

On two-sample mean tests under spiked covariances

Rui Wang^a, Xingzhong Xu^{a,b,*}

^aSchool of Mathematics and Statistics, Beijing Institute of Technology, Beijing 100081, China

^bBeijing Key Laboratory on MCAACI, Beijing Institute of Technology, Beijing 100081, China

Abstract

This paper considers testing the means of two p -variate normal samples in high dimensional settings. We show that under the null hypothesis, a necessary and sufficient condition for the asymptotic normality of Chen and Qin (2010)'s test statistic is that the eigenvalues of the covariance matrix are concentrated around their average. However, this condition is not satisfied when the variables are strongly correlated. To characterize the correlations between variables, we adopt a spiked covariance model. Under the spiked covariance model, we derive the asymptotic distribution of Chen and Qin (2010)'s test statistic and correct the critical value of their test statistic. The recently proposed random projection test procedures suggest that the power of tests may be boosted using the projected data. By maximizing an average signal to noise ratio, we find that the optimal projection subspace is the orthogonal complement of the principal subspace. We propose a new test procedure through the projection onto the estimated orthogonal complement of the principal subspace. The asymptotic normality of the new test statistic is proved and the asymptotic power function of the test is given. Theoretical and simulation results show that the new test outperforms the competing tests substantially under the spiked covariance model.

Keywords: high dimension, mean test, principal subspace, spiked covariance model

1. Introduction

Suppose $X_{k,1}, \dots, X_{k,n_k}$ are independent identically distributed (i.i.d.) p -dimensional normal random vectors with unknown mean vector μ_k and covariance matrix Σ , $k = 1, 2$. We consider the hypothesis testing problem

$$H_0 : \mu_1 = \mu_2 \quad \text{vs.} \quad H_1 : \mu_1 \neq \mu_2. \quad (1)$$

In this paper, the high dimensional setting is adopted, that is, the dimension p varies as n increases, where $n = n_1 + n_2 - 2$. Testing hypotheses (1) is important in many fields, including biology, finance and economics.

A classical test statistic for hypotheses (1) is Hotelling's T^2 test statistic $(\bar{X}_1 - \bar{X}_2)^T \mathbf{S}^{-1} (\bar{X}_1 - \bar{X}_2)$ where \bar{X}_1 and \bar{X}_2 are the two sample means and $\mathbf{S} = n^{-1} \sum_{k=1}^2 \sum_{i=1}^{n_k} (X_{k,i} - \bar{X}_k)(X_{k,i} - \bar{X}_k)^T$ is the pooled sample covariance matrix. However, Hotelling's test statistic is not defined when $p > n$. Moreover, Bai and Saranadasa (1996) showed that even if $p \leq n$, Hotelling's test suffers from low power when p is comparable to n . Perhaps, the main reason for the low power of Hotelling's test is that \mathbf{S} is a poor estimator of Σ when p is large compared with n . See Chen and Qin (2010) and the references therein. For testing hypotheses (1) in high dimensional settings, many test statistics are based on the estimation of $(\mu_1 - \mu_2)^T \mathbf{A}(\mu_1 - \mu_2)$ for a positive definite matrix \mathbf{A} . Bai and Saranadasa (1996) proposed a test based on

$$T_{BS} = \|\bar{X}_1 - \bar{X}_2\|^2 - \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \text{tr } \mathbf{S},$$

*Corresponding author

Email address: xuxz@bit.edu.cn (Xingzhong Xu)

an unbiased estimator of $\|\mu_1 - \mu_2\|^2$. Chen and Qin (2010) modified T_{BS} by removing terms $\sum_{i=1}^{n_k} X_{ki}^T X_{ki}$, $k = 1, 2$, and proposed a test based on

$$\begin{aligned} T_{CQ} &= \frac{\sum_{i \neq j}^{n_1} X_{1i}^T X_{1j}}{n_1(n_1 - 1)} + \frac{\sum_{i \neq j}^{n_2} X_{2i}^T X_{2j}}{n_2(n_2 - 1)} - 2 \frac{\sum_{i=1}^{n_1} \sum_{j=1}^{n_2} X_{1i}^T X_{2j}}{n_1 n_2} \\ &= \|\bar{X}_1 - \bar{X}_2\|^2 - \frac{1}{n_1} \text{tr } \mathbf{S}_1 - \frac{1}{n_2} \text{tr } \mathbf{S}_2, \end{aligned}$$

where $\mathbf{S}_k = (n_k - 1)^{-1} \sum_{i=1}^{n_k} (X_{k,i} - \bar{X}_k)(X_{k,i} - \bar{X}_k)^T$, $k = 1, 2$. As an estimator of $\|\mu_1 - \mu_2\|^2$, T_{CQ} is unbiased even if the covariances of the two populations are different. In contrast, T_{BS} is unbiased only when the covariances are the same or $n_1 = n_2$. Srivastava and Du (2008) proposed a test based on

$$T_{SD} = (\bar{X}_1 - \bar{X}_2)^T [\text{diag}(\mathbf{S})]^{-1} (\bar{X}_1 - \bar{X}_2),$$

where $\text{diag}(\mathbf{S})$ is a diagonal matrix with the same diagonal elements as \mathbf{S} 's.

In practice, it is often the case that the variables are strongly correlated. See, for example, Chen et al. (2011), Thulin (2014) and Ma et al. (2015). As noted by Ma et al. (2015), however, the tests of Bai and Saranadasa (1996), Srivastava and Du (2008) and Chen and Qin (2010) may not be valid when there are strong correlations between the variables. For example, the condition

$$\text{tr}(\boldsymbol{\Sigma}^4) = o(\text{tr}^2(\boldsymbol{\Sigma}^2)) \quad (2)$$

imposed by Chen and Qin (2010) is violated when $\boldsymbol{\Sigma}$ has a uniform correlation structure. More precisely, we suppose that $\boldsymbol{\Sigma} = (1 - \rho)\mathbf{I}_p + \rho\mathbf{1}_p\mathbf{1}_p^T$ where $0 < \rho < 1$, \mathbf{I}_p is the p dimensional identity matrix and $\mathbf{1}_p$ is the p dimensional vector with all the elements equal to one. In this case, $\boldsymbol{\Sigma}$ has eigenvalues $1 + \rho(p - 1)$ and $1 - \rho$ with multiplicities 1 and $p - 1$ respectively. Then (2) is violated since as $p \rightarrow \infty$, we have

$$\frac{\text{tr}(\boldsymbol{\Sigma}^4)}{\text{tr}^2(\boldsymbol{\Sigma}^2)} = \frac{(1 + \rho(p - 1))^4 + (1 - \rho)^4(p - 1)}{[(1 + \rho(p - 1))^2 + (1 - \rho)^2(p - 1)]^2} \rightarrow 1.$$

Note that under the uniform correlation structure, the largest eigenvalue of $\boldsymbol{\Sigma}$ is significantly larger than the rest of eigenvalues. This is a special case of the spiked covariance model

$$\boldsymbol{\Sigma} = \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^T + \sigma^2\mathbf{I}_p, \quad (3)$$

where $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_r)$, $\boldsymbol{\lambda}_1 \geq \dots \geq \boldsymbol{\lambda}_r > 0$, $r \geq 1$, \mathbf{V} is a $p \times r$ matrix with orthonormal columns and $\sigma^2 > 0$. The spiked covariance model (3) is adopted by many theoretical studies. See Cai et al. (2013), Birnbaum et al. (2013), Passemier et al. (2017) and the references therein. This model arises when variables are strongly correlated and the correlations are determined by a small number of factors. In Section 2, we will derive the asymptotic distribution of T_{CQ} under the spiked covariance model. Generally, T_{CQ} is asymptotically distributed as a weighted chi-squared random variable. In a special case, T_{CQ} is asymptotically distributed as the sum of a weighted chi-squared random variable and a normal random variable. We also correct the critical value of T_{CQ} under the spiked covariance model.

Recently, a class of test procedures are proposed through random projection. See, for example, Lopes et al. (2011), Thulin (2014) and Srivastava et al. (2016). The idea is to project data onto some random lower-dimensional subspaces, and then construct the test using the projected data. It has been shown that these procedures have substantially higher power than competing tests when the variables are correlated. This suggests that projecting data onto certain subspaces may lead to an improvement of the test procedures. Instead of using randomly chosen subspaces, we would like to find the optimal subspace. In Section 3, we will see that under the spiked covariance model, the optimal subspace is the orthogonal complement of the principal subspace. Motivated by this, we propose a new test procedure through projection onto the estimated orthogonal complement of the principal subspace. The asymptotic null distribution of our test

statistic is derived and the asymptotic power function is also given. Our analysis and simulations show that our test has very attractive power performance under the spiked covariance model.

The rest of the paper is organized as follows. In Section 2, we revisit Chen and Qin (2010)'s test. In Section 3, we propose a test procedure and exploit properties of the test. In Section 4, simulations are carried out and a real data example is given. Section 5 contains some discussion. The technical proofs are presented in Appendix.

2. Asymptotic properties of Chen and Qin (2010)'s test

Throughout the paper, we assume $p \rightarrow \infty$ as $n \rightarrow \infty$ and $n_1/n_2 \rightarrow c \in (0, +\infty)$, that is, we consider high dimensional and balanced data.

In Chen and Qin (2010), the asymptotic normality of T_{CQ} is derived under the condition (2). We shall show that under the null hypothesis, the condition (2) is essential for the asymptotic normality of T_{CQ} . We note that under the null hypothesis, T_{CQ} is a quadratic form of a standard normal random vector. To see this, let $Z_{k,i} = \Sigma^{-1/2} X_{k,i}$, $i = 1, \dots, n_k$, $k = 1, 2$. It can be seen that $Z_{k,i}$ is $N_p(0, \mathbf{I}_p)$ distributed under the null hypothesis. Write $Z = (Z_{1,1}^T, \dots, Z_{1,n_1}^T, Z_{2,1}^T, \dots, Z_{2,n_2}^T)^T$. Then $T_{CQ} = Z^T (\mathbf{B}_n \otimes \Sigma) Z$, where \otimes is the Kronecker product and

$$\mathbf{B}_n = \begin{pmatrix} \frac{1}{n_1(n_1-1)}(\mathbf{1}_{n_1}\mathbf{1}_{n_1}^T - \mathbf{I}_{n_1}) & -\frac{1}{n_1n_2}\mathbf{1}_{n_1}\mathbf{1}_{n_2}^T \\ -\frac{1}{n_1n_2}\mathbf{1}_{n_2}\mathbf{1}_{n_1}^T & \frac{1}{n_2(n_2-1)}(\mathbf{1}_{n_2}\mathbf{1}_{n_2}^T - \mathbf{I}_{n_2}) \end{pmatrix}.$$

Using characteristic function method, one can prove the following result which gives a necessary and sufficient condition for the asymptotic normality of the quadratic form of a standard normal random vector.

Lemma 1. *Suppose Y_n is a k_n dimensional standard normal random vector and \mathbf{A}_n is a $k_n \times k_n$ symmetric matrix. Then as $n \rightarrow \infty$, a necessary and sufficient condition for*

$$\frac{Y_n^T \mathbf{A}_n Y_n - \mathbb{E} Y_n^T \mathbf{A}_n Y_n}{[\text{Var}(Y_n^T \mathbf{A}_n Y_n)]^{1/2}} \xrightarrow{\mathcal{L}} N(0, 1) \quad (4)$$

is that

$$\frac{\lambda_1(\mathbf{A}_n^2)}{\text{tr}(\mathbf{A}_n^2)} \rightarrow 0, \quad (5)$$

where “ $\xrightarrow{\mathcal{L}}$ ” means convergence of a sequence of random variables in law and $\lambda_i(\cdot)$ means the i th largest eigenvalue.

To apply Lemma 1 to T_{CQ} , one needs to calculate the eigenvalues of $\mathbf{B}_n \otimes \Sigma$. Note that the eigenvalues of \mathbf{B}_n are $-1/n_1(n_1-1)$, $-1/n_2(n_2-1)$, $(n_1+n_2)/n_1n_2$ and 0 with multiplicities n_1-1 , n_2-1 , 1 and 1 respectively. Thus,

$$\text{tr}(\mathbf{B}_n \otimes \Sigma)^2 = \text{tr}(\mathbf{B}_n^2) \text{tr}(\Sigma^2) = \left(\frac{1}{n_1(n_1-1)} + \frac{1}{n_2(n_2-1)} + \frac{2}{n_1n_2} \right) \text{tr}(\Sigma^2),$$

and

$$\lambda_1((\mathbf{B}_n \otimes \Sigma)^2) = \lambda_1(\mathbf{B}_n^2)\lambda_1(\Sigma^2) = \left(\frac{1}{n_1} + \frac{1}{n_2} \right)^2 \lambda_1(\Sigma^2).$$

Because $n_1/n_2 \rightarrow c$, the condition

$$\frac{\lambda_1((\mathbf{B}_n \otimes \Sigma)^2)}{\text{tr}(\mathbf{B}_n \otimes \Sigma)^2} \rightarrow 0$$

is equivalent to $\lambda_1(\Sigma^2)/\text{tr} \Sigma^2 \rightarrow 0$. From

$$\frac{\lambda_1(\Sigma)^4}{(\sum_{i=1}^p \lambda_i(\Sigma)^2)^2} \leq \frac{\sum_{i=1}^p \lambda_i(\Sigma)^4}{(\sum_{i=1}^p \lambda_i(\Sigma)^2)^2} \leq \frac{\lambda_1(\Sigma)^2 \sum_{i=1}^p \lambda_i(\Sigma)^2}{(\sum_{i=1}^p \lambda_i(\Sigma)^2)^2} = \frac{\lambda_1(\Sigma)^2}{\sum_{i=1}^p \lambda_i(\Sigma)^2},$$

we can see that $\lambda_1^2(\Sigma)/\text{tr}(\Sigma^2) \rightarrow 0$ is equivalent to (2). Then Lemma 1 implies that under the null hypothesis, the condition (2) is a necessary and sufficient condition for

$$\frac{T_{CQ} - \mathbb{E} T_{CQ}}{[\text{Var}(T_{CQ})]^{1/2}} \xrightarrow{\mathcal{L}} N(0, 1).$$

The above result implies that Chen and Qin (2010)'s test procedure can be used only when the eigenvalues of Σ are concentrated around their average. In a class of applications, however, the correlations between variables are mainly driven by several common factors, and consequently, Σ has a few eigenvalues which are much larger than the others. See, for example, Jung and Marron (2009), Cai et al. (2013) and Wang and Fan (2017). To characterize such correlations between variables, we consider the spiked covariance model (3). For $p \geq q$, let $\mathbb{O}_{p \times q}$ denote the collection of $p \times q$ matrices with orthonormal columns. We make the following assumption for the covariance matrix Σ .

Assumption 1. *The covariance matrix Σ has structure $\Sigma = \mathbf{V}\Lambda\mathbf{V}^T + \sigma^2\mathbf{I}_p$, where $\mathbf{V} \in \mathbb{O}_{p \times r}$, r is a known number and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_r)$, $\lambda_1 \geq \dots \geq \lambda_r > 0$. As n, p tend to infinity, the parameters r, σ^2 are fixed and Λ satisfies*

$$\kappa p^\beta \geq \lambda_1 \geq \dots \geq \lambda_r \geq \kappa^{-1} p^\beta,$$

where $\kappa > 1$ and $\beta \geq 1/2$ are constants.

The covariance structure in Assumption 1 is commonly adopted in PCA study. See Cai et al. (2013), Birnbaum et al. (2013), Passemier et al. (2017) and the references therein. This covariance structure is also connected with the factor model. In fact, the model in Assumption 1 with $\beta = 1$ corresponds to the factor model in Ma et al. (2015) with homoscedastic noise.

In Assumption 1, the column space of \mathbf{V} is exactly the eigenspace of Σ associated with the r leading eigenvalues, and is therefore called principal subspace. Since the columns of \mathbf{V} are orthonormal, $\mathbf{V}\mathbf{V}^T$ is the orthogonal projection onto the principal subspace. Let $\tilde{\mathbf{V}}$ be a member of $\mathbb{O}_{p \times (p-r)}$ such that the columns of $\tilde{\mathbf{V}}$ are orthogonal to the columns of \mathbf{V} . Although such $\tilde{\mathbf{V}}$ is not unique, the orthogonal projection $\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T = \mathbf{I}_p - \mathbf{V}\mathbf{V}^T$ is unique and is equal to the orthogonal projection onto the orthogonal complement of the principal subspace.

For positive sequences $\{a_n\}$ and $\{b_n\}$, we write $a_n \asymp b_n$ to denote $a_n = O(b_n)$ and $b_n = O(a_n)$ as $n \rightarrow \infty$. Under Assumption 1, we have

$$\frac{\text{tr}(\Sigma^4)}{\text{tr}^2(\Sigma^2)} = \frac{\sum_{i=1}^r (\lambda_i + \sigma^2)^4 + (p-r)\sigma^8}{\left(\sum_{i=1}^r (\lambda_i + \sigma^2)^2 + (p-r)\sigma^4\right)^2} \asymp \frac{p^{4\beta} + p}{(p^{2\beta} + p)^2}.$$

The right hand side tends to 0 if and only if $\beta < 1/2$. Our previous arguments assert that the asymptotic distribution of T_{CQ} won't be normal for $\beta \geq 1/2$. To derive the asymptotic distribution of T_{CQ} for $\beta \geq 1/2$, note that the variation of T_{CQ} is mainly due to $\|\bar{X}_1 - \bar{X}_2\|^2$. Let $\tau = 1/n_1 + 1/n_2$. Under the null hypothesis, we have

$$\text{Var}(\|\bar{X}_1 - \bar{X}_2\|^2) = 2\tau^2 \text{tr}(\Sigma^2) = 2\tau^2 \sum_{i=1}^r (\lambda_i + \sigma^2)^2 + 2\tau^2(p-r)\sigma^4,$$

where the first term of the right hand side is of order $p^{2\beta}/n^2$ and the second term is of order p/n^2 . If $\beta = 1/2$, the two terms are of the same order. If $\beta > 1/2$, however, the second term is dominated by the first term. This implies that the asymptotic distributions of T_{CQ} are different for $\beta = 1/2$ and $\beta > 1/2$. Since the variance of $(\tau p^\beta)^{-1}\|\bar{X}_1 - \bar{X}_2\|^2$ is bounded away from 0 and $+\infty$ under the null hypothesis, we use τp^β to standardize T_{CQ} . The following two theorems give the asymptotic distributions of $(\tau p^\beta)^{-1}T_{CQ}$ for $\beta = 1/2$ and $\beta > 1/2$, respectively.

Theorem 1. *Under Assumption 1, suppose $\beta = 1/2$ and $\lambda_i/p^\beta \rightarrow \omega_i \in (0, +\infty)$, $i = 1, \dots, r$. Let Z_0, Z_1, \dots, Z_r be i.i.d. standard normal random variables, then we have*

(a) if $\mu_1 = \mu_2$, then

$$\frac{1}{\tau p^\beta} T_{CQ} \xrightarrow{w} \sqrt{2}\sigma^2 Z_0 + \sum_{i=1}^r \omega_i Z_i^2 - \sum_{i=1}^r \omega_i,$$

where " \xrightarrow{w} " denotes weak convergence;

(b) if $(\tau p^\beta)^{-1/2} (\mathbf{V}^T(\mu_1 - \mu_2))_i \rightarrow \zeta_i \in (-\infty, +\infty)$, $i = 1, \dots, r$, and $(\tau p^\beta)^{-1} \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 \rightarrow \zeta^* \in [0, +\infty)$, then

$$\frac{1}{\tau p^\beta} T_{CQ} \xrightarrow{w} \sqrt{2}\sigma^2 Z_0 + \sum_{i=1}^r (\sqrt{\omega_i} Z_i + \zeta_i)^2 + \zeta^* - \sum_{i=1}^r \omega_i.$$

Theorem 2. Under Assumption 1, suppose $\beta > 1/2$ and $\lambda_i/p^\beta \rightarrow \omega_i \in (0, +\infty)$, $i = 1, \dots, r$. Let Z_1, \dots, Z_r be i.i.d. standard normal random variables, then we have

(a) if $\mu_1 = \mu_2$, then

$$\frac{1}{\tau p^\beta} T_{CQ} \xrightarrow{w} \sum_{i=1}^r \omega_i Z_i^2 - \sum_{i=1}^r \omega_i;$$

(b) if $(\tau p^\beta)^{-1/2} (\mathbf{V}^T(\mu_1 - \mu_2))_i \rightarrow \zeta_i \in (-\infty, +\infty)$, $i = 1, \dots, r$, and $(\tau p^\beta)^{-1} \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 \rightarrow \zeta^* \in [0, +\infty)$, then

$$\frac{1}{\tau p^\beta} T_{CQ} \xrightarrow{w} \sum_{i=1}^r (\sqrt{\omega_i} Z_i + \zeta_i)^2 + \zeta^* - \sum_{i=1}^r \omega_i.$$

Remark 1. By the definitions of ζ_i and ζ^* , we have

$$\frac{1}{\tau p^\beta} \|\mu_1 - \mu_2\|^2 = \frac{1}{\tau p^\beta} \|\mathbf{V}^T(\mu_1 - \mu_2)\|^2 + \frac{1}{\tau p^\beta} \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 \rightarrow \sum_{i=1}^r \zeta_i^2 + \zeta^*.$$

Thus, $\sum_{i=1}^r \zeta_i^2$ and ζ^* characterize the signal strength in the principal subspace and the complement of the principal subspace, respectively. Under the conditions of Theorem 1 or Theorem 2, the following statements are equivalent:

- (1) $\zeta_1 = \dots = \zeta_r = \zeta^* = 0$.
- (2) $\|\mu_1 - \mu_2\|^2 = o(\tau p^\beta)$.
- (3) The asymptotic distributions of $(\tau p^\beta)^{-1} T_{CQ}$ are the same under the null hypothesis and the alternative hypothesis.
- (4) Any test procedure based on T_{CQ} has trivial power asymptotically.

It is implied by Theorem 1 and Theorem 2 that the original critical value of T_{CQ} can not be used when $\beta \geq 1/2$. Now we adjust the critical value of T_{CQ} such that the resulting test has correct level asymptotically. Consider the random variable $W = \sqrt{2p\sigma^2} Z_0 + \sum_{i=1}^r \lambda_i Z_i^2 - \sum_{i=1}^r \lambda_i$, where Z_0, Z_1, \dots, Z_r are i.i.d. $N(0, 1)$ random variables. Let $F(x; \lambda_1, \dots, \lambda_r, \sigma^2)$ be the cumulative distribution function of W . Under the conditions of Theorem 1 ($\beta = 1/2$), we have

$$\frac{W}{p^\beta} \xrightarrow{w} \sqrt{2}\sigma^2 Z_0 + \sum_{i=1}^r \omega_i Z_i^2 - \sum_{i=1}^r \omega_i.$$

Under the conditions of Theorem 2 ($\beta > 1/2$), we have

$$\frac{W}{p^\beta} \xrightarrow{w} \sum_{i=1}^r \omega_i Z_i^2 - \sum_{i=1}^r \omega_i.$$

Hence in both case, we have

$$\sup_{x \in \mathbb{R}} |\Pr\left(\frac{1}{\tau} T_{CQ} \leq x\right) - \Pr(W \leq x)| = o(1).$$

Thus, the test that rejects the null hypothesis if

$$\frac{1}{\tau} T_{CQ} > F^{-1}(1 - \alpha; \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_r, \sigma^2)$$

is asymptotically level α . However, the distribution $F(x; \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_r, \sigma^2)$ involves some unknown parameters. In order to consistently estimate $F(x; \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_r, \sigma^2)$, we need to estimate $\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_r$ and σ^2 . In Section 3, we will give their estimators $\hat{\boldsymbol{\lambda}}_1, \dots, \hat{\boldsymbol{\lambda}}_r$ and $\hat{\sigma}_*^2$, and Proposition 2 asserts that these estimators are ratio consistent. Now we propose a corrected T_{CQ} test which rejects the null hypothesis with α level of significance if

$$\tau^{-1} T_{CQ} > F^{-1}(1 - \alpha; \hat{\boldsymbol{\lambda}}_1, \dots, \hat{\boldsymbol{\lambda}}_r, \hat{\sigma}_*^2).$$

Then under the conditions of either Theorem 1 or Theorem 2, the corrected T_{CQ} test is asymptotically level α .

As we have seen in Remark 1, the corrected T_{CQ} test has trivial power α asymptotically if and only if

$$\|\mu_1 - \mu_2\|^2 = o(p^\beta/n). \quad (6)$$

Note that the trivial power region (6) becomes larger as β increases. This property is undesirable. In the best case ($\beta = 1/2$), the trivial power region of the corrected T_{CQ} test is $\|\mu_1 - \mu_2\|^2 = o(\sqrt{p}/n)$. In the next section, we will propose a new test. The trivial power region of the new test is $\|\mu_1 - \mu_2\|^2 = o(\sqrt{p}/n)$ which is independent of β .

3. A projection test

Recently, a class of test procedures have been proposed through random projection onto lower dimensional subspace. See, for example, Lopes et al. (2011), Thulin (2014) and Srivastava et al. (2016). It is known that random projection methods offer higher power when the variables are correlated. However, these test procedures are randomized, which is undesirable in practice. Then, is there an optimal projection which is nonrandomized?

For $\mathbf{O} \in \mathbb{O}_{p \times k}$ ($1 \leq k \leq p$), let

$$T(\mathbf{O}) = \|\mathbf{O}^T(\bar{X}_1 - \bar{X}_2)\|^2 - \frac{1}{n_1} \text{tr}(\mathbf{O}^T \mathbf{S}_1 \mathbf{O}) - \frac{1}{n_2} \text{tr}(\mathbf{O}^T \mathbf{S}_2 \mathbf{O}).$$

Then $T(\mathbf{O})$ is Chen and Qin (2010)'s statistic on the transformed data $\mathbf{O}^T X_{k,i}$. Under the condition (2), Chen and Qin (2010) proved that the asymptotic power of T_{CQ} under the local alternative is

$$\Phi\left(\Phi^{-1}(\alpha) + \frac{\|\mu_1 - \mu_2\|^2}{\sqrt{2\tau^2 \text{tr}(\boldsymbol{\Sigma}^2)}}\right),$$

where $\Phi(\cdot)$ is the cumulative distribution function of the standard normal random variable. Hence the power of T_{CQ} is related to $\|\mu_1 - \mu_2\|^2 / \sqrt{2\tau^2 \text{tr}(\boldsymbol{\Sigma}^2)}$, which may be viewed as a signal to noise ratio (SNR). Consequently, $\|\mathbf{O}^T(\mu_1 - \mu_2)\|^2 / \sqrt{2\tau^2 \text{tr}(\mathbf{O}^T \boldsymbol{\Sigma}^2 \mathbf{O})}$ measures the power of $T(\mathbf{O})$. Like Lopes et al. (2011), to consider an average-case scenario, we temporarily place a prior on $\mu_1 - \mu_2$. Suppose that the norm $\|\mu_1 - \mu_2\|$ is nonrandom while the orientation $\delta = (\mu_1 - \mu_2) / \|\mu_1 - \mu_2\|$ is from the uniform distribution on the unit sphere. In this case, an average SNR can be defined as

$$\mathbb{E}\left(\frac{\|\mathbf{O}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 \text{tr}(\mathbf{O}^T \boldsymbol{\Sigma}^2 \mathbf{O})}}\right) = \frac{k/p}{\sqrt{2\tau^2 \text{tr}(\mathbf{O}^T \boldsymbol{\Sigma}^2 \mathbf{O})}} \|\mu_1 - \mu_2\|^2. \quad (7)$$

The $T(\mathbf{O})$ that maximizes the average SNR can be expected to have the best average power behavior among $\{T(\mathbf{O}) : \mathbf{O} \in \mathbb{O}_{p \times k}, k \leq p\}$.

Next we maximize (7) under Assumption 1. Note that for fixed k , (7) is maximized when the columns of \mathbf{O} are equal to the last k eigenvectors of Σ . Thus, it remains to maximize

$$\frac{k/p}{\sqrt{2\tau^2 \sum_{i=p-k+1}^p \lambda_i^2(\Sigma)}} \quad (8)$$

over k . If $k \leq p-r$, (8) is equal to $\sqrt{k}/(\sqrt{2}\sigma^2\tau p)$ which is an increasing function of k . If $k > p-r$, we have

$$\begin{aligned} \frac{k/p}{\sqrt{2\tau^2 \sum_{i=p-k+1}^p \lambda_i^2(\Sigma)}} &\leq \frac{1}{\sqrt{2\tau^2(\kappa^{-2}p^{2\beta} + (p-r)\sigma^4)}} \\ &= \frac{p/(p-r)}{\sqrt{\kappa^{-2}p^{2\beta}/((p-r)\sigma^4) + 1}} \frac{(p-r)/p}{\sqrt{2\tau^2((p-r)\sigma^4)}}. \end{aligned}$$

Since $\beta \geq 1/2$, for sufficiently large p and all $k > p-r$, we have

$$\frac{k/p}{\sqrt{2\tau^2 \sum_{i=p-k+1}^p \lambda_i^2(\Sigma)}} < \frac{(p-r)/p}{\sqrt{2\tau^2((p-r)\sigma^4)}},$$

and (8) is maximized when $k = p-r$. Consequently, for sufficiently large p , (7) is maximized when $\mathbf{O} = \tilde{\mathbf{V}}$.

The above discussion motivates us to consider the variable

$$T_1 = \|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 - \frac{1}{n_1} \text{tr}(\tilde{\mathbf{V}}^T \mathbf{S}_1 \tilde{\mathbf{V}}) - \frac{1}{n_2} \text{tr}(\tilde{\mathbf{V}}^T \mathbf{S}_2 \tilde{\mathbf{V}}).$$

Note that based on the data $\tilde{\mathbf{V}}^T X_{ki}$, $i = 1, \dots, n_k$, $k = 1, 2$, the likelihood ratio test statistic for hypotheses (1) is equivalent to $\|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$. Hence T_1 can be regarded as a restricted likelihood ratio statistic. Below we show that T_1 is approximately a standardized chi-squared random variable.

Proposition 1. *Under Assumption 1, we have*

(a) *if $\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 = o(p/n)$,*

$$\frac{T_1 - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \stackrel{d}{=} \frac{\chi^2(p-r) - (p-r)}{\sqrt{2(p-r)}} + o_P(1), \quad (9)$$

where “ $\stackrel{d}{=}$ ” means having identical distribution;

(b) *if $p/n = o(\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2)$,*

$$\frac{T_1 - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{4\tau^2 \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 \sigma^2}} \xrightarrow{\mathcal{L}} N(0, 1). \quad (10)$$

Remark 2. Proposition 1 implies that under the null hypothesis,

$$\frac{T_1}{\sqrt{2\tau^2 p \sigma^4}} \xrightarrow{\mathcal{L}} N(0, 1).$$

However, compared with the standard normal distribution, the standardized chi-squared distribution is a better approximation of the distribution of the statistic. This is implied in the proof of Proposition 1.

Note that T_1 is dependent on the subspace $\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T$ which is typically unknown and thus needs to be estimated. Let $\hat{\mathbf{V}}$ and $\tilde{\mathbf{V}}$ denote the first r and last $p - r$ eigenvectors of \mathbf{S} , respectively. Anderson (1963) proved that the maximum likelihood estimator (MLE) of \mathbf{V} is $\hat{\mathbf{V}}$. This fact, together with the equalities $\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T = \mathbf{I}_p - \mathbf{V}\mathbf{V}^T$ and $\hat{\mathbf{V}}\hat{\mathbf{V}}^T = \mathbf{I}_p - \hat{\mathbf{V}}\hat{\mathbf{V}}^T$, implies that the MLE of $\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T$ is $\hat{\mathbf{V}}\hat{\mathbf{V}}^T$. Thus, as the main term of T_1 , $\|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$ can be estimated by $\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$. In T_1 , the centralization term $n_1^{-1} \text{tr}(\tilde{\mathbf{V}}^T \mathbf{S}_1 \tilde{\mathbf{V}}) + n_2^{-1} \text{tr}(\tilde{\mathbf{V}}^T \mathbf{S}_2 \tilde{\mathbf{V}})$ is an unbiased estimator of $\tau(p - r)\sigma^2$, such that $E T_1 = 0$ under the null hypothesis. However, compared with $\|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$, the centralization of $\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$ is more complicated.

Since $\hat{\mathbf{V}}$ is independent of \bar{X}_k , $k = 1, 2$, it is convenient to work with the conditional distribution of $\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$ given \mathbf{S} . Under the null hypothesis,

$$E[\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 | \mathbf{S}] = \tau \text{tr}(\hat{\mathbf{V}}^T \Sigma \hat{\mathbf{V}}) = \tau(p - r)\sigma^2 + \tau \text{tr}(\hat{\mathbf{V}}^T \mathbf{V} \Lambda \mathbf{V}^T \hat{\mathbf{V}}).$$

Note that

$$\text{tr}(\hat{\mathbf{V}}^T \mathbf{V} \Lambda \mathbf{V}^T \hat{\mathbf{V}}) = \text{tr}(\Lambda^{1/2} \mathbf{V}^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \mathbf{V} \Lambda^{1/2}) = \text{tr}(\Lambda^{1/2} (\mathbf{I}_r - \mathbf{V}^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \mathbf{V}) \Lambda^{1/2}) = \sum_{i=1}^r (1 - \sum_{l=1}^r (\hat{v}_l^T v_i)^2) \lambda_i,$$

where \hat{v}_i and v_i are the i th columns of $\hat{\mathbf{V}}$ and \mathbf{V} , respectively, $i = 1, \dots, r$. Under $p = O(n\lambda_r)$ and some other regular conditions, Theorem 3.2 in Wang and Fan (2017) asserts that

$$\hat{v}_j^T v_i = O_P(\epsilon_{j,n}), \quad 1 \leq i \neq j \leq r, \quad \text{and} \quad \hat{v}_i^T v_i = \frac{1}{\sqrt{1 + \frac{p}{n(\lambda_i + \sigma^2)} \sigma^2}} + O_P(\epsilon_{i,n}), \quad 1 \leq i \leq r, \quad (11)$$

where $\epsilon_{i,n}$ is a smaller order term, $i = 1, \dots, r$. This motivates us to approximate $\text{tr}(\hat{\mathbf{V}}^T \mathbf{V} \Lambda \mathbf{V}^T \hat{\mathbf{V}})$ by

$$\sum_{i=1}^r \left(1 - \frac{1}{1 + \frac{p}{n(\lambda_i + \sigma^2)} \sigma^2}\right) \lambda_i = \sum_{i=1}^r \frac{p\sigma^2}{n\lambda_i + (n+p)\sigma^2} \lambda_i.$$

Hence

$$E[\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 | \mathbf{S}] \approx \tau(p - r)\sigma^2 + \tau \sum_{i=1}^r \frac{p\sigma^2}{n\lambda_i + (n+p)\sigma^2} \lambda_i.$$

To centralize $\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$, we need to estimate $\lambda_1, \dots, \lambda_r$ and σ^2 . Anderson (1963) proved that the MLE of σ^2 is $\hat{\sigma}^2 = (p - r)^{-1} \sum_{i=r+1}^p \lambda_i(\mathbf{S})$ and the MLE of λ_i is $\lambda_i(\mathbf{S}) - \hat{\sigma}^2$, $i = 1, \dots, r$. However, Lemma 5 in Appendix B implies that $\hat{\sigma}^2$ is downward biased and $\lambda_i(\mathbf{S}) - \hat{\sigma}^2$ is upward biased. (See also Passemier et al. (2017) and Wang and Fan (2017).) Motivated by Lemma 5 in Appendix B, we propose the following bias-corrected estimators:

$$\hat{\sigma}_*^2 = \frac{n}{n - r} \hat{\sigma}^2, \quad \hat{\lambda}_i = \lambda_i(\mathbf{S}) - \frac{p + n - r}{n} \hat{\sigma}_*^2 \quad i = 1, \dots, r.$$

The following proposition gives the convergence rate of these estimators.

Proposition 2. *Under Assumption 1, we have*

$$\hat{\sigma}_*^2 = \sigma^2 + O_P\left(\max\left(\frac{1}{\sqrt{np}}, \frac{1}{p}\right)\right), \quad (12)$$

and

$$\frac{\hat{\lambda}_i}{\lambda_i} = 1 + O_P\left(\max\left(\frac{1}{\sqrt{n}}, \frac{1}{p^\beta}\right)\right). \quad (13)$$

Remark 3. Recently, Passemier et al. (2017) proposed a bias-corrected estimator of σ^2 :

$$\left(1 + \frac{1}{n-2} \left(r + \hat{\sigma}^2 \sum_{i=1}^r \frac{1}{\lambda_i}\right)\right) \hat{\sigma}^2.$$

In their paper, λ_i 's are fixed and known. This is different from Assumption 1 where λ_i 's are divergent and unknown.

Remark 4. Recently, Wang and Fan (2017) proposed an estimator of $\lambda_i(\Sigma) = \lambda_i + \sigma^2$, $i = 1, \dots, r$:

$$\max \left(\lambda_i(\mathbf{S}) - \frac{p}{n-2} \left(1 - \frac{p}{p-r} \frac{r}{n-2}\right)^{-1} \hat{\sigma}^2, 0 \right).$$

They showed that under $p > n - 2$, $p = O(n\lambda_r)$ and some other conditions,

$$\frac{1}{\lambda_i + \sigma^2} \max \left(\lambda_i(\mathbf{S}) - \frac{p}{n-2} \left(1 - \frac{p}{p-r} \frac{r}{n-2}\right)^{-1} \hat{\sigma}^2, 0 \right) = 1 + O_P \left(\frac{1}{\lambda_i} \sqrt{\frac{p}{n}} + \frac{1}{\sqrt{n}} \right).$$

Note that under Assumption 1 and $p > n - 2$, we have

$$\frac{1}{\lambda_i} \sqrt{\frac{p}{n}} + \frac{1}{\sqrt{n}} \asymp \frac{1}{\sqrt{n}} \asymp \max \left(\frac{1}{\sqrt{n}}, \frac{1}{p^\beta} \right), \quad i = 1, \dots, r.$$

In this case, Wang and Fan (2017)'s estimator and $\hat{\lambda}_i$ have the same convergence rate, although the estimands are slightly different. Compared with Wang and Fan (2017)'s result, Proposition 2 doesn't need the conditions $p > n - 2$ and $p = O(n\lambda_r)$.

Now we propose the following test statistic:

$$T_2 = \|\hat{\tilde{\mathbf{V}}}^T (\bar{X}_1 - \bar{X}_2)\|^2 - \tau(p-r)\hat{\sigma}_*^2 - \tau \sum_{i=1}^r \frac{p\hat{\sigma}_*^2}{n\hat{\lambda}_i + (n+p)\hat{\sigma}_*^2} \hat{\lambda}_i.$$

The following theorem establishes the asymptotic normality of T_2 .

Theorem 3. Under Assumption 1, suppose $p/n^2 \rightarrow 0$, we have

(a) if $(\mu_1 - \mu_2)^T \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T \hat{\Sigma} \hat{\tilde{\mathbf{V}}}^T (\mu_1 - \mu_2) = o_P(p/n)$,

$$\frac{T_2 - \|\hat{\tilde{\mathbf{V}}}^T (\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \xrightarrow{\mathcal{L}} N(0, 1);$$

(b) if $p/n = o_P((\mu_1 - \mu_2)^T \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T \hat{\Sigma} \hat{\tilde{\mathbf{V}}}^T (\mu_1 - \mu_2))$,

$$\frac{T_2 - \|\hat{\tilde{\mathbf{V}}}^T (\mu_1 - \mu_2)\|^2}{\sqrt{4\tau(\mu_1 - \mu_2)^T \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T \hat{\Sigma} \hat{\tilde{\mathbf{V}}}^T (\mu_1 - \mu_2)}} \xrightarrow{\mathcal{L}} N(0, 1);$$

(c) if $\|\mu_1 - \mu_2\|^2 = O(\sqrt{p}/n)$,

$$\frac{T_2 - \|\tilde{\mathbf{V}}^T (\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \xrightarrow{\mathcal{L}} N(0, 1).$$

Remark 5. The asymptotic normality of T_2 is closely related to the convergence rate of $\hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T$ to $\tilde{\mathbf{V}} \tilde{\mathbf{V}}^T$. Lemma 6 in Appendix B and the equality $\|\hat{\mathbf{V}} \hat{\mathbf{V}}^T - \mathbf{V} \mathbf{V}^T\|^2 = \|\hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T - \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T\|^2$ imply that $\|\hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T - \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T\|^2 = O_P(p/(p^\beta n))$. Hence $\hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T$ can consistently estimate $\tilde{\mathbf{V}} \tilde{\mathbf{V}}^T$ only if $p/(p^\beta n) \rightarrow 0$. Moreover, Cai et al. (2013)'s

Theorem 5 implies that no other estimator has faster convergence rate than $O_P(p/(p^\beta n))$. On the other hand, the asymptotic normality of T_2 requires the condition

$$p^{-1}(p^\beta \|\hat{\tilde{\mathbf{V}}}\hat{\tilde{\mathbf{V}}}^T - \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T\|^2 + 1)^2 \xrightarrow{P} 0,$$

which is equivalent to $\|\hat{\tilde{\mathbf{V}}}\hat{\tilde{\mathbf{V}}}^T - \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T\|^2 = o_P(\sqrt{p}/p^\beta)$. This is why Theorem 3 needs the condition $p/n^2 \rightarrow 0$. The proof of Theorem 3 implies that the asymptotic normality of T_2 is not valid if the condition $p/n^2 \rightarrow 0$ is violated.

Remark 6. In Theorem 3, the result (a) gives the asymptotic distribution of T_2 under the local alternative while the result (b) gives the asymptotic distribution of T_2 under a version of fixed alternatives. While both the condition and the conclusion of (a) involve the random quantity $\hat{\tilde{\mathbf{V}}}$, the result (c) only involves the nonrandom quantity $\tilde{\mathbf{V}}$.

By Proposition 2 and (a) of Theorem 3, the test that rejects the null hypothesis if $T_2/\sqrt{2\tau^2 p \hat{\sigma}_*^4} > \Phi^{-1}(1-\alpha)$ is asymptotically level α . However, using $\sqrt{2\tau^2 p \sigma^4}$ to standardize T_2 may not be accurate if n is small, and the asymptotic distribution $N(0, 1)$ may not be accurate if p is small. We would like to make further effort to construct a more accurate test. Note that

$$\begin{aligned} \text{Var}(\|\hat{\tilde{\mathbf{V}}}^T(\bar{X}_1 - \bar{X}_2)\|^2 | \mathbf{S}) &= 2\tau^2 (\text{tr}(\hat{\tilde{\mathbf{V}}}^T \mathbf{V} \Lambda \mathbf{V}^T \hat{\tilde{\mathbf{V}}})^2 + 2\sigma^2 \text{tr}(\hat{\tilde{\mathbf{V}}}^T \mathbf{V} \Lambda \mathbf{V}^T \hat{\tilde{\mathbf{V}}}) + \sigma^4(p-r)) \\ &= 2\tau^2 \left(\sum_{i=1}^r \left(1 - \sum_{l=1}^r (\hat{v}_l^T v_i)^2\right)^2 \lambda_i^2 + 2 \sum_{1 \leq i < j \leq r} \left(\sum_{l=1}^r (\hat{v}_l^T v_i)(\hat{v}_l^T v_j) \right)^2 \lambda_i \lambda_j + 2\sigma^2 \sum_{i=1}^r \left(1 - \sum_{l=1}^r (\hat{v}_l^T v_i)^2\right) \lambda_i + \sigma^4(p-r) \right) \\ &\approx 2\tau^2 \left(\sum_{i=1}^r \left(\frac{p\sigma^2}{n\lambda_i + (n+p)\sigma^2} \lambda_i \right)^2 + 2\sigma^2 \sum_{i=1}^r \frac{p\sigma^2}{n\lambda_i + (n+p)\sigma^2} \lambda_i + \sigma^4(p-r) \right), \end{aligned}$$

where in the last line, we use the approximation in (11). Hence we propose the following standardized statistic

$$Q = T_2 / \left(2\tau^2 \left(\sum_{i=1}^r \left(\frac{p\hat{\sigma}_*^2}{n\hat{\lambda}_i + (n+p)\hat{\sigma}_*^2} \hat{\lambda}_i \right)^2 + 2\hat{\sigma}_*^2 \sum_{i=1}^r \frac{p\hat{\sigma}_*^2}{n\hat{\lambda}_i + (n+p)\hat{\sigma}_*^2} \hat{\lambda}_i + \hat{\sigma}_*^4(p-r) \right) \right)^{1/2}.$$

In view of Proposition 1, we propose to reject the null hypothesis if

$$Q > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}},$$

where $\chi_{1-\alpha}^2(p-r)$ is the $1-\alpha$ quantile of a $\chi^2(p-r)$ random variable.

The following theorem gives the asymptotic power function of the proposed test. In particular, it shows that the proposed test is asymptotically level α .

Theorem 4. Under Assumption 1, suppose $p/n^2 \rightarrow 0$, we have

$$\Pr \left(Q > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \right) = \mathbb{E} \Phi \left(-\Phi^{-1}(1-\alpha) + \frac{\|\hat{\tilde{\mathbf{V}}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \right) + o(1). \quad (14)$$

If we further assume $\|\mu_1 - \mu_2\|^2 = O(\sqrt{p}/n)$, then

$$\Pr \left(Q > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \right) = \Phi \left(-\Phi^{-1}(1-\alpha) + \frac{\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \right) + o(1). \quad (15)$$

With the asymptotic power function of our test, we are in a position to provide a comparison with other tests. Again, to make an average-case comparison against other tests, we place a prior on $\mu_1 - \mu_2$. Suppose that the norm $\|\mu_1 - \mu_2\|$ is nonrandom while the orientation $\delta = (\mu_1 - \mu_2)/\|\mu_1 - \mu_2\|$ is uniformly distributed on the unit sphere. Then $\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 = \|\mu_1 - \mu_2\|^2 \|\tilde{\mathbf{V}}^T \delta\|^2$. Since $E \|\tilde{\mathbf{V}}^T \delta\|^2 = (p-r)/p$ and $\text{Var}(\|\tilde{\mathbf{V}}^T \delta\|^2) = 2r(p-r)/(p^2(p+2))$, we have $\|\tilde{\mathbf{V}}^T \delta\|^2 = 1 + O_P(1/p)$. In this case, the asymptotic power function of our test is equal to

$$\Phi\left(-\Phi^{-1}(1-\alpha) + \frac{\|\mu_1 - \mu_2\|^2}{\sqrt{2\tau^2 p \sigma^4}}\right).$$

So the trivial power region of our test is $\|\mu_1 - \mu_2\|^2 = o(\sqrt{p}/n)$.

Now we compare our test with the corrected T_{CQ} test. Recall that the trivial power region of the corrected T_{CQ} test is $\|\mu_1 - \mu_2\|^2 = o(p^\beta/n)$. Thus, when $\beta > 1/2$, the trivial power region of our test is smaller than that of the corrected T_{CQ} test. Consequently, our test is more powerful than the corrected T_{CQ} test in average.

We would also like to compare our test with T_{SD} . Srivastava and Du (2008) showed that the asymptotic power function of T_{SD} is

$$\Phi\left(-\Phi^{-1}(1-\alpha) + \frac{\|\Sigma_d^{-1/2}(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 \text{tr}(\mathbf{R}^2)}}\right),$$

where $\Sigma_d = \text{diag}(\Sigma)$ and $\mathbf{R} = \Sigma_d^{-1/2} \Sigma \Sigma_d^{-1/2}$ is the population correlation matrix. It is difficult to compare our test with T_{SD} under general covariance matrices. Here we consider two representative cases.

It is known that the power of T_{SD} is highest when the covariance matrix is diagonal. To compare our test with T_{SD} in this case, suppose that $\Sigma = \text{diag}(\lambda_1 + \sigma^2, \dots, \lambda_r + \sigma^2, \sigma^2, \dots, \sigma^2)$. Then $\mathbf{R} = \mathbf{I}_p$. Note that $\|\Sigma_d^{-1/2}(\mu_1 - \mu_2)\|^2 = \|\mu_1 - \mu_2\|^2 \|\Sigma^{-1/2}\delta\|^2$. We have

$$E \|\Sigma^{-1/2}\delta\|^2 = \frac{1}{p} \text{tr}(\Sigma^{-1}) = \sigma^{-2}(1 + o_P(1))$$

and

$$\text{Var}(\|\Sigma^{-1/2}\delta\|^2) = \frac{2}{p+2} \left(\frac{1}{p} \text{tr}(\Sigma^{-2}) - \left(\frac{1}{p} \text{tr}(\Sigma^{-1}) \right)^2 \right) = o\left(\frac{1}{p}\right).$$

It follows that $\|\Sigma^{-1/2}\delta\|^2 = \sigma^{-2}(1 + o_P(1))$. Hence the asymptotic power function of T_{SD} equals to

$$\Phi\left(-\Phi^{-1}(1-\alpha) + \frac{\|\mu_1 - \mu_2\|^2}{\sqrt{2\tau^2 p \sigma^4}}\right).$$

Thus, under the diagonal covariance matrix, the asymptotic power function of our test is the same as that of T_{SD} .

In the second case, we consider the uniform correlation structure $\Sigma = (1-\rho)\mathbf{I}_p + \rho\mathbf{1}_p\mathbf{1}_p^T$ ($0 < \rho < 1$) which is far away from a diagonal matrix. In this case, the diagonal entries of Σ are all equal to 1. We have

$$\frac{\|\Sigma_d^{-1/2}(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 \text{tr}(\mathbf{R}^2)}} = \frac{1}{\sqrt{\rho^2 p + (1-\rho^2)}} \frac{\|\mu_1 - \mu_2\|^2}{\sqrt{2\tau^2 p}}.$$

Then the asymptotic relative efficiency of our test with respect to T_{SD} is $\sqrt{\rho^2 p + (1-\rho^2)}$. Hence our test tends to be much more powerful than T_{SD} under the uniform correlation structure.

4. Numerical studies

4.1. Simulation results

In this section, we consider the simulation performance of the proposed test and compare it with several other alternatives: (1) The corrected T_{CQ} test procedure constructed in Section 2 (CCQ); (2) Ma

et al. (2015)'s test (FAST); (3) Chen and Qin (2010)'s test (CQ); (4) Srivastava and Du (2008)'s test (SD); (5) Lopes et al. (2011)'s test (LJW). The data generation mechanism is as follow. We randomly choose a $\mathbf{U} \in \mathbb{O}_{p \times p}$ from Haar invariant distribution. Let d_i equal to p^β plus a random error from $U(0, 1)$ (Uniform distribution on the interval $(0, 1)$), $i = 1, \dots, r$. Define the $p \times p$ diagonal matrix \mathbf{D} as $\mathbf{D} = \text{diag}(\sqrt{d_1}, \dots, \sqrt{d_r}, 1, \dots, 1)$. Then, we independently generate data by the formula

$$X_{k,i} = \mu_k + \mathbf{UDY}_{k,i} \quad i = 1, \dots, n_k \text{ and } k = 1, 2,$$

where $Y_{k,i}$ is a p dimensional random vector whose entries are i.i.d. random variables with common distribution F . We will consider three different distributions of F .

- Normal: $F \sim N(0, 1)$.
- Chi-squared: $F \sim (\chi^2(4) - 4)/\sqrt{8}$.
- Student's t : $F \sim t_4/\sqrt{2}$, where t_4 is a Student's t random variable with degree of freedom 4.

Throughout the simulations, we take the nominal level $\alpha = 0.05$ and $r = 2$.

The critical value of our test is determined by the quantile of $(\chi^2(p-r) - (p-r))/\sqrt{2(p-r)}$. Our test is asymptotically level α since the distribution of Q is asymptotically equal to that of $(\chi^2(p-r) - (p-r))/\sqrt{2(p-r)}$. Now we use Q-Q plot to compare these two distributions in the finite sample case. Figure 1 displays the Q-Q plots in different combinations of sample size and dimension. For $n_1 = n_2 = 50$, $p = 500$ or 800, the right tail of Q is a little heavier than the standardized chi-squared distribution. Otherwise, the distribution of Q can be well approximated by the standardized chi-squared distribution. Figure 2 displays the Q-Q plots in different combinations of F and β . It can be seen that the distribution of Q is very close to that of the standardized chi-squared distribution, even under non-normal distributions.

Next, we consider the simulation of the empirical level. Samples are repeatedly generated 2000 times to calculate empirical level. The results are listed in Tables 1-3. We can see that the empirical levels of the CQ test are larger than the nominal level in all cases. The empirical levels of the SD test are close to the nominal level when $\beta = 1/2$, but tend to be smaller than the nominal level as β increases. The empirical levels of the LJW test are very close to the nominal level. This is not surprising since the LJW test is exact under the normal distribution. The empirical levels of the FAST test are very close to the nominal level in most cases, but tend to be smaller than the nominal level when $n_1 = n_2 = 50$, $p = 800$ and $\beta = 1/2$. The empirical levels of the CCQ test are very close to the nominal level in all cases. The empirical levels of our test are a little inflated for $n_1 = n_2 = 50$, but converge to the nominal level as the sample size increases.

Now we consider the simulation of the empirical power. In view of Corollary 4, we define the SNR as $\text{SNR} = \|\mu_1 - \mu_2\|^2 / (\sigma^2 \sqrt{2\tau^2 p})$. We take $\mu_1 = \mathbf{0}_p$. The orientation of μ_2 is from the uniform distribution on the unit sphere. The norm of μ_2 is selected to make SNR equal to specific values. Samples are repeatedly generated 2000 times to calculate empirical power. The simulation results are illustrated in Figures 3 and 4. From the results, we can see that our test outperforms the other five tests substantially under the spiked covariance model.

4.2. Real data analysis

In this section, we study the practical problem considered in Ma et al. (2015). The task is to test whether Monday stock returns are equal to those of other trading days on average. Define an observation to be the log returns of stocks in a day. Hence p is the total number of stocks. Let sample 1 and sample 2 be the observations on Monday and the other trading days, respectively. We would like to test $H_0 : \mu_1 = \mu_2$ v.s. $H_1 : \mu_1 \neq \mu_2$. We collected the data of $p = 710$ stocks of China from 01/04/2013 to 12/31/2014. There are total $n_1 = 95$ Mondays and $n_2 = 388$ other trading days.

We assume $\Sigma_1 = \Sigma_2$. The first eigenvalue of \mathbf{S} is 0.14, which is significantly larger than the others. In fact, the second eigenvalue is 0.02. Hence there's clearly a spiked eigenvalue. We take $r = 1$ and perform our test. The p value is 0.149, which is obtained by permutation method with 1000 permutations. Hence, the null hypothesis can not be rejected for $\alpha = 0.05$. We draw the same conclusion as Ma et al. (2015).

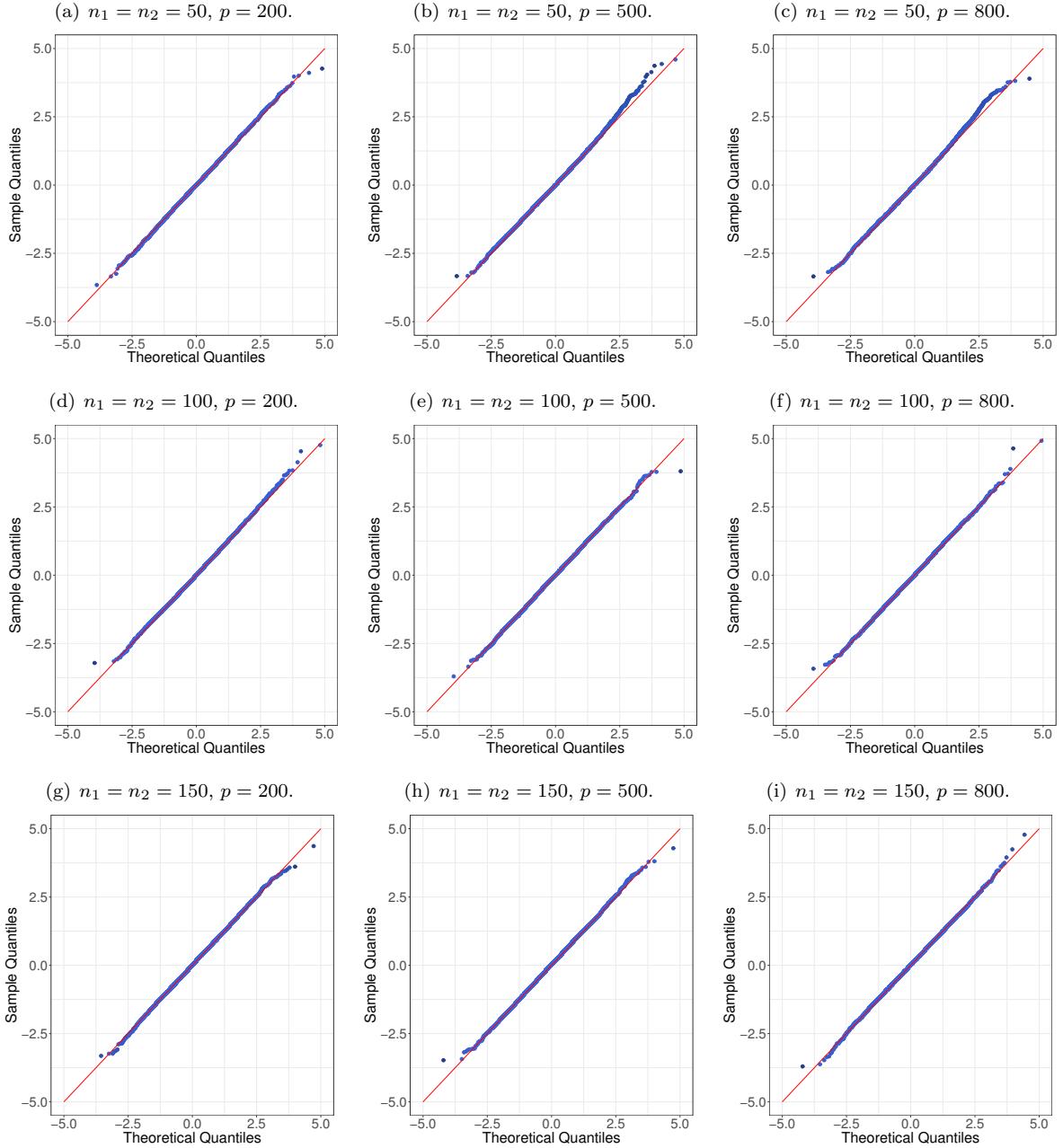


Figure 1: Q-Q plots of the empirical distribution of Q against that of $(\chi^2(p - r) - (p - r))/\sqrt{2(p - r)}$ based on 10000 independently generated Q . In all cases, F is normal and $\beta = 1$.

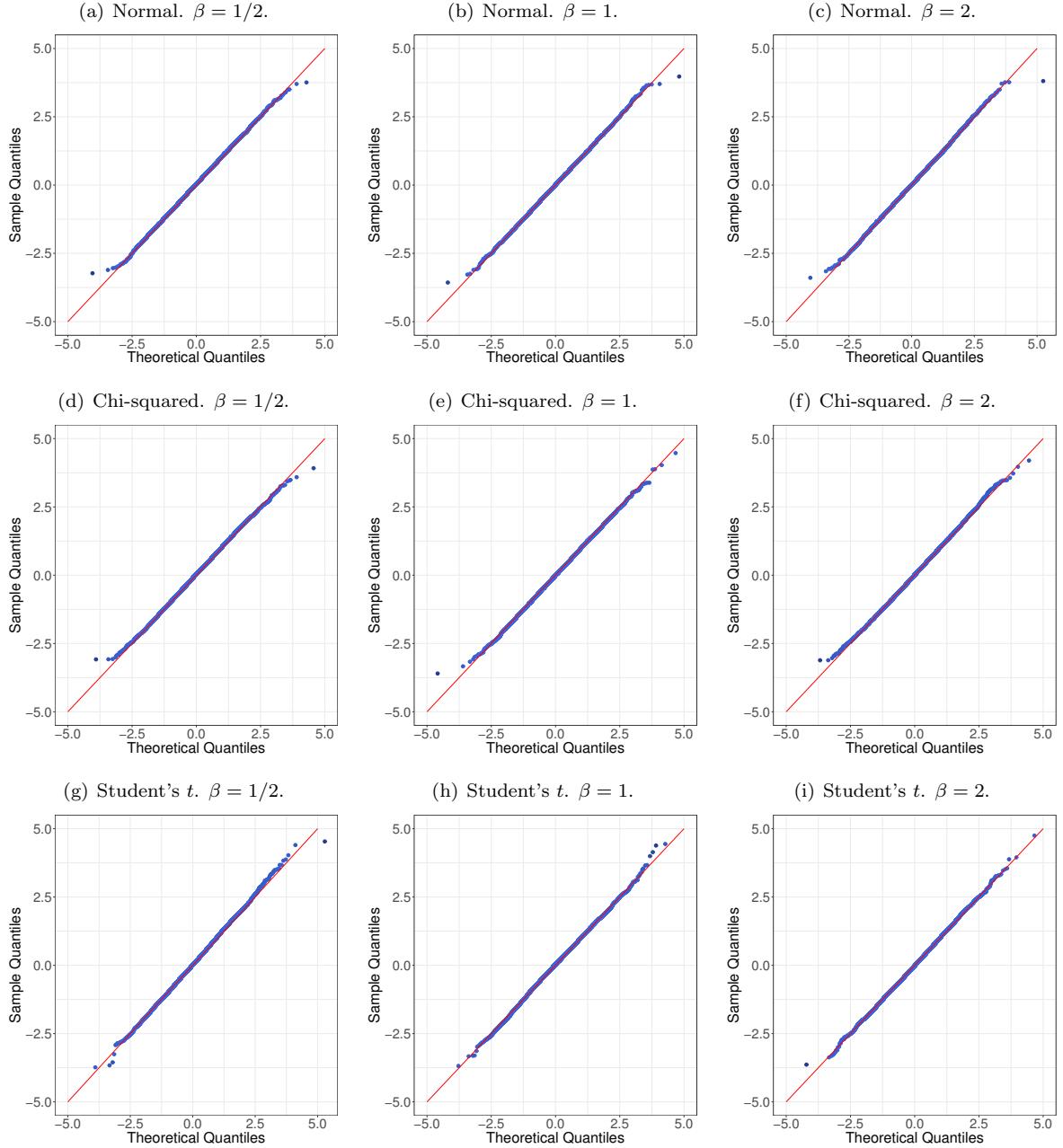


Figure 2: Q-Q plots of the empirical distribution of Q against that of $(\chi^2(p - r) - (p - r))/\sqrt{2(p - r)}$ based on 10000 independently generated Q . In all cases, $n_1 = n_2 = 100$ and $p = 500$.

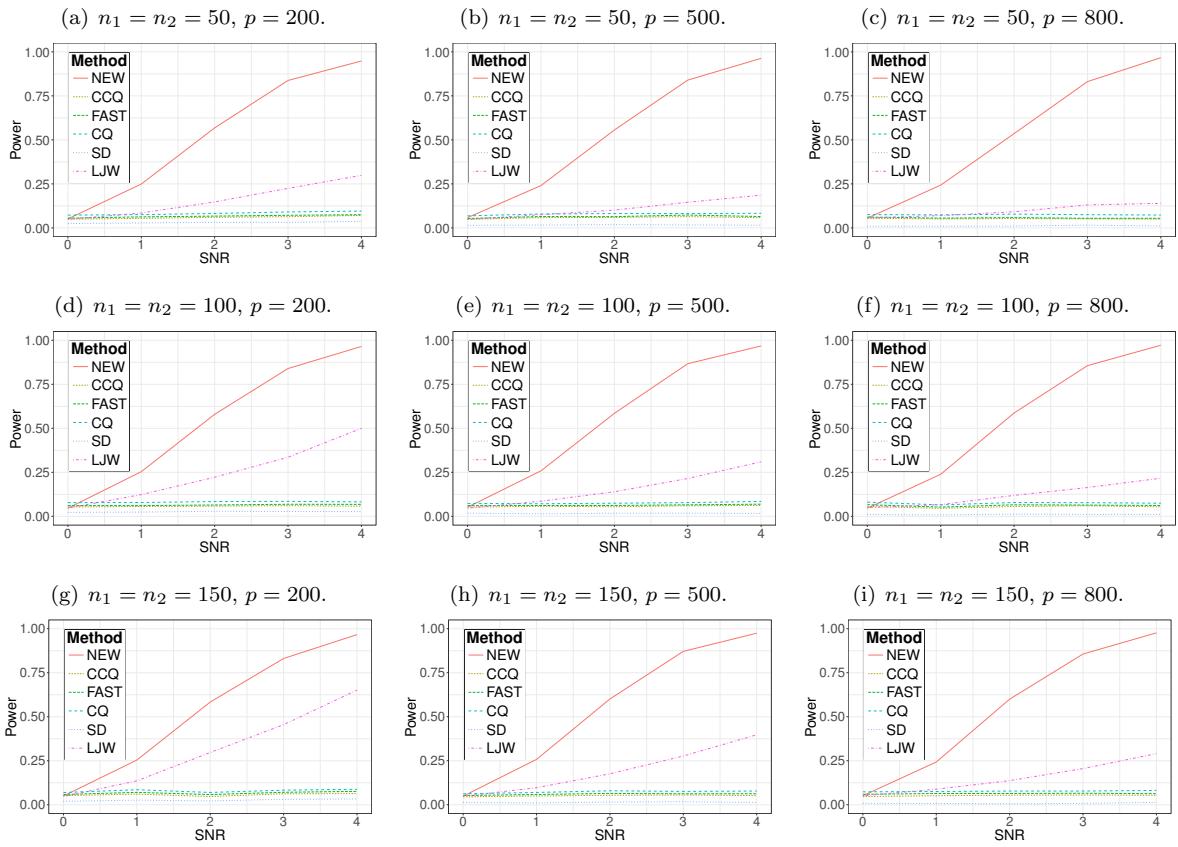


Figure 3: Empirical power of our test and competing tests. In all cases, F is normal and $\beta = 1$.

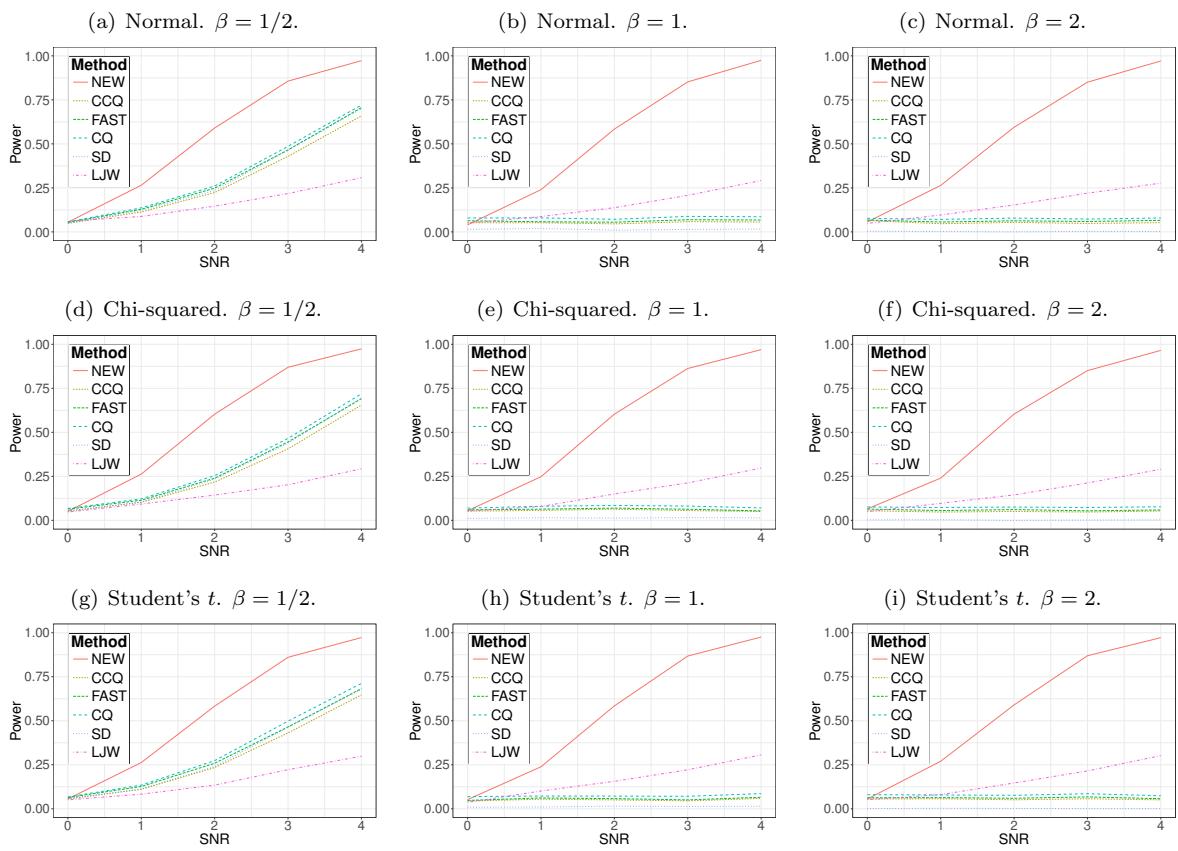


Figure 4: Empirical power of our test and competing tests. In all cases, $n_1 = n_2 = 100$ and $p = 500$.

Table 1: Empirical levels of our test and competing tests. $n_1 = n_2 = 50$.

p	Normal			Chi-squared			Student's t		
	200	500	800	200	500	800	200	500	800
$\beta = 1/2$									
NEW	0.0635	0.0565	0.0580	0.0610	0.0520	0.0695	0.0595	0.0640	0.0595
CCQ	0.0565	0.0535	0.0415	0.0520	0.0515	0.0480	0.0560	0.0420	0.0445
FAST	0.0525	0.0485	0.0290	0.0485	0.0395	0.0395	0.0560	0.0365	0.0325
CQ	0.0630	0.0645	0.0575	0.0675	0.0665	0.0750	0.0630	0.0725	0.0610
SD	0.0515	0.0500	0.0440	0.0525	0.0465	0.0580	0.0440	0.0515	0.0405
LJW	0.0535	0.0480	0.0455	0.0510	0.0575	0.0440	0.0440	0.0500	0.0505
$\beta = 1$									
NEW	0.0560	0.0585	0.0620	0.0625	0.0600	0.0505	0.0505	0.0585	0.0510
CCQ	0.0480	0.0445	0.0570	0.0490	0.0525	0.0525	0.0540	0.0400	0.0425
FAST	0.0505	0.0470	0.0590	0.0500	0.0550	0.0530	0.0560	0.0420	0.0435
CQ	0.0710	0.0785	0.0720	0.0725	0.0775	0.0730	0.0770	0.0715	0.0675
SD	0.0200	0.0165	0.0130	0.0225	0.0130	0.0100	0.0215	0.0075	0.0100
LJW	0.0470	0.0520	0.0460	0.0565	0.0470	0.0535	0.0485	0.0440	0.0540
$\beta = 2$									
New	0.0545	0.0535	0.0470	0.0550	0.0545	0.0585	0.0505	0.0585	0.0540
CCQ	0.0510	0.0505	0.0495	0.0480	0.0540	0.0530	0.0485	0.0410	0.0475
FAST	0.0540	0.0530	0.0510	0.0500	0.0565	0.0555	0.0520	0.0445	0.0495
CQ	0.0690	0.0760	0.0675	0.0745	0.0665	0.0630	0.0785	0.0745	0.0700
SD	0.0060	0.0010	0.0000	0.0045	0.0010	0.0010	0.0055	0.0010	0.0010
LJW	0.0435	0.0500	0.0535	0.0410	0.0450	0.0430	0.0530	0.0500	0.0515

5. Discussion

This paper is concerned with the problem of testing the equality of means under the spiked covariance model. We find that Chen and Qin (2010)'s test statistic is not asymptotically normal distributed under the spiked covariance model. Although a corrected T_{CQ} test can be defined, its power behavior is unsatisfactory. The recently proposed random projection test procedures suggest that the power of tests may be boosted using the projected data. By maximizing an average SNR, we find the optimal projection subspace is the orthogonal complement of the principal subspace. By projecting data onto the (estimated) optimal subspace, we proposed a new test. The asymptotic normality of the new test statistic is proved and the asymptotic power of the new test is given. Theoretical and simulation results show that the new test outperforms the competing tests substantially under the spiked covariance model.

In another paper, Zhao and Xu (2016) proved that their test statistic can be written in the form of projection. Their simulation results showed that their test performs well under strong correlations. Our work partially explains why their test performs well, although their projection is different from ours.

The spiked covariance model is an important correlation pattern and has been widely studied in terms of PCA. In PCA, authors focus on the principal subspace. As our work have shown, however, the complement of the principal subspace may be more useful in hypothesis testing.

In our paper, r is assumed to be known. If r is an unknown positive number, a consistent estimator of r is

$$\hat{r} = \operatorname{argmax}_{l \leq R} \frac{\lambda_l(\mathbf{S})}{\lambda_{l+1}(\mathbf{S})},$$

where R is a hyperparameter. See Ahn and Horenstein (2013) for detail.

The asymptotic normality of our test statistic relies on the assumption $\sqrt{p}/n \rightarrow 0$. In the situation of small n or very large p , the critical value of the new test can be determined by permutation method. It

Table 2: Empirical levels of our test and competing tests. $n_1 = n_2 = 100$.

p	Normal			Chi-squared			Student's t		
	200	500	800	200	500	800	200	500	800
$\beta = 1/2$									
NEW	0.0515	0.0470	0.0550	0.0595	0.0535	0.0535	0.0545	0.0565	0.0515
CCQ	0.0480	0.0475	0.0475	0.0600	0.0510	0.0380	0.0430	0.0430	0.0450
FAST	0.0510	0.0505	0.0475	0.0660	0.0555	0.0390	0.0480	0.0465	0.0460
CQ	0.0680	0.0610	0.0640	0.0650	0.0575	0.0630	0.0585	0.0705	0.0510
SD	0.0580	0.0520	0.0535	0.0560	0.0500	0.0485	0.0410	0.0615	0.0445
LJW	0.0440	0.0440	0.0460	0.0495	0.0470	0.0475	0.0440	0.0450	0.0505
$\beta = 1$									
NEW	0.0530	0.0500	0.0475	0.0445	0.0545	0.0510	0.0455	0.0475	0.0500
CCQ	0.0545	0.0480	0.0530	0.0500	0.0370	0.0400	0.0460	0.0405	0.0530
FAST	0.0550	0.0500	0.0550	0.0525	0.0415	0.0450	0.0465	0.0450	0.0555
CQ	0.0715	0.0665	0.0880	0.0655	0.0680	0.0730	0.0785	0.0655	0.0740
SD	0.0220	0.0135	0.0085	0.0205	0.0140	0.0095	0.0225	0.0115	0.0100
LJW	0.0525	0.0515	0.0485	0.0420	0.0475	0.0475	0.0485	0.0445	0.0465
$\beta = 2$									
NEW	0.0510	0.0500	0.0565	0.0470	0.0525	0.0445	0.0510	0.0465	0.0495
CCQ	0.0420	0.0605	0.0530	0.0545	0.0445	0.0485	0.0615	0.0490	0.0430
FAST	0.0460	0.0635	0.0580	0.0565	0.0475	0.0525	0.0640	0.0520	0.0455
CQ	0.0855	0.0775	0.0675	0.0615	0.0695	0.0755	0.0750	0.0740	0.0700
SD	0.0055	0.0025	0.0015	0.0030	0.0000	0.0005	0.0025	0.0005	0.0005
LJW	0.0480	0.0495	0.0545	0.0410	0.0455	0.0510	0.0535	0.0435	0.0550

remains a theoretical interest to study the asymptotic behavior of permutation based test in these situations.

The normality assumption is very important for this paper. In particular, our derivations and proofs heavily rely on the fact that \mathbf{S} is independent of \bar{X}_1, \bar{X}_2 . Nevertheless, the simulation results show that the performance of our test is satisfactory under certain non-normal distributions. It is an interesting but challenging problem to study the robustness of our test and extend the results to other distributions.

Acknowledgments

The authors thank the Editor-in-Chief, the Associate Editor and the referees for helpful comments and suggestions that improve the paper considerably. This work was supported by the National Natural Science Foundation of China under Grant No. 11471035, 11471030.

Appendix A Proofs of the results in Section 2

Proof of Lemma 1. By a standard orthogonal transformation, we can write

$$\frac{Y_n^T \mathbf{A}_n Y_n - \mathbb{E} Y_n^T \mathbf{A}_n Y_n}{[\text{Var}(Y_n^T \mathbf{A}_n Y_n)]^{1/2}} = \sum_{i=1}^{k_n} \frac{\lambda_i(\mathbf{A}_n)}{[2 \text{tr}(\mathbf{A}_n^2)]^{1/2}} (Z_{ni}^2 - 1), \quad (16)$$

where Z_{n1}, \dots, Z_{nk_n} are independent standard normal random variables.

Table 3: Empirical levels of our test and competing tests. $n_1 = n_2 = 150$.

p	Normal			Chi-squared			Student's t		
	200	500	800	200	500	800	200	500	800
$\beta = 1/2$									
NEW	0.0520	0.0495	0.0575	0.0560	0.0505	0.0540	0.0495	0.0540	0.0540
CCQ	0.0475	0.0545	0.0485	0.0450	0.0415	0.0480	0.0510	0.0570	0.0440
FAST	0.0510	0.0595	0.0510	0.0510	0.0440	0.0530	0.0625	0.0645	0.0495
CQ	0.0655	0.0780	0.0710	0.0720	0.0515	0.0635	0.0565	0.0725	0.0660
SD	0.0530	0.0675	0.0600	0.0585	0.0445	0.0585	0.0465	0.0600	0.0560
LJW	0.0540	0.0520	0.0535	0.0545	0.0555	0.0480	0.0470	0.0500	0.0605
$\beta = 1$									
NEW	0.0490	0.0520	0.0490	0.0515	0.0465	0.0425	0.0470	0.0520	0.0460
CCQ	0.0405	0.0505	0.0435	0.0470	0.0445	0.0515	0.0510	0.0490	0.0450
FAST	0.0430	0.0530	0.0455	0.0495	0.0480	0.0530	0.0545	0.0510	0.0470
CQ	0.0785	0.0690	0.0740	0.0720	0.0695	0.0715	0.0730	0.0795	0.0725
SD	0.0265	0.0115	0.0105	0.0220	0.0130	0.0090	0.0145	0.0145	0.0080
LJW	0.0430	0.0605	0.0460	0.0515	0.0580	0.0525	0.0580	0.0595	0.0510
$\beta = 2$									
NEW	0.0495	0.0495	0.0565	0.0510	0.0490	0.0560	0.0575	0.0520	0.0520
CCQ	0.0495	0.0405	0.0480	0.0460	0.0530	0.0545	0.0475	0.0395	0.0510
FAST	0.0525	0.0415	0.0505	0.0490	0.0560	0.0560	0.0490	0.0400	0.0555
CQ	0.0675	0.0705	0.0695	0.0765	0.0655	0.0775	0.0785	0.0720	0.0875
SD	0.0040	0.0020	0.0005	0.0030	0.0000	0.0000	0.0025	0.0010	0.0000
LJW	0.0490	0.0475	0.0490	0.0465	0.0535	0.0480	0.0490	0.0445	0.0490

If 5 holds, then

$$\begin{aligned}
 & \sum_{i=1}^{k_n} \mathbb{E} \left[\frac{\lambda_i^2(\mathbf{A}_n)}{2\text{tr}(\mathbf{A}_n^2)} (Z_{ni}^2 - 1)^2 \left\{ \frac{\lambda_i^2(\mathbf{A}_n)}{2\text{tr}(\mathbf{A}_n^2)} (Z_{ni}^2 - 1)^2 \geq \epsilon \right\} \right] \\
 & \leq \sum_{i=1}^{k_n} \frac{\lambda_i^2(\mathbf{A}_n)}{2\text{tr}(\mathbf{A}_n^2)} \mathbb{E} \left[(Z_{n1}^2 - 1)^2 \left\{ \frac{\lambda_1(\mathbf{A}_n^2)}{2\text{tr}(\mathbf{A}_n^2)} (Z_{n1}^2 - 1)^2 \geq \epsilon \right\} \right] \\
 & = \frac{1}{2} \mathbb{E} \left[(Z_{n1}^2 - 1)^2 \left\{ \frac{\lambda_1(\mathbf{A}_n^2)}{2\text{tr}(\mathbf{A}_n^2)} (Z_{n1}^2 - 1)^2 \geq \epsilon \right\} \right] \rightarrow 0.
 \end{aligned}$$

Hence 4 follows by Lindeberg's central limit theorem.

Conversely, if 4 holds, we will prove that there is a subsequence of $\{n\}$ along which 5 holds. Then 5 follows by a standard contradiction argument.

Denote $c_{ni} = \lambda_i(\mathbf{A}_n)/[2\text{tr}(\mathbf{A}_n^2)]^{1/2}$, $i = 1, \dots, k_n$. Since 4 holds, the characteristic function of $\sum_{i=1}^{k_n} c_{ni}(Z_{ni}^2 - 1)$ converges to $\exp(-t^2/2)$ for every t . Denote by $\log z$ ($z \in \mathbb{C}$) the principal branch of the complex loga-

rithm. For $t \in (-1/2, 1/2)$, we have

$$\begin{aligned} \mathbb{E} \left[\exp \left(it \sum_{j=1}^{k_n} c_{nj} (Z_{nj}^2 - 1) \right) \right] &= \exp \left(-i \left(\sum_{j=1}^{k_n} c_{nj} \right) t - \frac{1}{2} \sum_{j=1}^{k_n} \log(1 - i2c_{nj}t) \right) \\ &= \exp \left(-i \left(\sum_{j=1}^{k_n} c_{nj} \right) t + \frac{1}{2} \sum_{j=1}^{k_n} \sum_{l=1}^{+\infty} \frac{1}{l} (i2c_{nj}t)^l \right) = \exp \left(-i \left(\sum_{j=1}^{k_n} c_{nj} \right) t + \frac{1}{2} \sum_{l=1}^{+\infty} \left[\sum_{j=1}^{k_n} (c_{nj})^l \right] \frac{1}{l} (i2t)^l \right) \\ &= \exp \left(-\frac{1}{2} t^2 + \frac{1}{2} \sum_{l=3}^{+\infty} \left[\sum_{j=1}^{k_n} (c_{nj})^l \right] \frac{1}{l} (i2t)^l \right), \end{aligned}$$

where the second equality holds since $0 \leq c_{ni} \leq \sqrt{2}/2$ by definition. Let $b_{nl} = \sum_{j=1}^{k_n} (c_{nj})^l$, $n = 1, 2, \dots$ and $l = 3, 4, \dots$. Note that for $l \geq 3$, we have

$$|b_{nl}| = \left| \sum_{j=1}^{k_n} (c_{nj})^l \right| \leq \left| \sum_{j=1}^{k_n} (c_{nj})^2 \right| = 1/2.$$

By Helly's selection theorem, there's a subsequence of $\{n\}$ along which $\lim_{n \rightarrow \infty} b_{nl} = b_l$ exists for every l . For this subsequence, applying the dominated convergence theorem yields

$$\mathbb{E} \left[\exp \left(it \sum_{j=1}^{k_n} c_{nj} (Z_{nj}^2 - 1) \right) \right] \rightarrow \exp \left(-\frac{1}{2} t^2 + \frac{1}{2} \sum_{l=3}^{+\infty} b_l \frac{1}{l} (i2t)^l \right), \quad t \in \left(-\frac{1}{2}, \frac{1}{2} \right).$$

But the left hand side converges to $\exp(-t^2/2)$. It follows that

$$-\frac{1}{2} t^2 + \frac{1}{2} \sum_{l=3}^{+\infty} b_l \frac{1}{l} (i2t)^l = -\frac{1}{2} t^2 + 2\pi m i, \quad t \in \left(-\frac{1}{2}, \frac{1}{2} \right),$$

for some integer m . By the uniqueness of power series, we must have $m = 0$ and $b_l = 0$ for $l \geq 3$. Then 5 follows by noting that $b_{n4} \geq \max_j (c_{nj})^4$. \square

Proofs of Theorems 1 and 2. In both Theorems, (a) is a corollary of (b). We shall prove (b) of Theorems 1 and 2 simultaneously.

Note that $(n_k - 1)\mathbf{S}_k \sim \text{Wishart}_p(n_k - 1, \boldsymbol{\Sigma})$, $k = 1, 2$, where $\text{Wishart}_p(m, \Psi)$ is the p dimensional Wishart distribution with parameter Ψ and m degrees of freedom. We have

$$\mathbb{E} \left(\frac{1}{n_1} \text{tr } \mathbf{S}_1 + \frac{1}{n_2} \text{tr } \mathbf{S}_2 \right) = \tau \text{tr } \boldsymbol{\Sigma},$$

and

$$\begin{aligned} \text{Var} \left(\frac{1}{n_1} \text{tr } \mathbf{S}_1 + \frac{1}{n_2} \text{tr } \mathbf{S}_2 \right) &= \left(\frac{2}{n_1^2(n_1 - 1)} + \frac{2}{n_2^2(n_2 - 1)} \right) \text{tr } \boldsymbol{\Sigma}^2 \\ &= O \left(\frac{1}{n^3} (p^{2\beta} + p) \right) = O \left(\frac{p^{2\beta}}{n^3} \right). \end{aligned}$$

It follows that

$$\begin{aligned} \frac{1}{n_1} \text{tr } \mathbf{S}_1 + \frac{1}{n_2} \text{tr } \mathbf{S}_2 &= \tau \text{tr } \boldsymbol{\Sigma} + O_P \left(\frac{1}{n\sqrt{n}} p^\beta \right) \\ &= \tau \sum_{i=1}^r (\lambda_i + \sigma^2) + \tau(p - r)\sigma^2 + O_P \left(\frac{1}{n\sqrt{n}} p^\beta \right) \\ &= \tau p^\beta \sum_{i=1}^r \omega_i + \tau(p - r)\sigma^2 + o_P \left(\frac{1}{n} p^\beta \right). \end{aligned}$$

Thus,

$$\frac{1}{\tau p^\beta} \left(\frac{1}{n_1} \operatorname{tr} \mathbf{S}_1 + \frac{1}{n_2} \operatorname{tr} \mathbf{S}_2 \right) = \sum_{i=1}^r \omega_i + p^{1-\beta} \sigma^2 + o_P(1). \quad (17)$$

Next we deal with $\|\bar{X}_1 - \bar{X}_2\|^2$. Note that we have

$$\|\bar{X}_1 - \bar{X}_2\|^2 = \|\mathbf{V}^T(\bar{X}_1 - \bar{X}_2)\|^2 + \|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2.$$

These two terms are independent. For the first term, note that $\mathbf{V}^T(\bar{X}_1 - \bar{X}_2) \sim N_r(\mathbf{V}^T(\mu_1 - \mu_2), \tau(\Lambda + \sigma^2 \mathbf{I}_r))$, we have

$$\begin{aligned} \|\mathbf{V}^T(\bar{X}_1 - \bar{X}_2)\|^2 &\stackrel{d}{=} \sum_{i=1}^r \left(\sqrt{\tau(\lambda_i + \sigma^2)} Z_i + (\mathbf{V}^T(\mu_1 - \mu_2))_i \right)^2 \\ &= \tau p^\beta \sum_{i=1}^r \left(\sqrt{p^{-\beta}(\lambda_i + \sigma^2)} Z_i + \frac{1}{\sqrt{\tau p^\beta}} (\mathbf{V}^T(\mu_1 - \mu_2))_i \right)^2. \end{aligned}$$

By the assumptions of the theorem, we have that

$$\frac{1}{\tau p^\beta} \|\mathbf{V}^T(\bar{X}_1 - \bar{X}_2)\|^2 \xrightarrow{w} \sum_{i=1}^r (\sqrt{\omega_i} Z_i + \zeta_i)^2. \quad (18)$$

As for $\|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2$, we have that

$$\begin{aligned} \|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 &= \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2) + \tilde{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 \\ &= \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 + \|\tilde{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 + 2(\mu_1 - \mu_2)^T \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2)). \end{aligned}$$

Since $\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2) \sim N_{p-r}(\tilde{\mathbf{V}}^T(\mu_1 - \mu_2), \sigma^2 \tau \mathbf{I}_{p-r})$, by the central limit theorem, we have

$$\frac{\|\tilde{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 - \sigma^2 \tau(p-r)}{\sigma^2 \tau \sqrt{2(p-r)}} \xrightarrow{\mathcal{L}} N(0, 1).$$

For the intersection term, we have

$$\begin{aligned} 2(\mu_1 - \mu_2)^T \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2)) &\sim N(0, 4\sigma^2 \tau \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2) \\ &= O_P(\sqrt{\tau} \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|) = o_P(\tau p^\beta). \end{aligned}$$

It follows that

$$\frac{1}{\tau p^\beta} (\|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 - \sigma^2 \tau(p-r) - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2) \xrightarrow{\mathcal{L}} \sqrt{2} \sigma^2 \delta_{\{\frac{1}{2}\}}(\beta) Z_0, \quad (19)$$

where $\delta_{\frac{1}{2}}(\beta)$ equals 1 if $\beta = 1/2$ and equals 0 otherwise.

Combining (17) (18) and (19) leads to

$$\begin{aligned} \frac{1}{\tau p^\beta} T_{CQ} &= \frac{1}{\tau p^\beta} (\|\bar{X}_1 - \bar{X}_2\|^2 - \frac{1}{n_1} \operatorname{tr} \mathbf{S}_1 - \frac{1}{n_2} \operatorname{tr} \mathbf{S}_2) \\ &= \frac{1}{\tau p^\beta} \|\mathbf{V}^T(\bar{X}_1 - \bar{X}_2)\|^2 + \frac{1}{\tau p^\beta} (\|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 - \sigma^2 \tau(p-r) - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2) \\ &\quad - \frac{1}{\tau p^\beta} \left(\frac{1}{n_1} \operatorname{tr} \mathbf{S}_1 + \frac{1}{n_2} \operatorname{tr} \mathbf{S}_2 \right) + \frac{\sigma^2(p-r)}{p^\beta} + \frac{1}{\tau p^\beta} \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 \\ &= \sum_{i=1}^r (\sqrt{\omega_i} Z_i + \zeta_i)^2 + \sqrt{2} \sigma^2 \delta_{\{\frac{1}{2}\}}(\beta) Z_0 - \left(\sum_{i=1}^r \omega_i + p^{1-\beta} \sigma^2 \right) + \frac{\sigma^2(p-r)}{p^\beta} + \zeta^* + o_P(1) \\ &\xrightarrow{\mathcal{L}} \sum_{i=1}^r (\sqrt{\omega_i} Z_i + \zeta_i)^2 + \zeta^* + \sqrt{2} \sigma^2 \delta_{\{\frac{1}{2}\}}(\beta) Z_0 - \sum_{i=1}^r \omega_i. \end{aligned}$$

This implies the conclusions of Theorem 1 and Theorem 2. \square

Appendix B Proofs of the results in Section 3

Lemma 2 (Weyl's inequality). *Let \mathbf{A} and \mathbf{B} be two symmetric $n \times n$ matrices and $\mathbf{C} = \mathbf{A} + \mathbf{B}$. If $r + s - 1 \leq i \leq j + k - n$, we have*

$$\lambda_j(\mathbf{A}) + \lambda_k(\mathbf{B}) \leq \lambda_i(\mathbf{C}) \leq \lambda_r(\mathbf{A}) + \lambda_s(\mathbf{B}).$$

See, for example, Horn and Johnson (2012), Theorem 4.3.1.

Lemma 3 (Cai et al. (2015), Proposition 1). *Let \mathbf{A}_1 and \mathbf{A}_2 be $p \times p$ symmetric matrices. Let $r < p$ be arbitrary and let $\mathbf{V}_1, \mathbf{V}_2 \in \mathbb{O}_{p,r}$ be formed by the r leading singular vectors of \mathbf{A}_1 and \mathbf{A}_2 , respectively. Then*

$$\|\mathbf{A}_1 - \mathbf{A}_2\| \geq \frac{1}{2}(\lambda_r(\mathbf{A}_2) - \lambda_{r+1}(\mathbf{A}_2))\|\mathbf{V}_1\mathbf{V}_1^T - \mathbf{V}_2\mathbf{V}_2^T\|.$$

Lemma 4 (Davidson and Szarek (2001), Theorem II.7). *Let \mathbf{Z} be a $p \times n$ random matrix with i.i.d. $N(0, 1)$ entries. Then for any $t > 0$,*

$$\begin{aligned} \Pr(\sqrt{\lambda_1(\mathbf{Z}\mathbf{Z}^T)} > \sqrt{n} + \sqrt{p} + t) &\leq e^{-t^2/2}, \\ \Pr(\sqrt{\lambda_{\min(n,p)}(\mathbf{Z}\mathbf{Z}^T)} < \sqrt{n} - \sqrt{p} - t) &\leq e^{-t^2/2}. \end{aligned}$$

The following corollary of Lemma 4 is useful.

Corollary 1. *Suppose that \mathbf{W}_n is a $p \times p$ random matrix which is distributed as Wishart _{p} (n, \mathbf{I}_p). Then as $n, p \rightarrow \infty$,*

$$\left\| \frac{1}{n} \mathbf{W}_n - \mathbf{I}_p \right\| = O_P \left(\max \left(\sqrt{\frac{p}{n}}, \frac{p}{n} \right) \right).$$

Proof. Since the eigenvalues of $n^{-1}\mathbf{W}_n - \mathbf{I}_p$ are $n^{-1}\lambda_1(\mathbf{W}_n) - 1, \dots, n^{-1}\lambda_p(\mathbf{W}_n) - 1$, we have

$$\left\| \frac{1}{n} \mathbf{W}_n - \mathbf{I}_p \right\| = \max \left(\frac{1}{n} \lambda_1(\mathbf{W}_n) - 1, 1 - \frac{1}{n} \lambda_p(\mathbf{W}_n) \right).$$

This, together with the union bound, yields

$$\Pr \left(\left\| \frac{1}{n} \mathbf{W}_n - \mathbf{I}_p \right\| > 4 \left(\sqrt{\frac{p}{n}} + \frac{p}{n} \right) \right) \leq \Pr \left(\lambda_1(\mathbf{W}_n) > (\sqrt{n} + 2\sqrt{p})^2 \right) + \Pr \left(\lambda_p(\mathbf{W}_n) < n - 4\sqrt{np} - 4p \right).$$

For the first term, we have

$$\Pr \left(\lambda_1(\mathbf{W}_n) > (\sqrt{n} + 2\sqrt{p})^2 \right) = \Pr \left(\sqrt{\lambda_1(\mathbf{W}_n)} > \sqrt{n} + 2\sqrt{p} \right) \leq e^{-p/2},$$

where the last inequality follows from Lemma 4 with $t = \sqrt{p}$.

We now show that the second term is also bounded by $e^{-p/2}$. To see this, note that if $p > n/4$, then $n - 4\sqrt{np} - 4p \leq n - 4p < 0$. In this case, $\Pr \left(\lambda_p(\mathbf{W}_n) < n - 4\sqrt{np} - 4p \right) = 0$. If $p \leq n/4$, we have

$$\begin{aligned} \Pr \left(\lambda_p(\mathbf{W}_n) < n - 4\sqrt{np} - 4p \right) &\leq \Pr \left(\lambda_p(\mathbf{W}_n) < n - 4\sqrt{np} + 4p \right) \\ &= \Pr \left(\sqrt{\lambda_p(\mathbf{W}_n)} < \sqrt{n} - 2\sqrt{p} \right) \leq e^{-p/2}, \end{aligned}$$

where the last inequality follows from Lemma 4 with $t = \sqrt{p}$.

Now we conclude that

$$\Pr \left(\left\| \frac{1}{n} \mathbf{W}_n - \mathbf{I}_p \right\| > 4 \left(\sqrt{\frac{p}{n}} + \frac{p}{n} \right) \right) \leq 2e^{-p/2}.$$

Consequently,

$$\left\| \frac{1}{n} \mathbf{W}_n - \mathbf{I}_p \right\| = O_P \left(\sqrt{\frac{p}{n}} + \frac{p}{n} \right) = O_P \left(\max \left(\sqrt{\frac{p}{n}}, \frac{p}{n} \right) \right).$$

□

Lemma 5. Under Assumption 1, we have

$$\lambda_i(\mathbf{S}) = \boldsymbol{\lambda}_i + \frac{p+n-r}{n} \sigma^2 + O_P\left(\max\left(\frac{p^\beta}{\sqrt{n}}, 1\right)\right), \quad i = 1, \dots, r, \quad (20)$$

and

$$\hat{\sigma}^2 = \left(1 - \frac{r}{n}\right) \sigma^2 + O_P\left(\max\left(\frac{1}{\sqrt{np}}, \frac{1}{p}\right)\right). \quad (21)$$

Proof. Let $\boldsymbol{\Sigma} = \mathbf{U}\mathbf{E}\mathbf{U}^T$ denote the spectral decomposition of $\boldsymbol{\Sigma}$, where $\mathbf{U} = (\mathbf{V}, \tilde{\mathbf{V}})$ and $\mathbf{E} = \text{diag}(\boldsymbol{\lambda}_1 + \sigma^2, \dots, \boldsymbol{\lambda}_r + \sigma^2, \sigma^2, \dots, \sigma^2)$. Let \mathbf{Z} be a $p \times n$ random matrix with i.i.d. $N(0, 1)$ entries. Denote $\mathbf{Z} = (\mathbf{Z}_{(1)}^T, \mathbf{Z}_{(2)}^T)^T$, where $\mathbf{Z}_{(1)}$ and $\mathbf{Z}_{(2)}$ are the first r rows and last $p-r$ rows of \mathbf{Z} . Then the sample covariance matrix \mathbf{S} has the same distribution as the random matrix $n^{-1}\mathbf{U}\mathbf{E}^{1/2}\mathbf{Z}\mathbf{Z}^T\mathbf{E}^{1/2}\mathbf{U}^T$. So we have $\lambda_i(\mathbf{S}) \stackrel{d}{=} n^{-1}\lambda_i(\mathbf{Z}^T\mathbf{E}\mathbf{Z})$, $i = 1, \dots, r$ and $\text{tr}(\mathbf{S}) \stackrel{d}{=} n^{-1}\text{tr}(\mathbf{Z}^T\mathbf{E}\mathbf{Z})$.

To prove (20), we only need to consider $n^{-1}\lambda_i(\mathbf{Z}^T\mathbf{E}\mathbf{Z})$, $i = 1, \dots, r$. Note that $\mathbf{Z}^T\mathbf{E}\mathbf{Z} = \mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)\mathbf{Z}_{(1)} + \sigma^2\mathbf{Z}_{(2)}^T\mathbf{Z}_{(2)}$. From this and Weyl's inequality, for $i = 1, \dots, r$, we have

$$\begin{aligned} & \left| \frac{1}{n}\lambda_i(\mathbf{Z}^T\mathbf{E}\mathbf{Z}) - \frac{1}{n}\lambda_i(\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)\mathbf{Z}_{(1)}) - \frac{p-r}{n}\sigma^2 \right| \\ &= \left| \lambda_i\left(\frac{1}{n}\mathbf{Z}^T\mathbf{E}\mathbf{Z}\right) - \lambda_i\left(\frac{1}{n}\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)\mathbf{Z}_{(1)} + \frac{p-r}{n}\sigma^2\mathbf{I}_n\right) \right| \\ &\leq \left\| \frac{1}{n}\mathbf{Z}^T\mathbf{E}\mathbf{Z} - \frac{1}{n}\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)\mathbf{Z}_{(1)} - \frac{p-r}{n}\sigma^2\mathbf{I}_n \right\| \\ &= \left\| \frac{1}{n}\sigma^2\mathbf{Z}_{(2)}^T\mathbf{Z}_{(2)} - \frac{p-r}{n}\sigma^2\mathbf{I}_n \right\| \\ &= \frac{p-r}{n}\sigma^2 \left\| \frac{1}{p-r}\mathbf{Z}_{(2)}^T\mathbf{Z}_{(2)} - \mathbf{I}_n \right\|. \end{aligned}$$

But Corollary 1 implies that

$$\left\| \frac{1}{p-r}\mathbf{Z}_{(2)}^T\mathbf{Z}_{(2)} - \mathbf{I}_n \right\| = O_P\left(\max\left(\sqrt{\frac{n}{p-r}}, \frac{n}{p-r}\right)\right).$$

Thus, for $i = 1, \dots, r$,

$$\frac{1}{n}\lambda_i(\mathbf{Z}^T\mathbf{E}\mathbf{Z}) = \frac{1}{n}\lambda_i(\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)\mathbf{Z}_{(1)}) + \frac{p-r}{n}\sigma^2 + O_P\left(\max\left(\sqrt{\frac{p}{n}}, 1\right)\right). \quad (22)$$

Next we deal with $n^{-1}\lambda_i(\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)\mathbf{Z}_{(1)})$. For $i = 1, \dots, r$, we have

$$\begin{aligned} & \left| \frac{1}{n}\lambda_i(\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)\mathbf{Z}_{(1)}) - (\boldsymbol{\lambda}_i + \sigma^2) \right| \\ &= \left| \lambda_i\left(\frac{1}{n}(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)^{1/2}\right) - \lambda_i(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r) \right| \\ &\leq \left\| \frac{1}{n}(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)^{1/2} - (\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r) \right\| \\ &= \left\| (\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)^{1/2} \left(\frac{1}{n}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T - \mathbf{I}_r \right) (\boldsymbol{\Lambda} + \sigma^2\mathbf{I}_r)^{1/2} \right\| \\ &\leq (\boldsymbol{\lambda}_1 + \sigma^2) \left\| \frac{1}{n}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T - \mathbf{I}_r \right\| \\ &= O_P\left(\frac{\boldsymbol{\lambda}_1}{\sqrt{n}}\right), \end{aligned}$$

where the first inequality follows from Weyl's inequality and the last equality follows from Corollary 1. This, together with (22), leads to

$$\frac{1}{n} \lambda_i(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) = \boldsymbol{\lambda}_i + \frac{p+n-r}{n} \sigma^2 + O_P\left(\max\left(\sqrt{\frac{p}{n}}, \frac{\boldsymbol{\lambda}_1}{\sqrt{n}}, 1\right)\right), \quad i = 1, \dots, r.$$

Then (20) follows from $\boldsymbol{\lambda}_1 \asymp p^\beta$ and $\beta \geq 1/2$.

Now we prove (21). We note that

$$\hat{\sigma}^2 = \frac{1}{p-r} \left(\text{tr}(\mathbf{S}) - \sum_{i=1}^r \lambda_i(\mathbf{S}) \right) \stackrel{d}{=} \frac{1}{n(p-r)} \left(\text{tr}(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) - \sum_{i=1}^r \lambda_i(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) \right).$$

For the term $\text{tr}(\mathbf{Z}^T \mathbf{E} \mathbf{Z})$, we have $\text{tr}(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) = \text{tr}(\mathbf{Z}_{(1)}^T (\boldsymbol{\Lambda} + \sigma^2 \mathbf{I}_r) \mathbf{Z}_{(1)}) + \sigma^2 \text{tr}(\mathbf{Z}_{(2)}^T \mathbf{Z}_{(2)})$. Since $\text{tr}(\mathbf{Z}_{(2)}^T \mathbf{Z}_{(2)}) \sim \chi^2(n(p-r))$, by the central limit theorem, we have

$$\text{tr}(\mathbf{Z}_{(2)}^T \mathbf{Z}_{(2)}) = n(p-r) + O_P(\sqrt{np}).$$

It follows that

$$\text{tr}(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) = \text{tr}(\mathbf{Z}_{(1)}^T (\boldsymbol{\Lambda} + \sigma^2 \mathbf{I}_r) \mathbf{Z}_{(1)}) + n(p-r)\sigma^2 + O_P(\sqrt{np}). \quad (23)$$

In view of (22), we have

$$\sum_{i=1}^r \lambda_i(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) = \text{tr}(\mathbf{Z}_{(1)}^T (\boldsymbol{\Lambda} + \sigma^2 \mathbf{I}_r) \mathbf{Z}_{(1)}) + r(p-r)\sigma^2 + O_P\left(\max(\sqrt{np}, n)\right).$$

This, together with (23), yields

$$\frac{1}{n(p-r)} \left(\text{tr}(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) - \sum_{i=1}^r \lambda_i(\mathbf{Z}^T \mathbf{E} \mathbf{Z}) \right) = \left(1 - \frac{r}{n}\right) \sigma^2 + O_P\left(\max\left(\frac{1}{\sqrt{np}}, \frac{1}{p}\right)\right).$$

□

Lemma 6. Under Assumption 1, we have

$$\|\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \mathbf{V}\mathbf{V}^T\|^2 = O_P\left(\frac{p}{p^\beta n}\right).$$

Remark 7. Cai et al. (2013), Theorem 5 asserts that under certain conditions,

$$\mathbb{E} \|\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \mathbf{V}\mathbf{V}^T\|_F^2 = O\left(\frac{p}{p^\beta n}\right),$$

where $\|\cdot\|_F$ is the Frobenius norm. Moreover, they proved that the convergence rate $p/(p^\beta n)$ is in fact minimax optimal. However, Cai et al. (2013), Theorem 5 needs the condition $\log p = O(n)$ which is unwanted. This condition is used to control the expectation and hence can be dropped in Lemma 6.

Proof. Since $\|\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \mathbf{V}\mathbf{V}^T\|^2 \leq 1$, the conclusion is trivial when $p/(p^\beta n)$ is unbounded. So without loss of generality, we assume $p/(p^\beta n) = O(1)$. Define \mathbf{U} , \mathbf{E} , \mathbf{Z} , $\mathbf{Z}_{(1)}$ and $\mathbf{Z}_{(2)}$ as in the proof of Lemma 5. Without loss of generality, we assume that $\mathbf{S} = n^{-1} \mathbf{U} \mathbf{E}^{1/2} \mathbf{Z} \mathbf{Z}^T \mathbf{E}^{1/2} \mathbf{U}^T$. Similar to the proof of Cai et al. (2013), Theorem 5, we define

$$\mathbf{S}_0 = \frac{1}{n} \mathbf{V} (\boldsymbol{\Lambda} + \sigma^2 \mathbf{I}_r)^{1/2} \mathbf{Z}_{(1)} \mathbf{Z}_{(1)}^T (\boldsymbol{\Lambda} + \sigma^2 \mathbf{I}_r)^{1/2} \mathbf{V}^T + \sigma^2 \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T.$$

It can be seen that the set of eigenvalues of \mathbf{S}_0 is the union of the nonzero eigenvalues of the matrix $n^{-1}\mathbf{V}(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{V}^T$ and $\sigma^2\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T$. All the nonzero eigenvalues of $\sigma^2\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T$ are σ^2 . We note that with probability 1, the first matrix is of rank r . Define the event

$$A = \left\{ \frac{1}{n}\lambda_r(\mathbf{V}(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{V}^T) > \sigma^2 \right\}.$$

On the event A , the eigenspace of \mathbf{S}_0 associated with the r leading eigenvalues is exactly $\mathbf{V}\mathbf{V}^T$, although the individual columns of \mathbf{V} need not be the leading eigenvectors of \mathbf{S}_0 . Applying Lemma 3 to \mathbf{S} and \mathbf{S}_0 yields

$$\|\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \mathbf{V}\mathbf{V}^T\| \mathbf{1}_{\{A\}} \leq \frac{2}{\lambda_r(\mathbf{S}_0) - \lambda_{r+1}(\mathbf{S}_0)} \|\mathbf{S} - \mathbf{S}_0\| \mathbf{1}_{\{A\}}. \quad (24)$$

Note that

$$\lambda_r\left(\frac{1}{n}\mathbf{V}(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{V}^T\right) = \frac{1}{n}\lambda_r((\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}).$$

By Weyl's inequality,

$$\begin{aligned} & \left| \frac{1}{n}\lambda_r((\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}) - (\lambda_r + \sigma^2) \right| \\ &= \left| \frac{1}{n}\lambda_r((\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}) - \lambda_r(\Lambda + \sigma^2\mathbf{I}_r) \right| \\ &\leq \left\| (\Lambda + \sigma^2\mathbf{I}_r)^{1/2} \left(\frac{1}{n}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T - \mathbf{I}_r \right) (\Lambda + \sigma^2\mathbf{I}_r)^{1/2} \right\| \\ &\leq (\lambda_1 + \sigma^2) \left\| \frac{1}{n}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T - \mathbf{I}_r \right\| \\ &= O_P\left(\frac{\lambda_1 + \sigma^2}{\sqrt{n}}\right), \end{aligned}$$

where the last equality follows from Corollary 1. Hence

$$\frac{1}{n}\lambda_r(\mathbf{V}(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{V}^T) = (\lambda_1 + \sigma^2)\left(1 + O_P\left(\frac{1}{\sqrt{n}}\right)\right).$$

It follows that $P(A) \rightarrow 1$. In view of (24), it is sufficient to show that

$$\frac{2}{\lambda_r(\mathbf{S}_0) - \lambda_{r+1}(\mathbf{S}_0)} \|\mathbf{S} - \mathbf{S}_0\| \mathbf{1}_{\{A\}} = O_P\left(\sqrt{\frac{p}{p^\beta n}}\right).$$

Note that on event A , $\lambda_{r+1}(\mathbf{S}_0) = \sigma^2$ and

$$\lambda_r(\mathbf{S}_0) = \frac{1}{n}\lambda_r(\mathbf{V}(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{V}^T).$$

Hence

$$\begin{aligned} & \frac{2}{\lambda_r(\mathbf{S}_0) - \lambda_{r+1}(\mathbf{S}_0)} \mathbf{1}_{\{A\}} = \frac{2}{\frac{1}{n}\lambda_r(\mathbf{V}(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{V}^T) - \sigma^2} \mathbf{1}_{\{A\}} \\ &= \frac{2}{\lambda_r\left(1 + O_P\left(\frac{1}{\sqrt{n}}\right)\right)} \mathbf{1}_{\{A\}} = \frac{2}{\lambda_r} (1 + o_P(1)). \end{aligned}$$

Next we bound $\|\mathbf{S} - \mathbf{S}_0\|$. We have

$$\begin{aligned}\|\mathbf{S} - \mathbf{S}_0\| &= \|(\mathbf{V}\mathbf{V}^T + \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T)(\mathbf{S} - \mathbf{S}_0)(\mathbf{V}\mathbf{V}^T + \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T)\| \\ &\leq \|\mathbf{V}\mathbf{V}^T(\mathbf{S} - \mathbf{S}_0)\mathbf{V}\mathbf{V}^T\| + 2\|\mathbf{V}\mathbf{V}^T(\mathbf{S} - \mathbf{S}_0)\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T\| + \|\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T(\mathbf{S} - \mathbf{S}_0)\tilde{\mathbf{V}}\tilde{\mathbf{V}}^T\| \\ &\leq \|\mathbf{V}^T(\mathbf{S} - \mathbf{S}_0)\mathbf{V}\| + 2\|\mathbf{V}^T(\mathbf{S} - \mathbf{S}_0)\tilde{\mathbf{V}}\| + \|\tilde{\mathbf{V}}^T(\mathbf{S} - \mathbf{S}_0)\tilde{\mathbf{V}}\| \\ &= 2\left\|\frac{\sigma}{n}(\Lambda + \sigma^2\mathbf{I}_r)^{1/2}\mathbf{Z}_{(1)}\mathbf{Z}_{(2)}^T\right\| + \sigma^2\left\|\frac{1}{n}\mathbf{Z}_{(2)}\mathbf{Z}_{(2)}^T - \mathbf{I}_{p-r}\right\| \\ &\leq \frac{2\sqrt{(\lambda_1 + \sigma^2)\sigma^2}}{n}\|\mathbf{Z}_{(1)}\mathbf{Z}_{(2)}^T\| + \sigma^2\left\|\frac{1}{n}\mathbf{Z}_{(2)}\mathbf{Z}_{(2)}^T - \mathbf{I}_{p-r}\right\|.\end{aligned}$$

By Corollary 1, we have $\|n^{-1}\mathbf{Z}_{(2)}\mathbf{Z}_{(2)}^T - \mathbf{I}_{p-r}\| = O_p(\max(\sqrt{p/n}, p/n))$. By the independence of $\mathbf{Z}_{(1)}$ and $\mathbf{Z}_{(2)}$, we have

$$\begin{aligned}\mathbb{E}\|\mathbf{Z}_{(1)}\mathbf{Z}_{(2)}^T\|^2 &\leq \mathbb{E}\|\mathbf{Z}_{(1)}\mathbf{Z}_{(2)}^T\|_F^2 = \mathbb{E}[\text{tr}(\mathbf{Z}_{(1)}\mathbf{Z}_{(2)}^T\mathbf{Z}_{(2)}\mathbf{Z}_{(1)}^T)] \\ &= \mathbb{E}\mathbb{E}[\text{tr}(\mathbf{Z}_{(1)}\mathbf{Z}_{(2)}^T\mathbf{Z}_{(2)}\mathbf{Z}_{(1)}^T)|\mathbf{Z}_{(1)}] = (p-r)\mathbb{E}[\text{tr}(\mathbf{Z}_{(1)}\mathbf{Z}_{(1)}^T)] = rn(p-r).\end{aligned}$$

Hence $\|\mathbf{Z}_{(1)}\mathbf{Z}_{(2)}^T\| = O_P(\sqrt{np})$. Combining these bounds leads to

$$\|\mathbf{S} - \mathbf{S}_0\| = O_P\left(\sqrt{\frac{\lambda_1 p}{n}}\right) + O_P\left(\max\left(\sqrt{\frac{p}{n}}, \frac{p}{n}\right)\right) = O_P\left(\sqrt{\frac{p^\beta p}{n}}\right) + O_P\left(\frac{p}{n}\right).$$

Thus,

$$\frac{2}{\lambda_r(\mathbf{S}_0) - \lambda_{r+1}(\mathbf{S}_0)}\|\mathbf{S} - \mathbf{S}_0\| = O_P\left(\sqrt{\frac{p}{p^\beta n}}\right) + O_P\left(\frac{p}{p^\beta n}\right) = O_P\left(\sqrt{\frac{p}{p^\beta n}}\right),$$

where the last equality holds since we have assumed $p/(p^\beta n) = O(1)$. This completes the proof. \square

Proof of Proposition 1. Note that

$$\frac{n_k - 1}{\sigma^2} \text{tr}(\tilde{\mathbf{V}}^T \mathbf{S}_k \tilde{\mathbf{V}}) \sim \chi^2((p-r)(n_k - 1)), \quad k = 1, 2.$$

It follows from the central limit theorem that

$$\frac{1}{\sigma^2} \sqrt{\frac{n_k - 1}{2(p-r)}} \left(\text{tr}(\tilde{\mathbf{V}}^T \mathbf{S}_k \tilde{\mathbf{V}}) - (p-r)\sigma^2 \right) \xrightarrow{\mathcal{L}} N(0, 1), \quad k = 1, 2.$$

Then we have

$$\frac{1}{n_k} \text{tr}(\tilde{\mathbf{V}}^T \mathbf{S}_k \tilde{\mathbf{V}}) = \frac{p-r}{n_k} \sigma^2 + O_P\left(\frac{\sqrt{p}}{n\sqrt{n}}\right), \quad k = 1, 2.$$

Thus,

$$\begin{aligned}T_1 - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 &= \|\tilde{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 - \tau(p-r)\sigma^2 + O_P\left(\frac{\sqrt{p}}{n\sqrt{n}}\right) \\ &= T_1^{(1)} + T_1^{(2)} + O_P\left(\frac{\sqrt{p}}{n\sqrt{n}}\right),\end{aligned}$$

where

$$\begin{aligned}T_1^{(1)} &= \|\tilde{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 - \tau(p-r)\sigma^2, \\ T_1^{(2)} &= 2(\mu_1 - \mu_2)^T \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T ((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2)).\end{aligned}$$

It follows from the fact $\|\tilde{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 \sim \sigma^2 \tau \chi^2(p-r)$ that

$$\frac{T_1^{(1)}}{\sqrt{2\tau^2(p-r)\sigma^4}} \sim \frac{\chi^2(p-r) - (p-r)}{\sqrt{2(p-r)}}.$$

For $T_1^{(2)}$, we have

$$\frac{T_1^{(2)}}{\sqrt{4\tau\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2\sigma^2}} \sim N(0, 1).$$

Thus, if $\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 = o(p/n)$, then $T_1^{(1)}$ dominates $T_1^{(2)}$ and (9) follows from Slutsky's theorem. On the other hand, if $p/n = o(\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2)$, then $T_1^{(2)}$ dominates $T_1^{(1)}$ and (10) follows from Slutsky's theorem. \square

Proof of Proposition 2. The conclusion (12) is a direct corollary of (20) in Lemma 5. By (21) and (12), for $i = 1, \dots, r$, we have

$$\begin{aligned} \hat{\boldsymbol{\lambda}}_i &= \lambda_i(\mathbf{S}) - \frac{p+n-r}{n} \hat{\sigma}_*^2 \\ &= \boldsymbol{\lambda}_i + \frac{p+n-r}{n} \sigma^2 + O_P\left(\max\left(\frac{p^\beta}{\sqrt{n}}, 1\right)\right) - \frac{p+n-r}{n} \sigma^2 - O_P\left(\frac{n+p}{n} \max\left(\frac{1}{\sqrt{np}}, \frac{1}{p}\right)\right) \\ &= \boldsymbol{\lambda}_i + O_P\left(\max\left(\frac{p^\beta}{\sqrt{n}}, 1\right)\right), \end{aligned}$$

which proves (13). \square

Proof of Theorem 3. Proposition 2 implies that $\hat{\sigma}_*^2 \xrightarrow{P} \sigma^2$ and $\hat{\boldsymbol{\lambda}}_i/\boldsymbol{\lambda}_i \xrightarrow{P} 1$. Hence

$$\tau \sum_{i=1}^r \frac{p\hat{\sigma}_*^2}{n\hat{\boldsymbol{\lambda}}_i + (n+p)\hat{\sigma}_*^2} \hat{\boldsymbol{\lambda}}_i = \tau \frac{p}{n} \sum_{i=1}^r \frac{\hat{\sigma}_*^2}{1 + \frac{(n+p)}{n\hat{\boldsymbol{\lambda}}_i} \hat{\sigma}_*^2} = \tau \frac{p}{n} \sum_{i=1}^r \frac{\sigma^2(1+o_P(1))}{1 + \frac{(n+p)}{n\boldsymbol{\lambda}_i} \sigma^2(1+o_P(1))}.$$

Since $p/n^2 \rightarrow 0$ implies $p/(n\boldsymbol{\lambda}_i) \rightarrow 0$, we have

$$\tau \sum_{i=1}^r \frac{p\hat{\sigma}_*^2}{n\hat{\boldsymbol{\lambda}}_i + (n+p)\hat{\sigma}_*^2} \hat{\boldsymbol{\lambda}}_i = r\tau \frac{p}{n} \sigma^2(1+o_P(1)) = o_P(\sqrt{\tau^2 p}).$$

In view of (12), we have

$$\begin{aligned} \tau(p-r)\hat{\sigma}_*^2 &= \tau(p-r)\sigma^2 + O_P\left(\frac{p}{n} \max\left(\frac{1}{\sqrt{np}}, \frac{1}{p}\right)\right) \\ &= \tau(p-r)\sigma^2 + O_P\left(\frac{\sqrt{p}}{n} \max\left(\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{p}}\right)\right) = \tau(p-r)\sigma^2 + o_P(\sqrt{\tau^2 p}). \end{aligned}$$

Thus,

$$\begin{aligned} T_2 - \|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 &= \|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 - \|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 - \tau(p-r)\sigma^2 + o_P(\sqrt{\tau^2 p}) \\ &= P_1 + P_2 + o_P(\sqrt{\tau^2 p}), \end{aligned}$$

where

$$\begin{aligned} P_1 &= \|\hat{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 - \tau(p-r)\sigma^2, \\ P_2 &= 2(\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T ((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2)). \end{aligned}$$

It can be seen that

$$\frac{P_2}{\sqrt{4\tau(\mu_1 - \mu_2)^T \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T (\mu_1 - \mu_2)}} \sim N(0, 1).$$

Now we derive the asymptotic distribution of P_1 . To make clear the mode of convergence, we need a metric for weak convergence. For two distribution function F and G , the Levy metric ρ of F and G is defined as

$$\rho(F, G) = \inf\{\epsilon : F(x - \epsilon) - \epsilon \leq G(x) \leq F(x + \epsilon) + \epsilon \text{ for all } x\}.$$

It's well known that $\rho(F_n, F) \rightarrow 0$ if and only if $F_n \xrightarrow{\mathcal{L}} F$.

Note that \bar{X}_1 , \bar{X}_2 , and \mathbf{S} are mutually independent and $\hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T$ only depends on \mathbf{S} . Then the conditional distribution of $\hat{\tilde{\mathbf{V}}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))$ given \mathbf{S} is $N(0, \tau \hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})$. Thus,

$$\tau^{-1} \|\hat{\tilde{\mathbf{V}}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 \stackrel{p-r}{=} \sum_{i=1}^{p-r} \lambda_i(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}}) \xi_i^2, \quad (25)$$

where ξ_1, \dots, ξ_{p-r} are i.i.d. standard normal random variables which are independent of $\hat{\tilde{\mathbf{V}}}$. In view of Lemma 1, the asymptotic distribution of P_1 relies on the asymptotic behavior of $\lambda_i(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})$, $i = 1, \dots, p-r$. Note that

$$\lambda_1(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}}) = \lambda_1(\hat{\tilde{\mathbf{V}}}^T(\mathbf{V} \Lambda \mathbf{V}^T + \sigma^2 \mathbf{I}_p) \hat{\tilde{\mathbf{V}}}) \leq \kappa p^\beta \lambda_1(\hat{\tilde{\mathbf{V}}}^T \mathbf{V} \mathbf{V}^T \hat{\tilde{\mathbf{V}}}) + \sigma^2 = \kappa p^\beta \|\mathbf{V} \mathbf{V}^T - \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T\|^2 + \sigma^2. \quad (26)$$

On the other hand, for $i = r+1, \dots, p-r$, we have

$$\lambda_i(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}}) = \lambda_i(\hat{\tilde{\mathbf{V}}}^T \mathbf{V} \Lambda \mathbf{V}^T \hat{\tilde{\mathbf{V}}}) + \sigma^2 = \sigma^2, \quad (27)$$

where the last equality holds since $\text{Rank}(\hat{\tilde{\mathbf{V}}}^T \mathbf{V} \Lambda \mathbf{V}^T \hat{\tilde{\mathbf{V}}}) \leq \text{Rank}(\mathbf{V}) = r$. It follows from (26) and (27) that

$$\begin{aligned} \text{tr}(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})^2 &= O_P\left(\left(p^\beta \|\mathbf{V} \mathbf{V}^T - \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T\|^2 + 1\right)^2\right) + (p-2r)\sigma^4 \\ &= O_P\left(\left(\frac{p}{n} + 1\right)^2\right) + (p-2r)\sigma^4 = p\sigma^4(1 + o_P(1)). \end{aligned} \quad (28)$$

This, combined with (26), yields

$$\frac{\lambda_1^2(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})}{\text{tr}(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})^2} = O_P\left(p^{-1}(p^\beta \|\mathbf{V} \mathbf{V}^T - \hat{\tilde{\mathbf{V}}} \hat{\tilde{\mathbf{V}}}^T\|^2 + 1)^2\right) = O_P\left(\frac{(p/n + 1)^2}{p}\right) = o_P(1).$$

Then for every subsequence $\{n(k)\}$ of $\{n\}$, there's a further subsequence $\{n(k(l))\}$ along which

$$\frac{\lambda_1^2(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})}{\text{tr}(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})^2} \xrightarrow{a.s.} 0.$$

This, together with (25) and Lemma 1, implies that along $\{n(k(l))\}$ we have

$$\rho(\mathcal{L}(Y_n | \mathbf{S}), N(0, 1)) \xrightarrow{a.s.} 0, \quad (29)$$

where

$$Y_n = \frac{\|\hat{\tilde{\mathbf{V}}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 - \tau \text{tr}(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})}{\sqrt{2\tau^2 \text{tr}(\hat{\tilde{\mathbf{V}}}^T \Sigma \hat{\tilde{\mathbf{V}}})^2}},$$

and $\mathcal{L}(Y_n | \mathbf{S})$ is the conditional distribution of Y_n given \mathbf{S} . By the definition of weak convergence, (29) implies that for every continuous bounded function $f(\cdot)$, $E[f(Y_n) | \mathbf{S}] \xrightarrow{a.s.} E[f(\xi^*)]$ along $\{n(k(l))\}$, where ξ^* is a standard normal random variable. By the dominated convergence theorem, $E[f(Y_n)] \rightarrow E[f(\xi^*)]$ along $\{n(k(l))\}$. This implies that $Y_n \xrightarrow{\mathcal{L}} N(0, 1)$ along $\{n(k(l))\}$. Thus, for every subsequence of $\{n\}$, there is a further subsequence along which $Y_n \xrightarrow{\mathcal{L}} N(0, 1)$. This means $Y_n \xrightarrow{\mathcal{L}} N(0, 1)$, or

$$\frac{\|\hat{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 - \tau \text{tr}(\hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}})}{\sqrt{2\tau^2 \text{tr}(\hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}})^2}} \xrightarrow{\mathcal{L}} N(0, 1).$$

By (26) and (27), we have

$$\begin{aligned} \text{tr}(\hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}) &= \sum_{i=1}^r \lambda_i(\hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}) + \sum_{i=r+1}^{p-r} \lambda_i(\hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}) \\ &= O_P\left(\frac{p}{n} + 1\right) + (p - 2r)\sigma^2 = (p - r)\sigma^2 + o_P(\sqrt{p}). \end{aligned} \quad (30)$$

It follows from (28), (30) and Slutsky's theorem that

$$\frac{P_1}{\sqrt{2\tau^2 p \sigma^4}} = \frac{\|\hat{\mathbf{V}}^T((\bar{X}_1 - \mu_1) - (\bar{X}_2 - \mu_2))\|^2 - \sigma^2 \tau(p - r)}{\sqrt{2\tau^2 p \sigma^4}} \xrightarrow{\mathcal{L}} N(0, 1).$$

If $(\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}^T (\mu_1 - \mu_2) = o_P(p/n)$,

$$\frac{T_2 - \|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} = \frac{P_1}{\sqrt{2\tau^2 p \sigma^4}} + o_P(1) \xrightarrow{\mathcal{L}} N(0, 1).$$

On the other hand, if $p/n = o_P((\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}^T (\mu_1 - \mu_2))$,

$$\frac{T_2 - \|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{4\tau(\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}^T (\mu_1 - \mu_2)}} = \frac{P_2}{\sqrt{4\tau(\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}^T (\mu_1 - \mu_2)}} + o_P(1) \xrightarrow{\mathcal{L}} N(0, 1).$$

Hence (a) and (b) hold.

Now we prove (c). Note that

$$\begin{aligned} (\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}^T (\mu_1 - \mu_2) &\leq \lambda_1(\hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}})(\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T (\mu_1 - \mu_2) \\ &\leq (\kappa p^\beta \lambda_1(\hat{\mathbf{V}}^T \mathbf{V} \mathbf{V}^T \hat{\mathbf{V}}) + \sigma^2) \|\mu_1 - \mu_2\|^2. \end{aligned}$$

But

$$\lambda_1(\hat{\mathbf{V}}^T \mathbf{V} \mathbf{V}^T \hat{\mathbf{V}}) = \|\mathbf{V}^T \hat{\mathbf{V}}\|^2 = \|\mathbf{V} \mathbf{V}^T - \hat{\mathbf{V}} \hat{\mathbf{V}}^T\|^2 = O_P\left(\frac{p}{p^\beta n}\right),$$

where the second equality follows from Golub and Van Loan (2013), Theorem 2.5.1, and the last equality follows from Lemma 6. Hence if $\|\mu_1 - \mu_2\|^2 = O(\sqrt{p}/n)$, we have

$$(\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \hat{\Sigma} \hat{\mathbf{V}}^T (\mu_1 - \mu_2) = O_P\left(\frac{p}{n} + 1\right) \|\mu_1 - \mu_2\|^2 = o_P\left(\frac{p}{n}\right).$$

Then the condition of (a) satisfies and we have

$$\frac{T_2 - \|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \xrightarrow{\mathcal{L}} N(0, 1).$$

It remains to show $|\|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2| = o(\sqrt{p}/n)$. We have

$$\begin{aligned} & |\|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2 - \|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2| = |(\mu_1 - \mu_2)^T (\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T)(\mu_1 - \mu_2)| \\ & \leq \|\mu_1 - \mu_2\|^2 \|\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T\| = \|\mu_1 - \mu_2\|^2 \|\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \mathbf{V}\mathbf{V}^T\| = O\left(\frac{\sqrt{p}}{n}\right) o_P(1) = o_P\left(\frac{\sqrt{p}}{n}\right). \end{aligned}$$

This completes the proof. \square

Proof of Theorem 4. Define the event

$$A = \left\{ (\mu_1 - \mu_2)^T (\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \Sigma \hat{\mathbf{V}}\hat{\mathbf{V}}^T)(\mu_1 - \mu_2) \leq \frac{p}{n} \sqrt{\frac{\sqrt{p}}{n} + \frac{1}{\sqrt{p}}} \right\}.$$

It can be seen that the condition of (a) of Theorem 3 holds on event A . To deal with the case where the condition of (a) of Theorem 3 is not satisfied, we would like to shrink μ_1 and μ_2 on A^c . Define the shrinkage factor h as

$$h = \begin{cases} 1 & \text{on } A, \\ \left(\frac{\frac{p}{n} \sqrt{\frac{\sqrt{p}}{n} + \frac{1}{\sqrt{p}}}}{(\mu_1 - \mu_2)^T (\hat{\mathbf{V}}\hat{\mathbf{V}}^T - \Sigma \hat{\mathbf{V}}\hat{\mathbf{V}}^T)(\mu_1 - \mu_2)} \right)^{1/2} & \text{on } A^c. \end{cases}$$

Then $h \in (0, 1]$ is a random variable which only depends on \mathbf{S} .

Let $\bar{X}_k^* = \bar{X}_k + (h-1)\mu_k$, $k = 1, 2$. Then $\bar{X}_k^* | \mathbf{S} \sim N(h\mu_k, n_k^{-1}\Sigma)$, $k = 1, 2$. Define

$$\begin{aligned} T_2^* &= \|\hat{\mathbf{V}}^T(\bar{X}_1^* - \bar{X}_2^*)\|^2 - \tau(p-r)\hat{\sigma}_*^2 - \tau \sum_{i=1}^r \frac{p\hat{\sigma}_*^2}{n\hat{\lambda}_i + (n+p)\hat{\sigma}_*^2} \hat{\lambda}_i, \\ Q^* &= T_2^*/\left(2\tau^2 \left(\sum_{i=1}^r \left(\frac{p\hat{\sigma}_*^2}{n\hat{\lambda}_i + (n+p)\hat{\sigma}_*^2} \hat{\lambda}_i \right)^2 + 2\hat{\sigma}_*^2 \sum_{i=1}^r \frac{p\hat{\sigma}_*^2}{n\hat{\lambda}_i + (n+p)\hat{\sigma}_*^2} \hat{\lambda}_i + \hat{\sigma}_*^4(p-r) \right)\right)^{1/2}. \end{aligned}$$

By Proposition 2,

$$\frac{p\hat{\sigma}_*^2}{n\hat{\lambda}_i + (n+p)\hat{\sigma}_*^2} \hat{\lambda}_i \leq \frac{p}{n} \hat{\sigma}_*^2 = o_P(p).$$

Then we have

$$Q^* = \frac{T_2^*}{\sqrt{2\tau^2 p \sigma^4}} (1 + o_P(1)), \quad (31)$$

where the term $o_P(1)$ only depends on \mathbf{S} . Note that for T_2^* , the condition of (a) of Theorem 3 holds. Then similar to the proof of (a) of Theorem 3, we can show that

$$\rho \left(\mathcal{L} \left(\frac{T_2^* - \|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \middle| \mathbf{S} \right), N(0, 1) \right) \xrightarrow{P} 0. \quad (32)$$

Hence for every subsequence $\{n(k)\}$ of $\{n\}$, there's a further subsequence $\{n(k(l))\}$ along which (31) and (32) hold almost surely. Thus, along $\{n(k(l))\}$ we have almost surely that

$$\begin{aligned} & \Pr \left(Q^* > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \middle| \mathbf{S} \right) = \Pr \left(\frac{T_2^*}{\sqrt{2\tau^2 p \sigma^4}} > \Phi^{-1}(1-\alpha) + o(1) \middle| \mathbf{S} \right) \\ & = \Pr \left(\frac{T_2^* - \|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} > \Phi^{-1}(1-\alpha) - \frac{\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} + o(1) \middle| \mathbf{S} \right) \\ & = \Phi \left(-\Phi^{-1}(1-\alpha) + \frac{\|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \right) + o(1). \end{aligned}$$

That is to say,

$$\Pr \left(Q^* > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \middle| \mathbf{S} \right) = \Phi \left(-\Phi^{-1}(1-\alpha) + \frac{\|\tilde{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \right) + o_P(1). \quad (33)$$

On the event A , we have $Q = Q^*$. On the event A^c , we have

$$\frac{p}{n} \sqrt{\frac{\sqrt{p}}{n} + \frac{1}{\sqrt{p}}} < (\mu_1 - \mu_2)^T \hat{\mathbf{V}} \hat{\mathbf{V}}^T \Sigma \hat{\mathbf{V}} \hat{\mathbf{V}}^T (\mu_1 - \mu_2) \leq (O_P(\frac{p}{n}) + \sigma^2) \|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2.$$

This implies that

$$\frac{\|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{p}/n} \geq \frac{\sqrt{p} \sqrt{\frac{\sqrt{p}}{n} + \frac{1}{\sqrt{p}}}}{O_P(\frac{p}{n}) + \sigma^2} = \frac{\sqrt{\frac{\sqrt{p}}{n} + \frac{1}{\sqrt{p}}}}{O_P(\frac{\sqrt{p}}{n} + \frac{1}{\sqrt{p}})} = \frac{1}{o_P(1)} \rightarrow \infty.$$

So (33) implies that

$$\Pr \left(Q^* > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \middle| \mathbf{S} \right) \mathbf{1}_{\{A^c\}} = \mathbf{1}_{\{A^c\}} + o_P(1).$$

By induction on p , one can show that the conditional distribution $\mathcal{L}(\|\hat{\mathbf{V}}^T(\bar{X}_1 - \bar{X}_2)\|^2 | \mathbf{S})$ is stochastically larger than the conditional distribution $\mathcal{L}(\|\hat{\mathbf{V}}^T(\bar{X}_1^* - \bar{X}_2^*)\|^2 | \mathbf{S})$. Hence we have

$$\begin{aligned} \Pr \left(Q > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \middle| \mathbf{S} \right) \mathbf{1}_{\{A^c\}} &\geq \Pr \left(Q^* > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \middle| \mathbf{S} \right) \mathbf{1}_{\{A^c\}} \\ &= \mathbf{1}_{\{A^c\}} + o_P(1). \end{aligned}$$

Thus,

$$\begin{aligned} \Pr \left(Q > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \middle| \mathbf{S} \right) &= \Pr \left(Q^* > \frac{\chi_{1-\alpha}^2(p-r) - (p-r)}{\sqrt{2(p-r)}} \middle| \mathbf{S} \right) + o_P(1) \\ &= \Phi \left(-\Phi^{-1}(1-\alpha) + \frac{\|\hat{\mathbf{V}}^T(\mu_1 - \mu_2)\|^2}{\sqrt{2\tau^2 p \sigma^4}} \right) + o_P(1). \end{aligned}$$

Note that the term $o_P(1)$ satisfies $|o_P(1)| \leq 2$. Hence (14) follows from the dominated convergence theorem.

Finally, (15) follows directly from the fact

$$Q = \frac{T_2}{\sqrt{2\tau^2 p \sigma^4}} (1 + o_P(1))$$

and (c) of Theorem 3. This completes the proof. □

References

- Ahn, S.C., Horenstein, A.R., 2013. Eigenvalue ratio test for the number of factors. *Econometrica* 81, 1203–1227. doi:10.3982/ecta8968.
- Anderson, T.W., 1963. Asymptotic theory for principal component analysis. *The Annals of Mathematical Statistics* 34, 122–148.
- Bai, Z., Saranadasa, H., 1996. Effect of high dimension: by an example of a two sample problem. *Statistica Sinica* 6, 311–329.
- Birnbaum, A., Johnstone, I.M., Nadler, B., Paul, D., 2013. Minimax bounds for sparse PCA with noisy high-dimensional data. *The Annals of Statistics* 41, 1055–1084. doi:10.1214/12-aos1014.
- Cai, T., Ma, Z., Wu, Y., 2015. Optimal estimation and rank detection for sparse spiked covariance matrices. *Probability Theory & Related Fields* 161, 781–815.

- Cai, T.T., Ma, Z., Wu, Y., 2013. Sparse PCA: Optimal rates and adaptive estimation. *The Annals of Statistics* 41, 3074–3110. doi:10.1214/13-aos1178.
- Chen, L.S., Paul, D., Prentice, R.L., Wang, P., 2011. A regularized hotelling'sT2test for pathway analysis in proteomic studies. *Journal of the American Statistical Association* 106, 1345–1360. doi:10.1198/jasa.2011.ap10599.
- Chen, S.X., Qin, Y.L., 2010. A two-sample test for high-dimensional data with applications to gene-set testing. *The Annals of Statistics* 38, 808–835. doi:10.1214/09-aos716.
- Davidson, K.R., Szarek, S.J., 2001. *Handbook of the Geometry of Banach Spaces*. volume 1. North-Holland, Amsterdam. doi:10.1016/S1874-5849(01)80010-3.
- Golub, G.H., Van Loan, C.F., 2013. *Matrix Computations*. Fourth ed., The Johns Hopkins University Press.
- Horn, R.A., Johnson, C.R., 2012. *Matrix Analysis*. 2nd ed., Cambridge University Press, New York.
- Jung, S., Marron, J.S., 2009. PCA consistency in high dimension, low sample size context. *The Annals of Statistics* 37, 4104–4130. doi:10.1214/09-aos709.
- Lopes, M., Jacob, L., Wainwright, M.J., 2011. A more powerful two-sample test in high dimensions using random projection, in: *Advances in Neural Information Processing Systems* 24. Curran Associates, Inc., pp. 1206–1214.
- Ma, Y., Lan, W., Wang, H., 2015. A high dimensional two-sample test under a low dimensional factor structure. *Journal of Multivariate Analysis* 140, 162–170. doi:10.1016/j.jmva.2015.05.005.
- Passemier, D., Li, Z., Yao, J., 2017. On estimation of the noise variance in high dimensional probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79, 51–67. doi:10.1111/rssb.12153.
- Srivastava, M.S., Du, M., 2008. A test for the mean vector with fewer observations than the dimension. *Journal of Multivariate Analysis* 99, 386–402. doi:10.1016/j.jmva.2006.11.002.
- Srivastava, R., Li, P., Ruppert, D., 2016. RAPTT: An exact two-sample test in high dimensions using random projections. *Journal of Computational and Graphical Statistics* 25, 954–970. doi:10.1080/10618600.2015.1062771.
- Thulin, M., 2014. A high-dimensional two-sample test for the mean using random subspaces. *Computational Statistics & Data Analysis* 74, 26–38. doi:10.1016/j.csda.2013.12.003.
- Wang, W., Fan, J., 2017. Asymptotics of empirical eigenstructure for high dimensional spiked covariance. *The Annals of Statistics* 45, 1342–1374. doi:10.1214/16-aos1487.
- Zhao, J., Xu, X., 2016. A generalized likelihood ratio test for normal mean when p is greater than n . *Computational Statistics & Data Analysis* 99, 91–104. doi:10.1016/j.csda.2016.01.006.