# SUPPLEMENT TO "TESTING HIGH DIMENSIONAL COVARIANCE MATRICES, WITH APPLICATION TO DETECTING SCHIZOPHRENIA RISK GENES"

By Lingxue Zhu[*], Jing Lei[*], Bernie Devlin[†] and Kathryn Roeder[*]

*Carnegie Mellon University[*] and University of Pittsburgh[†]*

This document provides supplementary material to the article "Testing High Dimensional Covariance Matrices, with Application to Detecting Schizophrenia Risk Genes" written by the same authors.

**S1. Simulations.** In this section, we present the remaining simulation results for comparing `sLED` with other existing methods, including `Sfrob` (Schott, 2007), `Ustat` (Li and Chen, 2012), `Max` (Cai, Liu and Xia, 2013), `MBoot` (Chang et al., 2016), and `RProj` (Wu and Li, 2015).

The samples are generated by $X_i = \Sigma_1^{1/2} Z_i$ for $i = 1, \cdots, n$, and $Y_l = \Sigma_2^{1/2} Z_{n+l}$ for $l = 1, \cdots, m$, where $\{Z_i\}_{i=1,n+m}$ are independent $p$-dimensional random variables with *i.i.d.* coordinates $Z_{ij}$, $j = 1, \cdots, p$. We let $\Sigma_2 = \Sigma_1$ under $H_0$ and $\Sigma_2 = \Sigma_1 + D$ under $H_1$. For the different choices of $\Sigma_1$ and $D$, please refer to the main manuscript (Zhu et al., 2017). We consider the following four distributions for $Z_{ij}$:

1. Standard Normal $N(0, 1)$, which leads to multinomial Gaussian samples $X$ and $Y$.
2. Centralized Gamma distribution with $\alpha = 4, \beta = 0.5$ (i.e., the theoretical expectation $\alpha\beta = 2$ is subtracted from $\Gamma(4, 0.5)$ samples).
3. $t$-distribution with degrees of freedom 12 ($t(12)$).
4. Centralized Negative Binomial distribution with mean $\mu = 2$ and dispersion parameter $\phi = 2$ (i.e., the theoretical expectation $\mu = 2$ is subtracted from NB$(2, 2)$ samples).

We compare the empirical power among different testing procedures, where 100 permutations are used to compute the $p$-values for all methods, except for `MBoot` where 100 bootstrap repetitions are used. Table S1 summarizes the empirical power under different covariance structures and differential matrices when $Z_{ij}$'s are sampled from $t$-distribution and centralized NB$(2, 2)$. The results for standard Normal and centralized Gamma distributions are presented in the main manuscript. The smoothing parameter for sLED is set to be $\sqrt{R} = 0.3\sqrt{p}$, and 100 random projections are

used for `Rproj`. We also examine the sensitivity of sLED to the smoothing parameter in Figure S1, where $c$ is varied among $\{0.10, 0.12, \cdots, 0.30\}$ (recall that $\sqrt{R} = c\sqrt{p}$). We see that `sLED` achieves superior power to other approaches under most scenarios, and the results remain robust to many choices of $c$'s. Finally, Figure S2 shows the distribution of the empirical size in 100 repetitions under all these scenarios, including 4 different choices of $\Sigma_1$, 3 different choices of $p$, and 4 different choices of the distribution of $Z_{ij}$. We see that for all of the 6 testing procedures, the empirical size is comparable and controlled around $\alpha = 0.05$. In addition, the empirical size of sLED is robust to the choice of $c$. We point out that only a small number of permutations is used in this simulation. In practice, more permutations will be conducted and the size will be more properly controlled.

**S2. Proofs for consistency.** In this section, we prove Theorems 1 to 3 for the asymptotic power of sLED.

*Notation.* For a set $\mathcal{A}$, let $|\mathcal{A}|$ be its cardinality, and $\mathcal{A}^c$ be its complement. For $Z = (Z_1, \cdots Z_N) = (X_1, \cdots, X_n, Y_1, \cdots, Y_m)$, we denote $Z_{ki}$ to be the $i$-th coordinate of the $k$-th sample $Z_k$, and

$$\hat{\Sigma} = \frac{1}{N} \sum_{k=1}^{N} Z_k Z_k^T \,,\ \ \bar{Z} = \frac{1}{N} \sum_{k=1}^{N} Z_k = (\bar{Z}_1, \cdots, \bar{Z}_p)^T \,,$$

$$m_z = ||Z||_\infty \,,\ \ \overline{m}_z = ||\bar{Z}||_\infty \,,\ \ \overline{m}_z^{(2q)} = \max_{1 \le i,j \le p} \frac{1}{N} \sum_{k=1}^{N} Z_{ki}^q Z_{kj}^q \,,\ q = 1, 2 \,.$$

PROOF OF THEOREM 1. By Theorem 2 and Lemma 2, there exist some constants $C', C''$ depending on $(\underline{c}, \bar{c}, \nu^2, \delta)$, such that if $(n, p)$ are sufficiently large, with probability at least $1 - \delta$,

$$||\hat{D}^*||_\infty \le C'\sqrt{\frac{\log p}{n}} \,,\ \ ||\hat{D} - D||_\infty \le C''\sqrt{\frac{\log p}{n}} \,.$$

Then we apply Theorem 3 on both $\hat{D}, -\hat{D}$ and $\hat{D}^*, -\hat{D}^*$, and this together with assumption (A4) imply the desired conclusion with $C = C' + C''$. □

PROOF OF THEOREM 2. First, note that for $\forall \epsilon > 0$,

$$(S2.1) \quad \mathbb{P}\left( ||\hat{D}^*||_\infty > \epsilon \right) \le \mathbb{P}\left( ||\hat{\Sigma}_1^* - \hat{\Sigma}||_\infty > \frac{\epsilon}{2} \right) + \mathbb{P}\left( ||\hat{\Sigma}_2^* - \hat{\Sigma}||_\infty > \frac{\epsilon}{2} \right) \,.$$

Now for any $\delta > 0$ and constants $C_1, C_2$, define

$$\mathcal{A} = \left\{ Z : m_z \le C_2\sqrt{\log\left(\frac{C_1 np}{\delta}\right)} \,,\ \overline{m}_z \le C_2\sqrt{\frac{\log(C_1 p/\delta)}{n}} \,,\ \overline{m}_z^{(2q)} \le C_2 \,,\ q = 1, 2 \right\} \,.$$

TABLE S1

*Empirical power in 100 repetitions, where $n = m = 100$, nominal level $\alpha = 0.05$, and $Z_{ij}$'s are sampled from centralized Negative Binomial $(2, 2)$ (top) and t-distribution with degrees of freedom 12 (bottom). Under each scenario, the largest power is highlighted.*

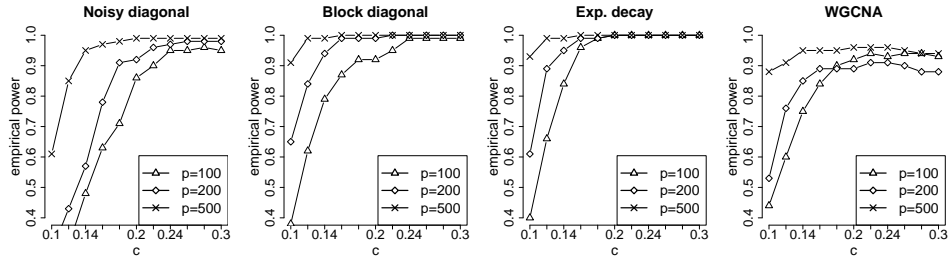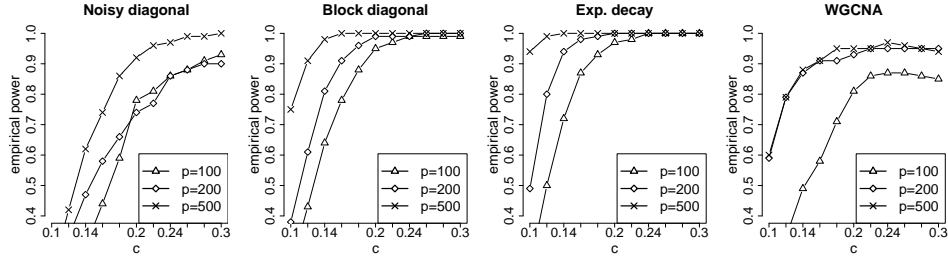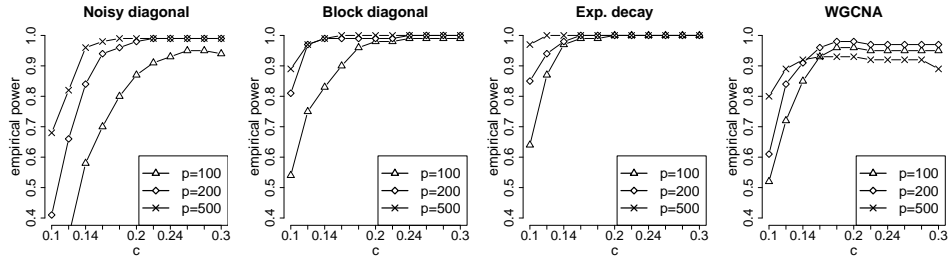| D | $\Sigma_1$ | Noisy diagonal | | | Block diagonal | | | Exp. decay | | | WGCNA | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | p | 100 | 200 | 500 | 100 | 200 | 500 | 100 | 200 | 500 | 100 | 200 | 500 |
| | | Centralized Negative Binomial | | | | | | | | | | | |
| **Block** | Max | 0.59 | 0.18 | 0.11 | 0.87 | 0.69 | 0.25 | **1.00** | 0.94 | 0.50 | 0.80 | 0.84 | 0.28 |
| | MBoot | 0.47 | 0.12 | 0.07 | 0.80 | 0.60 | 0.16 | 0.99 | 0.82 | 0.33 | 0.72 | 0.73 | 0.17 |
| | Ustat | 0.69 | 0.62 | 0.68 | 0.89 | 0.93 | 0.94 | 0.99 | 0.99 | **1.00** | 0.57 | 0.78 | 0.74 |
| | Sfrob | 0.66 | 0.63 | 0.69 | 0.91 | 0.91 | 0.94 | 0.98 | 0.99 | **1.00** | 0.58 | 0.79 | 0.80 |
| | RProj | 0.11 | 0.13 | 0.06 | 0.18 | 0.13 | 0.11 | 0.21 | 0.21 | 0.10 | 0.15 | 0.10 | 0.16 |
| | sLED | **0.91** | **0.90** | **0.99** | **0.99** | **1.00** | **1.00** | **1.00** | **1.00** | **1.00** | **0.86** | **0.95** | **0.95** |
| **Spiked** | Max | 0.04 | 0.07 | 0.06 | 0.60 | 0.20 | 0.10 | 0.93 | 0.82 | 0.25 | 0.93 | 0.41 | 0.06 |
| | MBoot | 0.03 | 0.04 | 0.03 | 0.51 | 0.20 | 0.02 | 0.92 | 0.76 | 0.18 | 0.89 | 0.33 | 0.05 |
| | Ustat | **0.25** | **0.12** | 0.03 | 0.82 | 0.35 | 0.13 | 0.98 | 0.94 | 0.72 | 0.36 | 0.12 | 0.03 |
| | Sfrob | **0.25** | 0.11 | 0.03 | 0.85 | 0.35 | 0.12 | 0.98 | 0.96 | 0.69 | 0.40 | 0.12 | 0.06 |
| | RProj | 0.08 | 0.07 | 0.02 | 0.31 | 0.17 | 0.10 | 0.35 | 0.17 | 0.06 | 0.58 | 0.17 | **0.13** |
| | sLED | 0.22 | 0.04 | **0.08** | **0.97** | **0.68** | **0.16** | **0.99** | **1.00** | **1.00** | **0.96** | **0.48** | **0.13** |
| | | T-distribution | | | | | | | | | | | |
| **Block** | Max | 0.22 | 0.15 | 0.12 | 0.88 | 0.73 | 0.23 | **1.00** | 0.85 | 0.27 | **0.97** | 0.40 | 0.15 |
| | MBoot | 0.23 | 0.14 | 0.13 | 0.88 | 0.68 | 0.20 | **1.00** | 0.84 | 0.23 | 0.91 | 0.37 | 0.13 |
| | Ustat | 0.53 | 0.63 | 0.74 | 0.97 | 0.93 | 0.96 | **1.00** | 0.99 | **1.00** | 0.75 | 0.67 | 0.84 |
| | Sfrob | 0.52 | 0.63 | 0.71 | 0.97 | 0.93 | 0.97 | **1.00** | 0.99 | 0.99 | 0.74 | 0.67 | 0.76 |
| | RProj | 0.08 | 0.13 | 0.11 | 0.28 | 0.12 | 0.02 | 0.28 | 0.19 | 0.08 | 0.17 | 0.13 | 0.06 |
| | sLED | **0.95** | **0.99** | **0.99** | **0.99** | **1.00** | **1.00** | **1.00** | **1.00** | **1.00** | 0.95 | **0.97** | **0.92** |
| **Spiked** | Max | 0.13 | 0.08 | 0.05 | 0.78 | 0.50 | 0.08 | 0.96 | 0.79 | 0.18 | 0.85 | 0.27 | 0.10 |
| | MBoot | 0.12 | 0.07 | 0.03 | 0.78 | 0.49 | 0.07 | 0.97 | 0.77 | 0.11 | 0.83 | 0.32 | 0.13 |
| | Ustat | 0.14 | 0.07 | **0.10** | 0.79 | 0.36 | 0.07 | **1.00** | 0.91 | 0.68 | 0.30 | 0.06 | 0.05 |
| | Sfrob | 0.12 | **0.10** | **0.10** | 0.80 | 0.36 | 0.07 | **1.00** | 0.91 | 0.66 | 0.29 | 0.09 | 0.05 |
| | RProj | 0.07 | 0.06 | 0.06 | 0.34 | 0.20 | 0.08 | 0.36 | 0.16 | 0.07 | 0.57 | 0.27 | **0.14** |
| | sLED | **0.40** | 0.09 | 0.03 | **0.96** | **0.76** | **0.14** | **1.00** | **1.00** | **1.00** | **0.95** | **0.55** | 0.08 |

(a) $Z_{ij} \sim$ centralized $\Gamma(4, 0.5)$.



(b) $Z_{ij} \sim$ centralized $\mathrm{NB}(2, 2)$.



(c) $Z_{ij} \sim t(12)$.

FIG. S1. *Empirical power of sLED in 100 repetitions using different smoothing parameters* $\sqrt{R} = c\sqrt{p}$ *for* $c \in \{0.10, 0.12, \cdots, 0.30\}$, *where* $D$ *has sparse block difference and* $Z_{ij}$'s *are sampled from different distributions.*

(a) 6 testing procedures
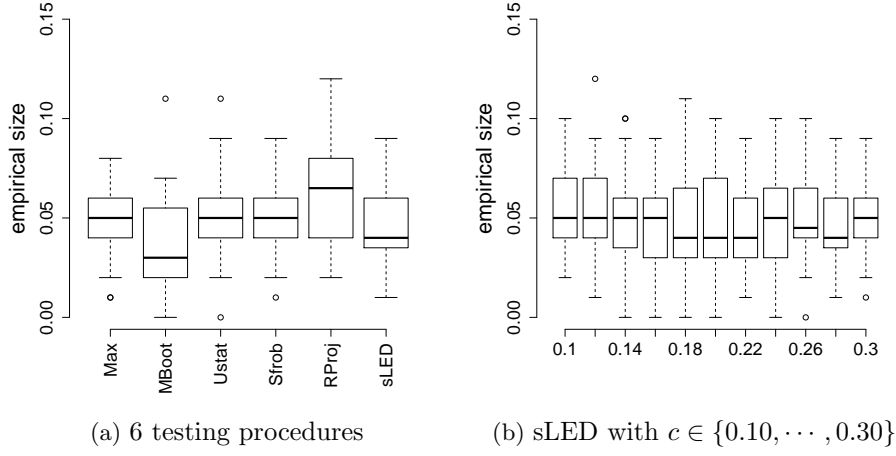
(b) sLED with $c \in \{0.10, \cdots, 0.30\}$

FIG. S2. **(a)** *Boxplot of the empirical size in 100 repetitions of the 6 testing procedures under different scenarios, where the smoothing parameter for sLED is chosen to be $c = 0.3$.* **(b)** *Boxplot of the empirical size in 100 repetitions of sLED under different scenarios, using different smoothing parameters $c \in \{0.10, 0.12, \cdots, 0.30\}$.*

By Lemma 2, there exist some constants $C_1$, $C_2$ depending on $(\underline{c}, \bar{c}, \nu^2)$, such that if $(n, p)$ are sufficiently large, $\mathbb{P}(Z \notin \mathcal{A}) \leq \delta/4$. Therefore, in order to show that

$$\mathbb{P}\left(||\hat{\Sigma}_1^* - \hat{\Sigma}||_\infty > \frac{\epsilon}{2}\right) \leq \mathbb{P}\left(||\hat{\Sigma}_1^* - \hat{\Sigma}||_\infty > \frac{\epsilon}{2}\Big| Z \in \mathcal{A}\right) + \mathbb{P}(Z \notin \mathcal{A}) \leq \frac{\delta}{2},$$

it suffices to show that given any $Z \in \mathcal{A}$, the conditional probability satisfies

$$(S2.2) \qquad \mathbb{P}_Z\left(||\hat{\Sigma}_1^* - \hat{\Sigma}||_\infty > \frac{\epsilon}{2}\right) \leq \frac{\delta}{4}.$$

For any $1 \leq i, j \leq p$, we first bound the $(i,j)$-th entry:

$$\mathbb{P}_Z\left(|\hat{\Sigma}_{1,ij}^* - \hat{\Sigma}_{ij}| > \frac{\epsilon}{2}\right) \leq \underbrace{\mathbb{P}_Z\left(\left|\frac{1}{n}\sum_{k=1}^n Z_{ki}^* Z_{kj}^* - \frac{1}{N}\sum_{k=1}^N Z_{ki} Z_{kj}\right| > \frac{\epsilon}{4}\right)}_{\Delta_1} +$$

$$+ \underbrace{\mathbb{P}_Z\left(\left|\bar{X}_i^* \bar{X}_j^* - \bar{Z}_i \bar{Z}_j\right| > \frac{\epsilon}{4}\right)}_{\Delta_2},$$

where $\bar{X}_i^* = \frac{1}{n}\sum_{k=1}^n Z_{ki}^*$. Now we bound $\Delta_1$ and $\Delta_2$ separately.

(i) $\Delta_1$: Note that for any $(k, i, j)$,

$$\left| Z^*_{ki} Z^*_{kj} \right| \le (m_z)^2, \ \mathrm{var}_Z \left( Z^*_{ki} Z^*_{kj} \right) \le \frac{1}{N} \sum_{l=1}^{N} Z^2_{li} Z^2_{lj} \le \overline{m}_z^{(4)}.$$

By Lemma 1, there exists a constant $C'_2$ depending on $(C_2, \nu^2)$, such that if $(n, p)$ are sufficiently large,

(S2.3) $$\Delta_1 \le 2 \exp \left\{ -\frac{n\epsilon^2/C'_2}{1 + \log(C_1 np/\delta)\epsilon} \right\}.$$

(ii) $\Delta_2$: Note that

$$\bar{X}^*_i \bar{X}^*_j - \bar{Z}_i \bar{Z}_j = (\bar{X}^*_i - \bar{Z}_i)(\bar{X}^*_j - \bar{Z}_j) + \bar{Z}_j(\bar{X}^*_i - \bar{Z}_i) + \bar{Z}_i(\bar{X}^*_j - \bar{Z}_j),$$

and for any $(k, i, j)$,

$$|\bar{Z}_i| \le \overline{m}_z, \ |Z^*_{ki}| \le m_z, \ \mathrm{var}_Z(Z^*_{ki}) \le \frac{1}{N} \sum_{l=1}^{N} Z^2_{li} \le \overline{m}_z^{(2)}.$$

Therefore,

$$\Delta_2 \le 2 \max_i \left[ \mathbb{P}_Z \left( \left| \frac{1}{n} \sum_{k=1}^{n} Z^*_{ki} - \bar{Z}_i \right| > \sqrt{\frac{\epsilon}{8}} \right) + \mathbb{P}_Z \left( \left| \frac{1}{n} \sum_{k=1}^{n} Z^*_{ki} - \bar{Z}_i \right| > \frac{\epsilon}{16\overline{m}_z} \right) \right].$$

Applying Lemma 1 on both terms, we know that there exists a constant $C''_2$ depending on $(C_2, \nu^2)$, such that if $(n, p)$ are sufficiently large,

(S2.4)
$$\Delta_2 \le 4 \exp \left\{ -\frac{n\epsilon/C''_2}{1 + \sqrt{\log(C_1 np/\delta)}\sqrt{\epsilon}} \right\} +$$
$$4 \exp \left\{ -\frac{n\epsilon^2/C''_2}{\frac{\log(C_1 p/\delta)}{n} + \sqrt{\frac{\log(C_1 p/\delta)\log(C_1 np/\delta)}{n}}\epsilon} \right\}.$$

Combining the results in (S2.3) and (S2.4), and note that $(\log p)^3 = O(n)$ by assumption (A3), we have $\Delta_1, \Delta_2 \le \frac{\delta}{8} p^{-2}$ if $(n, p)$ are sufficiently large, as long as

$$\epsilon \ge C' \sqrt{\frac{\log(C_1 p^2/\delta)}{n}}$$

for some constant $C'$ depending on $C'_2$ and $C''_2$. Finally, (S2.2) follows from a union bound over $1 \le i, j \le p$. Similar statement also holds for $||\hat{\Sigma}^*_2 - \hat{\Sigma}||_\infty$ with sample size $m$, and the final result follows from (S2.1) and the fact that $\underline{c} n \le m \le \bar{c} n$.                                        □

PROOF OF THEOREM 3. (i) Note that a feasible solution of (2.12) or (2.14) in the main manuscript always satisfies $||H||_1 \leq R$, where $H = vv^T$ if using (2.14). Then the result directly follows from the Hölder's inequality:

$$\text{tr}\left(\hat{D}H\right) \leq ||\hat{D}||_\infty ||H||_1 \,.$$

(ii) Let $v^*$ be the $R$-sparse leading eigenvector of $D$, then $||v^*||_2 = 1$ and $||v^*(v^*)^T||_1 = ||v^*||_1^2 \leq ||v^*||_0 = R$, so $v^*(v^*)^T$ is feasible for (2.12) and (2.14). The result follows from

$$\tilde{\lambda}_1^R(\hat{D}) - \lambda_1^R(D) \geq (v^*)^T \hat{D} v^* - (v^*)^T D v^*$$

and $\left|(v^*)^T (\hat{D} - D) v^*\right| \leq ||\hat{D} - D||_\infty ||v^*(v^*)^T||_1.$

□

**S3. Lemmas.** In this section, we state and prove the lemmas that are used in Section S2.

LEMMA 1 (Bernstein inequality for sampling without replacement). *Let* $\mathcal{Z} = \{z_1, ..., z_N\}$ *be a finite set containing* $N$ *real numbers, and* $(z_1^*, ..., z_n^*)$ *be i.i.d. random variables that are drawn without replacement from* $\mathcal{Z}$*. Let*

$$\bar{z} = \max_{1 \leq i \leq N} |z_i| \,, \ \mu_z = \frac{1}{N} \sum_{i=1}^N z_i \,, \ \sigma_z^2 = \frac{1}{N} \sum_{i=1}^N (z_i - \mu_z)^2 \,,$$

*then for any* $\epsilon > 0$,

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n z_i^* - \mu_z\right| \geq \epsilon\right) \leq 2 \exp\left\{-\frac{n\epsilon^2}{2\sigma_z^2 + \frac{4}{3}\bar{z}\epsilon}\right\} \,.$$

*As a consequence, for any* $t > 0$,

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n z_i^* - \mu_z\right| > \frac{4\bar{z}}{3} \frac{t}{n} + \sqrt{2\sigma_z^2 \frac{t}{n}}\right) \leq 2e^{-t} \,.$$

PROOF. See Proposition 1.4 in Bardenet et al. (2015). □

LEMMA 2 (Sub-gaussian tail bound). *Under assumptions (A1)-(A2), for* $\forall \delta > 0$*, there exist constants* $C_1, C_2$ *depending on* $(\underline{c}, \bar{c}, \nu^2)$*, such that if* $(n, p)$ *are sufficiently large, with probability at least* $1 - \delta$,

(i) $||\hat{\Sigma}_q - \Sigma_q||_\infty \le C_2 \sqrt{\frac{\log(C_1 p^2/\delta)}{N}}$ *for* $q = 1, 2$. *As a consequence,*

$$||\hat{D} - D||_\infty \le 2C_2 \sqrt{\frac{\log(C_1 p^2/\delta)}{N}} \, .$$

(ii) $\overline{m}_z \le C_2 \sqrt{\frac{\log(C_1 p/\delta)}{N}}$. *This together with (i) imply that*

$$\overline{m}_z^{(2)} \le 2\nu^2 + 2C_2 \sqrt{\frac{\log(C_1 p^2/\delta)}{N}} \, .$$

(iii) $m_z \le C_2 \sqrt{\log(C_1 N p/\delta)}$.
(iv) $\overline{m}_z^{(4)} \le C_2 \left[ 1 + \frac{\log(C_1 p^2/\delta)}{N} \right]$.

PROOF.     (i) See for example, Lemma 12 in Yuan (2010).
(ii) The first part is standard Hoeffding's bound on $\frac{1}{N} \sum_{k=1}^N Z_{ki}$, with a
     union bound over $1 \le i \le p$. The second part follows from

$$\overline{m}_z^{(2)} \le \max\{||\hat{\Sigma}_1||_\infty, ||\hat{\Sigma}_2||_\infty\} + (\overline{m}_z)^2 \, .$$

(iii) By Markov inequality, $\forall \epsilon, t > 0$,

$$\mathbb{P}\left( \max_{k,i} Z_{ki} > \epsilon \right) \le e^{-t\epsilon} \mathbb{E}\left[ e^{t \max_{k,i} Z_{ki}} \right] = e^{-t\epsilon} \mathbb{E}\left[ \max_{k,i} e^{tZ_{ki}} \right]$$

$$\le e^{-t\epsilon} \sum_{k=1}^N \sum_{i=1}^p \mathbb{E}\left[ e^{tZ_{ki}} \right] \le Np \cdot e^{-t\epsilon + \frac{t^2 \nu^2}{2}} \, .$$

Finally, take $t = \frac{\epsilon}{\nu^2}$, and note that similar arguments hold for $-Z_{ki}$.
(iv) For any given $(i, j)$, let $W_k = Z_{ki}^2 Z_{kj}^2$, and define its cumulant gener-
     ating function
$$\Psi_k(\theta) = \log \mathbb{E}\left[ e^{\theta(W_k - \mathbb{E}(W_k))} \right] \, .$$

Note that $\Psi_1 = \cdots = \Psi_n$ and $\Psi_{n+1} = \cdots = \Psi_{n+m}$. By Markov
inequality, for any $t, \theta > 0$,

(S3.1)     $$\mathbb{P}\left( \left| \frac{1}{n} \sum_{k=1}^n W_k - \mathbb{E}(W_1) \right| > t \right) \le 2 \exp\{-n\theta t + n\Psi(\theta)\} \, ,$$

where $\Psi(\theta) = \max\{\Psi_1(\theta), \Psi_{n+1}(\theta)\}$ is an upper bound of the cumulant
generating functions. Since $Z_{ki}, Z_{kj}$ are sub-gaussian, there exists a

small constant $\theta_0 \neq 0$, such that $\Psi(\theta_0) < \infty$. Plugging in $\theta_0$ to (S3.1), we know that with probability at least $1 - \frac{\delta}{2}p^{-2}$,

$$\frac{1}{n}\sum_{k=1}^{n} W_k - \mathbb{E}(W_1) \leq \frac{\log(4p^2/\delta)}{n\theta_0} + \frac{\Psi(\theta_0)}{\theta_0}.$$

The same arguments also hold for $\frac{1}{m}\sum_{k=(n+1)}^{n+m} W_k - \mathbb{E}(W_{n+1})$. Then the final result follows from a union bound over $(i, j)$ and the fact that $\mathbb{E}(W_k) \leq C\nu^4$ for some constant $C$.

$\square$

## References.

BARDENET, R., MAILLARD, O.-A. et al. (2015). Concentration inequalities for sampling without replacement. *Bernoulli* **21** 1361–1385.

CAI, T., LIU, W. and XIA, Y. (2013). Two-sample covariance matrix testing and support recovery in high-dimensional and sparse settings. *Journal of the American Statistical Association* **108** 265–277.

CHANG, J., ZHOU, W., ZHOU, W.-X. and WANG, L. (2016). Comparing Large Covariance Matrices under Weak Conditions on the Dependence Structure and its Application to Gene Clustering. *arXiv preprint arXiv:1505.04493v3*.

LI, J. and CHEN, S. X. (2012). Two sample tests for high-dimensional covariance matrices. *The Annals of Statistics* **40** 908–940.

SCHOTT, J. R. (2007). A test for the equality of covariance matrices when the dimension is large relative to the sample sizes. *Computational Statistics & Data Analysis* **51** 6535–6542.

WU, T.-L. and LI, P. (2015). Tests for High-Dimensional Covariance Matrices Using Random Matrix Projection. *arXiv preprint arXiv:1511.01611*.

YUAN, M. (2010). High dimensional inverse covariance matrix estimation via linear programming. *The Journal of Machine Learning Research* **11** 2261–2286.

ZHU, L., LEI, J., DEVLIN, B. and ROEDER, K. (2017). Testing High Dimensional Covariance Matrices, with Application to Detecting Schizophrenia Risk Genes. *The Annals of Applied Statistics*.

DEPARTMENT OF STATISTICS
CARNEGIE MELLON UNIVERSITY
5000 FORBES AVENUE
PITTSBURGH, PENNSYLVANIA 15213
USA
E-MAIL: lzhu@cmu.edu
      jinglei@andrew.cmu.edu
      roeder@andrew.cmu.edu

DEPARTMENT OF PSYCHIATRY AND HUMAN GENETICS
UNIVERSITY OF PITTSBURGH SCHOOL OF MEDICINE
3811 O'HARA STREET
PITTSBURGH, PENNSYLVANIA 15213
USA
E-MAIL: devlinbj@upmc.edu