$$\begin{split} \bar{\mu}_n &= n \left(p \left(\frac{5}{2} - n \right) \log \left(1 - \frac{1}{n} \right) - \left(p - n + \frac{3}{2} \right) \log \left(1 - \frac{p}{n} \right) \\ &+ \frac{1}{3} \left(\sum_{j=1}^p \frac{1}{(n-j)^2} - \frac{p}{(n-1)^2} \right) \right), \end{split}$$

and $f_n = \frac{1}{2}p(p-1)$. Define Z_n as in (12). Then, under the null hypothesis in (14), we have

$$\lim_{n \to \infty} \sup_{-\infty < x < \infty} |P(Z_n \le x) - P(\chi_{f_n}^2 \le x)| = 0.$$

In addition, if $\lim_{n\to\infty} p_n = \infty$, we have

$$\frac{-2\log\Lambda_n - \bar{\mu}_n}{\bar{\sigma}_n} \xrightarrow{d} N(0, 1)$$

3 Simulation study

In this section, we will have some simulation studies to compare the performance of three different approaches to the likelihood ratio test statistics and to compare our adjusted likelihood test statistics with some other approaches in the literature.

3.1 Comparisons of the likelihood ratio tests under different approaches

In this subsection, we will compare the accuracy of the three different approaches to the likelihood ratio test statistic $-2 \log \Lambda_n$ under the null hypothesis of (1) and the performance of $-2 \log \Lambda_n$ under different normalizations proposed by Jiang and Yang (2013), Jiang and Qi (2015a) and approaches in the present paper in terms of sizes and powers.

First, we carry out a finite-sample simulation study to compare the performance of the three approximation methods including the classical chi-square approximation (9), the normal approximation (3), and adjusted chi-square approximation (13). For the three methods, we demonstrate how well the proposed limiting distributions fit the histograms of the three test statistics based on 10,000 random samples of size n = 101. From (21) in Lemma 1, the moment-generating function of log W_n is distribution-free under the null hypothesis in (1). Since the three test statistics are functions of A_n and hence they are also functions of log W_n from (2), they are distribution-free under the null hypothesis in (1). Therefore, the underlying distribution in our study is assumed to be a multivariate normal distribution with independent standard normal components.

The simulation study consists of two cases. In the first case, we test the independence of k = 3 block vectors and the ratio of their dimensions is kept fixed as $q_1:q_2:q_3 =$ 2:2:1 with $p = q_1 + q_2 + q_3 = 5$, 30, 60, and 100, respectively, and the sample size n = 101 is fixed. Figure 1 contains the histograms for the three test statistics considered in Sect. 2 including classical chi-square approach "Chisq" given in (9), normal approximation "CLT" given in (3), and adjusted likelihood ratio approach "ALRT" in (12) and (13). The second case is the test of complete independence in (14) with dimension parameter p = 5, 30, 60, and 100 and a fixed sample size n = 101, and corresponding histograms for the three test statistics are included in Fig. 2.

From both Figs. 1 and 2, one can easily conclude the following common features:

- (i) The classical chi-square approximation works very well for small *p*, and it becomes worse with the increase of *p* and eventually departs completely from the histograms of the test statistic.
- (ii) The normal approximation shows some lack of fit to the histograms for small p such as p = 5, and the approximation is getting better with the increase of p. For p = 100, the normal distribution is slightly away from the histograms of the test statistic. In the range of $2 \le p \le n 1 = 100$, p = 100 represents the extremal case when one can apply the likelihood ratio method, and the sample covariance matrices in this case are nearly singular. In this case, the convergence rate in (3) is slow. The fit improves with the increase of n.
- (iii). The adjusted likelihood ratio test (ALRT) statistic given by (12) works very well in the entire range $2 \le p \le n 1 = 100$, that is, for smaller *p*, ALRT behaves like the classical chi-square approximation, while for large *p*, the ALRT performs equally well as the normal approximation.

In summary, the adjusted likelihood ratio test (ALRT) outperforms over the classical chi-square approximation and the normal approximation. By using the ALRT, one does not have to differentiate whether a value of p is small or moderately large. For very large p, the ALRT and the normal approximation work equally well.

Now we compare the performance of $-2 \log A_n$ under different normalizations proposed by Jiang and Yang (2013), Jiang and Qi (2015a) and approaches in the present paper in terms of sizes and powers. Jiang and Yang (2013) and Jiang and Qi (2015a) use the same normalization constants. Under the null hypothesis in (1), they show the following central limit theorem under different constraints on q_1, \ldots, q_k with fixed $k \ge 2$:

$$T_1 := \frac{-2\log \Lambda_n - m_n}{\tau_n} \xrightarrow{d} N(0, 1) \text{ as } n \to \infty,$$
(16)

where

$$m_{n} = n \left(r_{n-1}^{2} \left(p - n + \frac{3}{2} \right) - \sum_{i=1}^{k} r_{n-1,i}^{2} \left(q_{i} - n + \frac{3}{2} \right) \right),$$

$$\tau_{n}^{2} = 2n^{2} \left(r_{n-1}^{2} - \sum_{i=1}^{k} r_{n-1,i}^{2} \right),$$

 $r_x = (-\log(1-\frac{p}{x}))^{1/2}$ for x > p, and $r_{x,i} = (-\log(1-\frac{q_i}{x}))^{1/2}$ for $x > q_i$ and $1 \le i \le k$. In fact, Jiang and Yang (2013) prove (16) under assumption that

25 30 35

n = 101, p = 5

9

0.08

Density 0.06

0.04

0.02

0.0

0 5 10

15 20 ALRT

n = 101 , p = 30



n = 101 , p = 100

0.005

0.004

Density 0.003

0.002

0.001

000





3600



n = 101 , p = 100

Fig. 1 Test of independence in (1) with k = 3: histograms for three test statistics including classical chi-square approach "Chisq" given in (9), normal approximation "CLT" given in (3), and adjusted likelihood ratio approach "ALRT" in (12) and (13)





 $q_i/n \rightarrow y_i \in (0, 1)$ for $1 \le i \le k$ and Jiang and Qi (2015a) assume only that q_i 's are of the same order and $\min_{1\le i\le k} q_i \rightarrow \infty$ as $n \rightarrow \infty$.

By using T_1 in (16), one rejects the null hypothesis in (1) at level α if $T_1 > m_n + \tau_n z_\alpha$.

As some empirical evidence has shown, the classical chi-square approach to the likelihood ratio test statistics is improper to our current settings. Our comparison of the likelihood ratio test statistics under different normalizations will focus on the test statistics T_0 , Z_n , and T_1 , defined in Eqs. (3), (12), and (16), respectively.

In our simulation, we generate 10,000 random samples of size *n* from a multivariate normal distribution $N_p(0, \Sigma_{\delta})$, where $\Sigma_{\delta} = (1-\delta)\mathbf{I}_p + \delta \mathbf{J}_p$, where \mathbf{I}_p denotes a $p \times p$ identity matrix, \mathbf{J}_p denotes a $p \times p$ matrix with all entries equal to 1, and $\delta \in [0, 1)$ is a constant.

First, we set k = 3. With different choices of (q_1, q_2, q_3) and n, we estimate the sizes (when $\delta = 0$) and powers (when $\delta = 0.1$ and 0.2) for the three test statistics based on the 10000 random samples. The comparison results are given in Table 1.

Then, we consider the case when k is large. In the study, we assume p is an even integer and set k = p/2. We discuss two different cases as follows.

Case 1 (Balanced case) $q_1 = \cdots = q_k = 2;$

Case 2 (Unbalanced case) $q_1 = k + 1$, $q_2 = \cdots = q_k = 1$.

Again the sizes and powers are estimated from 10,000 random samples of size *n* from $N_p(0, \Sigma_{\delta})$ with different choices of *p* and *n* under each of the two cases above. The simulation results are given in Table 2.

From Tables 1 and 2, we can conclude that test statistic Z_n gives the most accurate sizes (type I errors), especially when dimension parameter p is not large. Both T_0 and T_1 are also working very well in terms of sizes when p is reasonably large. T_0 and T_1 have larger powers than Z_n especially when p is small, and this can be explained by the difference among sizes of the three test statistics. Note that all three test statistics are linear transformations of likelihood ratio test statistic $-2 \log A_n$. Since T_0 and T_1 have larger and less accurate sizes, their rejection regions are larger than those of Z_n and naturally their powers are larger than those of Z_n . When p increases, the size of Z_n is getting closer to that of T_0 and T_1 and the powers of all test statistics are comparable. From Tables 1 and 2, we note that the performance of T_0 and T_1 is quite similar in terms of both sizes (type I errors) and powers.

3.2 Comparisons of adjusted log-likelihood ratio test statistic and other methods

In this subsection, we plan to compare our adjusted log-likelihood ratio test statistic, i.e., Z_n in (12) with other three test statistics, including two trace criterion test statistics by Jiang et al. (2013) and Li et al. (2017) and Schott type statistics by Bao et al. (2017). Jiang et al. (2013) and Bao et al. (2017) propose test statistics for test (1) for any fixed $k \ge 2$ while Li et al. (2017) consider test (1) for k = 2 only.

Bao et al. (2017) compare numerically the performance of test statistics in Jiang et al. (2013), Jiang and Yang (2013), and Bao et al. (2017) when k = 3. They conclude that Schott type statistics are very robust, and when the sample size n and the total dimension p are large, all three test statistics perform very satisfactorily in terms of the empirical sizes, except the likelihood ratio test statistics when the total dimension p is

ble 1 Comparison	s of the likelihood r	atio test statistics u	nder different	normalizations.	The sizes	and power	rs are estima	ted based o	n 10,000 sin	nulations, an	id the type I
ors for all tests are	set to be 0.05										

	20 m 20 m 20	0000								
(q_1, q_2, q_3)	u	Size $(\delta = 0)$			Power ($\delta =$	0.1)		Power ($\delta = 0$	0.2)	
		T_0	Z_n	T_1	T_0	Z_n	T_1	T_0	Z_n	T_1
(2, 2, 1)	20	0.0730	0.0529	0.0717	0.1158	0.0911	0.1151	0.2648	0.2233	0.2637
	50	0.0653	0.0481	0.0652	0.2233	0.1848	0.2230	0.6570	0.6109	0.6569
	100	0.0699	0.0533	0.0698	0.4294	0.3722	0.4293	0.9471	0.9327	0.9471
(10, 10, 5)	50	0.0599	0.0545	0.0591	0.3049	0.2869	0.3015	0.7478	0.7327	0.7457
	100	0.0561	0.0505	0.0559	0.7912	0.7783	0.7907	0.7912	0.7783	0.7907
	150	0.0553	0.0497	0.0552	0.9701	0.9679	0.9700	1.0000	1.0000	1.0000
(30, 30, 10)	100	0.0531	0.0510	0.0520	0.5097	0.5004	0.5044	0.8920	0.8882	0.8901
	150	0.0520	0.0508	0.0519	0.9154	0.9126	0.9152	6666.0	7666.0	0.9999
(40, 20, 20)	100	0.0567	0.0542	0.0545	0.4413	0.4343	0.4347	0.7980	0.7937	0.7942

simulations, and the type	
imated based on 10,000	
sizes and powers are est	
nt normalizations. The	
statistics under differe	
the likelihood ratio test	be 0.05
ble 2 Comparisons of	ors for all tests are set to

	01 all web alv or										
d	Case no.	u	Size $(\delta = 0)$			Power ($\delta =$	0.1)		Power ($\delta =$	0.2)	
			T_0	Z_n	T_1	T_0	Z_n	T_1	T_0	Z_n	T_1
10	1	20	0.0659	0.0537	0.0627	0.1389	0.1200	0.1348	0.3630	0.3328	0.3560
	1	50	0.0635	0.0534	0.0631	0.3721	0.3429	0.3713	0.8981	0.8838	0.8980
20	1	50	0.0584	0.0530	0.0579	0.5360	0.5152	0.5333	0.9761	0.9735	0.9755
	1	100	0.0589	0.0533	0.0587	0.9456	0.9409	0.9456	1.0000	1.0000	1.0000
40	1	50	0.0598	0.0568	0.0562	0.5468	0.5372	0.5346	0.9769	0.9760	0.9759
	1	100	0.0563	0.0539	0.0557	0.9928	0.9922	0.9926	1.0000	1.0000	1.0000
10	2	20	0.0690	0.0565	0.0670	0.1246	0.1051	0.1215	0.2847	0.2496	0.2774
	2	50	0.0644	0.0541	0.0640	0.3008	0.2683	0.3002	0.8140	0.7895	0.8134
20	2	50	0.0591	0.0525	0.0583	0.3993	0.3790	0.3974	0.9205	0.9137	0.9196
	2	100	0.0611	0.0544	0.0609	0.8615	0.8475	0.8611	1.0000	0.9999	1.0000
40	2	50	0.0597	0.0561	0.0561	0.3643	0.3535	0.3533	0.8657	0.8586	0.8583
	2	100	0.0584	0.0547	0.0581	0.9415	0.9382	0.9413	1.0000	1.0000	1.0000

too close to the sample size *n*. In terms of the empirical powers, Schott type statistics are very competitive among the three test statistics in most cases.

Since Li et al. (2017)'s test statistics apply to the case k = 2 only, for convenience, we will compare the aforementioned test statistics in case k = 2 and compare them with the test statistics in the present paper.

Jiang et al. (2013)'s large-dimensional trace criterion test statistic (T_2) is defined as

$$L_n = \mathbf{tr}(A_{21}A_{11}^{-1}A_{12}A_{22}^{-1}).$$

where **tr**(*A*) denotes the trace of matrix *A*. If $r_{n1} := q_2/q_1 \to r_1 \in (0, \infty), r_{n2} := q_2/(n - 1 - q_2) \to r_2 \in (0, \infty), q_2 < n$, then

$$T_2 := \frac{L_n - a_n}{\sqrt{b_n}} \stackrel{d}{\to} N(0, 1) \text{ as } n \to \infty$$
(17)

under the null hypothesis in (1), where

$$b_n = \frac{2h_n^2 r_{n1}^2 r_{n2}^2}{(r_{n1} + r_{n2})^2}, \ a_n = \frac{q_2 r_{n2}}{r_{n1} + r_{n2}}, \ h_n = \sqrt{r_{n1} + r_{n2} - r_{n1} r_{n2}}.$$

It is easy to verify that $a_n = q_1 q_2 / (n-1)$ and $b_n = 2q_1 q_2 (n-1-q_1)(n-1-q_2) / (n-1)^4$.

In case k = 2, the Schott type statistics proposed by Bao et al. (2017) reduce to

$$\operatorname{tr}\left(A_{22}^{-1/2}A_{21}A_{11}^{-1}A_{12}A_{22}^{-1/2}\right) = \operatorname{tr}\left(A_{21}A_{11}^{-1}A_{12}A_{22}^{-1}\right),$$

which is equal to L_n . Theorem 3.1 in Bao et al. (2017) leads to (17) in this case.

Define for i, j = 1, 2

$$\gamma_{ij} = \frac{1}{(n-2)(n+1)} \left(\mathbf{tr}(A_{ij}A_{ji}) - \frac{1}{n-1} \mathbf{tr}(A_{ii}) \mathbf{tr}(A_{jj}) \right).$$

The trace criterion test statistic by Li et al. (2017) is defined as γ_{12} . Under the null hypothesis in (1) for k = 2, it is proved in Li et al. (2017) that

$$T_3 := \sqrt{\frac{(n-2)(n+1)}{2}} \frac{\gamma_{12}}{\sqrt{\gamma_{11}\gamma_{22}}} \xrightarrow{d} N(0,1) \text{ as } n \to \infty$$
(18)

if $p = q_1 + q_2 \rightarrow \infty$ as $n \rightarrow \infty$ and

$$0 < \lim_{n \to \infty} \frac{1}{p} \operatorname{tr}(\boldsymbol{\Sigma}^{\mathbf{i}}) < \infty \quad \text{for } i = 1, 2, 4.$$
⁽¹⁹⁾

Note that the test statistic T_3 is not distribution-free under the null hypothesis of (1), and the asymptotic normality (18) is proved under condition (19); however, condition

 $p = q_1 + q_2 \rightarrow \infty$ is less restrictive than that required for other statistics mentioned above.

Given a size $\alpha \in (0, 1)$, test T_2 (or T_3) rejects the null hypothesis if $T_2 > z_{\alpha}$ (or $T_3 > z_{\alpha}$).

Now we compare the performance of test statistics T_2 and T_3 and our adjusted loglikelihood ratio test statistic Z_n . In our simulation study, we assume $q_1 > q_2 \ge 1$ and $p = q_1 + q_2 < n$. Our samples are generated from the populations similar to those in Jiang et al. (2013). Let $\mathbf{z} = (z_1, \ldots, z_p)'$ be a random vector whose components are independent normal random variables with mean 0 and variance 1.

Model 1 $\mathbf{x} = (x_1, ..., x_p)'$, where $x_i = (1 + c)z_i$ for $i = 1, ..., p_1, x_{p_1+j} = z_{p_1+j} + cz_j$ for $j = 1, ..., p_2$, and c is a constant;

Model 2 **x** = $(x_1, ..., x_p)'$, where $x_i = (1 + c)z_i$ for $i = 1, ..., p_1, x_{p_1+j} = z_{p_1+j} + cz_j$ for $j = 1, ..., p_2 - 1, x_p = p^{-1/4}z_p$, and c is a constant.

With a selection of (q_1, q_2, n, c) , we generate 10,000 random samples of size *n* from each model above and then estimate the sizes of the tests (when c = 0) or the powers of the tests (when $c \neq 0$). The size α is set to be 0.05 in the simulation.

Tables 3 and 4 present results for the numerical comparisons on the three test statistics. From the two tables, empirical sizes for three tests are close to the nominal level 0.05 in most cases under both Model 1 and Model 2. In view of empirical powers, Z_n and T_2 are comparable in most cases while T_3 is better than both Z_n and T_2 under Model 1. Under Model 2, Z_n has a slightly larger power than T_2 in most cases and both are significantly better than T_3 .

4 Proofs

4.1 Some Lemmas

For two sequences of numbers $\{a_n\}$ and $\{b_n\}$, the notation $a_n = O(b_n)$ as $n \to \infty$ means $\limsup_{n\to\infty} \frac{a_n}{b_n} < \infty$, and $a_n = o(b_n)$ as $n \to \infty$ means $\lim_{n\to\infty} \frac{a_n}{b_n} = 0$.

Throughout the paper, $\Gamma(x)$ stands for the Gamma function, given by

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} \mathrm{d}t \quad x > 0.$$

Define the multivariate Gamma function by

$$\Gamma_p(x) := \pi^{p(p-1)/4} \prod_{i=1}^p \Gamma\left(x - \frac{1}{2}(i-1)\right)$$
(20)

with x > (p - 1)/2. See p. 62 in Muirhead (1982).

We will need the following formula for the moments of W_n , where W_n is defined in Eq. (2).

Lemma 1 (*Theorem 11.2.3 from Muirhead (1982)*) Let $p = \sum_{i=1}^{k} q_i$ and W_n be Wilks' likelihood ratio statistic defined as (2). Then, under (1),

Com

Table 3Coare set to be	mparisons on 0.05	size and power u	nder Model 1. T	he sizes and pow	ers are estimated	l based on 10,000) simulations und	ler Model 1, and	the type I errors	for all tests
(q_1, q_2)	u	Size $(c = 0)$			Power ($c =$	0.1)		Power ($c =$	0.2)	
		Z_n	T_2	T_3	Z_n	T_2	T_3	Z_n	T_2	T_3
(6, 4)	20	0.0570	0.0518	0.0587	0.0666	0.0595	0.0746	0660.0	0.0940	0.1345
	50	0.0525	0.0567	0.0612	0.0851	0.0932	0.1060	0.2385	0.2594	0.2973
	100	0.0523	0.0621	0.0612	0.1398	0.1616	0.1685	0.5578	0.5942	0.6173
(18, 12)	50	0.0588	0.0523	0.0543	0.0791	0.0744	0.0951	0.1860	0.1888	0.3065
	100	0.0481	0.0496	0.0529	0.1190	0.1235	0.1509	0.5622	0.5799	0.6830
	150	0.0470	0.0488	0.0533	0.1840	0.1918	0.2185	0.8561	0.8663	0.9074
(36, 24)	100	0.0519	0.0524	0.0548	0.1050	0.1052	0.1477	0.4225	0.4536	0.6923
	150	0.0515	0.0504	0.0507	0.1647	0.1677	0.2192	0.8139	0.8312	0.9305
	200	0.0528	0.0531	0.0544	0.2503	0.2545	0.3109	0.9701	0.9732	0.9922
(60, 40)	150	0.0484	0.0482	0.0523	0.1339	0.1370	0.2201	0.6536	0.7049	0.9403
	200	0.0477	0.0473	0.0516	0.2112	0.2138	0.3110	0.9370	0.9486	0.9953
	300	0.0514	0.0510	0.0514	0.4073	0.4116	0.5225	1.0000	1.0000	1.0000

	cn.									
(q_1, q_2)	и	Size $(c = 0)$			Power $(c = 0)$	(1)		Power $(c = 0)$.(2)	
		Z_n	T_2	T_3	\mathbf{Z}_n	T_2	T_3	Z_n	T_2	T_3
(6, 4)	20	0.0570	0.0518	0.0601	0.0903	0.0836	0.0802	0.2168	0.2024	0.1464
	50	0.0525	0.0567	0.0627	0.1879	0.2028	0.1128	0.7439	0.7474	0.3421
	100	0.0523	0.0621	0.0605	0.4344	0.4667	0.1818	0.9944	0.9949	0.7049
(18, 12)	50	0.0588	0.0523	0.0541	0.1513	0.1391	0.0991	0.5664	0.4817	0.3172
	100	0.0481	0.0496	0.0526	0.4172	0.3888	0.1558	0.9967	0.9858	0.7095
	150	0.0470	0.0488	0.0511	0.7130	0.6808	0.2337	1.0000	1.0000	0.9262
(36, 24)	100	0.0519	0.0524	0.0564	0.2901	0.2596	0.1476	0.9434	0.8424	0.7062
	150	0.0515	0.0504	0.0511	0.6009	0.5238	0.2221	1.0000	0.9979	0.9400
	200	0.0528	0.0531	0.0548	0.8461	0.7755	0.3161	1.0000	1.0000	0.9940
(60, 40)	150	0.0484	0.0482	0.0515	0.4087	0.3467	0.2217	0.9927	0.9511	0.9438
	200	0.0477	0.0473	0.0494	0.7179	0.6008	0.3160	1.0000	0.9997	0.9958
	300	0.0514	0.0510	0.0524	0.9800	0.9344	0.5331	1.0000	1.0000	1.0000

$$\mathbb{E}(W_n^t) = \frac{\Gamma_p\left(\frac{n-1}{2}+t\right)}{\Gamma_p\left(\frac{n-1}{2}\right)} \prod_{i=1}^k \frac{\Gamma_{q_i}\left(\frac{n-1}{2}\right)}{\Gamma_{q_i}\left(\frac{n-1}{2}+t\right)}$$
(21)

for any t > (p - n)/2, where $\Gamma_p(x)$ is defined as (4.1).

To prove our main result in Theorem 1, we will establish the central limit theorem for log W_n by showing that the moment-generating function for properly normalized log W_n converges to that of a normal distribution. Since $\mathbb{E}(e^{t \log W_n}) = \mathbb{E}(W_n^t)$, Lemma 1 bridges the moment-generating functions for log W_n and the multivariate Gamma functions defined in (20). Before we proceed to prove Theorem 1, we will present several lemmas, some of which involve deliberate expansions of the Gamma functions.

Define

$$\xi(x) = -2(\log(1-x) + x), \quad x \in [0, 1).$$
(22)

 $\xi(x)$ is nonnegative in its domain. By the definition of $\xi(x)$, $\xi(0) = 0$, and $\xi'(x) = \frac{2x}{1-x}$. We have

$$\xi(x) = \xi(x) - \xi(0)$$
$$= \int_0^x \xi'(t) dt$$
$$= 2 \int_0^x \frac{t}{1-t} dt.$$

Substitute t = ux; then, dt = xdu. Therefore,

$$\xi(x) = 2 \int_0^1 \frac{ux^2}{1 - ux} du = 2x^2 \int_0^1 \frac{u}{1 - ux} du, \quad x \in [0, 1).$$
(23)

We also define

$$\eta(x) = \frac{\xi(x)}{x^2} = 2 \int_0^1 \frac{u}{1 - ux} \mathrm{d}u, \quad x \in [0, 1).$$

Then,

 $\eta(x) \ge 1$ is increasing in [0,1), and $\lim_{x \uparrow 1} \eta(x) = \infty$. (24)

One can easily verify

$$\frac{n}{q(n-q)} \le 2, \quad \text{for } 1 \le q < n \tag{25}$$

which together with (24) yields

$$\max_{1 \le i < n} \frac{\frac{1}{8\sqrt{\xi\left(\frac{i}{n}\right)}}}{\frac{n-i}{4}} = \frac{1}{2} \max_{1 \le i < n} \frac{\frac{n}{i(n-i)}}{\sqrt{\eta\left(\frac{i}{n}\right)}} \le 1.$$
(26)

Lemma 2 Let r_n be any sequence of positive integers and satisfies $r_n \to \infty$ and $r_n/n \to 0$ as $n \to \infty$. Then,

$$\lim_{n \to \infty} \max_{r_n \le q < n} \frac{\frac{q}{n(n-q)}}{\xi\left(\frac{q}{n}\right)} = 0,$$
(27)

and

$$\lim_{n \to \infty} \max_{r_n \le q < n} \frac{\frac{1}{(n-q)^2}}{\xi\left(\frac{q}{n}\right)} = 0.$$
 (28)

Moreover, we haven for any $\varepsilon > 0$

$$\lim_{n \to \infty} \max_{r_n \le q < n} \frac{\left(\frac{q}{n(n-q)}\right)^{1+\varepsilon}}{\xi\left(\frac{q}{n}\right)} = 0.$$
(29)

Proof It is easy to see that

$$\max_{r_n \le q \le n - r_n} \frac{n}{q(n-q)} = \frac{n}{r_n(n-r_n)} \to 0 \quad \text{as } n \to \infty.$$
(30)

Note that

$$\frac{\frac{q}{n(n-q)}}{\xi\left(\frac{q}{n}\right)} = \frac{\frac{n}{q(n-q)}}{\eta\left(\frac{q}{n}\right)}.$$

Then, it follows from (30) and (24) that

$$\max_{r_n \leq q \leq n-r_n} \frac{\frac{q}{n(n-q)}}{\xi\left(\frac{q}{n}\right)} = \max_{r_n \leq q \leq n-r_n} \frac{\frac{n}{q(n-q)}}{\eta\left(\frac{q}{n}\right)} \leq \frac{n}{r_n(n-r_n)} \to 0,$$

and from (25) that

$$\max_{n-r_n < q < n} \frac{\frac{q}{n(n-q)}}{\xi(\frac{q}{n})} \le \max_{n-r_n < q < n} \frac{2}{\eta(\frac{q}{n})} \le \frac{2}{\eta\left(1 - \frac{r_n}{n}\right)} \to 0$$

since $\eta(x) \to \infty$ as $x \uparrow 1$. Therefore, we obtain (27). Similarly, we can show (28). Finally, (29) follows from (27) and (25). This completes the proof of the lemma. \Box

Lemma 3 Let $p = p_n$ satisfy p < n and $p \to \infty$ as $n \to \infty$. Assume $k = k_n$ is a sequence of positive integers and q_1, \ldots, q_k are positive integers such that $p = \sum_{i=1}^k q_i \cdot \sigma_n^2$ is defined as in (6). Then,

$$\sigma_n^2 = n^2 \left(\xi\left(\frac{p}{n}\right) - \sum_{i=1}^k \xi\left(\frac{q_i}{n}\right) \right),\tag{31}$$

and

$$n^{2}\left(1-\sum_{i=1}^{k}\left(\frac{q_{i}}{p}\right)^{2}\right)\xi\left(\frac{p}{n}\right)\leq\sigma_{n}^{2}\leq n^{2}\xi\left(\frac{p}{n}\right).$$
(32)

Furthermore, if for some $\delta \in (0, 1)$, $\max_{1 \le i \le k} q_i \le \delta p$ for all large n, we have

$$n^{2}(1-\delta)\xi\left(\frac{p}{n}\right) \le \sigma_{n}^{2} \le n^{2}\xi\left(\frac{p}{n}\right).$$
(33)

Proof (31) is trivial by using the Definition (22) and the fact that $p = \sum_{i=1}^{k} q_i$. We will prove (32) and (33).

Since $\xi(x)$ is nonnegative in its domain, it follows that $\sigma_n^2 \le n^2 \xi(\frac{p}{n})$, that is, the first half of (32) is true. It follows from (23) that

$$\begin{split} \xi\left(\frac{p}{n}\right) &- \sum_{i=1}^{k} \xi\left(\frac{q_{i}}{n}\right) = 2\left(\frac{p}{n}\right)^{2} \int_{0}^{1} \frac{u}{1 - \frac{pu}{n}} du - \sum_{i=1}^{k} 2\left(\frac{q_{i}}{n}\right)^{2} \int_{0}^{1} \frac{u}{1 - \frac{q_{i}u}{n}} du \\ &\geq \left(\left(\frac{p}{n}\right)^{2} - \sum_{i=1}^{k} \left(\frac{q_{i}}{n}\right)^{2}\right) \times 2 \int_{0}^{1} \frac{u}{1 - \frac{pu}{n}} du \\ &= 2 \int_{0}^{1} \frac{u\left(\frac{p}{n}\right)^{2}}{1 - \frac{pu}{n}} (1 - \sum_{i=1}^{k} \left(\frac{q_{i}}{p}\right)^{2}) du \\ &= \left(1 - \sum_{i=1}^{k} \left(\frac{q_{i}}{p}\right)^{2}\right) \times 2 \int_{0}^{1} \frac{u\left(\frac{p}{n}\right)^{2}}{1 - \frac{pu}{n}} du \\ &= \left(1 - \sum_{i=1}^{k} \left(\frac{q_{i}}{p}\right)^{2}\right) \xi\left(\frac{p}{n}\right), \end{split}$$

which, together with (31), yields (32).

Since $\max_{1 \le i \le k} q_i \le \delta p$ and $\sum_{1 \le i \le k} q_i = p$,

$$\sum_{i=1}^k \left(\frac{q_i}{p}\right)^2 \le \frac{\max_{1\le i\le k} q_i}{p} \sum_{i=1}^k \frac{q_i}{p} = \frac{\max_{1\le i\le k} q_i}{p} \le \delta,$$

we can easily conclude (33) from (32). This completes the proof of the lemma.

Lemma 4 Let $\delta \in (0, 1)$ be any given number. Then,

$$\log \frac{\Gamma(x+b)}{\Gamma(x)} = (x+b)\log(x+b) - x\log x - b$$
$$-\left(\frac{1}{2x} + \frac{1}{12x^2}\right)b + \frac{b^2}{4x^2} + O\left(\frac{|b|^3 + 1}{x^3}\right)$$
(34)

931

holds uniformly on $b \in [-\delta x, \delta x]$ *as* $x \to \infty$ *. Furthermore,*

$$\log\left(\frac{\Gamma(x+b)}{\Gamma(x)} \cdot \frac{\Gamma(z)}{\Gamma(z+b)}\right) = b\left(\log x - \log z - \left(\frac{1}{2x} - \frac{1}{2z}\right) - \left(\frac{1}{12x^2} - \frac{1}{12z^2}\right)\right) + b^2\left(\frac{1}{2x} - \frac{1}{2z} + \frac{1}{4x^2} - \frac{1}{4z^2}\right) + O\left(\frac{|b|^3(x-z)}{xz^2} + \frac{|b|^3 + 1}{z^3}\right)$$
(35)

uniformly over $b \in [-\delta z, \delta z]$ and $x \ge z$ as $z \to \infty$.

Proof A similar expansion to (34) has been proved in Lemma A.1 in Jiang and Qi (2015b). We need to use the ;well-known Stirling formula (see, e.g., Ahlfors (1979)):

$$\log \Gamma(x) = \left(x - \frac{1}{2}\right) \log(x) - x + \frac{1}{2} \log(2\pi) + \frac{1}{12x} + O\left(\frac{1}{x^3}\right)$$
(36)

as $x \to \infty$. For any fixed $\delta \in (0, 1)$, we have that

$$\log \Gamma(x+b) - \log \Gamma(x)$$

$$= (x+b)\log(x+b) - x\log x - b - \frac{1}{2}\log\left(1+\frac{b}{x}\right)$$

$$+ \frac{1}{12}\left(\frac{1}{x+b} - \frac{1}{x}\right) + O\left(\frac{1}{x^3}\right)$$

uniformly on $b \in [-\delta x, \delta x]$ as $x \to \infty$. Then, (34) follows from the facts that

$$\log\left(1+\frac{b}{x}\right) = \frac{b}{x} - \frac{b^2}{2x^2} + O\left(\frac{|b|^3}{x^3}\right)$$

and

$$\frac{1}{x+b} - \frac{1}{x} = -\frac{b}{x^2} + O\left(\frac{b^2}{x^3}\right) = -\frac{b}{x^2} + O\left(\frac{|b|^3 + 1}{x^3}\right)$$

uniformly on $b \in [-\delta x, \delta x]$ as $x \to \infty$.

To prove (35), we have

$$c(x,b) := (x+b)\log(x+b) - x\log x - b - \left(\frac{1}{2x} + \frac{1}{12x^2}\right)b + \frac{b^2}{4x^2}$$
$$= \int_0^b \frac{d}{dt}((x+t)\log(x+t))dt - b - \left(\frac{1}{2x} + \frac{1}{12x^2}\right)b + \frac{b^2}{4x^2}$$

$$= \int_{0}^{b} \log(x+t) dt - \left(\frac{1}{2x} + \frac{1}{12x^{2}}\right) b + \frac{b^{2}}{4x^{2}}$$

$$= b \int_{0}^{1} \log(x+bv) dv - \left(\frac{1}{2x} + \frac{1}{12x^{2}}\right) b + \frac{b^{2}}{4x^{2}}$$

$$= b \log x + \frac{b^{2}}{2x} - \left(\frac{1}{2x} + \frac{1}{12x^{2}}\right) b + \frac{b^{2}}{4x^{2}}$$

$$+ b \int_{0}^{1} \left(\log\left(1 + \frac{bv}{x}\right) - \frac{bv}{x}\right) dv$$

$$= b \log x - \left(\frac{1}{2x} + \frac{1}{12x^{2}}\right) b + \left(\frac{1}{2x} + \frac{1}{4x^{2}}\right) b^{2}$$

$$- b \int_{0}^{1} \int_{0}^{bv/x} \frac{s}{1+s} ds dv$$

$$= b \log x - \left(\frac{1}{2x} + \frac{1}{12x^{2}}\right) b + \left(\frac{1}{2x} + \frac{1}{4x^{2}}\right) b^{2}$$

$$- b \int_{0}^{1} v^{2} \int_{0}^{b/x} \frac{u}{1+uv} du dv.$$

Then, we obtain from (34) that uniformly over $b \in [-\delta z, \delta z]$ and $x \ge z$

$$\log\left(\frac{\Gamma(x+b)}{\Gamma(x)} \cdot \frac{\Gamma(z)}{\Gamma(z+b)}\right)$$

= $c(x,b) - c(z,b) + O\left(\frac{|b|^3 + 1}{z^3}\right)$
= $b(\log x - \log z) - \left(\frac{1}{2x} + \frac{1}{12x^2}\right)b + \left(\frac{1}{2x} + \frac{1}{4x^2}\right)b^2 + \left(\frac{1}{2z} + \frac{1}{12z^2}\right)b$
 $- \left(\frac{1}{2z} + \frac{1}{4z^2}\right)b^2 - b\int_0^1 v^2 \int_{b/z}^{b/x} \frac{u}{1+uv} du dv + O\left(\frac{|b|^2 + 1}{z^3}\right),$

and (35) follows since

$$\left| b \int_0^1 v^2 \int_{b/z}^{b/x} \frac{u}{1+uv} \mathrm{d}u \mathrm{d}v \right| \le \frac{b^2}{(1-\delta)z} \int_0^1 v^2 \left| \int_{b/z}^{b/x} \mathrm{d}u \right| \mathrm{d}v = \frac{|b|^3 (x-z)}{3(1-\delta)xz^2}.$$

This completes the proof of the lemma.

Lemma 5 As $n \to \infty$,

$$\sum_{i=1}^{q} \left(\frac{1}{n-i} - \frac{1}{n-1} \right) = -\log\left(1 - \frac{q}{n}\right) - \frac{q}{n} + O\left(\frac{q}{n(n-q)}\right)$$
(37)

and

$$\sum_{i=1}^{q} \left(\log(n-1) - \log(n-i) \right) = \left(n - q - \frac{1}{2} \right) \log \left(1 - \frac{q}{n} \right) + \frac{(n-1)q}{n} + O\left(\frac{q}{n(n-q)} \right)$$
(38)

hold uniformly on $1 \le q < n$.

Proof By the partial sum of harmonic series,

$$\sum_{i=1}^{k} \frac{1}{i} = \log k + \gamma + \frac{1}{2k} - \psi(k),$$

where γ is the Euler–Mascheroni constant and $0 < \psi(k) < \frac{2}{k(k+1)}$. See, for example, Young (1991). Rewrite the above equation as

$$\sum_{i=1}^{k-1} \frac{1}{i} = \log k + \gamma - \frac{1}{2k} - \psi(k).$$
(39)

By applying (39), we have

$$\sum_{i=1}^{q} \left(\frac{1}{n-i} - \frac{1}{n-1} \right)$$

= $\sum_{i=1}^{n-1} \frac{1}{i} - \sum_{i=1}^{n-q-1} \frac{1}{i} - \frac{q}{n-1}$
= $-\log\left(1 - \frac{q}{n}\right) - \frac{q}{n} + \frac{q}{2n(n-q)} - \frac{q}{n(n-1)} + \psi(n-q) - \psi(n).$

By noting that $|\psi(n-q) - \psi(n)| \le \frac{2}{(n-q)(n-q+1)}$, to show (37), it suffices to show $\frac{1}{(n-q)(n-q+1)} = O(\frac{q}{n(n-q)})$. In fact, we have from (25) that

$$\frac{1}{(n-q)(n-q+1)} = \frac{n}{q(n-q)} \frac{q}{n(n-q+1)} < \frac{2q}{n(n-q)}.$$

This finishes the proof of (37).

To show (38), we apply the Stirling formula (36). By setting x = n and x = n - q, respectively, and then taking the difference, we have as $n \to \infty$,

$$\log \Gamma(n) - \log \Gamma(n-q) = \left(n - \frac{1}{2}\right) \log n - \left(n - q - \frac{1}{2}\right) \log(n-q) - q$$

$$+ \frac{1}{12} \left(\frac{1}{n} - \frac{1}{n-q} \right) + O\left(\frac{1}{(n-q)^2} \right)$$

$$= \left(n - \frac{1}{2} \right) \log n - \left(n - q - \frac{1}{2} \right) \log(n-q) - q - \frac{q}{12n(n-q)}$$

$$+ O\left(\frac{1}{(n-q)^2} \right)$$

$$= \left(n - \frac{1}{2} \right) \log n - \left(n - q - \frac{1}{2} \right) \log(n-q) - q + O\left(\frac{q}{n(n-q)} \right)$$

In the last step, we have used inequality (25) to get

$$\frac{1}{(n-q)^2} = \frac{n}{q(n-q)} \frac{q}{n(n-q)} \leq \frac{2q}{n(n-q)}$$

Since $\Gamma(m) = (m-1)!$ for any integer $m \ge 1$, we have

$$\begin{split} &\sum_{i=1}^{q} (\log(n-1) - \log(n-i)) \\ &= q \log(n-1) - (\log \Gamma(n) - \log \Gamma(n-q)) \\ &= q \log(n-1) - \left(n - \frac{1}{2}\right) \log n + \left(n - q - \frac{1}{2}\right) \log(n-q) + q \\ &+ O\left(\frac{q}{n(n-q)}\right) \\ &= \left(n - q - \frac{1}{2}\right) \log\left(1 - \frac{q}{n}\right) + q + q \log\left(1 - \frac{1}{n}\right) + O\left(\frac{q}{n(n-q)}\right) \\ &= \left(n - q - \frac{1}{2}\right) \log\left(1 - \frac{q}{n}\right) + q \left(1 - \frac{1}{n}\right) + O\left(\frac{q}{n(n-q)}\right), \end{split}$$

proving (38).

Lemma 6 Assume $\xi(x)$ is defined as in (22). Then, there exists a sequence of integers $\{s_n\}$ satisfying $s_n \uparrow \infty$ and $s_n = O(\log n)$ such that as $n \to \infty$

$$\log\left(\left(\frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n-1}{2}+t)}\right)^{q} \frac{\Gamma_{q}\left(\frac{n-1}{2}+t\right)}{\Gamma_{q}\left(\frac{n-1}{2}\right)}\right)$$
(40)
$$= t\left(\left(q-n+\frac{3}{2}\right)\log\left(1-\frac{q}{n}\right) - \frac{n-2}{n}q + \frac{1}{3}\frac{q}{(n-1)^{2}} - \frac{1}{3}b(n,q)\right) + \frac{t^{2}}{2}\left(\xi\left(\frac{q}{n}\right) + 2b(n,q) - \frac{2q}{(n-1)^{2}}\right) + O\left(\frac{q^{2}|t|^{3}}{n^{2}(n-q)} + \frac{q(t^{2}+1)}{n(n-q)}\right)$$

Deringer

holds uniformly over $|t| \leq \frac{n-q}{4}$ and $1 \leq q \leq n - s_n$, and

$$\log\left(\left(\frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n-1}{2}+t)}\right)^{q} \frac{\Gamma_{q}\left(\frac{n-1}{2}+t\right)}{\Gamma_{q}\left(\frac{n-1}{2}\right)}\right)$$
(41)
$$= t\left(\left(q-n+\frac{3}{2}\right)\log(1-\frac{q}{n}) - \frac{n-2}{n}q + \frac{1}{3}\frac{q}{(n-1)^{2}} - \frac{1}{3}b(n,q)\right) + \frac{t^{2}}{2}\left(\xi\left(\frac{q}{n}\right) + 2b(n,q) - \frac{2q}{(n-1)^{2}}\right) + O\left(\frac{|t|^{3}+1}{s_{n}} + \frac{(|t|+t^{2})(\xi(1-\frac{s_{n}}{n}))^{1/2}}{s_{n}^{2}}\right)$$

holds uniformly over $|t| \leq \frac{n-q}{4}$ and $n - s_n < q < n$.

Proof Let $\Gamma'(x)$ denote the first derivative of $\Gamma(x)$. Define

$$c_{n,i} = \frac{1}{\sqrt{\xi(1-\frac{i}{n})}} \left(\sup_{1/4 \le x \le i+1} \frac{|\Gamma'(x)|}{\Gamma(x)} + i + 4 \right)$$
(42)

for $1 \le i < n$. Observe that for each fixed k, $\lim_{n\to\infty} c_{n,k} = 0$. For each $k \ge 1$, there exists a positive integer $j_k \ge 2$ such that $c_{j,k} < \frac{1}{k^3}$ for all $j \ge j_k$. We can choose $j_k \ge 2j_{k-1}$ for each $k \ge 2$. From this, we can conclude that $j_k \ge 2^k$ for all $k \ge 1$. Define $s_n = k$ if $j_k \le n < j_{k+1}$. Then, $s_n \to \infty$ as $n \to \infty$, and $s_n = O(\log n)$ since $2^{s_n} \le n$. Moreover,

$$c_{n,s_n} \le \frac{1}{s_n^3}, \ s_n \le \frac{\left(\xi\left(1 - \frac{s_n}{n}\right)\right)^{1/2}}{s_n^3} \text{ and } \frac{\left(\xi\left(1 - \frac{s_n}{n}\right)\right)^{1/2}}{s_n^2} \ge 1$$
 (43)

for all large n.

We first apply (35) in Lemma 4 to get

$$\begin{split} \log\left(\frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} \cdot \frac{\Gamma\left(\frac{n-i}{2}\right)}{\Gamma\left(\frac{n-i}{2}+t\right)}\right) \\ &= t\left(\left(\log\frac{n-1}{2} - \log\frac{n-i}{2}\right) - \left(\frac{1}{n-1} - \frac{1}{n-i}\right) - \left(\frac{1}{3(n-1)^2} - \frac{1}{3(n-i)^2}\right)\right) \\ &+ t^2\left(\left(\frac{1}{n-1} - \frac{1}{n-i}\right) + \left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right)\right) + O\left(\frac{|t|^3(i-1)}{n(n-i)^2} + \frac{|t|^3+1}{(n-i)^3}\right) \\ &= t\left((\log(n-1) - \log(n-i)) - \left(\frac{1}{n-1} - \frac{1}{n-i}\right) - \frac{1}{3}\left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right)\right) \\ &+ t^2\left(\left(\frac{1}{n-1} - \frac{1}{n-i}\right) + \left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right)\right) + O\left(\frac{|t|^3(i-1)}{n(n-i)^2} + \frac{|t|^3+1}{(n-i)^3}\right) \\ \end{split}$$

uniformly over $|t| \le \frac{n-i}{4}$ and $1 \le i \le n - s_n$ as $n \to \infty$. Rewrite the above equation as

$$\log\left(\frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} \cdot \frac{\Gamma\left(\frac{n-i}{2}\right)}{\Gamma\left(\frac{n-i}{2}+t\right)}\right)$$

= $t\left((\log(n-1) - \log(n-i)) - \left(\frac{1}{n-1} - \frac{1}{n-i}\right) - \frac{1}{3}\left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right)\right)$
+ $t^2\left(\left(\frac{1}{n-1} - \frac{1}{n-i}\right) + \left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right)\right) + d_{n,i}(t)$ (44)

for $|t| \leq \frac{n-i}{4}$ and $1 \leq i \leq n - s_n$, where

$$|d_{n,i}(t)| \le K \left(\frac{|t|^3(i-1)}{n(n-i)^2} + \frac{|t|^3 + 1}{(n-i)^3} \right) \quad \text{for } |t| \le \frac{n-i}{4}$$
(45)

for all $1 \le i \le n - s_n$ for some constant K > 0.

From the definition of $\Gamma_p(x)$ in (20), we have

$$\log\left(\left(\frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n-1}{2}+t)}\right)^{q} \frac{\Gamma_{q}\left(\frac{n-1}{2}+t\right)}{\Gamma_{q}\left(\frac{n-1}{2}\right)}\right) = -\sum_{i=1}^{q} \log\left(\frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} \frac{\Gamma\left(\frac{n-i}{2}\right)}{\Gamma\left(\frac{n-i}{2}+t\right)}\right).$$
(46)

Note that the range of t in (44) and (45) depends on i for $1 \le i \le n - s_n$. For each $1 \le q \le n - s_n$, we limit t in a common range $|t| \le \frac{n-q}{4}$ for the first q terms with i = 1, ..., q. First, note that

$$\sum_{i=1}^{q} \frac{1}{(n-i)^2} \le \sum_{i=1}^{q} \frac{1}{(n-i)(n-i-1)}$$
$$\le \sum_{i=1}^{q} \left(\frac{1}{(n-i-1)} - \frac{1}{n-i}\right)$$
$$= \frac{1}{n-q-1} - \frac{1}{n-1}$$
$$= \frac{q}{(n-1)(n-q-1)}$$
$$\le \frac{2q}{n(n-q)}$$

for all large n. From (45), we have

$$\sum_{i=1}^{q} |d_{n,i}(t)| \le K \sum_{i=1}^{q} \left(\frac{|t|^{3}(i-1)}{n(n-i)^{2}} + \frac{|t|^{3}+1}{(n-i)^{3}} \right)$$
$$\le K \sum_{i=1}^{q} \left(\frac{|t|^{3}(i-1)}{n(n-i)^{2}} + \frac{|t|^{2}+1}{(n-i)^{2}} \right)$$

$$\leq K \left(\frac{q}{n}|t|^{3} + t^{2} + 1\right) \sum_{i=1}^{q} \frac{1}{(n-i)^{2}}$$
$$\leq 2K \left(\frac{q}{n}|t|^{3} + t^{2} + 1\right) \frac{2q}{n(n-q)}$$
(47)

for all large n, which coupled with Eq. (44), Definition (7), and Lemma 5 yields that

$$\begin{split} &\sum_{i=1}^{q} \log \left(\frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} \cdot \frac{\Gamma\left(\frac{n-i}{2}\right)}{\Gamma\left(\frac{n-i}{2}+t\right)} \right) \\ &= t \sum_{i=1}^{q} \left((\log(n-1) - \log(n-i)) - \left(\frac{1}{n-1} - \frac{1}{n-i}\right) - \frac{1}{3} \left(\frac{1}{(n-1)^2} - \frac{1}{(n-1)^2}\right) \right) \\ &- \frac{1}{(n-i)^2} \right) + t^2 \sum_{i=1}^{q} \left(\left(\frac{1}{n-1} - \frac{1}{n-i}\right) + \left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right) \right) \\ &+ O\left(\frac{q^2 |t|^3}{n^2 (n-q)} + \frac{q(t^2+1)}{n (n-q)}\right) \\ &= t \left(\left(n-q-\frac{3}{2}\right) \log\left(1-\frac{q}{n}\right) + \frac{(n-2)q}{n} - \frac{1}{3} \frac{q}{(n-1)^2} + \frac{1}{3} b(n,q) \right) \\ &- \frac{t^2}{2} \left(\xi\left(\frac{q}{n}\right) + 2b(n,q) - \frac{2q}{(n-1)^2} \right) + O\left(\frac{q^2 |t|^3}{n^2 (n-q)} + \frac{q(t^2+|t|+1)}{n (n-q)} \right) \\ &= t \left(\left(n-q-\frac{3}{2}\right) \log\left(1-\frac{q}{n}\right) + \frac{(n-2)q}{n} - \frac{1}{3} \frac{q}{(n-1)^2} + \frac{1}{3} b(n,q) \right) \\ &- \frac{t^2}{2} \left(\xi\left(\frac{q}{n}\right) + 2b(n,q) - \frac{2q}{(n-1)^2} \right) + O\left(\frac{q^2 |t|^3}{n^2 (n-q)} + \frac{q(t^2+1)}{n (n-q)} \right) \end{split}$$

holds uniformly over $|t| \le \frac{n-q}{4}$ and $1 \le q \le n - s_n$. In the last step, we drop the |t| term inside the big "O" since $|t| \le \frac{t^2+1}{2}$. This proves (40) by using (46). Next, we will prove (41). First, we show that

$$\log\left(\frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} \cdot \frac{\Gamma\left(\frac{n-i}{2}\right)}{\Gamma\left(\frac{n-i}{2}+t\right)}\right)$$
(48)
= $t\left((\log(n-1) - \log(n-i)) - \left(\frac{1}{n-1} - \frac{1}{n-i}\right) - \frac{1}{3}\left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right)\right)$
+ $t^2\left(\left(\frac{1}{n-1} - \frac{1}{n-i}\right) + \left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right)\right)$
+ $O\left(\frac{|t|\left(\xi\left(1 - \frac{s_n}{n}\right)\right)^{1/2}}{s_n^3} + t^2 + \frac{1}{n}\right)$

uniformly over $|t| \le \frac{n-i}{4}$ and $n - s_n < i < n$ as $n \to \infty$. Using (42) and (43), we conclude that

$$\begin{aligned} \left| \log \Gamma\left(\frac{n-i}{2}+t\right) - \log \Gamma\left(\frac{n-i}{2}\right) \right| &= \left| \int_{\frac{n-i}{2}}^{\frac{n-i}{2}+t} \frac{\Gamma'(x)}{\Gamma(x)} dx \right| \\ &\leq \left| t \right| \sup_{1/4 \leq x \leq s_n} \frac{\left| \Gamma'(x) \right|}{\Gamma(x)} \\ &\leq \left| t \right| \left(\xi \left(1-\frac{s_n}{n}\right) \right)^{1/2} c_{n,s_n} \\ &\leq \frac{\left| t \right| \left(\xi \left(1-\frac{s_n}{n}\right) \right)^{1/2}}{s_n^3} \end{aligned}$$

for $|t| \leq \frac{n-i}{4}$ and $n - s_n \leq i < n$.

We can use (34) and estimate c(x, b) defined in the proof of Lemma 4 to derive that

$$\log \frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} = t\left(\log \frac{n-1}{2} - \frac{1}{n-1} - \frac{1}{3(n-1)^2}\right) + t^2\left(\frac{1}{n-1} + \frac{1}{(n-1)^2}\right) + O\left(\frac{|t|^3}{n^2} + \frac{1}{n^3}\right)$$

uniformly over $|t| \leq \frac{n-1}{4}$ as $n \to \infty$. Since for any $n - s_n \leq i < n$, any $|t| \leq \frac{n-i}{4}$ is also bounded by $\frac{n-1}{4}$, the above approximation holds true uniformly for $|t| \leq \frac{n-i}{4}$ and $n-s_n \le i < n \text{ as } n \to \infty$. Also note that uniformly over $|t| \le \frac{n-i}{4}$ and $n-s_n \le i < n$

$$\left| t \left(\log \frac{n-i}{2} - \frac{1}{n-i} - \frac{1}{3(n-i)^2} \right) + t^2 \left(\frac{1}{n-i} + \frac{1}{(n-i)^2} \right) \right|$$

= $O(s_n |t| + t^2) = O\left(\frac{|t| (\xi(1 - \frac{s_n}{n}))^{1/2}}{s_n^3} + t^2 \right)$

for all large n. We have used the second inequality in (43). Immediately, (48) follows from the above three estimates.

Now rewrite (48) as

$$\log\left(\frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} \cdot \frac{\Gamma\left(\frac{n-i}{2}\right)}{\Gamma\left(\frac{n-i}{2}+t\right)}\right)$$
(49)

$$= t \Big((\log(n-1) - \log(n-i)) - \Big(\frac{1}{n-1} - \frac{1}{n-i}\Big) - \frac{1}{3}\Big(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\Big) \Big) \\ + t^2 \Big(\Big(\frac{1}{n-1} - \frac{1}{n-i}\Big) + \Big(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\Big)\Big) + d_{n,i}(t)$$

for $|t| \leq \frac{n-i}{4}$ and $n - s_n < i < n$, where

$$|d_{n,i}(t)| = O\left(\frac{|t|\left(\xi\left(1-\frac{s_n}{n}\right)\right)^{1/2}}{s_n^3} + t^2 + \frac{1}{n}\right)$$
(50)

holds uniformly over $|t| \leq \frac{n-i}{4}$ and $n - s_n < i < n$ as $n \to \infty$.

For any $n - s_n < q < n$, limit *t* in the common range of $|t| \le \frac{n-q}{4}$ so that we can apply both (45) and (50) for different *i* with $1 \le i \le q$ and estimate $\sum_{i=1}^{q} |d_{n,i}(t)|$. For the sum of the first $n - s_n$ terms, we apply (45) and get

$$\sum_{i=1}^{n-s_n} |d_{n,i}(t)| \le 2K \left(\frac{n-s_n}{n} |t|^3 + t^2 + 1 \right) \frac{2(n-s_n)}{ns_n} \le \frac{8K(|t|^3 + 1)}{s_n}$$

which is essentially the estimation in (47) with the choice $q = n - s_n$, which together with (50) and (43) yields

$$\sum_{i=1}^{q} |d_{n,i}(t)| = \sum_{i=1}^{n-s_n} |d_{n,i}(t)| + \sum_{i=n-s_n+1}^{q} |d_{n,i}(t)|$$

$$\leq \frac{8K(|t|^3 + 1)}{s_n} + O\left(\frac{|t|(\xi(1 - \frac{s_n}{n}))^{1/2}}{s_n^2} + s_n t^2 + \frac{s_n}{n}\right)$$

$$= O\left(\frac{|t|^3 + 1}{s_n} + \frac{(|t| + t^2)(\xi(1 - \frac{s_n}{n}))^{1/2}}{s_n^2}\right)$$

uniformly over $|t| \leq \frac{n-q}{4}$ and $n - s_n < q < n$.

Therefore, it follows from the above estimate and Lemma 5 that

$$\begin{split} &\sum_{i=1}^{q} \log \left(\frac{\Gamma\left(\frac{n-1}{2}+t\right)}{\Gamma\left(\frac{n-1}{2}\right)} \cdot \frac{\Gamma\left(\frac{n-i}{2}\right)}{\Gamma\left(\frac{n-i}{2}+t\right)} \right) \\ &= t \sum_{i=1}^{q} \left((\log(n-1) - \log(n-i)) - \left(\frac{1}{n-1} - \frac{1}{n-i}\right) \right) \\ &- \frac{1}{3} \left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right) \right) \\ &+ t^2 \sum_{i=1}^{q} \left(\left(\frac{1}{n-1} - \frac{1}{n-i}\right) + \left(\frac{1}{(n-1)^2} - \frac{1}{(n-i)^2}\right) \right) + \sum_{i=1}^{q} d_{n,i}(t) \\ &= t \left(\left(n-q-\frac{3}{2}\right) \log \left(1-\frac{q}{n}\right) + \frac{(n-2)q}{n} - \frac{1}{3} \frac{q}{(n-1)^2} + \frac{1}{3} b(n,q) \right) \\ &- \frac{t^2}{2} \left(\xi\left(\frac{q}{n}\right) + 2b(n,q) - \frac{2q}{(n-1)^2} \right) + O\left(\frac{|t|^3+1}{s_n} + \frac{(|t|+t^2) \left(\xi\left(1-\frac{s_n}{n}\right)\right)^{1/2}}{s_n^2} \right) \end{split}$$

uniformly over $|t| \le \frac{n-q}{4}$ and $n - s_n < q < n$. This together with (46) proves (41). The proof of the lemma is completed.

4.2 Proof of Theorem 1

We will first show (8), that is,

$$\frac{\bar{\mu}_n - \mu_n}{\bar{\sigma}_n} \to 0 \quad \text{and} \quad \frac{\sigma_n}{\bar{\sigma}_n} \to 1 \quad \text{as } n \to \infty.$$

In fact, it suffices to show that

$$\frac{n(b(n, p) - \sum_{i=1}^{k} b(n, q_i))}{\sigma_n} \to 0 \quad \text{as } n \to \infty.$$
(51)

Since $0 \le b(n, p) - \sum_{i=1}^{k} b(n, q_i) \le \sum_{j=n-p}^{\infty} \frac{1}{j^2} = O(\frac{1}{n-p})$, we have

$$\frac{b(n, p) - \sum_{i=1}^{k} b(n, q_i)}{\sqrt{\xi\left(\frac{p}{n}\right)}} = O\left(\frac{\frac{1}{n-p}}{\sqrt{\xi\left(\frac{p}{n}\right)}}\right) = O\left(\sqrt{\frac{\frac{1}{(n-p)^2}}{\xi\left(\frac{p}{n}\right)}}\right) \to 0$$

from (28) in Lemma 2. Moreover, since we have

$$\frac{1}{\sqrt{\xi\left(\frac{p}{n}\right)}} \le \frac{n}{\sigma_n} \le \frac{1}{\sqrt{1-\delta}} \frac{1}{\sqrt{\xi\left(\frac{p}{n}\right)}}$$
(52)

from (32) in Lemma 3, (51) is obtained.

Set $V_n = -2 \log \Lambda_n$. Then, it follows from (2) that $V_n = -n \log W_n$, and $-(V_n - \bar{\mu}_n)/\bar{\sigma}_n = \frac{n}{\bar{\sigma}_n} \log W_n + \frac{\bar{\mu}_n}{\bar{\sigma}_n}$. To show (3), it suffices to show the moment-generating function of $-(V_n - \bar{\mu}_n)/\bar{\sigma}_n$ converges to that of the standard normal in a neighborhood of zero, that is, for some $\delta_1 > 0$

$$\mathbb{E}\exp(-\frac{V_n-\bar{\mu}_n}{\bar{\sigma}_n}s) = \mathbb{E}(W_n^{ns/\bar{\sigma}_n})e^{\bar{\mu}_n s/\bar{\sigma}_n} \to e^{\frac{s^2}{2}}, \quad |s| \le \delta_1$$

or equivalently

$$\log \mathbb{E}(W_n^{ns/\bar{\sigma}_n}) + \frac{\bar{\mu}_n s}{\bar{\sigma}_n} \to \frac{s^2}{2}, \quad |s| \le \delta_1.$$
(53)

In the proof below, we will choose $\delta_1 = \frac{\sqrt{1-\delta}}{8}$ and assume that $|s| \le \delta_1$. Set $t = \frac{ns}{\sigma_n}$. Since

$$\frac{1}{\sqrt{\xi(\frac{p}{n})}} \le 2(n-p)$$

from (26), we have $|t| = \frac{n|s|}{\bar{\sigma}_n} \le \frac{n|s|}{\sigma_n} \le \frac{n-p}{4}$. Thus, we can apply (21) in Lemma 1. It is trivial that $|t| \le \frac{n-q_i}{4}$ for all $0 \le i \le k$, and we can apply Lemma 6 with q = p, and $q = q_i$ uniformly over $1 \le i \le k$. Define

$$\begin{split} R(q) &= \log\left(\left(\frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n-1}{2}+t\right)}\right)^{q} \frac{\Gamma_{q}\left(\frac{n-1}{2}+t\right)}{\Gamma_{q}\left(\frac{n-1}{2}\right)}\right) \\ &- t\left(\left(q-n+\frac{3}{2}\right)\log\left(1-\frac{q}{n}\right) - \frac{n-2}{n}q + \frac{1}{3}\frac{q}{(n-1)^{2}} - \frac{1}{3}b(n,q)\right) \\ &- \frac{t^{2}}{2}\left(\xi\left(\frac{q}{n}\right) + 2b(n,q) - \frac{2q}{(n-1)^{2}}\right). \end{split}$$

Then, it follows from (21), (5), and (31) that

$$\log \mathbb{E}(W_n^t) = \log \frac{\Gamma_p\left(\frac{n-1}{2} + t\right)}{\Gamma_p\left(\frac{n-1}{2}\right)} - \sum_{i=1}^k \log \frac{\Gamma_{q_i}\left(\frac{n-1}{2} + t\right)}{\Gamma_{q_i}\left(\frac{n-1}{2}\right)} \\ = \log \left(\left(\frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n-1}{2} + t\right)} \right)^p \frac{\Gamma_p\left(\frac{n-1}{2} + t\right)}{\Gamma_p\left(\frac{n-1}{2}\right)} \right) \\ - \sum_{i=1}^k \log \left(\left(\frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n-1}{2} + t\right)} \right)^{q_i} \frac{\Gamma_{q_i}\left(\frac{n-1}{2} + t\right)}{\Gamma_{q_i}\left(\frac{n-1}{2}\right)} \right) \\ = -\frac{t\bar{\mu}_n}{n} + \frac{t^2\bar{\sigma}_n^2}{2n^2} + R(p) - \sum_{i=1}^k R(q_i) \\ = -\frac{\bar{\mu}_n s}{\bar{\sigma}_n} + \frac{s^2}{2} + R(p) - \sum_{i=1}^k R(q_i).$$

Therefore, in order to prove (53), it suffices to show that

$$\lim_{n \to \infty} \sum_{i=1}^{k} R(q_i) = 0 \quad \text{and} \quad \lim_{n \to \infty} R(p) = 0.$$
(54)

Let s_n be a sequence of positive integers defined in Lemma 6, satisfying $s_n \uparrow \infty$ and $s_n = O(\log n)$ as $n \to \infty$. Since $\max_{1 \le i \le k} q_i \le \delta n \le n - s_n$ for all large *n*, we have from (40) that

$$\sum_{1 \le i \le k} |R(q_i)| = O\left(\sum_{i=1}^k \frac{q_i^2 |t|^3}{n^2 (n-q_i)} + \sum_{i=1}^k \frac{q_i (t^2+1)}{n(n-q_i)}\right)$$
$$= O\left(\frac{p|t|^3}{n^2 \max(n-p, s_n)} \sum_{i=1}^k q_i + \frac{t^2+1}{n \max(n-p, s_n)} \sum_{i=1}^k q_i\right)$$

$$= O\left(\frac{p^2|t|^3}{n^2 \max(n-p,s_n)} + \frac{p(t^2+1)}{n \max(n-p,s_n)}\right)$$

= $O\left(\frac{p^2|t|^3}{n^2 \max(n-p,s_n)} + \frac{pt^2}{n \max(n-p,s_n)} + \frac{p}{n \max(n-p,s_n)}\right).$

In the above estimation, we have used the following facts:

$$\max_{1 \le i \le k} q_i \le \min(p, n - s_n), \quad \sum_{i=1}^k q_i = p, \quad \min_{1 \le i \le k} (n - q_i) \ge \max(n - p, s_n).$$

Set $r_n = \min(p_n, s_n)$. Then, $r_n \to \infty$ and $r_n/n \to 0$ as $n \to \infty$. From (52) and (24),

$$|t| = \frac{n}{\tilde{\sigma}_n} |s| \le \frac{n}{\sigma_n} |s| \le \frac{1}{\sqrt{1-\delta}} \frac{1}{\sqrt{\xi\left(\frac{p}{n}\right)}} \frac{\sqrt{1-\delta}}{8} = \frac{1}{8} \frac{1}{\sqrt{\eta\left(\frac{p}{n}\right)}} \frac{n}{p} \le \frac{n}{8p},$$

and thus we have

$$\frac{p^{2}|t|^{3}}{n^{2}\max(n-p,s_{n})} \leq \frac{p^{2}}{n^{2}}\frac{n^{3}}{8^{3}p^{3}}\frac{1}{\max(n-p,s_{n})}$$

$$= \frac{1}{8^{3}}\frac{n}{p\max(n-p,s_{n})}(I_{(p\leq n-s_{n})} + I_{(p>n-s_{n})})$$

$$= \frac{1}{8^{3}}\left(\frac{n}{p(n-p)}I_{(p\leq n-s_{n})} + \frac{n}{ps_{n}}I_{(p>n-s_{n})}\right)$$

$$\leq \frac{1}{8^{3}}\left(\frac{n}{p(n-p)}I_{(p\leq n-s_{n})} + \frac{n}{(n-s_{n})s_{n}}I_{(p>n-s_{n})}\right)$$

$$\leq \frac{2}{8^{3}}\max_{r_{n}\leq q\leq n-r_{n}}\frac{n}{q(n-q)}$$

$$\to 0$$

as $n \to \infty$ from (30). Similarly, we have

$$\frac{pt^2}{n\max(n-p,s_n)} \le \frac{2}{8^2} \max_{r_n \le q \le n-r_n} \frac{n}{q(n-q)} \to 0$$

and

$$\frac{p}{n\max(n-p,s_n)} \le \frac{1}{s_n} \to 0.$$

Hence, we have shown that $\lim_{n\to\infty} \sum_{1\leq i\leq k} |R(q_i)| = 0$. Similarly, we can show that $\lim_{n\to\infty} R(p)I(p\leq n-s_n) = 0$.

To finish the proof of (54), we need to show that $\lim_{n\to\infty} R(p)I(n-s_n .$ $In fact, if <math>n - s_n for large$ *n*, we have

$$|t| = \frac{n}{\bar{\sigma}_n} |s| \le \frac{n}{\sigma_n} |s| \le \frac{|s|}{(1-\delta)^{1/2}} \frac{1}{\sqrt{\xi\left(\frac{p}{n}\right)}} \le \frac{|s|}{(1-\delta)^{1/2}} \frac{1}{\sqrt{\xi\left(1-\frac{s_n}{n}\right)}} \to 0$$

and thus

$$R(p) = O\left(\frac{|t|^3 + 1}{s_n} + \frac{(|t| + t^2)(\xi(1 - \frac{s_n}{n}))^{1/2}}{s_n^2}\right) = O\left(\frac{1}{s_n}\right) \to 0$$

as $n \to \infty$. Consequently, we have proved (54). The proof of the theorem is completed.

4.3 Proof of Theorem 2

To prove (13), it suffices to show that for any subsequence $\{n'\}$ of $\{n\}$, there is a further subsequence $\{n''\}$ such that (13) holds along $\{n''\}$. The subsequence $\{n''\}$ can be selected in a way that both the limits of $k_{n''}$ and $p_{n''}$ exist, and both the limits can be infinity. For the sake of simplicity, we can assume both the limits of k_n and p_n exist along the entire sequence and prove (13) holds. Note that if the limit of a sequence of integers is finite, the sequence takes a constant value ultimately. We will show (13) under each of the following two assumptions:

Case 1 $p_n = p$ and $k_n = k$ for all large *n*, where both *p* and *k* are fixed integers; *Case 2* $\lim_{n\to\infty} p_n = \infty$.

Proof for Case 1. We can assume q_1, \ldots, q_k are fixed integers. Otherwise, we can use subsequential limit argument since q_1, \ldots, q_k are bounded by p and their subsequential limits always exist. Thus, under case 1, (9) holds. Since ρ defined in (11) converges to one, we have $-2 \log \Lambda_n$ converges in distribution to a chi-square distribution with f degrees of freedom. Review Z_n in (12). To prove (13), it suffices to verify that

$$\lim_{n \to \infty} \frac{2f_n}{\bar{\sigma}_n^2} = 1 \text{ and } \lim_{n \to \infty} \left(f_n - \bar{\mu}_n \sqrt{\frac{2f_n}{\bar{\sigma}_n^2}} \right) = 0.$$
(55)

Note that $b(n, p) - \sum_{i=1}^{k} b(n, q_i) = o(\frac{1}{n^2})$. Then, by using Taylor's expansion, we have from (6) that

$$\begin{split} \bar{\sigma}_n^2 &= \sigma_n^2 + o(1) \\ &= 2n^2 \Big(\frac{p}{n} + \frac{1}{2} \Big(\frac{p}{n}\Big)^2 - \sum_{i=1}^k \Big(\frac{q_i}{n} + \frac{1}{2} \Big(\frac{q_i}{n}\Big)^2\Big) + O\Big(\frac{1}{n^3}\Big)\Big) + o(1) \\ &= p^2 - \sum_{i=1}^k q_i^2 + o(1) \\ &= 2f_n + o(1), \end{split}$$

🖉 Springer

which implies the first limit in Eq. (55). Similarly, we have from (5) that

$$\begin{split} \bar{\mu}_n &= \mu_n + o\left(\frac{1}{n}\right) \\ &= n\Big(\sum_{i=1}^k \left(q_i - n + \frac{3}{2}\right)\Big(-\frac{q_i}{n} - \frac{1}{2}\Big(\frac{q_i}{n}\Big)^2 + O\Big(\frac{1}{n^3}\Big)\Big) - \Big(p - n + \frac{3}{2}\Big) \\ &\quad \Big(\frac{p}{n} + \frac{1}{2}\Big(\frac{p}{n}\Big)^2 + O\Big(\frac{1}{n^3}\Big)\Big)\Big) + o\left(\frac{1}{n}\right) \\ &= \frac{1}{2}\left(p^2 - \sum_{i=1}^k q_i^2\right) + O\left(\frac{1}{n}\right) \\ &= f_n + O\left(\frac{1}{n}\right), \end{split}$$

which yields that $f_n - \bar{\mu}_n \sqrt{\frac{2f_n}{\bar{\sigma}_n^2}} = f_n - f_n(1 + o(1)) = o(1)$ since f_n is bounded. This proves the second limit in (55).

Proof for Case 2. Note that (13) is equivalent to

$$\lim_{n \to \infty} \sup_{-\infty < x < \infty} \left| P\left(\frac{Z_n - f_n}{\sqrt{2f_n}} \le x\right) - P\left(\frac{\chi_{f_n}^2 - f_n}{\sqrt{2f_n}} \le x\right) \right| = 0.$$
(56)

Since $\chi_{f_n}^2$ can be written as a sum of f_n independent chi-squared random variables each having one degree of freedom, it follows from the central limit theorem that

$$\lim_{n \to \infty} \sup_{-\infty < x < \infty} \left| P\left(\frac{\chi_{f_n}^2 - f_n}{\sqrt{2f_n}} \le x\right) - \Phi(x) \right| = 0.$$

Therefore, in order to show (56), we only need to show that

$$\frac{Z_n - f_n}{\sqrt{2f_n}}$$
 converges in distribution to $N(0, 1)$,

which follows from Theorem 1 since

$$\frac{Z_n - f_n}{\sqrt{2f_n}} = \frac{-2\log\Lambda_n - \bar{\mu}_n}{\bar{\sigma}_n}$$

This completes the proof of Theorem 2.

Acknowledgements The authors would like to thank three anonymous reviewers for their constructive comments and suggestions that have led to improvements in the paper. Qi's research was supported by NSF Grant DMS-1005345, and Wang's research was supported by NSFC Grant No. 11671021, NSFC Grant No. 11471222 and Foundation of Beijing Education Bureau Grant No. 201510028002.

References

- Ahlfors, L. V. (1979). Complex analysis: An introduction to the theory of analytic functions of one complex variable (3rd ed.). New York: McGraw-Hill.
- Bai, Z., Jiang, D., Yao, J., Zheng, S. (2009). Corrections to LRT on large dimensional covariance matrix by RMT. Annals of Statistics, 37(6B), 3822–3840.
- Bao, Z., Hu, J., Pan, G., Zhou, W. (2017). Test of independence for high-dimensional random vectors based on freeness in block correlation matrices. *Electronic Journal of Statistics*, 11(1), 1527–1548.
- Chen, S. X., Zhang, L. X., Zhong, P. S. (2010). Tests for high-dimensional covariance matrices. *Journal of the American Statistical Association*, 105(490), 810–819.
- Jiang, D., Jiang, T., Yang, F. (2012). Likelihood ratio tests for covariance matrices of high-dimensional normal distributions. *Journal of Statistical Planning and Inference*, 142(8), 2241–2256.
- Jiang, D., Bai, Z., Zheng, S. (2013). Testing the independence of sets of large-dimensional variables. Science China Mathematics, 56(1), 135–147.
- Jiang, T., Qi, Y. (2015a). Likelihood ratio tests for high-dimensional normal distributions. Scandinavian Journal of Statistics, 42(4), 988–1009.
- Jiang, T., Qi, Y. (2015b). Supplement to "Likelihood ratio tests for high-dimensional normal Distributions". http://www.stat.umn.edu/~tjiang/papers/SJSJQ.pdf.
- Jiang, T., Yang, F. (2013). Central limit theorems for classical likelihood ratio tests for high-dimensional normal distributions. Annals of Statistics, 41(4), 2029–2074.
- Ledoit, O., Wolf, M. (2002). Some hypothesis tests for the covariance matrix when the dimension is large compared to the sample size. *Annals of Statistics*, 30(4), 1081–1102.
- Li, W., Chen, J., Yao, J. (2017). Testing the independence of two random vectors where only one dimension is large. *Statistics*, 51(1), 141–153.
- Muirhead, R. J. (1982). Aspects of multivariate statistical theory. New York: Wiley.
- Schott, J. R. (2001). Some tests for the equality of covariance matrices. Journal of Statistical Planning and Inference, 94(1), 25–36.
- Schott, J. R. (2005). Testing for complete independence in high dimensions. Biometrika, 92(4), 951–956.
- Schott, J. R. (2007). A test for the equality of covariance matrices when the dimension is large relative to the sample sizes. *Computational Statistics and Data Analysis*, 51(12), 6535–6542.
- Srivastava, M. S., Reid, N. (2012). Testing the structure of the covariance matrix with fewer observations than the dimension. *Journal of Multivariate Analysis*, 112(C), 156–171.
- Wilks, S. S. (1935). On the independence of k sets of normally distributed statistical variables. *Econometrica*, 3(3), 309–326.
- Young, R. M. (1991). 75.9 Euler's Constant. Mathematical Gazette, 75(472), 187–190.