# Supplement to "Joint mean and covariance estimation with unreplicated matrix-variate data"

Michael Hornstein, Roger Fan, Kerby Shedden, Shuheng Zhou Department of Statistics, University of Michigan

# Outline

We provide additional simulation and data analysis results in Section A and B. We state some preliminary results and notation in Section C. We prove Theorem 1 in Section D and Corollary 2 in Section D.2. We prove Theorem 3 in Section E, with additional lemmas proved in Section F. We prove entrywise convergence of the sample correlation matrices for Algorithm 1 in Section G. We prove Theorem 4 in Section H, and we prove additional lemmas used in the proof of Theorem 4 in Section I. In Section J we provide additional comparisons between our method and some related methods on both simulated and real data.

# A Additional simulation results

Figure S1 demonstrates the effect of mean structure on covariance estimation. As expected, when there is no mean structure Gemini performs competitively. As more mean structure is added, however, its performance quickly decays to be worse than Algorithm 2. This also provides evidence that the plug-in estimator  $\hat{\tau}_{init}$  used in Algorithm 2 is appropriately selecting genes to group center, as when there are no or very few differentially expressed genes Algorithm 2 is still never worse than Gemini. Algorithm 1 does not perform as well as Algorithm 2 but still tends to eventually outperform Gemini as more mean structure is added. As the sample size increases, the difference between Algorithm 2 and Algorithm 1 decreases as the added noise from group centering becomes less of a factor. We still recommend using Algorithm 2 in most realistic scenarios, but this reinforces our theoretical finding that the two algorithms have the same error rates.



Figure S1: Performance of Gemini, Algorithm 1, and Algorithm 2 for estimating B under different mean and covariance structures. As the sample size increases, we can see that Algorithm 1 improves relative to Gemini and begins to catch up to Algorithm 2. Gemini's performance always degrades as the true differences grow or more differentially expressed genes are added, while Algorithm 1 and 2 are stable. We set  $B^{-1}$  as Erdős-Rényi (ER) or star-block with blocks of size 4 (SB). All plots use A from an AR1(0.8) model with m = 2000 and are averaged over 200 replications. In the left plot the first 50 genes are differentially expressed at the level specified on the x-axis. As indicated, the three groups of lines correspond to n = 20, 40, and 80. In the right two columns there are m1 number of genes with exponentially decaying true differences between groups, scaled so that the largest difference is 5 (resulting in an average difference of approximately 1).

## **B** Additional data analysis

As discussed in Section 3.1, it is particularly important that the design effect is accurately estimated in order for the test statistics to be properly calibrated. The first plot of Figure S2a displays the sensitivity of the estimated design effect (21) for Algorithm 2 to the GLasso penalty parameter and the number of group centered columns. In the case that all columns are group centered, Algorithm 2 reduces to Algorithm 1. If we group center all genes, the estimated design effect is sensitive to the penalty parameter, but if we group center a small proportion of genes, it is less sensitive to the penalty parameter. This is further evidence that it may be advantageous to avoid over-centering the data when the true mean difference



(a) The first plot displays the estimated design effect vs. the penalty multiplier for Algorithm 2. Each curve corresponds to a different number of columns being group centered. As the curves go from top to bottom, the number of group centered columns increases from 10 to 2000. The second plot shows a quantile plot of test statistics from the data vs. simulated test statistics; in the simulation, the population person-person covariance matrix is  $\hat{B}$ , as estimated from the UC data.



(b) Quantile plot and inverse covariance graphs. The first two plots correspond to  $\lambda = 0.4$  and 128 group centered genes. The third plot corresponds to  $\lambda = 0.5$  and 128 group centered genes. Green circles correspond to twins with UC, orange circles to twins without UC. Twins are aligned vertically.

Figure S2

vector  $\gamma$  may be sparse. The second plot of Figure S2a shows a quantile plot comparing the distribution of test statistics from the UC data to test statistics from a simulation whose population correlation structure is matched to the UC data. The quantile plot reveals that we can reproduce the pattern of overdispersion in the test statistics using simulated data having person-person as well as gene-gene correlations. Such correlations therefore provide a possible explanation for the overdispersion of the test statistics.

Figure S2b displays a quantile plot and inverse covariance graph for  $\lambda = 0.4$  and 128 group centered genes. Under these settings the test statistics appear correctly calibrated, coinciding with the central portion of the reference line. Furthermore, the inverse covariance graph is sparse (38 edges). In the inverse covariance graph, there are more edges between subjects with UC than between the healthy subjects, which could be explained by the existence of subtypes of UC inducing correlations between subsets of subjects. The third plot of Figure S2b displays a sparser inverse covariance graph, corresponding to a larger penalty  $\lambda = 0.5$ . There are three edges between twin pairs, and there are more edges between subjects with UC than between those without UC.

#### **B.1** Stability simulation

Table S1 shows the results from a simulation analogous to Table 2, demonstrating stability across iterations of the procedure. Iteration 1 begins by group centering 1280 genes and this number is halved in each successive iteration. We can see from the table that the gene rankings generated by Algorithm 2 are robust to misspecifying the number of differentially expressed genes. When the number of group centered genes is 160 or below (iterations 4 through 8), the commonly selected genes among the top 20 genes are stable. Furthermore, the true positives remain stable as we decrease the amount of genes centered, while the false positives decrease.

Table S1: Number of genes in common among genes ranked in the top 20 when different numbers of genes are group centered. This simulation is analogous to Table 2. Note that the maximum possible value for any entry of the table is 20; if entry (i, j) is 20, then iterations i and j selected the same top twenty genes. The first 10 genes have a difference of 1.5 and the second 10 have a difference of 1. All remaining genes have a true mean difference of zero. We use B as estimated from the UC data, and A is from an AR1(0.8) model. These simulations have n = 20 individuals and 2000 genes and are averaged over 200 replications. The last two rows display the average number of true and false positives among the genes ranked in the top 20 of each iteration.

	1	2	3	4	5	6	7	8
1	20.0	17.6	15.8	14.8	14.3	14.0	14.0	13.9
2	17.6	20.0	17.9	16.8	16.2	15.9	15.8	15.8
3	15.8	17.9	20.0	18.7	18.1	17.8	17.7	17.6
4	14.8	16.8	18.7	20.0	19.3	19.0	18.9	18.8
5	14.3	16.2	18.1	19.3	20.0	19.6	19.5	19.4
6	14.0	15.9	17.8	19.0	19.6	20.0	19.8	19.7
7	14.0	15.8	17.7	18.9	19.5	19.8	20.0	19.8
8	13.9	15.8	17.6	18.8	19.4	19.7	19.8	20.0
TP	12.7	14.3	15.6	16.4	16.7	16.8	16.8	16.8
FP	7.3	5.7	4.4	3.6	3.3	3.2	3.2	3.2

## C Preliminary results

In this section, we refresh notation and introduce propositions that are shared in the proofs of the theorems. For convenience, we first restate some notation.

$$D = \begin{bmatrix} 1_{n_1} & 0\\ 0 & 1_{n_2} \end{bmatrix} \in \mathbb{R}^{n \times 2}$$
(S1)

$$\Omega = (D^T B^{-1} D)^{-1} \text{ and } \Omega_{n,m} = (D^T B^{-1}_{n,m} D)^{-1}$$
(S2)

$$\Delta = B_{n,m}^{-1} - B^{-1} \tag{S3}$$

$$\widehat{\beta}(\widehat{B}^{-1}) = (D^T \widehat{B}^{-1} D)^{-1} D^T \widehat{B}^{-1} X \in \mathbb{R}^{2 \times m}$$
(S4)

When D has the form (S1), the singular values are  $\sigma_{\max}(D) = \sqrt{n_{\max}}$  and  $\sigma_{\min}(D) = \sqrt{n_{\min}}$ . The condition number is  $\kappa(D) = \sigma_{\max}(D)/\sigma_{\min}(D) = \sqrt{n_{\text{ratio}}}$  where  $n_{\text{ratio}} = \max(n_1, n_2)/\min(n_1, n_2)$ . We first state some convenient notation and bounds.

$$r_a := a_{\max}/a_{\min}$$
 and  $r_b := b_{\max}/b_{\min};$ 

$$1/\varphi_{\min}(A) = \|A^{-1}\|_2 \leq \|\rho(A)^{-1}\|_2/a_{\min} = \frac{1}{a_{\min}\varphi_{\min}(\rho(A))},$$
(S5)

$$1/\varphi_{\min}(B) = \|B^{-1}\|_2 \leq \|\rho(B)^{-1}\|_2/b_{\min} = \frac{1}{b_{\min}\varphi_{\min}(\rho(B))},$$
(S6)

$$1/\varphi_{\min}(\rho(A)) = \|\rho(A)^{-1}\|_2 \leqslant a_{\max} \|A^{-1}\|_2,$$
(S7)

$$1/\varphi_{\min}(\rho(B)) = \|\rho(B)^{-1}\|_2 \le b_{\max}\|B^{-1}\|_2$$
(S8)

$$||A||_{2} \leq a_{\max} ||\rho(A)||_{2}, \quad ||B||_{2} \leq b_{\max} ||\rho(B)||_{2}, \tag{S9}$$

$$\|\rho(A)\|_2 \leq \|A\|_2/a_{\min}, \text{ and } \|\rho(B)\|_2 \leq \|B\|_2/b_{\min}.$$
 (S10)

The eigenvalues of the correlation matrices satisfy

$$0 < \varphi_{\min}(\rho(A)) \leq 1 \leq \varphi_{\max}(\rho(A)) \text{ and } 0 < \varphi_{\min}(\rho(B)) \leq 1 \leq \varphi_{\max}(\rho(B)).$$
(S11)

In the remainder of this section, we state preliminary results and highlight important intermediate steps that are used in the proofs of Theorems 1 and 3. First we state propositions used in mean estimation for Theorems 1 and 3.

#### C.1 Propositions

We now state propositions used in the proofs of Lemmas S5 and S6. We defer the proof of Proposition S1 to Section D.5.

**Proposition S1.** For  $\Omega$  as defined in (S2) and some design matrix D,

$$\|\Omega\|_2 \leq \|B\|_2 / \sigma_{\min}^2(D)$$

In the case that D is defined as in (S1), we have  $\|\Omega\|_2 \leq \|B\|_2/n_{\min}$ . Furthermore,

$$\lambda_{\min}(\Omega) \ge \frac{\lambda_{\min}(B)}{n_{\max}}.$$
 (S12)

We state the following perturbation bound.

**Theorem S2** (Golub & Van Loan, Theorem 2.3.4). If A is invertible and  $||A^{-1}E||_p < 1$ , then A + E is invertible and

$$\|(A+E)^{-1} - A^{-1}\|_p \leq \frac{\|E\|_p \|A^{-1}\|_p^2}{1 - \|A^{-1}E\|_p} \leq \frac{\|E\|_p \|A^{-1}\|_p^2}{1 - \|A^{-1}\|_p \|E\|_p}$$

In Proposition S3, we provide auxiliary upper bounds that depend on  $\|\Delta\|_2$ ,  $\|B\|_2$ ,  $\kappa(D)$ , and  $\sigma_{\min}(D)$ . We defer the proof of Proposition S3 to the end of this section, for clarity of presentation.

**Proposition S3.** Let  $\Delta = B_{n,m}^{-1} - B^{-1}$ .

$$\delta_0(\Delta) := \|\Omega_{n,m} - \Omega\|_2 \leqslant \frac{1}{\sigma_{\min}^2(D)} \frac{\|B\|_2^2 \|\Delta\|_2}{1/\kappa^2(D) - \|B\|_2 \|\Delta\|_2}$$
(S13)

$$\delta_1(\Delta) := \left\| \Omega D^T \Delta \right\|_2 \leqslant \sigma_{\max}(D) \|B\|_2 \|\Delta\|_2 / \sigma_{\min}^2(D) = \frac{\sqrt{n_{\max}}}{n_{\min}} \|B\|_2 \|\Delta\|_2.$$
(S14)

If  $||(D^T B^{-1}D)^{-1}D^T \Delta D||_2 < 1$ , then

$$\delta_2(\Delta) := \left\| (\Omega_{n,m} - \Omega) D^T \Delta \right\|_2 \leqslant \frac{\kappa(D)}{\sigma_{\min}(D)} \frac{\|B\|_2^2 \|\Delta\|_2^2}{1/\kappa^2(D) - \|B\|_2 \|\Delta\|_2}$$
(S15)

$$\delta_3(\Delta) := \left\| \left(\Omega_{n,m} - \Omega\right) D^T B^{-1} \right\|_2 \leq \frac{\kappa(D)}{\sigma_{\min}(D)} \frac{\|B\|_2^2 \|B^{-1}\|_2 \|\Delta\|_2}{1/\kappa^2(D) - \|B\|_2 \|_2 \|\Delta\|_2}$$
(S16)

The following proposition is a corollary of Proposition S3.

**Proposition S4.** When D has the form (S1), and  $\Omega$  is as defined in (S2),

$$\delta_{0}(\Delta) = \|\Omega_{n,m} - \Omega\|_{2} \leqslant \frac{1}{n_{\min}} \frac{\|B\|_{2}^{2} \|\Delta\|_{2}}{1/n_{ratio} - \|B\|_{2} \|\Delta\|_{2}}$$
  

$$\delta_{1}(\Delta) = \|\Omega D^{T} \Delta\|_{2} \leqslant \frac{\sqrt{n_{ratio}}}{\sqrt{n_{\min}}} \|B\|_{2} \|\Delta\|_{2}$$
  

$$\delta_{2}(\Delta) = \|(\Omega_{n,m} - \Omega) D^{T} \Delta\|_{2} \leqslant \frac{\sqrt{n_{ratio}}}{\sqrt{n_{\min}}} \frac{\|B\|_{2}^{2} \|\Delta\|_{2}^{2}}{1/n_{ratio} - \|B\|_{2} \|\Delta\|_{2}}$$

Let K be defined as in Theorem 1. We express the entrywise rates of convergence of the sample correlation matrices  $\widehat{\Gamma}(B)$  and  $\widehat{\Gamma}(A)$ , respectively, in terms of the following quantities:

$$\widetilde{\alpha} = C_A K \frac{\log^{1/2}(m)}{\sqrt{m}} \left( 1 + \frac{\|B\|_1}{n} \right) + \frac{\|B\|_1}{n_{\min}} \text{ and } \widetilde{\eta} = C_B K \frac{\log^{1/2}(m \vee n)}{\sqrt{n}} + \frac{\|B\|_1}{n}.$$
 (S17)

## D Proof of Theorem 1 and Corollary 2

#### D.1 Proof of Theorem 1

Let  $B_{n,m} \in \mathbb{R}^{n \times n}$  denote a fixed positive definite matrix. Let D be as defined as in (4). Define  $\Delta_{n,m} = B_{n,m}^{-1} - B^{-1}$  and

$$\Omega = (D^T B^{-1} D)^{-1} \text{ and } \Omega_{n,m} = (D^T B^{-1}_{n,m} D)^{-1}.$$
 (S18)

Note that we can decompose the error for all j as

$$\|\widehat{\beta}_{j}(B_{n,m}^{-1}) - \beta_{j}^{*}\|_{2} \leq \|\widehat{\beta}_{j}(B^{-1}) - \beta_{j}^{*}\|_{2} + \|\widehat{\beta}_{j}(B_{n,m}^{-1}) - \widehat{\beta}_{j}(B^{-1})\|_{2} =: \mathbf{I} + \mathbf{II}.$$
(S19)

We will use the following lemmas, which are proved in subsections D.4 and D.3, to bound these two terms on the right-hand side, respectively.

**Lemma S5.** Let  $\mathcal{E}_2$  denote the event

$$\mathcal{E}_{2} = \left\{ \|\widehat{\beta}_{j}(B^{-1}) - \beta_{j}^{*}\|_{2} \leqslant s_{n,m} \right\}, \quad with \quad s_{n,m} = C_{3}d^{1/2}\sqrt{\frac{\log(m)\|B\|_{2}}{n_{\min}}}.$$
 (S20)

Then  $P(\mathcal{E}_2) \ge 1 - 2/m^d$ .

**Lemma S6.** Let  $B_{n,m} \in \mathbb{R}^{n \times n}$  denote a fixed matrix such that  $B_{n,m} > 0$ . Let  $X_j \in \mathbb{R}^n$  denote the jth column of X, where X is a realization of model (2). Let  $\mathcal{E}_3$  denote the event

$$\mathcal{E}_{3} = \left\{ \|\widehat{\beta}_{j}(B_{n,m}^{-1}) - \widehat{\beta}_{j}(B^{-1})\|_{2} \leq t_{n,m} \right\}, \quad with \quad t_{n,m} = \widetilde{C}n_{\min}^{-1/2} \|\Delta_{n,m}\|_{2}.$$
(S21)

for some absolute constant  $\widetilde{C}$ . Then  $P(\mathcal{E}_3) \ge 1 - 2/m^d$ .

The proof of (18) follows from the union bound  $P(\mathcal{E}_2 \cap \mathcal{E}_3) \ge 1 - P(\mathcal{E}_2) - P(\mathcal{E}_3) \ge 1 - 4/m^d$ . Next we prove (20). Let  $r_{n,m} = s_{n,m} + t_{n,m}$ , as defined in (18). Let  $\delta = (1, -1) \in \mathbb{R}^2$ . Then

$$\left|\hat{\gamma}_{j}(B_{n,m}^{-1}) - \gamma_{j}\right| = \left|\delta^{T}\left(\hat{\beta}_{j}(B_{n,m}^{-1}) - \beta_{j}^{*}\right)\right| \leq \|\delta\|_{2}\|\hat{\beta}_{j}(B_{n,m}^{-1}) - \beta_{j}^{*}\|_{2} = \sqrt{2}\|\hat{\beta}_{j}(B_{n,m}^{-1}) - \beta_{j}^{*}\|_{2},$$

where we used the Cauchy-Schwarz inequality. Hence if  $\|\widehat{\beta}_j(B_{n,m}^{-1}) - \beta_j\|_2 \leq r_{n,m}$ , it follows

that  $|\hat{\gamma}_j(B_{n,m}^{-1}) - \gamma_j| \leq \sqrt{2}r_{n,m}$ . The result holds by applying a union bound over the variables  $j = 1, \ldots, m$ .  $\Box$ 

This completes the proof of Theorem 1.

#### D.2 Proof of Corollary 2 and Corollary 5

First note that by Proposition S4,

$$\begin{aligned} \left| \delta^{T} (D^{T} \hat{B}^{-1} D)^{-1} \delta - \delta^{T} (D^{T} B^{-1} D)^{-1} \delta \right| &= \left| \delta^{T} \left( (D^{T} \hat{B}^{-1} D)^{-1} - (D^{T} B^{-1} D)^{-1} \right) \delta \right| \\ &\leq \left\| \delta \right\|_{2}^{2} \left\| (D^{T} \hat{B}^{-1} D)^{-1} - (D^{T} B^{-1} D)^{-1} \right\|_{2} \\ &= 2 \left\| (D^{T} \hat{B}^{-1} D)^{-1} - (D^{T} B^{-1} D)^{-1} \right\|_{2} \\ &\leq 2 \frac{\left\| B \right\|_{2}^{2} \left\| \Delta \right\|_{2}}{n_{\min}}. \end{aligned}$$
(S22)

Note that by Proposition S1,

$$|\delta^T \Omega \delta| \ge \frac{\lambda_{\min}(B)}{n_{\max}}.$$
(S23)

Corollary 2 follows from (S22) and (S23), which provide an upper bound on the numerator and lower bound on the denominator, respectively.

Corollary 5 holds because by (28) of Theorem 4,

$$\left|\delta^{T}\left(\widehat{\Omega}-\Omega\right)\delta\right| \leq 2\frac{\|B\|_{2}^{2}}{n_{\min}}\left(\frac{C'\lambda_{A}\sqrt{|B^{-1}|_{0,\text{off}} \vee 1}}{b_{\min}\varphi_{\min}^{2}(\rho(B))}\right) \leq 2C'\frac{\kappa(B)}{n_{\min}}\lambda_{A}\sqrt{|B^{-1}|_{0,\text{off}} \vee 1} \quad (S24)$$

#### D.3 Proof of Lemma S5

First, we show that

$$\|\Omega^{1/2}\|_F + d^{1/2} K^2 \sqrt{\log(m)} \|\Omega\|_2^{1/2} / \sqrt{c} \leqslant s_{n,m},$$
(S25)

with  $s_{n,m}$  as defined in (19). Because  $\|\Omega^{1/2}\|_F \leq \sqrt{2}\|\Omega^{1/2}\|_2$ , it follows that

$$\begin{split} \|\Omega^{1/2}\|_{F} + d^{1/2}K^{2}\sqrt{\log(m)}\|\Omega\|_{2}^{1/2}/\sqrt{c} &\leq \left(\sqrt{2} + d^{1/2}K^{2}\sqrt{\log(m)}/\sqrt{c}\right)\|\Omega\|_{2}^{1/2} \\ &\leq C_{3}d^{1/2}\sqrt{\log(m)}\|\Omega\|_{2}^{1/2} \leqslant C_{3}d^{1/2}\sqrt{\frac{\log(m)\|B\|_{2}}{n_{\min}}}, \end{split}$$

where the last step follows from Proposition S1. Next, we express  $\hat{\beta}_j(B^{-1}) - \beta_j^*$  as

$$\hat{\beta}_j(B^{-1}) - \beta_j^* = \Omega^{1/2} \eta_j, \text{ where } \eta_j = \Omega^{-1/2} \left( \hat{\beta}_j(B^{-1}) - \beta_j^* \right).$$

By the bound (S25), event  $\mathcal{E}_2^c$  implies  $\{\|\Omega^{1/2}\eta_j\|_2 > \|\Omega^{1/2}\|_F + d^{1/2}K^2\sqrt{\log(m)}\|\Omega\|_2^{1/2}/\sqrt{c}\}$ . Therefore,

$$\begin{split} P\left(\|\Omega\eta_{j}\|_{2} \ge s_{n,m}\right) &\leqslant P\left(\|\Omega\eta_{j}\|_{2} > \|\Omega^{1/2}\|_{F} + d^{1/2}K^{2}\sqrt{\log(m)}\|\Omega\|_{2}^{1/2}/\sqrt{c}\right) \\ &\leqslant P\left(\left|\|\Omega^{1/2}\eta_{j}\|_{2} - \|\Omega^{1/2}\|_{F}\right| > d^{1/2}K^{2}\sqrt{\log(m)}\|\Omega\|_{2}^{1/2}/\sqrt{c}\right) \\ &\leqslant 2\exp\left(\frac{-c\left(d^{1/2}K^{2}\sqrt{\log(m)}\|\Omega\|_{2}^{1/2}/\sqrt{c}\right)^{2}}{K^{4}\|\Omega^{1/2}\|_{2}^{2}}\right) \\ &= 2\exp\left(\frac{-d\log(m)\|\Omega\|_{2}}{\|\Omega^{1/2}\|_{2}^{2}}\right) = 2\exp\left(-d\log(m)\right) = 2/m^{d}. \end{split}$$

#### D.4 Proof of Lemma S6

The proof will proceed in the following steps. First, we show that  $\hat{\beta}_j(B_{n,m}^{-1}) - \hat{\beta}_j(B^{-1})$  can be expressed as  $VZ_j$ , where

$$V = \left(\Omega_{n,m} D^T B_{n,m}^{-1} - \Omega D^T B^{-1}\right) B^{1/2} \in \mathbb{R}^{2 \times m}$$

is a fixed matrix, and  $Z_j = B^{-1/2}X_j$ . Second, we show that

$$\|V\|_F + d^{1/2}K^2 \log^{1/2}(m) \|V\|_2 / \sqrt{c} \leq \widetilde{C} n_{\min}^{-1/2} \|\Delta\|_2.$$

Third, we use the first and second steps combined with the Hanson-Wright inequality to show that with high probability,  $\|VZ_j\|_2$  is at most  $\tilde{C}n_{\min}^{-1/2}\|\Delta\|_2$ .

For the first step of the proof, let  $Z_j = B^{-1/2}X_j$ , and note that  $\hat{\beta}_j(B_{n,m}^{-1}) - \hat{\beta}_j(B^{-1}) = VZ_j$ , where  $V \in \mathbb{R}^{2 \times m}$  is a fixed matrix, because

$$\hat{\beta}_{j}(B_{n,m}^{-1}) - \hat{\beta}_{j}(B^{-1}) = \left[ (D^{T}B_{n,m}^{-1}D)^{-1}D^{T}B_{n,m}^{-1} - \Omega D^{T}B^{-1} \right] B^{1/2}(B^{-1/2}X_{j})$$
$$= \left[ (D^{T}B_{n,m}^{-1}D)^{-1}D^{T}B_{n,m}^{-1} - \Omega D^{T}B^{-1} \right] B^{1/2}Z_{j}.$$

For the second step of the proof, we show that  $\|V\|_F + d^{1/2}K^2 \log^{1/2}(m)\|V\|_2/\sqrt{c} \leq \widetilde{C}n_{\min}^{-1/2}\|\Delta\|_2$ . First we obtain an upper bound on V. By the triangle inequality,

$$\begin{split} \|\Omega_{n,m}D^{T}B_{n,m}^{-1} - \Omega D^{T}B^{-1}\|_{2} &= \left\|\Omega_{n,m}D^{T}B_{n,m}^{-1} - \Omega D^{T}B^{-1}\right\|_{2} \\ &\leqslant \left\|\left(\Omega_{n,m} - \Omega\right)D^{T}(B_{n,m}^{-1} - B^{-1})\right\|_{2} + \left\|\left(\Omega_{n,m} - \Omega\right)D^{T}B^{-1}\right\|_{2} + \left\|\Omega D^{T}\Delta\right\|_{2} \\ &= \delta_{2}(\Delta) + \delta_{3}(\Delta) + \delta_{1}(\Delta). \end{split}$$

We bound each of the three terms using Proposition S3,

$$\delta_{2}(\Delta) = \left\| (\Omega_{n,m} - \Omega) D^{T} \Delta \right\|_{2} \leqslant \frac{\sqrt{n_{\text{ratio}}}}{\sqrt{n_{\min}}} \frac{\|B\|_{2}^{2} \|\Delta\|_{2}^{2}}{1/n_{\text{ratio}} - \|B\|_{2} \|\Delta\|_{2}}$$
  
$$\delta_{3}(\Delta) = \left\| (\Omega_{n,m} - \Omega) D^{T} B^{-1} \right\|_{2} \leqslant \frac{\sqrt{n_{\text{ratio}}}}{\sqrt{n_{\min}}} \frac{\|B\|_{2}^{2} \|B^{-1}\|_{2} \|\Delta\|_{2}}{1/n_{\text{ratio}} - \|B\|_{2} \|2\|\Delta\|_{2}}$$
  
$$\delta_{1}(\Delta) = \left\| \Omega D^{T} \Delta \right\|_{2} \leqslant \frac{\sqrt{n_{\text{ratio}}}}{\sqrt{n_{\min}}} \|B\|_{2} \|\Delta\|_{2}.$$

Applying the above bounds yields

$$\begin{split} \|V\|_{2} &\leqslant \frac{\sqrt{n_{\text{ratio}}}}{\sqrt{n_{\min}}} \|\Delta\|_{2} \|B\|_{2}^{1/2} \left( \frac{\|B\|_{2}^{2} \|\Delta\|_{2}}{1/\kappa^{2}(D) - \|B\|_{2} \|\Delta\|_{2}} + \frac{\|B\|_{2}^{2} \|B^{-1}\|_{2}}{1/\kappa^{2}(D) - \|B\|_{2} \|2\|\Delta\|_{2}} + \|B\|_{2} \right) \\ &\leqslant \quad \widetilde{C}n_{\min}^{-1/2} \|\Delta\|_{2}. \end{split}$$

For the third step of the proof, we use the Hanson-Wright inequality to bound  $\|VZ_j\|_2$ :

$$P\left(\|VZ_{j}\|_{2} > \widetilde{C}n_{\min}^{-1/2}\|\Delta\|_{2}\right) \leq P\left(\|VZ_{j}\|_{2} > \|V\|_{F} + d^{1/2}K^{2}\log^{1/2}(m)\|V\|_{2}/\sqrt{c}\right)$$

$$= P\left(\|VZ_{j}\|_{2} - \|V\|_{F} > d^{1/2}K^{2}\log^{1/2}(m)\|V\|_{2}/\sqrt{c}\right)$$

$$\leq P\left(\|\|VZ_{j}\|_{2} - \|V\|_{F}\| > d^{1/2}K^{2}\log^{1/2}(m)\|V\|_{2}/\sqrt{c}\right)$$

$$\leq 2\exp\left(-\frac{c\left(d^{1/2}K^{2}\log^{1/2}(m)\|V\|_{2}/\sqrt{c}\right)^{2}}{K^{4}\|V\|_{2}^{2}}\right) \qquad (\text{Hanson-Wright inequality})$$

$$= 2\exp\left(-d\log(m)\right) = 2/m^{d}.$$

#### D.5 Proof of Proposition S1

Let  $D = U\Psi V^T$  be the singular value decomposition of D, with  $U \in \mathbb{R}^{n \times 2}$ ,  $\Psi \in \mathbb{R}^{2 \times 2}$ , and  $V \in \mathbb{R}^{2 \times 2}$ . Then  $(D^T B^{-1} D)^{-1} = (V\Psi U^T B^{-1} U\Psi V^T)^{-1} = V\Psi^{-1} (U^T B^{-1} U)^{-1} \Psi^{-1} V^T$ . Thus

$$\begin{split} \| (D^T B^{-1} D)^{-1} \|_2 &= \| \Psi^{-1} (U^T B^{-1} U)^{-1} \Psi^{-1} \|_2 \qquad \text{(because } V \text{ is square, orthonormal)} \\ &\leq \| \Psi^{-1} \|_2^2 \| (U^T B^{-1} U)^{-1} \|_2 \qquad \text{(sub-multiplicative property)} \\ &= \sigma_{\max}^2 (\Psi^{-1}) \| (U^T B^{-1} U)^{-1} \|_2 \\ &= \| (U^T B^{-1} U)^{-1} \|_2 / \sigma_{\min}^2 (\Psi) = \| (U^T B^{-1} U)^{-1} \|_2 / \sigma_{\min}^2 (D), \end{split}$$

where  $\sigma_{\min}(D) = \sigma_{\min}(\Psi)$ , because  $\Psi$  is the diagonal matrix of singular values of D. Next, note that  $\|(U^T B^{-1} U)^{-1}\|_2 = 1/\varphi_{\min}(U^T B^{-1} U)$  and

$$\varphi_{\min}(U^T B^{-1} U) = \min_{\eta \in \mathbb{R}^2} \eta^T U^T B^{-1} U \eta / \eta^T \eta.$$

We perform the change of variables  $\gamma = U\eta$ , under which  $\eta^T \eta = \gamma^T U^T U \gamma = \gamma^T \gamma$  (that is, U preserves the length of  $\eta$  because the columns of U are orthonormal). Hence

$$\varphi_{\min}(U^T B^{-1} U) = \min_{\gamma \in \operatorname{col}(U), \gamma \neq 0} \gamma^T B^{-1} \gamma / \gamma^T \gamma$$
$$\geqslant \min_{\gamma \neq 0} \gamma^T B^{-1} \gamma / \gamma^T \gamma$$
$$= \varphi_{\min}(B^{-1}) = 1 / \|B\|_2.$$

We have shown that  $1/\varphi_{\min}(U^T B^{-1}U) \leq ||B||_2$ , which implies that

$$||(U^T B^{-1} U)^{-1}||_2 \le ||B||_2$$

Therefore

$$||(D^T B^{-1} D)^{-1}||_2 \leq ||B||_2 / \sigma_{\min}^2(D).$$

In the special case of the two-group design matrix,  $\sigma_{\min}^2(D) = n_{\min}$ , so  $\|(D^T B^{-1} D)^{-1}\|_2 \leq \|B\|_2 / n_{\min}$ .

The proof of (S12) is as follows:

$$\lambda_{\min}(\Omega) = \frac{1}{\lambda_{\max}\left(\Omega^{-1}\right)} = \frac{1}{\lambda_{\max}\left(D^T B^{-1} D\right)} \ge \frac{1}{\|D\|_2^2 \lambda_{\max}(B^{-1})} = \frac{\lambda_{\min}(B)}{\|D\|_2^2} = \frac{\lambda_{\min}(B)}{n_{\max}}$$

#### D.6 Proof of Proposition S3

By the definitions of  $\Omega_{n,m}$  in (S2) and  $\Delta = B_{n,m}^{-1} - B^{-1}$ , we have by Theorem S2

$$\begin{split} \|\Omega_{n,m} - \Omega\|_{2} &= \|(D^{T}B_{n,m}D)^{-1} - \Omega\|_{2} \\ &= \left\| \left( D^{T}B_{n,m}^{-1}D - D^{T}B^{-1}D + D^{T}B^{-1}D \right)^{-1} - \Omega \right\|_{2} \\ &= \left\| \left( D^{T}B^{-1}D + D^{T}\Delta D \right)^{-1} - \Omega \right\|_{2} \\ &\leqslant \frac{\|D^{T}\Delta D\|_{2}\|\Omega\|_{2}^{2}}{1 - \|\Omega\|_{2}\|D^{T}\Delta D\|_{2}} \quad \text{(by Theorem S2)} \\ &\leqslant \frac{(\sigma_{\max}^{2}(D)/\sigma_{\min}^{4}(D)) \|B\|_{2}^{2}\|\Delta\|_{2}}{1 - \kappa^{2}(D)\|B\|_{2}\|\Delta\|_{2}}. \end{split}$$

In the last step we apply Proposition S1. Thus

$$\begin{aligned} \|\Omega_{n,m} - \Omega\|_{2} &\leqslant \frac{1}{\sigma_{\min}^{2}(D)} \frac{\kappa^{2}(D) \|B\|_{2}^{2} \|\Delta\|_{2}}{1 - \kappa^{2}(D) \|B\|_{2} \|\Delta\|_{2}} \\ &= \frac{1}{\sigma_{\min}^{2}(D)} \frac{\|B\|_{2}^{2} \|\Delta\|_{2}}{(1/\kappa^{2}(D)) - \|B\|_{2} \|\Delta\|_{2}}. \end{aligned}$$

We prove (S14) using the submultiplicative property of the operator norm and Proposition S1:

$$\|\Omega D^T \Delta\|_2 \leq \frac{\|B\|_2}{\sigma_{\min}^2(D)} \sigma_{\max}(D) \|\Delta\|_2 = \frac{\kappa(D)}{\sigma_{\min}(D)} \|B\|_2 \|\Delta\|_2.$$

We prove (S15), as follows:

$$\begin{split} \left\| \left(\Omega_{n,m} - \Omega\right) D^{T} \Delta \right\|_{2} &\leq \left\|\Omega_{n,m} - \Omega\right\|_{2} \left\| D^{T} \right\|_{2} \|\Delta\|_{2} \\ &\leq \left[ \frac{1}{\sigma_{\min}^{2}(D)} \frac{\|B\|_{2}^{2} \|\Delta\|_{2}}{(1/\kappa^{2}(D)) - \|B\|_{2} \|\Delta\|_{2}} \right] \sigma_{\max}(D) \|\Delta\|_{2} \qquad \text{(by Proposition S3)} \\ &= \frac{\kappa(D)}{\sigma_{\min}(D)} \frac{\|B\|_{2}^{2} \|\Delta\|_{2}^{2}}{(1/\kappa^{2}(D)) - \|B\|_{2} \|\Delta\|_{2}}. \end{split}$$

The proof of (S16) is analogous.  $\Box$ 

## E Proof of Theorem 3

Note that the proof in the current Section follows exactly the same steps as the proof of Theorems 3.1 and 3.2 in Zhou (2014a). Theorem 3 **Part II** is proved in Section E.2. To prove Theorem 3 **Part I**, we first state Lemma S7, which establishes rates of convergence for estimating  $A^{-1}$  and  $B^{-1}$  in the operator and the Frobenius norm. We then state the auxiliary Lemma S8, which is identical to that for Theorems 11.1 and 11.2 of Zhou (2014a), except that we plug in  $\tilde{\alpha}$  and  $\tilde{\eta}$  as defined in (S17). Putting these results together proves Theorem 3, **Part I**. We prove these auxiliary results in Section F.

Let  $\mathcal{X}_0$  denote the event

$$\forall i, j \qquad \left| \frac{(e_i - p_i)^T X X^T (e_j - p_j)}{\operatorname{tr}(A^*) \sqrt{b_{ii}^* b_{jj}^*}} - \rho_{ij}(B) \right| \leqslant \widetilde{\alpha}$$
(S26)

$$\forall i, j \qquad \left| \frac{X_i^T (I - P_2) X_j}{\operatorname{tr}(B^*) \sqrt{a_{ii}^* a_{jj}^*}} - \rho_{ij}(A) \right| \le \widetilde{\eta},$$
(S27)

with  $\mathcal{X}_0(B)$  and  $\mathcal{X}_0(A)$  denoting the events defined by equations (S26) and (S27), respectively. Let  $\tilde{\alpha}$  and  $\tilde{\eta}$  be as defined in (S17). On event  $\mathcal{X}_0(A)$ , for all j,  $\hat{\Gamma}_{jj}(A) = \rho_{jj}(A) = 1$  and

$$\max_{j,k,j\neq k} |\widehat{\Gamma}_{jk}(A) - \rho_{jk}(A)| \leq \frac{2\widetilde{\eta}}{1 - \widetilde{\eta}}$$
(S28)

On event  $\mathcal{X}_0(B)$ , for all j,  $\widehat{\Gamma}_{jj}(B) = \rho_{jj}(B) = 1$  and

$$\max_{j,k,j+k} |\widehat{\Gamma}_{jk}(B) - \rho_{jk}(B)| \leq \frac{2\widetilde{\alpha}}{1 - \widetilde{\alpha}}.$$
(S29)

**Lemma S7.** Suppose (A1) and (A2) hold. Let  $\widehat{W}_1$  and  $\widehat{W}_2$  be as defined in (10). Let  $\widehat{A}_{\rho}$ and  $\widehat{B}_{\rho}$  be as defined in (8a) and (8b). For some absolute constants 18 < C, C' < 36, the following events hold with probability at least  $1 - 2/(n \vee m)^2$ ,

$$\delta_{A,2} := \|\widehat{W}_1 \widehat{A}_\rho \widehat{W}_1 / \operatorname{tr}(B) - A\|_2 \leqslant C a_{\max} \kappa(\rho(A))^2 \lambda_B \sqrt{|A^{-1}|_{0,off} \vee 1}$$
(S30)

$$\delta_{B,2} := \|\widehat{W}_2 \widehat{B}_\rho \widehat{W}_2 / \operatorname{tr}(A) - B\|_2 \leqslant C' b_{\max} \kappa(\rho(B))^2 \lambda_A \sqrt{|B^{-1}|_{0,off} \vee 1}$$
(S31)

$$\delta_{A,F} := \|\widehat{W}_1 \widehat{A}_\rho \widehat{W}_1 / \operatorname{tr}(B) - A\|_F \leqslant C a_{\max} \kappa(\rho(A))^2 \lambda_B \sqrt{|A^{-1}|_{0,off} \vee m}$$
(S32)

$$\delta_{B,F} := \|\widehat{W}_2 \widehat{B}_\rho \widehat{W}_2 / \operatorname{tr}(A) - B\|_F \leqslant C' b_{\max} \kappa(\rho(B))^2 \lambda_A \sqrt{|B^{-1}|_{0,off}} \lor n;$$
(S33)

and for some 10 < C, C' < 19,

$$\begin{split} \delta_{A,2}^{-} &:= \left\| \operatorname{tr}(B) \left( \widehat{W}_{1} \widehat{A}_{\rho} \widehat{W}_{1} \right)^{-1} - A^{-1} \right\|_{2} \leqslant \frac{C \lambda_{B} \sqrt{|A^{-1}|_{0,off} \vee 1}}{a_{\min} \varphi_{\min}^{2}(\rho(A))} \\ \delta_{B,2}^{-} &:= \left\| \operatorname{tr}(A) \left( \widehat{W}_{2} \widehat{B}_{\rho} \widehat{W}_{2} \right)^{-1} - B^{-1} \right\|_{2} \leqslant \frac{C' \lambda_{A} \sqrt{|B^{-1}|_{0,off} \vee 1}}{b_{\min} \varphi_{\min}^{2}(\rho(B))} \\ \delta_{A,F}^{-} &:= \left\| \operatorname{tr}(B) \left( \widehat{W}_{1} \widehat{A}_{\rho} \widehat{W}_{1} \right)^{-1} - A^{-1} \right\|_{F} \leqslant \frac{C \lambda_{B} \sqrt{|A^{-1}|_{0,off} \vee m}}{a_{\min} \varphi_{\min}^{2}(\rho(A))} \\ \delta_{B,F}^{-} &:= \left\| \operatorname{tr}(A) \left( \widehat{W}_{2} \widehat{B}_{\rho} \widehat{W}_{2} \right)^{-1} - B^{-1} \right\|_{F} \leqslant \frac{C' \lambda_{A} \sqrt{|B^{-1}|_{0,off} \vee m}}{b_{\min} \varphi_{\min}^{2}(\rho(B))} \end{split}$$

Lemma S8 follows from Theorems 11.1 and 11.2 of Zhou (2014a,b), where we now plug in  $\tilde{\alpha}$  and  $\tilde{\eta}$  as defined in (S17). For completeness, we provide a sketch in Section F.2.

**Lemma S8.** Suppose (A1) and (A2) hold. For  $\varepsilon_1, \varepsilon_2 \in (0, 1)$ , let

$$\lambda_A = \widetilde{\eta} / \varepsilon_1, \quad \lambda_B = \widetilde{\alpha} / \varepsilon_2,$$

for  $\tilde{\alpha}$ ,  $\tilde{\eta}$  as defined in (S17), and suppose  $\lambda_A, \lambda_B < 1$ . Then on event  $\mathcal{X}_0$ , for 18 < C, C' < 36,

$$\begin{split} \|\widehat{A \otimes B} - A \otimes B\|_{2} &\leq \frac{\lambda_{A} \wedge \lambda_{B}}{2} \|A\|_{2} \|B\|_{2} + C\lambda_{B} a_{\max} \|B\|_{2} \kappa(\rho(A))^{2} \sqrt{|A^{-1}|_{0,off} \vee 1} \\ &+ C' \lambda_{A} b_{\max} \|A\|_{2} \kappa(\rho(B))^{2} \sqrt{|B^{-1}|_{0,off} \vee 1} \\ &+ 2 \left[ C' \lambda_{A} b_{\max} \kappa(\rho(B))^{2} \sqrt{|B^{-1}|_{0,off} \vee 1} \right] \left[ C\lambda_{B} a_{\max} \kappa(\rho(A))^{2} \sqrt{|A^{-1}|_{0,off} \vee 1} \right], \end{split}$$

and for 10 < C, C' < 19,

$$\begin{split} \widehat{\|A \otimes B^{-1} - A^{-1} \otimes B^{-1}\|_{2}} &\leqslant \frac{\lambda_{A} \wedge \lambda_{B}}{3} \|A^{-1}\|_{2} \|B^{-1}\|_{2} + C\lambda_{B} \|B^{-1}\|_{2} \frac{\sqrt{|A^{-1}|_{0,off} \vee 1}}{a_{\min}\varphi_{\min}^{2}(\rho(A))} \\ &+ C'\lambda_{A} \|A^{-1}\|_{2} \frac{\sqrt{|B^{-1}|_{0,off} \vee 1}}{b_{\min}\varphi_{\min}^{2}(\rho(B))} + \frac{3}{2} \left[ C\lambda_{B} \frac{\sqrt{|A^{-1}|_{0,off} \vee 1}}{a_{\min}\varphi_{\min}^{2}(\rho(A))} \right] \left[ C'\lambda_{A} \frac{\sqrt{|B^{-1}|_{0,off} \vee 1}}{b_{\min}\varphi_{\min}^{2}(\rho(B))} \right] \end{split}$$

;

For 18 < C, C' < 36,

$$\begin{split} \widehat{\|A \otimes B} - A \otimes B\|_{F} &\leqslant \frac{\lambda_{A} \wedge \lambda_{B}}{2} \|A\|_{F} \|B\|_{F} + C\lambda_{B} a_{\max} \|B\|_{F} \kappa(\rho(A))^{2} \sqrt{|A^{-1}|_{0,off} \vee m} \\ &+ C' \lambda_{A} b_{\max} \|A\|_{F} \kappa(\rho(B))^{2} \sqrt{|B^{-1}|_{0,off} \vee n} \\ &+ 2 \left[ C' \lambda_{A} b_{\max} \kappa(\rho(B))^{2} \sqrt{|B^{-1}|_{0,off} \vee n} \right] \left[ C \lambda_{B} a_{\max} \kappa(\rho(A))^{2} \sqrt{|A^{-1}|_{0,off} \vee m} \right], \end{split}$$

and for 10 < C, C' < 19,

$$\begin{split} \|\widehat{A \otimes B}^{-1} - A^{-1} \otimes B^{-1}\|_{F} &\leqslant \frac{\lambda_{A} \wedge \lambda_{B}}{3} \|A^{-1}\|_{2} \|B^{-1}\|_{F} + C\lambda_{B} \|B^{-1}\|_{F} \frac{\sqrt{|A^{-1}|_{0,off} \vee m}}{a_{\min}\varphi_{\min}^{2}(\rho(A))} \\ &+ C'\lambda_{A} \|A^{-1}\|_{F} \frac{\sqrt{|B^{-1}|_{0,off} \vee n}}{b_{\min}\varphi_{\min}^{2}(\rho(B))} + \frac{7}{5} \left[ C\lambda_{B} \frac{\sqrt{|A^{-1}|_{0,off} \vee m}}{a_{\min}\varphi_{\min}^{2}(\rho(A))} \right] \left[ C'\lambda_{A} \frac{\sqrt{|B^{-1}|_{0,off} \vee n}}{b_{\min}\varphi_{\min}^{2}(\rho(B))} \right]. \end{split}$$

## E.1 Proof of Theorem 3, Part I

We state additional helpful bounds:

$$(a_{\min} \lor \varphi_{\min}(A))\sqrt{m} \leqslant ||A||_F = \left(\sum_{i=1}^m \varphi_i^2(A)\right)^{1/2} \leqslant \sqrt{m} ||A||_2,$$
 (S34)

$$(b_{\min} \vee \varphi_{\min}(B))\sqrt{n} \leqslant \|B\|_F = \left(\sum_{i=1}^m \varphi_i^2(B)\right)^{1/2} \leqslant \sqrt{n} \|B\|_2, \tag{S35}$$

$$\sqrt{m}/a_{\max} = \left(\frac{1}{a_{\max}} \vee \frac{1}{\varphi_{\max}(A)}\right)\sqrt{m} \leqslant \|A^{-1}\|_F \leqslant \sqrt{m}\|A^{-1}\|_2, \tag{S36}$$

and

$$\sqrt{n}/b_{\max} = \left(\frac{1}{b_{\max}} \vee \frac{1}{\varphi_{\max}(B)}\right)\sqrt{n} \le \|B^{-1}\|_F \le \sqrt{n}\|B^{-1}\|_2.$$
(S37)

Proof of Theorem 3, Part I. We plug in bounds as in (S9) and (S10) into Lemma

S8 to obtain under (A1) and (A2),  $\left\|\widehat{A \otimes B} - A \otimes B\right\|_2 \leq \|A\|_2 \|B\|_2 \delta$ , where

$$\begin{split} \delta &= \frac{\lambda_A \wedge \lambda_B}{2} + \frac{Cr_a \kappa(\rho(A))}{\varphi_{\min}(\rho(A))} \lambda_B \sqrt{|A^{-1}|_{0,\text{off}} \vee 1} + \frac{C'r_b \kappa(\rho(B))}{\varphi_{\min}(\rho(B))} \lambda_A \sqrt{|B^{-1}|_{0,\text{off}} \vee 1} \\ &+ 2 \left[ \frac{Cr_a \kappa(\rho(A))}{\varphi_{\min}(\rho(A))} \lambda_B \sqrt{|A^{-1}|_{0,\text{off}} \vee 1} \right] \left[ \frac{C'r_b \kappa(\rho(B))}{\varphi_{\min}(\rho(B))} \lambda_A \sqrt{|B^{-1}|_{0,\text{off}} \vee 1} \right] \\ &= \frac{\lambda_A \wedge \lambda_B}{2} + \log^{1/2} (m \vee n) \left( \sqrt{\frac{|A^{-1}|_{0,\text{off}} \vee 1}{m}} + \sqrt{\frac{|B^{-1}|_{0,\text{off}} \vee 1}{n}} \right) + o(1). \end{split}$$

For the inverse, we plug in bounds as in (S7) and (S8) into Lemma S8 to obtain under (A1) and (A2),  $\left\|\widehat{A \otimes B}^{-1} - A^{-1} \otimes B^{-1}\right\|_2 \leq \|A^{-1}\|_2 \|B^{-1}\|_2 \delta'$ , where

$$\begin{split} \delta' &= \frac{\lambda_A \wedge \lambda_B}{3} + \frac{Cr_a \lambda_B \sqrt{|A^{-1}|_{0,\text{off}} \vee 1}}{\varphi_{\min}(\rho(A))} + \frac{C'r_b \lambda_A \sqrt{|B^{-1}|_{0,\text{off}} \vee 1}}{\varphi_{\min}(\rho(B))} \\ &+ \frac{3}{2} \left[ \frac{Cr_a \lambda_B \sqrt{|A^{-1}|_{0,\text{off}} \vee 1}}{\varphi_{\min}(\rho(A))} \right] \left[ \frac{C'r_b \lambda_A \sqrt{|B^{-1}|_{0,\text{off}} \vee 1}}{\varphi_{\min}(\rho(B))} \right] \\ &\approx \frac{\lambda_A \wedge \lambda_B}{3} + \log^{1/2} (m \vee n) \left( \sqrt{\frac{|A^{-1}|_{0,\text{off}} \vee 1}{m}} + \sqrt{\frac{|B^{-1}|_{0,\text{off}} \vee 1}{n}} \right) + o(1). \end{split}$$

The bounds in the Frobenius norm are proved in a similar manner; see Zhou (2014a) to finish.  $\Box$ 

#### E.2 Proof of Theorem 3, Part II

Let  $\widehat{B}^{-1} = \widehat{W}_2 \widehat{B}_{\rho} \widehat{W}_2$ . Let  $\widehat{\Delta} = \widehat{B}^{-1} - B^{-1}$ . Let  $\mathcal{E}_0(B)$  denote the event given by equations (S34) and (S34), which we know has probability at least  $1 - 2/(n \vee m)^2$  from Lemma S7, and define the event

$$\mathcal{E}_4 = \left\{ \| \hat{\beta}_j(\hat{B}^{-1}) - \beta_j^* \|_2 \leqslant s_{n,m} + t'_{n,m} \right\},\tag{S38}$$

where  $s_{n,m}$  is as defined in (19) and

$$t'_{n,m} := C\lambda_A \sqrt{\frac{n_{\text{ratio}} \left( |B_0^{-1}|_{0,\text{off}} \vee 1 \right)}{n_{\min}}}.$$
 (S39)

Under  $\mathcal{E}_0(B)$ , we see that

$$\|\widehat{\Delta}\|_2 \leqslant \frac{C'\lambda_A \sqrt{|B^{-1}|_{0,\text{off}} \vee 1}}{b_{\min}\varphi_{\min}^2(\rho(B))} = o(1).$$
(S40)

Using Proposition S1 and the fact that  $||D||_2 = \sqrt{n_{\text{max}}}$ , we get that

$$\|\Omega D^T \widehat{\Delta} D\|_2 \leqslant n_{\text{ratio}} \|B\|_2 \|\widehat{\Delta}\|_2, \tag{S41}$$

From (S40) we know that  $\|\widehat{\Delta}\|_2 \leq 1/(n_{\text{ratio}}\|B\|_2)$ , which we can plug into (S41) to show that  $\|\Omega D^T \widehat{\Delta} D\|_2 < 1$ . This implies that  $\widetilde{C} n_{\min}^{-1/2} \|\widehat{\Delta}\|_2 \leq t'_{n,m}$ . Therefore, we can apply Theorem 1 to get that the conditional probability of  $\mathcal{E}_4$  given  $\mathcal{E}_0(B)$  is at least  $1 - 4/(n \vee m)^2$ .

We can then bound the unconditional probability,

$$P(\mathcal{E}_4^c) \leq P(\mathcal{E}_4^c \mid \mathcal{E}_0(B)) P(\mathcal{E}_0(B)) + P(\mathcal{E}_0(B)^c)$$
  
$$\leq P(\mathcal{E}_4^c \mid \mathcal{E}_0(B)) + P(\mathcal{E}_0(B)^c)$$
  
$$\leq \frac{4}{(n \lor m)^2} + \frac{2}{(n \lor m)^2}.$$

## F More proofs for Theorem 3

The proof of Lemma S7 appears in Section F.1. The proofs of auxiliary lemmas appear in Section F.2.

#### F.1 Proof of Lemma S7

In order to prove Lemma S7, we need Theorem S9, which shows explicit non-asymptotic convergence rates in the Frobenius norm for estimating  $\rho(A)$ ,  $\rho(B)$ , and their inverses. Theorem S9 follows from the standard proof; see Rothman et al. (2008); Zhou et al. (2011) We also need Proposition S11 and Lemma S10, which are stated below and proved in Section F.2.

**Theorem S9.** Suppose that (A2) holds. Let  $\hat{A}_{\rho}$  and  $\hat{B}_{\rho}$  be the unique minimizers defined by (8a) and (8b) with sample correlation matrices  $\hat{\Gamma}(A)$  and  $\hat{\Gamma}(B)$  as their input. Suppose that event  $\mathcal{X}_0$  holds, with

$$\widetilde{\eta} \sqrt{|A^{-1}|_{0,off} \vee 1} = o(1) \quad and \quad \widetilde{\alpha} \sqrt{|B^{-1}|_{0,off} \vee 1} = o(1).$$
Set for some  $0 < \epsilon, \varepsilon < 1, \ \lambda_B = \widetilde{\alpha}/\varepsilon \ and \ \lambda_A = \widetilde{\eta}/\epsilon.$ 
(S42)

Then on event  $\mathcal{X}_0$ , we have for 9 < C < 18

$$\begin{aligned} \left\| \widehat{A}_{\rho} - \rho(A) \right\|_{2} &\leq \left\| \widehat{A}_{\rho} - \rho(A) \right\|_{F} \leq C\kappa(\rho(A))^{2} \lambda_{B} \sqrt{|A^{-1}|_{0,off} \vee 1}, \\ \left\| \widehat{B}_{\rho} - \rho(B) \right\|_{2} &\leq \left\| \widehat{B}_{\rho} - \rho(B) \right\|_{F} \leq C\kappa(\rho(B))^{2} \lambda_{A} \sqrt{|B^{-1}|_{0,off} \vee 1}, \end{aligned}$$

and

$$\left\|\hat{A}_{\rho}^{-1} - \rho(A)^{-1}\right\|_{2} \leqslant \left\|\hat{A}_{\rho}^{-1} - \rho(A)^{-1}\right\|_{F} < \frac{C\lambda_{B}\sqrt{|A^{-1}|_{0,off} \vee 1}}{2\varphi_{\min}^{2}(\rho(A))},\tag{S43}$$

$$\left\|\widehat{B}_{\rho}^{-1} - \rho(B)^{-1}\right\|_{2} \leq \left\|\widehat{B}_{\rho}^{-1} - \rho(B)^{-1}\right\|_{F} \leq \frac{C\lambda_{A}\sqrt{|B^{-1}|_{0,off} \vee 1}}{2\varphi_{\min}^{2}(\rho(B))}.$$
 (S44)

We now state an auxiliary result, Lemma S10, where we prove a bound on the error in the diagonal entries of the covariance matrices, and on their reciprocals. The following Lemma provides bounds analogous to those in Claim 15.1 Zhou (2014a,b).

**Lemma S10.** Let  $\widehat{W}_1$  and  $\widehat{W}_2$  be as defined in (10). Let  $W_1 = \sqrt{\operatorname{tr}(B)} \operatorname{diag}(A)^{1/2}$  and  $W_2 = \sqrt{\operatorname{tr}(A)} \operatorname{diag}(B)^{1/2}$ . Suppose event  $\mathcal{X}_0$  holds, as defined in (S26), (S27). For  $\eta' := \frac{\tilde{\eta}}{\sqrt{1-\tilde{\eta}}} \leq \frac{\lambda_B}{6}$  and  $\alpha' := \frac{\tilde{\alpha}}{\sqrt{1-\tilde{\alpha}}} \leq \frac{\lambda_A}{6}$ ,

$$\begin{aligned} \left\|\widehat{W}_{1}-W_{1}\right\|_{2} &\leqslant \widetilde{\eta}\sqrt{\operatorname{tr}\left(B\right)}\sqrt{a_{\max}}, \qquad \left\|\widehat{W}_{1}^{-1}-W_{1}^{-1}\right\|_{2} &\leqslant \frac{\widetilde{\eta}}{1-\widetilde{\eta}}/\sqrt{\operatorname{tr}\left(B\right)}\sqrt{a_{\min}}, \\ \left\|\widehat{W}_{2}-W_{2}\right\|_{2} &\leqslant \widetilde{\alpha}\sqrt{\operatorname{tr}\left(A\right)}\sqrt{b_{\max}}, \quad and \quad \left\|\widehat{W}_{2}^{-1}-W_{2}^{-1}\right\|_{2} &\leqslant \frac{\widetilde{\alpha}}{1-\widetilde{\alpha}}/\sqrt{\operatorname{tr}\left(A\right)}\sqrt{b_{\min}}. \end{aligned}$$

**Proposition S11.** (Zhou, 2014a). Let  $\widehat{W}$  and W be diagonal positive definite matrices. Let

 $\hat{\Psi}$  and  $\Psi$  be symmetric positive definite matrices. Then

$$\begin{split} \left\|\widehat{W}\widehat{\Psi}\widehat{W} - W\Psi W\right\|_{2} &\leq \left(\left\|\widehat{W} - W\right\|_{2} + \|W\|_{2}\right)^{2} \left\|\widehat{\Psi} - \Psi\right\|_{2} \\ &+ \left\|\widehat{W} - W\right\|_{2} \left(\left\|\widehat{W} - W\right\|_{2} + 2\right) \|\Psi\|_{2} \\ \left\|\widehat{W}\widehat{\Psi}\widehat{W} - W\Psi W\right\|_{F} &\leq \left(\left\|\widehat{W} - W\right\|_{2} + \|W\|_{2}\right)^{2} \left\|\widehat{\Psi} - \Psi\right\|_{F} \\ &+ \left\|\widehat{W} - W\right\|_{2} \left(\left\|\widehat{W} - W\right\|_{2} + 2\right) \|\Psi\|_{F}. \end{split}$$

*Proof* of Lemma S7. Assume that event  $\mathcal{X}_0$  holds. The proof follows exactly that of Lemma 15.3 in Zhou (2014a,b), in view of Theorem S9, Lemma S10 and Proposition 15.2 from Zhou (2014a,b), which is restated immediately above in Proposition S11.  $\Box$ 

It remains to prove Lemma S10.

*Proof* of Lemma S10. Suppose that event  $\mathcal{X}_0$  holds. Then

$$\max_{i=1,\dots,m} \left| \frac{\sqrt{X_i^T (I - P_2) X_i}}{\sqrt{a_{ii} \operatorname{tr}(B)}} - 1 \right| \leq \left( 1 - \sqrt{1 - \widetilde{\eta}} \right) \bigvee \left( \sqrt{1 + \widetilde{\eta}} - 1 \right) \leq \widetilde{\eta}.$$

Thus for all i,

$$\frac{1}{\sqrt{1+\widetilde{\eta}}} \leqslant \frac{\sqrt{a_{ii}\operatorname{tr}(B)}}{\sqrt{X_i^T(I-P_2)X_i}} \leqslant \frac{1}{\sqrt{1-\widetilde{\eta}}},$$

 $\mathbf{SO}$ 

$$\left|\frac{\sqrt{a_{ii}\operatorname{tr}(B)}}{\sqrt{X_i^T(I-P_2)X_i}} - 1\right| \leqslant \left(\frac{1-\sqrt{1-\widetilde{\eta}}}{\sqrt{1-\widetilde{\eta}}}\right) \bigvee \left(\frac{\sqrt{1+\widetilde{\eta}}-1}{\sqrt{1+\widetilde{\eta}}}\right) \leqslant \frac{\widetilde{\eta}}{\sqrt{1-\widetilde{\eta}}}$$

#### F.2 Proof of Lemma S8

In order to prove Lemma S8, we state Lemma S12, Lemma S13, and Proposition S14. Let  $\|\cdot\|$  denote a matrix norm such that  $\|A \otimes B\| = \|A\| \|B\|$ . Let

$$\Delta := \widehat{W}_1 \widehat{A}_{\rho} \widehat{W}_1 \otimes \widehat{W}_2 \widehat{B}_{\rho} \widehat{W}_2 / \operatorname{tr}(A) \operatorname{tr}(B) - A \otimes B, \qquad (S45)$$

$$\Delta' := \operatorname{tr}(A) \operatorname{tr}(B) \left(\widehat{W}_1 \widehat{A}_\rho \widehat{W}_1\right)^{-1} \otimes \left(\widehat{W}_2 \widehat{B}_\rho \widehat{W}_2\right)^{-1} - A^{-1} \otimes B^{-1}.$$
(S46)

Lemma S12 is identical to Lemma 15.5 of Zhou (2014a), except that we now plug in quantities  $\tilde{\alpha}$  and  $\tilde{\eta}$  as defined in (S17). Likewise, Proposition S14 is analogous to (20) in Theorem 4.1 of Zhou (2014a), except that we now use the centered data matrix  $(I - P_2)X$ , together with the rates  $\tilde{\alpha}$ ,  $\tilde{\eta}$ .

**Lemma S12.** Let  $\widehat{A \otimes B}$  be as in (11). Then for  $\Sigma = A \otimes B$ ,

$$\left\|\widehat{A\otimes B}^{-1} - \Sigma^{-1}\right\| \leq (\widetilde{\alpha} \wedge \widetilde{\eta}) \|A^{-1}\| \|B^{-1}\| + (1 + \widetilde{\alpha} \wedge \widetilde{\eta}) \|\Delta'\|$$
(S47)

$$\left\|\widehat{A\otimes B} - \Sigma\right\| \leq \frac{\lambda_A \wedge \lambda_B}{2} \|A\| \|B\| + (1 + \frac{\lambda_A \wedge \lambda_B}{2}) \|\Delta\|.$$
(S48)

Lemma S13 is a helpful bound on the difference of Kronecker products.

**Lemma S13.** (*Zhou*, 2014*a*). For matrices  $A_1$  and  $B_1$ , let  $\Delta_A := A_1 - A$  and  $\Delta_B := B_1 - B$ . Then

$$||A_1 \otimes B_1 - A \otimes B|| \le ||\Delta_A|| ||B|| + ||\Delta_B|| ||A|| + ||\Delta_A|| ||\Delta_B||.$$

**Proposition S14.** Under the event  $\mathcal{X}_0$ , as defined in as defined in (S26), (S27),

$$\left| \| (I - P_2) X \|_F^2 - \operatorname{tr}(A) \operatorname{tr}(B) \right| \leq (\widetilde{\alpha} \wedge \widetilde{\eta}) \operatorname{tr}(A) \operatorname{tr}(B).$$

*Proof* of Lemma S8. Assume that event  $\mathcal{X}_0$  as defined in (S26), (S27) holds. The proof follows exactly the steps in Theorems 11.1 and 11.2 in Supplementary Material of Zhou (2014a,b).  $\Box$ 

*Proof* of Lemma S12. By the triangle inequality and the sub-multiplicativity of the norm  $\|\cdot\|$ , with  $\Delta$  and  $\Delta'$  as defined in (S45) and (S46),

$$\operatorname{tr}(A)\operatorname{tr}(B)\left\|\left(\widehat{W}_{1}^{-1}\widehat{A}_{\rho}^{-1}\widehat{W}_{1}^{-1}\right)\otimes\left(\widehat{W}_{2}^{-1}\widehat{B}_{\rho}^{-1}\widehat{W}_{2}^{-1}\right)\right\| \leq \|A^{-1}\|\|B^{-1}\| + \|\Delta'\| \tag{S49}$$

$$\left\| \left( \widehat{W}_1 \widehat{A}_{\rho} \widehat{W}_1 \right) \otimes \left( \widehat{W}_2 \widehat{B}_{\rho} \widehat{W}_2 \right) / \operatorname{tr}(A) \operatorname{tr}(B) \right\| \leq \|A\| \|B\| + \|\Delta\|.$$
(S50)

Following proof of Lemma 15.5 Zhou (2014a,b), we have by definition of  $\Delta'$ , and Proposition S14, and (S49),

$$\left\|\widehat{A \otimes B}^{-1} - A^{-1} \otimes B^{-1}\right\| \leq (\widetilde{\alpha} \wedge \widetilde{\eta}) \left(\|A^{-1}\|\|B^{-1}\| + \|\Delta'\|\right) + \|\Delta'\|.$$

By Proposition S14, we have for  $\lambda_A \ge 3\widetilde{\alpha}$ ,  $\lambda_B \ge 3\widetilde{\eta}$ , where  $\widetilde{\alpha} \wedge \widetilde{\eta} \le \frac{\lambda_A \wedge \lambda_B}{3}$ ,

$$\left|\frac{1}{\|(I-P_2)X\|_F^2} - \frac{1}{\operatorname{tr}(A)\operatorname{tr}(B)}\right| = \left|\frac{\|(I-P_2)X\|_F^2 - \operatorname{tr}(A)\operatorname{tr}(B)}{\|(I-P_2)X\|_F^2\operatorname{tr}(A)\operatorname{tr}(B)}\right|$$
  

$$\leq \left|\frac{\widetilde{\alpha} \wedge \widetilde{\eta}}{\|(I-P_2)X\|_F^2}\right| \leq \frac{\widetilde{\alpha} \wedge \widetilde{\eta}}{\operatorname{tr}(A)\operatorname{tr}(B)(1-\widetilde{\alpha} \wedge \widetilde{\eta})}$$
  
thus  $\left|\frac{\operatorname{tr}(A)\operatorname{tr}(B)}{\|(I-P_2)X\|_F^2} - 1\right| \leq \frac{\widetilde{\alpha} \wedge \widetilde{\eta}}{1-\widetilde{\alpha} \wedge \widetilde{\eta}} \leq \frac{\lambda_A \wedge \lambda_B}{2}.$  (S51)

By the triangle inequality, the definition of  $\Delta$  in (S45), and (S50) and (S51),

$$\left\|\widehat{A\otimes B} - A\otimes B\right\| \leq \frac{\lambda_A + \lambda_B}{2} \|A\| \|B\| + (1 + \frac{\lambda_A + \lambda_B}{2}) \|\Delta\|;$$

See the proof of Lemma 15.5 Zhou (2014a,b).  $\Box$ 

*Proof* of Proposition S14. Suppose event  $\mathcal{X}_0$  holds. Note that

$$E[\|(I - P_2)X\|_F^2] = \operatorname{tr}((I - P_2)E[XX^T](I - P_2)) = \operatorname{tr}(A)\operatorname{tr}(\widetilde{B})$$

Decomposing by columns, we obtain the inequality,

$$\left| \| (I - P_2) X \|_F^2 - \operatorname{tr}(A) \operatorname{tr}(B) \right| = \left| \sum_{j=1}^m \| (I - P_2) X_j \|_2^2 - a_{jj} \operatorname{tr}(B) \right|$$
  
$$\leqslant \sum_{j=1}^m \left| X_j^T (I - P_2) X_j - a_{jj} \operatorname{tr}(B) \right| \leqslant \sum_{j=1}^m \widetilde{\eta}_{jj} a_{jj} \operatorname{tr}(B) \leqslant \widetilde{\eta} \operatorname{tr}(A) \operatorname{tr}(B).$$

Decomposing by rows, we obtain the inequality,

$$\left| \| (I - P_2) X \|_F^2 - \operatorname{tr}(A) \operatorname{tr}(B) \right| = \left| \sum_{i=1}^n \| (e_i - p_i)^T X \|_2^2 - b_{ii} \operatorname{tr}(A) \right|$$
  
$$\leq \sum_{i=1}^n \left| (e_i - p_i)^T X X^T (e_i - p_i) - b_{ii} \operatorname{tr}(A) \right| \leq \sum_{i=1}^n \widetilde{\alpha}_{ii} b_{ii} \operatorname{tr}(A) \leq \widetilde{\alpha} \operatorname{tr}(A) \operatorname{tr}(B).$$

Therefore  $|||(I - P_2)X||_F^2 - \operatorname{tr}(A)\operatorname{tr}(B)| \leq (\widetilde{\alpha} \wedge \widetilde{\eta})\operatorname{tr}(A)\operatorname{tr}(B).$ 

## G Entrywise convergence of sample correlations

In this section we prove entrywise rates of convergence for the sample correlation matrices in Theorem S15. The theorem applies to the Kronecker product model,  $\text{Cov}(\text{vec}(X)) = A^* \otimes B^*$ , where for identifiability we define the sample covariance matrices as

$$A^* = \frac{m}{\operatorname{tr}(A)}A$$
 and  $B^* = \frac{\operatorname{tr}(A)}{m}B$ ,

with the scaling chosen so that  $A^*$  has trace m. Let  $\rho(A) \in \mathbb{R}^{m \times m}$  and  $\rho(B) \in \mathbb{R}^{n \times n}$  denote the correlation matrices corresponding to covariance matrices  $A^*$  and  $B^*$ , respectively. Assume that that the mean of X satisfies the two-group model (4). Let  $P_2$  be as defined in (13). The matrix  $I - P_2$  is a projection matrix of rank n - 2 that performs within-group centering. The sample covariance matrices are defined as

$$S(B^*) = \frac{1}{m} \sum_{j=1}^{m} (I - P_2) X_j X_j^T (I - P_2), \qquad (S52)$$

$$S(A^*) = X^T (I - P_2) X/n,$$
 (S53)

where  $S(B^*)$  has null space of dimension two.

**Theorem S15.** Consider a data generating random matrix as in (2). Let C be some absolute constant. Let  $\tilde{\alpha}$  and  $\tilde{\eta}$  be as defined in (S17). Let  $m \lor n \ge 2$ . Then with probability at least  $1 - \frac{3}{(m \lor n)^2}$ , for  $\tilde{\alpha}, \tilde{\eta} < 1/3$ , and  $\hat{\Gamma}(A)$  and  $\hat{\Gamma}(B)$  as in (7),

$$\begin{aligned} \forall i \neq j, \ \left| \widehat{\Gamma}_{ij}(B) - \rho_{ij}(B) \right| &\leq \frac{\widetilde{\alpha}}{1 - \widetilde{\alpha}} + \left| \rho_{ij}(B) \right| \frac{\widetilde{\alpha}}{1 - \widetilde{\alpha}} \leq 3\widetilde{\alpha}, \\ \forall i \neq j, \ \left| \widehat{\Gamma}_{ij}(A) - \rho_{ij}(A) \right| &\leq \frac{\widetilde{\eta}}{1 - \widetilde{\eta}} + \left| \rho_{ij}(A) \right| \frac{\widetilde{\eta}}{1 - \widetilde{\eta}} \leq 3\widetilde{\eta}. \end{aligned}$$

We state three results used in the proof of Theorem S15: Proposition S16 provides an entrywise rate of convergence of  $S(B^*)$ , Proposition S17 provides an entrywise rate of convergence of  $S(A^*)$ , and Lemma S18 states that these entrywise rates imply  $\mathcal{X}_0$ . Let

$$\widetilde{B} := (I - P_2)B^*(I - P_2) = \operatorname{Cov}((I - P_2)X_j),$$
(S54)

where  $X_j$  is the *j*th column of X. Let  $\tilde{b}_{ij}$  denote the (i, j)th entry of  $\tilde{B}$ .

**Proposition S16.** Let d > 2. Then with probability at least  $1 - 2/m^{d-2}$ ,

$$\forall i, j \left| S_{ij}(B^*) - b_{ij}^* \right| \le \phi_{B,ij},\tag{S55}$$

with

$$\phi_{B,ij} = C \frac{\log^{1/2}(m)}{\sqrt{m}} \frac{\|A^*\|_F}{\sqrt{m}} \sqrt{\tilde{b}_{ii}\tilde{b}_{jj}} + \frac{3\|B^*\|_1}{n_{\min}}.$$
(S56)

**Proposition S17.** Let d > 2. Then with probability at least  $1 - 2/n^{d-2}$ ,

$$\forall i, j \left| S_{ij}(A^*) - a_{ij}^* \operatorname{tr}(B^*) / n \right| > \phi_{A,ij},$$
(S57)

with

$$\phi_{A,ij} = (a_{ij}^*/n) \left| \operatorname{tr}\left(\tilde{B}\right) - \operatorname{tr}\left(B^*\right) \right| + d^{1/2} K \log^{1/2} (n \vee m) (1/n) \sqrt{a_{ij}^{*2} + a_{ii}^* a_{jj}^*} \|\tilde{B}\|_F.$$
(S58)

**Lemma S18.** Suppose that (A2) holds and that  $m \lor n \ge 2$ . The event (S57) defined in Proposition S17 implies that  $\mathcal{X}_0(A)$  holds. Similarly, the event (S55) defined in Proposition S16 implies  $\mathcal{X}_0(B)$ . Hence  $\mathbb{P}(\mathcal{X}_0) \ge 1 - \frac{3}{(m \lor n)^2}$ .

Proposition S16 is proved in section G.1. Proposition S17 is proved in section G.2. Lemma S18 is proved in section G.3. Note that Lemma S18 follows from Propositions S16 and S17. We now prove Theorem S15, which follows from Lemma S18.

Proof of Theorem S15. Let  $q_i$  denote the *i*th column of  $I - P_2$ , so that  $q_i^T X X^T q_j$ is the (i, j)th entry of  $(I - P_2) X X^T (I - P_2)$ . Under  $\mathcal{X}_0(B)$ , the sample correlation  $\widehat{\Gamma}(B)$  satisfies the following bound:

$$\begin{aligned} \left| \widehat{\Gamma}_{ij}(B) - \rho_{ij}(B) \right| &= \left| \frac{q_i^T X X^T q_j}{\sqrt{q_i^T X X^T q_i} \sqrt{q_j^T X X^T q_j}} - \rho_{ij}(B) \right| \\ &= \left| \frac{q_i^T X X^T q_i / (\operatorname{tr}(A^*) \sqrt{b_{ii}^* b_{jj}^*})}{\sqrt{q_i^T X X^T q_i / (b_{ii}^* \operatorname{tr}(A^*))} \sqrt{q_j^T X X^T q_j / (b_{jj}^* \operatorname{tr}(A^*))}} - \rho_{ij}(B) \right| \\ &\leqslant \left| \frac{q_i^T X X^T q_i / (\operatorname{tr}(A^*) \sqrt{b_{ii}^* b_{jj}^*}) - \rho_{ij}(B)}{\sqrt{q_i^T X X^T q_i / (b_{ii}^* \operatorname{tr}(A^*))} \sqrt{q_j^T X X^T q_j / (b_{jj}^* \operatorname{tr}(A^*))}} \right| \\ &+ \left| \frac{\rho_{ij}(B)}{\sqrt{q_i^T X X^T q_i / (b_{ii}^* \operatorname{tr}(A^*))} \sqrt{q_j^T X X^T q_j / (b_{jj}^* \operatorname{tr}(A^*))}} - \rho_{ij}(B) \right| \\ &\leqslant \frac{\widetilde{\alpha}}{1 - \widetilde{\alpha}} + |\rho_{ij}(B)| \left| \frac{1}{1 - \widetilde{\alpha}} - 1 \right| \\ &\leqslant 3\widetilde{\alpha}, \end{aligned}$$

where the first inequality holds by  $\mathcal{X}_0(B)$  and the second inequality holds for  $\tilde{\alpha} \leq 1/3$ . Similarly, under  $\mathcal{X}_0(A)$  we obtain an entrywise bound on the sample correlation  $\widehat{\Gamma}(A)$ :

$$\begin{aligned} \left| \widehat{\Gamma}_{ij}(A) - \rho_{ij}(A) \right| &= \left| \frac{X_i^T(I - P_2)X_j}{\sqrt{X_i^T(I - P_2)X_i}\sqrt{X_j^T(I - P_2)X_j}} - \rho_{ij}(A) \right| \\ &= \left| \frac{X_i^T(I - P_2)X_j/\left(\operatorname{tr}(B^*)\sqrt{a_{ii}^*a_{jj}^*}\right)}{\sqrt{X_i^T(I - P_2)X_j/\left(a_{ij}^*\operatorname{tr}(B^*)\right)}\sqrt{X_j^T(I - P_2)X_j/\left(a_{jj}^*\operatorname{tr}(B^*)\right)}} - \rho_{ij}(A) \right| \\ &\leqslant \left| \frac{X_i^T(I - P_2)X_i/\left(a_{ii}^*\operatorname{tr}(B^*)\right)\sqrt{X_j^T(I - P_2)X_j/\left(a_{jj}^*\operatorname{tr}(B^*)\right)}}{\sqrt{X_i^T(I - P_2)X_i/\left(a_{ii}^*\operatorname{tr}(B^*)\right)}\sqrt{X_j^T(I - P_2)X_j/\left(a_{jj}^*\operatorname{tr}(B^*)\right)}} \right| \\ &+ \left| \frac{\rho_{ij}(A)}{\sqrt{X_i^T(I - P_2)X_i/\left(a_{ii}^*\operatorname{tr}(B^*)\right)}\sqrt{X_j^T(I - P_2)X_j/\left(a_{jj}^*\operatorname{tr}(B^*)\right)}} - \rho_{ij}(A) \right| \\ &\leqslant \left| \frac{\widetilde{\eta}}{1 - \widetilde{\eta}} + |\rho_{ij}(A)| \left| \frac{1}{1 - \widetilde{\eta}} - 1 \right| \leqslant 3\widetilde{\eta}, \end{aligned}$$

where the first inequality holds by  $\mathcal{X}_0(A)$ , and the second inequality holds for  $\tilde{\eta} < 1/3$ .

By Lemma S18, the event  $\mathcal{X}_0 = \mathcal{X}_0(B) \cap \mathcal{X}_0(A)$  holds with probability at least  $1 - 3/(n \vee m)^2$ , which completes the proof.  $\Box$ 

#### G.1 Proof of Proposition S16

We first present Lemma S19 and Lemma S20, which decompose the rate of convergence into a bias term and a variance term, respectively. We then combine the rates for the bias and variance terms to prove the entrywise rate of convergence for the sample covariance. Define

$$\mathcal{B}(B^*) := E[S(B^*)] - B^* \quad \text{and} \tag{S59}$$

$$\sigma(B^*) := S(B^*) - E[S(B^*)].$$
(S60)

We state maximum entrywise bounds on  $\mathcal{B}(B^*)$  and  $\sigma(B^*)$  in Lemma S19 and Lemma S20, respectively. Proofs for these lemmas are provided in Section G.4 and G.5 respectively.

**Lemma S19.** For  $\mathcal{B}(B^*)$  as defined in (S59),

$$\|\mathcal{B}(B^*)\|_{\max} \leq \frac{3\|B^*\|_1}{n_{\min}}.$$
 (S61)

**Lemma S20.** Let  $\sigma(B^*)$  be as defined in (S60). With probability at least  $1 - 2/m^d$ ,

$$|\sigma_{ij}(B^*)| = \left| S_{ij}(B^*) - b_{ij}^* \right| < C \log^{1/2}(m) \frac{\|A^*\|_F}{\operatorname{tr}(A^*)} \sqrt{\widetilde{b}_{ii}\widetilde{b}_{jj}}.$$

We now prove the entrywise rate of convergence for the sample covariance  $S(B^*)$ .

*Proof* of Proposition S16. By the triangle inequality,

$$\begin{aligned} |S_{ij}(B^*) - b_{ij}^*| &\leq |S_{ij}(B^*) - E[S_{ij}(B^*)]| + |E[S_{ij}(B^*)] - b_{ij}^*| \\ &= |\mathcal{B}_{ij}(B^*)| + |\sigma_{ij}(B^*)| \\ &\leq \phi_{B,ij}, \end{aligned}$$

where the last step follows from Lemmas S19 and S20.  $\Box$ 

**Remark.** Note that the first term of (S56) is of order  $\log^{1/2}(m)/\sqrt{m}$ , and the second term is of order  $||B^*||_1/n_{\min}$ .

#### G.2 Proof of Proposition S17

We express the (i, j)th entry of  $S(A^*)$  as a quadratic form in order to apply the Hanson-Wright inequality to obtain an entrywise large deviation bound. Without loss of generality, let i = 1, j = 2. The (1, 2) entry of  $S(A^*)$  can be expressed as a quadratic form, as follows,

$$S_{12}(A^*) = X_1^T (I - P_2) X_2 / n$$
  
=  $(1/2) \begin{bmatrix} X_1^T & X_2^T \end{bmatrix} \begin{bmatrix} 0 & (I - P_2) \\ (I - P_2) & 0 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} / n$   
=  $(1/2) \begin{bmatrix} X_1^T & X_2^T \end{bmatrix} \left( \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \otimes (I - P_2) \right) \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} / n$ 

We decorrelate the random vector  $(X_1, X_2) \in \mathbb{R}^{2n}$  so that we can apply the Hanson-Wright inequality. The covariance matrix used for decorrelation is

$$\operatorname{Cov}\left(\begin{bmatrix} X_1 \\ X_2 \end{bmatrix}\right) = \begin{bmatrix} a_{11}^* & a_{12}^* \\ a_{21}^* & a_{22}^* \end{bmatrix} \otimes B^* =: A_{\{1,2\}}^* \otimes B^*,$$

with

$$A_{\{1,2\}}^* = \begin{bmatrix} a_{11}^* & a_{12}^* \\ a_{21}^* & a_{22}^* \end{bmatrix} \in \mathbb{R}^{2 \times 2}.$$

Decorrelating the quadratic form yields

$$S_{12}(A^*) = Z^T \Phi Z,$$

where  $Z \in \mathbb{R}^{2n}$ , with E[Z] = 0 and  $Cov(Z) = I_{2n \times 2n}$ , and

$$\Phi = (1/2n) \left( (A_{\{1,2\}}^*)^{1/2} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} (A_{\{1,2\}}^*)^{1/2} \right) \otimes B^{1/2} (I - P_2) B^{1/2}.$$
(S62)

To apply the Hanson-Wright inequality, we first find the trace and Frobenius norm of  $\Phi$ .

For the trace, note that

$$\operatorname{tr}\left(\left(A_{\{1,2\}}^{*}\right)^{1/2}\begin{bmatrix}0&1\\1&0\end{bmatrix}\left(A_{\{1,2\}}^{*}\right)^{1/2}\right) = \operatorname{tr}\left(\begin{bmatrix}0&1\\1&0\end{bmatrix}A_{\{1,2\}}^{*}\right) = 2a_{12}^{*}.$$
 (S63)

For the Frobenius norm, note that

$$\left\| (A_{\{1,2\}}^*)^{1/2} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} (A_{\{1,2\}}^*)^{1/2} \right\|_F^2 = \operatorname{tr} \left( \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} A_{\{1,2\}}^* \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} A_{\{1,2\}}^* \right)$$
$$= \operatorname{tr} \left( \begin{bmatrix} a_{12}^{*2} + a_{11}^* a_{22}^* & 2a_{12}^* a_{22}^* \\ 2a_{12}^* a_{22}^* & a_{12}^{*2} + a_{11}^* a_{22}^* \end{bmatrix} \right)$$
$$= 2a_{12}^{*2} + 2a_{11}^* a_{22}^*,$$

Therefore the trace of  $\Phi$  is

$$\operatorname{tr}\left(\Phi\right) = a_{12}^{*}\operatorname{tr}\left(\widetilde{B}\right)/n,\tag{S64}$$

and the Frobenius norm of  $\Phi$  is

$$\|\Phi\|_F = (1/n)\sqrt{a_{12}^{*2} + a_{11}^* a_{22}^*} \|\widetilde{B}\|_F.$$
(S65)

Applying the Hanson-Wright inequality yields

$$P\left(|S_{12}(A^*) - a_{12}^* \operatorname{tr}(B^*)/n| > \phi_{A,12}\right)$$
  

$$\leq P\left(\left|S_{12}(A^*) - a_{12}^* \operatorname{tr}\left(\widetilde{B}\right)/n\right| + (a_{12}^*/n) \left|\operatorname{tr}\left(\widetilde{B}\right) - \operatorname{tr}(B^*)\right| > \phi_{A,12}\right)$$
  

$$= P\left(\left|S_{12}(A) - a_{12}^* \operatorname{tr}\left(\widetilde{B}\right)/n\right| > d^{1/2} K \log^{1/2} (n \lor m) \|\Phi\|_F\right)$$
  

$$\leq 2/(n \lor m)^d.$$

By the union bound,

$$P(\forall i, j | S_{ij}(A^*) - a_{ij} \operatorname{tr} (B^*) / n | < \phi_{A,ij})$$
  

$$\geq 1 - \sum_{i=1}^{m} \sum_{j=1}^{m} P(|S_{ij}(A^*) - a_{ij} \operatorname{tr} (B^*) / n | > \phi_{A,ij})$$
  

$$\geq 1 - 2m^2 / (n \lor m)^d \geq 2 / (n \lor m)^{d-2}.$$

## G.3 Proof of Lemma S18

For the event (S55) from Proposition S16,

$$\left|S_{ij}(B^*) - b_{ij}^*\right| < \phi_{B,ij} = K^2 d \frac{\log^{1/2}(m)}{\sqrt{m}} C_A \sqrt{\widetilde{b}_{ii}\widetilde{b}_{jj}} + \left|b_{ij}^* - \widetilde{b}_{ij}\right|,$$

dividing by  $\sqrt{b_{ii}^* b_{jj}^*}$  yields

$$\left|\frac{q_i X X^T q_j}{\operatorname{tr}(A^*) \sqrt{b_{ii}^* b_{jj}^*}} - \rho_{ij}(B)\right| < K^2 dC_A \frac{\log^{1/2}(m)}{\sqrt{m}} \sqrt{\frac{\widetilde{b}_{ii} \widetilde{b}_{jj}}{b_{ii}^* b_{jj}^*}} + \frac{\left|b_{ij} - \widetilde{b}_{ij}\right|}{\sqrt{b_{ii}^* b_{jj}^*}}.$$
(S66)

By Lemma S19,

$$\widetilde{b}_{ij} = b_{ij} \left[ 1 + O\left(\frac{\|B\|_1}{n}\right) \right],$$

so the right-hand side of (S66) is less than or equal to  $\tilde{\alpha}$ . Hence event (S55) implies  $\mathcal{X}_0(B)$ . Therefore, we know that  $P(\mathcal{X}_0(B)) \ge 1 - 2/m^{d-2}$ .

Similarly, event (S57) in Proposition S17:

$$\begin{aligned} \left| S_{ij}(A^*) - a_{ij}^* \operatorname{tr}(B^*) / n \right| &< \phi_{A,ij} \\ &= (a_{ij}^* / n) \left| \operatorname{tr}\left(\widetilde{B}\right) - \operatorname{tr}(B) \right| + d^{1/2} K \log^{1/2} (n \lor m) (1/n) \sqrt{a_{ij}^{*2} + a_{ii}^* a_{jj}^*} \|\widetilde{B}\|_F, \end{aligned}$$

implies that

$$\begin{aligned} \left| \frac{X_{j}^{T}(I-P_{2})X_{t}}{\operatorname{tr}(B^{*})\sqrt{a_{jj}^{*}a_{tt}^{*}}} - \rho_{jt}(A) \right| \\ &< |\rho_{jt}(A)| \frac{\left| \operatorname{tr}\left(\widetilde{B}\right) - \operatorname{tr}\left(B^{*}\right) \right|}{\operatorname{tr}(B^{*})} + d^{1/2}K \log^{1/2}(n \vee m)\sqrt{\rho_{jt}(A)^{2} + 1} \frac{\|\widetilde{B}\|_{F}}{\operatorname{tr}(B^{*})} \\ &= |\rho_{jt}(A)| \frac{\left| \operatorname{tr}\left(\widetilde{B}\right) - \operatorname{tr}\left(B^{*}\right) \right|}{\operatorname{tr}(B^{*})} + d^{1/2}K C_{B} \frac{\|\widetilde{B}\|_{F}}{\|B^{*}\|_{F}} \sqrt{\rho_{jt}(A)^{2} + 1} \frac{\log^{1/2}(n \vee m)}{\sqrt{n}} \\ &\leqslant \widetilde{\eta}, \end{aligned}$$

which is the event  $\mathcal{X}_0(A)$ . Therefore, we get that  $P(\mathcal{X}_0(A)) \ge 1 - 2/(n \lor m)^d$ .

We can obtain the  $P(\mathcal{X}_0)$  by using a union bound put together  $P(\mathcal{X}_0(B))$  and  $P(\mathcal{X}_0(A))$ , completing the proof.  $\Box$ 

#### G.4 Proof of Lemma S19

Recall that  $\widetilde{B} = (I - P_2)B^*(I - P_2)$ . The matrix  $\widetilde{B} - B^*$  can be expressed as

$$\widetilde{B} - B^* = (I - P_2)B^*(I - P_2) - B^* = -P_2B^* - B^*P_2 + P_2B^*P_2.$$

By the triangle inequality,  $\|\tilde{B} - B^*\|_{\max} \leq \|P_2B^*\|_{\max} + \|B^*P_2\|_{\max} + \|P_2B^*P_2\|_{\max}$ . We bound each term on the right-hand side.

First we bound  $||P_2B^*||_{\text{max}}$  and  $||B^*P_2||_{\text{max}}$ . Let  $p_i$  denote the *i*th column of  $P_2$ . The (i, j)th entry satisfies

$$|p_i^T b_j^*| \le ||B^* p_i||_{\infty} \le ||B^*||_{\infty} ||p_i||_{\infty} = ||B^*||_1 ||p_i||_{\infty} = ||B^*||_1 / n_{\min},$$

so  $||P_2B^*||_{\max} \leq ||B^*||_1/n_{\min}$ . Because  $P_2$  and  $B^*$  are symmetric,  $||P_2B^*||_{\max} = ||B^*P_2||_{\max}$ .

We now bound  $||P_2B^*P_2||_{\text{max}}$ . Let  $B^{1/2}$  denote the symmetric square root of  $B^*$ . We can express  $p_i^T B^* p_j$  as an inner product  $(B^{1/2} p_i)^T (B^{1/2} p_j)$ , so

$$|(P_2B^*P_2)_{ij}| = |(B^{1/2}p_i)^T(B^{1/2}p_j)| \le (p_i^TB^*p_i)^{1/2} (p_j^TB^*p_j)^{1/2}$$
(S67)

$$\leq \|p_i\|_2 \|p_j\|_2 \|B\|_2 \leq \|B^*\|_2 / n_{\min}, \tag{S68}$$

where (S67) follows from the Cauchy Schwarz inequality, and (S68) holds because

$$||p_i||_2 = \begin{cases} 1/\sqrt{n_1} & \text{if } i \in \{1, \dots, n_1\} \\ 1/\sqrt{n_2} & \text{if } i \in \{n_1 + 1, \dots, n\}. \end{cases}$$

## G.5 Proof of Lemma S20

Let  $B^{1/2}$  denote the symmetric square root of  $B^*$ . Let  $Z_j = (a_{jj}^* B^*)^{-1/2} X_j$ . We express  $S_{ij}(B^*)$  as a quadratic form in order to use the Hanson-Wright inequality to prove a large deviation bound. That is, we show that  $S_{ij}(B^*) = \operatorname{vec}(Z)^T \Phi^{ij} \operatorname{vec}(Z)$ , with

$$\Phi^{ij} = (1/m)A^* \otimes B^{1/2}(e_j - p_j)(e_i - p_i)^T B^{1/2}.$$
 (S69)

We express  $S_{ij}(B^*)$  as a quadratic form, as follows:

$$S_{ij}(B^*) = \frac{1}{m} \sum_{k=1}^m (e_i - p_i)^T X_k X_k^T (e_j - p_j) = \frac{1}{m} \sum_{k=1}^m \operatorname{tr} \left[ (e_i - p_i)^T X_k X_k^T (e_j - p_j) \right]$$
  
$$= \frac{1}{m} \sum_{k=1}^m X_k^T (e_j - p_j) (e_i - p_i)^T X_k$$
  
$$= \frac{1}{m} \operatorname{vec}(X)^T \left( I_{m \times m} \otimes (e_j - p_j) (e_i - p_i)^T \right) \operatorname{vec}(X)$$
  
$$= \operatorname{vec}(Z)^T \Phi^{ij} \operatorname{vec}(Z)$$

where

$$\operatorname{tr}(\Phi^{ij}) = \operatorname{tr}(B^{1/2}(e_j - p_j)(e_i - p_i)^T B^{1/2}) = (e_i - p_i)^T B^*(e_j - p_j) = \widetilde{b}_{ij},$$
(S70)

$$\|\Phi^{ij}\|_F = \frac{1}{m} \|A^*\|_F \|B^{1/2}(e_j - p_j)(e_i - p_i)^T B^{1/2}\|_F$$
(S71)

$$= \frac{1}{m} \|A^*\|_F \left( (e_i - p_i)^T B^* (e_i - p_i) \right)^{1/2} \left( (e_j - p_j)^T B^* (e_j - p_j) \right)^{1/2} = \frac{1}{m} \|A^*\|_F \sqrt{\tilde{b}_{ii} \tilde{b}_{jj}}.$$

Therefore, we get that

$$\begin{split} P\left(\forall i, j \ \left| S_{ij}(B^*) - \widetilde{b}_{ij} \right| &\leq K^2 d \log^{1/2}(m) \|\Phi^{ij}\|_F / c' \right) \\ &= P\left( \forall i, j \ \left| \operatorname{vec}(Z)^T \Phi^{ij} \operatorname{vec}(Z) - \operatorname{tr}\left(\Phi^{ij}\right) \right| &\leq K^2 d \log^{1/2}(m) \|\Phi^{ij}\|_F / c' \right) \\ &\geqslant 1 - 2m^2 \exp\left( -c \min\left( d^2 \log(m) / c'^2, \frac{d \log^{1/2}(m) \|\Phi^{ij}\|_F / c'}{\|\Phi^{ij}\|_2} \right) \right) \\ &\geqslant 1 - 2/m^{d-2}. \end{split}$$

If the event  $\left\{ \forall i, j \left| S_{ij}(B^*) - \widetilde{b}_{ij} \right| \leq K^2 d \log^{1/2}(m) \|\Phi^{ij}\|_F / c' \right\}$  holds, it follows that

$$\left|S_{ij}(B^*) - b_{ij}^*\right| \le \left|S_{ij}(B^*) - \widetilde{b}_{ij}\right| + |b_{ij}^* - \widetilde{b}_{ij}| \le K^2 d \log^{1/2}(m) \|\Phi^{ij}\|_F / c' + |b_{ij} - \widetilde{b}_{ij}|.$$

The Lemma is thus proved.  $\hfill\square$ 

# H Proof of Theorem 4

# H.1 Notation

Notation	Meaning			
Mean structure				
$\mu \in \mathbb{R}^m$	Vector of grand means of each gene			
$\gamma \in \mathbb{R}^m$	Vector of mean differences for each gene			
$\nu = \frac{1}{2} \begin{bmatrix} \frac{1}{n_1} 1_{n_1}^T & \frac{1}{n_2} 1_{n_2}^T \end{bmatrix}^T \in \mathbb{R}^n$	Inner product with $\nu$ computes global mean			
Outcome of model selection step				
$J_0 \subset \{1, 2, \dots, m\}$	Indices selected for group centering			
$J_1 \subset \{1, 2, \dots, m\}$	Indices selected for global centering			
Sizes of gene subsets				
$m_0 =  J_0 $	Number of group centered genes			
$m_1 =  J_1 $	Number of globally centered genes			
Projection matrices				
$P_1 = 1_n \nu^T$	Projection matrix that performs global centering			
$P_2$ (as in (S81))	Projection matrix that performs group centering			
Sample covariance matrices				
$S(B, J_0, J_1) = \frac{m_1}{m} S_1(B) + \frac{m_0}{m} S_2(B)$	Model selection sample covariance matrix			
$S_1(B, J_1) = \frac{1}{m_1} \sum_{j \in J_1} (I - P_1) X_j X_j^T (I - P_1)$	Globally centered sample covariance matrix			
$S_2(B, J_0) = \frac{1}{m_0} \sum_{j \in J_0} (I - P_2) X_j X_j^T (I - P_2)$	Group centered sample covariance matrix			
Decomposition of $S(B, J_0, J_1)$				
$S_{\rm I} = S(B, J_0, J_1) - \mathbb{E} \left[ S(B, J_0, J_1) \right]$	Bias			
$S_{\rm II} = \frac{1}{m} (I - P_1) M_{J_1} M_{J_1}^T (I - P_1)$	False negatives (deterministic)			
$S_{\text{III}} = \frac{1}{m}(I - P_1)M_{J_1}\varepsilon^T(I - P_1)$	False negatives (random)			
$S_{\rm IV} = m^{-1}(I - P_2)\varepsilon_{J_0}\varepsilon_{J_0}^T(I - P_2) +$	True negatives			
$m^{-1}(I-P_1)\varepsilon_{J_1}\varepsilon_{J_1}^T(I-P_1)$				

#### H.2 Two-Group Model and Centering

We begin by introducing some relevant notation for the two-group model and centering. Define the group membership vector  $\delta_n \in \mathbb{R}^n$  as

$$\delta_n := \begin{bmatrix} 1_{n_1}^T & -1_{n_2}^T \end{bmatrix}^T \in \mathbb{R}^n.$$
(S72)

In the two-group model, the mean matrix M can be expressed as

$$M = 1_n \mu^T + (1/2)\delta_n \gamma^T,$$
 (S73)

where  $\mu \in \mathbb{R}^m$  is a vector of grand means, and  $\gamma \in \mathbb{R}^m$  is the vector of mean differences. According to (S73), the (i, j)th entry of M can be expressed as

$$m_{ij} = \begin{cases} \mu_j + \gamma_j/2 & \text{if sample } i \text{ is in group one} \\ \mu_j - \gamma_j/2 & \text{if sample } i \text{ is in group two.} \end{cases}$$
(S74)

Define the vector  $\nu \in \mathbb{R}^n$  as

$$\nu = \frac{1}{2} \begin{bmatrix} \frac{1}{n_1} 1_{n_1}^T & \frac{1}{n_2} 1_{n_2}^T \end{bmatrix}^T \in \mathbb{R}^n,$$
(S75)

so that for the *j*th column of the data matrix  $X_j \in \mathbb{R}^n$ ,

$$\mathbb{E}\left(\nu^{T}X_{j}\right) = \frac{1}{2}\mathbb{E}\left(\frac{1}{n_{1}}\sum_{k=1}^{n_{1}}X_{jk} + \frac{1}{n_{2}}\sum_{k=n_{1}+1}^{n}X_{jk}\right) = \mu_{j}.$$
(S76)

Note that

$$\nu^T \mathbf{1}_n = (1/2)(1+1) = 1, \text{ and } \nu^T \delta_n = (1/2)(1-1) = 0.$$
 (S77)

Next we define a projection matrix that performs global centering. Define the non-orthogonal projection matrix

$$P_1 := \mathbf{1}_n \nu^T \in \mathbb{R}^{n \times n}.$$
(S78)

Applying the projection matrix to the mean matrix yields

$$P_1 M = 1_n \nu^T \left( 1_n \mu^T + (1/2)\delta_n \gamma^T \right) = 1_n \mu^T + (1/2)(\nu^T \delta_n) 1_n \gamma^T = 1_n \mu^T,$$
(S79)

with residuals

$$(I - P_1)M = M - P_1M = M - 1_n\mu^T = (1/2)\delta_n\gamma^T.$$
 (S80)

Define

$$P_2 = \begin{bmatrix} n_1^{-1} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^T & \\ & n_2^{-1} \mathbf{1}_{n_2} \mathbf{1}_{n_2}^T \end{bmatrix}.$$
 (S81)

Note that  $P_2 1_n = 1_n$  and  $P_2 \delta_n = \delta_n$ , so

$$P_2 M = P_2 1_n \mu^T + (1/2) P_2 \delta_n \gamma^T = 1_n \mu^T + (1/2) \delta_n \gamma^T = M,$$
(S82)

and therefore  $(I - P_2)M = 0$ .

Define

$$\check{B} = (I - P_1)B(I - P_1) = \left(\check{b}_{ij}\right) \tag{S83}$$

$$\widetilde{B} = (I - P_2)B(I - P_2) = \left(\widetilde{b}_{ij}\right)$$
(S84)

$$\breve{B} = (I - P_1)B(I - P_2) = \left(\breve{b}_{ij}\right).$$
(S85)

Let  $\check{b}_{\max}$ ,  $\tilde{b}_{\max}$ , and  $\check{b}_{\max}$  denote the maximum diagonal entries of  $\check{B}$ ,  $\tilde{B}$ , and  $\check{B}$ , respectively.

#### H.3 Model Selection Centering

For a subset  $J \subset \{1, \ldots, m\}$ , let  $X_J$  denote the submatrix of X consisting of columns indexed by J. For the fixed sets of genes  $J_0$  and  $J_1$ , define the sample covariance

$$S(B, J_0, J_1) = m^{-1} \sum_{k \in J_0} (I - P_2) X_k X_k^T (I - P_2)^T + m^{-1} \sum_{k \in J_1} (I - P_1) X_k X_k^T (I - P_1)^T =: I + II.$$
(S86)

Note that  $\mathbb{E}[S(B, J_0, J_1)] = B^{\sharp}$ , with

$$B^{\sharp} = \frac{\operatorname{tr}(A_{J_0})}{m}(I - P_2)B(I - P_2) + \frac{\operatorname{tr}(A_{J_1})}{m}(I - P_1)B(I - P_1).$$
(S87)

Define the sample correlation matrix,

$$\widehat{\Gamma}_{ij}(B) = \frac{(S(B, J_0, J_1))_{ij}}{\sqrt{(S(B, J_0, J_1))_{ii}(S(B, J_0, J_1))_{jj}}}.$$
(S88)

The baseline Gemini estimators Zhou (2014a) are then defined as follows, using a pair of penalized estimators for the correlation matrices  $\rho(A) = (a_{ij}/\sqrt{a_{ii}a_{jj}})$  and  $\rho(B) = (b_{ij}/\sqrt{b_{ii}b_{jj}})$ :

$$\widehat{A}_{\rho} = \arg\min_{A_{\rho}>0} \left\{ \operatorname{tr}\left(\widehat{\Gamma}(A)A_{\rho}^{-1}\right) + \log|A_{\rho}| + \lambda_{B}|A_{\rho}^{-1}|_{1,\mathrm{off}} \right\},$$
(S89a)

$$\widehat{B}_{\rho} = \arg\min_{B_{\rho} > 0} \left\{ \operatorname{tr} \left( \widehat{\Gamma}(B) B_{\rho}^{-1} \right) + \log |B_{\rho}| + \lambda_A |B_{\rho}^{-1}|_{1, \operatorname{off}} \right\}.$$
(S89b)

We will focus on  $\hat{B}_{\rho}$  using the input as defined in (S88).

The proof proceeds as follows. Lemma S22, the equivalent of Proposition S16 for Algorithm 1, establishes entry-wise convergence rates of the sample covariance matrix for fixed sets of group and globally centered genes. We use this to prove Theorem S21 below in Section H.4 and to prove Theorem 4 in Section H.5.

#### H.4 Convergence for fixed gene sets

We first state a standalone result, Theorem S21, which provides rates of convergence when  $S(B, J_0, J_1)$  as in (S86) is calculated using fixed sets of group centered and globally centered genes,  $J_0$  and  $J_1$ , respectively. This result shows how the algorithm used in the preliminary step to choose which genes to group center can be decoupled from the rest of the estimation procedure. The proof is presented below in Section H.4.2.

**Theorem S21.** Suppose that (A1), (A2'), and (A3) hold. Let  $J_0$  and  $J_1$  denote sets such that  $J_0 \cap J_1 = \emptyset$  and  $J_0 \cup J_1 = \{1, \ldots, m\}$ . Let  $m_0 = |J_0|$  and  $m_1 = |J_1|$  denote the sizes of the sets. Let  $\tau_{global} > 0$  satisfy

$$\max_{j \in J_1} |\gamma_j| \leqslant \tau_{global},\tag{S90}$$

for  $\tau_{global} = C\sqrt{\log(m)} \| (D^T B^{-1} D)^{-1} \|_2^{1/2} \approx \sqrt{\frac{\log(m)}{n}}.$ 

Consider the data as generated from model (S73) with  $\varepsilon = B^{1/2}ZA^{1/2}$ , where  $A \in \mathbb{R}^{m \times m}$ and  $B \in \mathbb{R}^{n \times n}$  are positive definite matrices, and Z is an  $n \times m$  random matrix as defined in Theorem 1. Let  $\lambda_A$  denote the penalty parameter for estimating B. Suppose the penalty parameter  $\lambda_A$  in (S89b) satisfies

$$\lambda_A \ge C'' \left[ C_A K \frac{\log^{1/2}(m \lor n)}{\sqrt{m}} + \frac{\|B\|_1}{n_{\min}} \right].$$
(S91)

where C'' is an absolute constant.

(I) Let  $\mathcal{E}_4(J_0, J_1)$  be the event such that

$$\left\| \operatorname{tr} \left( A \right) \left( \widehat{W}_2 \widehat{B}_{\rho} \widehat{W}_2 \right)^{-1} - B^{-1} \right\|_2 \leqslant \frac{C' \lambda_A \sqrt{|B^{-1}|_{0, \text{off}} \vee 1}}{b_{\min} \varphi_{\min}^2(\rho(B))}.$$
(S92)

Then  $P(\mathcal{E}_4(J_0, J_1)) \ge 1 - C/m^d$ .

(II) With probability at least  $1 - C'/m^d$ , for all j,

$$\|\widehat{\beta}_{j}(\widehat{B}^{-1}) - \beta_{j}^{*}\|_{2} \leq C_{1}\lambda_{A}\sqrt{\frac{n_{ratio}\left(|B^{-1}|_{0,off} \vee 1\right)}{n_{\min}}} + C_{2}\sqrt{\log(m)}\|(D^{T}B^{-1}D)^{-1}\|_{2}^{1/2}.$$
 (S93)

#### H.4.1 Decomposition of sample covariance matrix

The error in the sample covariance  $S(B, J_0, J_1)$  can be decomposed as

$$S(B, J_0, J_1) - B = [B^{\sharp} - B] + [S(B, J_0, J_1) - B^{\sharp}],$$
(S94)

where the first term corresponds to bias and the second term to variance. We now further decompose the variance term. The first term of  $S(B, J_0, J_1)$  in (S86) can be decomposed as,

$$I = m^{-1}(I - P_2)X_{J_0}X_{J_0}^T(I - P_2)$$
  
=  $m^{-1}(I - P_2)(M_{J_0} + \varepsilon_{J_0})(M_{J_0} + \varepsilon_{J_0})^T(I - P_2)$   
=  $m^{-1}(I - P_2)\varepsilon_{J_0}\varepsilon_{J_0}^T(I - P_2) + m^{-1}(I - P_2)M_{J_0}\varepsilon_{J_0}^T(I - P_2)$   
+  $m^{-1}(I - P_2)\varepsilon_{J_0}M_{J_0}^T(I - P_2) + m^{-1}(I - P_2)M_{J_0}M_{J_0}^T(I - P_2),$  (S95)

and the second term can be decomposed analogously, as

$$II = m^{-1}(I - P_1)\varepsilon_{J_1}\varepsilon_{J_1}^T(I - P_1) + m^{-1}(I - P_1)M_{J_1}\varepsilon_{J_1}^T(I - P_1) + m^{-1}(I - P_1)\varepsilon_{J_1}M_{J_1}^T(I - P_1) + m^{-1}(I - P_1)M_{J_1}M_{J_1}^T(I - P_1).$$
(S96)

By the above decompositions, it follows that  $S(B, J_0, J_1)$  can be expressed as

$$S(B, J_0, J_1) = S_{\rm II} + S_{\rm III} + S_{\rm III}^T + S_{\rm IV},$$
(S97)

with

$$S_{\rm II} = m^{-1}(I - P_2)M_{J_0}M_{J_0}^T(I - P_2) + m^{-1}(I - P_1)M_{J_1}M_{J_1}^T(I - P_1).$$
 (S98)

$$S_{\rm III} = m^{-1}(I - P_2)M_{J_0}\varepsilon_{J_0}^T(I - P_2) + m^{-1}(I - P_1)M_{J_1}\varepsilon_{J_1}^T(I - P_1)$$
(S99)

$$S_{\rm IV} = m^{-1}(I - P_2)\varepsilon_{J_0}\varepsilon_{J_0}^T(I - P_2) + m^{-1}(I - P_1)\varepsilon_{J_1}\varepsilon_{J_1}^T(I - P_1)$$
(S100)

For each of  $S_{\rm II}$ ,  $S_{\rm III}$ , and  $S_{\rm IV}$ , the first term comes from (S96) and the second term comes from (S97).

The terms  $S_{\text{II}}$  and  $S_{\text{III}}$  can be simplified, as follows. Because  $(I - P_2)M_{J_0} = 0$ , it follows that the first term of  $S_{\text{II}}$  is zero:

$$m^{-1}(I - P_2)M_{J_0}M_{J_0}^T(I - P_2) = 0.$$

and the first term of  $S_{\rm III}$  is also zero,

$$m^{-1}(I - P_2)M_{J_0}\varepsilon_{J_0}^T(I - P_2) = 0,$$

Therefore the terms  $S_{\rm II}$  and  $S_{\rm III}$  are equal to

$$S_{\rm II} = m^{-1} (I - P_1) M_{J_1} M_{J_1}^T (I - P_1), \qquad (S101)$$

$$S_{\rm III} = m^{-1}(I - P_1)M_{J_1}\varepsilon_{J_1}^T(I - P_1).$$
(S102)

Let  $S_{\rm I} = B^{\sharp} - B$ . We have thus decomposed the error in the sample covariance as

$$S(B, J_0, J_1) - B = \underbrace{S_{\mathrm{I}}}_{\mathrm{bias}} + \underbrace{\left[\left(S_{\mathrm{IV}} - B^{\sharp}\right) + S_{\mathrm{III}} + S_{\mathrm{II}}\right]}_{\mathrm{variance}}.$$
 (S103)

In Lemma S23, we provide an error bound for each term in the decomposition (S104).

We next state Lemma S22, which establishes the maximum of entry-wise errors for estimating B using the sample covariance for fixed gene sets as defined in (S104). Lemma S22 is used in the proof of Theorem S21. Following, we state Lemma S23, which is used in the proof of Lemma S22.

**Lemma S22.** Suppose the conditions of Theorem S21 hold. Let  $\mathcal{E}_6(J_0, J_1)$  denote the event

$$\mathcal{E}_{6}(J_{0}, J_{1}) = \left\{ \|S(B, J_{0}, J_{1}) - B\|_{\infty} \leq C_{A} K \frac{\log^{1/2}(m \vee n)}{\sqrt{m}} + \frac{\|B\|_{1}}{n_{\min}} \right\}.$$
 (S104)

Then  $\mathcal{E}_6(J_0, J_1)$  holds with probability at least  $1 - \frac{8}{(m \vee n)^2}$ .

**Lemma S23.** Let the model selection-based sample covariance  $S(B, J_0, J_1)$  be as defined in (S86), where  $J_1$  and  $J_0$  are fixed sets of variables that are globally centered, and group centered, respectively. Let  $m_0 = |J_0|$  and  $m_1 = |J_1|$ . Define the rates

$$r_1 = \frac{3 \|B\|_1}{n_{\min}},\tag{S105}$$

$$r_2 = (4m)^{-1} \|\gamma_{J_1}\|_2^2, \qquad (S106)$$

$$r_3 = C_3 d^{1/2} K^2 \log^{1/2}(m) m^{-1} \left(\gamma_{J_1}^T A_{J_1} \gamma_{J_1}\right)^{1/2} \check{b}_{\max}^{1/2}, \tag{S107}$$

$$r_4 = C_4 d^{1/2} K \log^{1/2}(m) m^{-1} \|A\|_F \|B\|_2.$$
(S108)

(I) Deterministically,

$$\left\|B^{\sharp} - B\right\|_{\infty} \leqslant r_1 \quad and \quad \left\|S_{\mathrm{II}}\right\|_{\infty} \leqslant r_2.$$
(S109)

(II) Define the events

$$\mathcal{E}_{\mathrm{I}} = \left\{ \left\| S_{\mathrm{IV}} - B^{\sharp} \right\|_{\infty} \leqslant r_{4} \right\} \quad and \quad \mathcal{E}_{\mathrm{II}} = \left\{ \left\| S_{\mathrm{III}} \right\|_{\infty} \leqslant r_{3} \right\}.$$
(S110)

Then  $\mathcal{E}_{I}$  and  $\mathcal{E}_{II}$  occur with probability at least  $1-2/m^{d}$ .

Lemmas S22 and S23 are proved in Section I. We analyze term  $S_{\rm I}$  in Section I.2, term  $S_{\rm II}$  in Section I.3, term  $S_{\rm III}$  in Section I.4, and term  $S_{\rm IV}$  in Section I.5.

#### H.4.2 Proof of Theorem S21

Let us first define the event  $\mathcal{E}_{global}$ , that is, the GLS error based on the true  $B^{-1}$  is small:

$$\mathcal{E}_{\text{global}} = \left\{ \left\| \hat{\gamma}(B^{-1}) - \gamma \right\|_{\infty} < \sqrt{\log(m)} \| (D^T B^{-1} D)^{-1} \|_2^{1/2} \right\}.$$
 (S111)

Let  $\mathcal{E}_4(J_0, J_1)$  be defined as in (S93), denoting small operator norm error in estimating  $B^{-1}$ :

$$\mathcal{E}_4(J_0, J_1) = \left\{ \left\| \operatorname{tr} \left( A \right) \left( \widehat{W}_2 \widehat{B}_\rho \widehat{W}_2 \right)^{-1} - B^{-1} \right\|_2 \leqslant \frac{C' \lambda_A \sqrt{|B^{-1}|_{0, \operatorname{off}} \vee 1}}{b_{\min} \varphi_{\min}^2(\rho(B))} \right\}.$$
 (S112)

Note that  $\mathcal{E}_4(J_0, J_1)$  holds deterministically under event  $\mathcal{E}_6(J_0, J_1)$  as defined in (S105) of Lemma S22.

Define the event bounding the perturbation in mean estimation due to error in estimating  $B^{-1}$ :

$$\mathcal{E}_{5}(J_{0}, J_{1}) = \left\{ \left\| \widehat{\gamma}(\widehat{B}^{-1}) - \widehat{\gamma}(B^{-1}) \right\|_{\infty} < C n_{\min}^{-1/2} \left\| \widehat{B}^{-1} - B^{-1} \right\|_{2} \right\}.$$
 (S113)

Conditional on a fixed matrix  $\hat{B}^{-1}$  that satisfies  $\mathcal{E}_4(J_0, J_1)$ , event  $\mathcal{E}_5(J_0, J_1)$  holds with probability at least  $1 - C/m^d$ , by Lemma S6 (used in the proof of Theorem 1).

The overall rate of convergence follows by applying the union bound to the events  $\mathcal{E}_{global} \cap \mathcal{E}_4(J_0, J_1) \cap \mathcal{E}_5(J_0, J_1)$ , as follows:

$$\begin{split} P(\mathcal{E}_{global}^{c} \cup \mathcal{E}_{4}(J_{0}, J_{1})^{c} \cup \mathcal{E}_{5}(J_{0}, J_{1})^{c}) \\ &\leqslant P(\mathcal{E}_{global}^{c}) + P(\mathcal{E}_{4}(J_{0}, J_{1})^{c}) + P(\mathcal{E}_{5}(J_{0}, J_{1})^{c} \mid \mathcal{E}_{4}(J_{0}, J_{1})) P(\mathcal{E}_{4}(J_{0}, J_{1})) \\ &+ P(\mathcal{E}_{5}(J_{0}, J_{1})^{c} \mid \mathcal{E}_{4}(J_{0}, J_{1})^{c}) P(\mathcal{E}_{4}(J_{0}, J_{1})^{c}) \\ &\leqslant P(\mathcal{E}_{global}^{c}) + P(\mathcal{E}_{4}(J_{0}, J_{1})^{c}) + P(\mathcal{E}_{4}(J_{0}, J_{1})^{c}) + P(\mathcal{E}_{5}(J_{0}, J_{1})^{c} \mid \mathcal{E}_{4}(J_{0}, J_{1})) \\ &= P(\mathcal{E}_{global}^{c}) + 2P(\mathcal{E}_{4}(J_{0}, J_{1})^{c}) + P(\mathcal{E}_{5}(J_{0}, J_{1})^{c} \mid \mathcal{E}_{4}(J_{0}, J_{1})), \end{split}$$

where  $P(\mathcal{E}_{global}^{c})$  and  $P(\mathcal{E}_{5}(J_{0}, J_{1})^{c} | \mathcal{E}_{4}(J_{0}, J_{1}))$  are bounded in Theorem 1, and  $P(\mathcal{E}_{4}(J_{0}, J_{1})^{c})$  has high probability under Lemma S22.

#### H.5 Proof of Theorem 4

Let  $\hat{\gamma}^{\text{init}}$  denote the output from Algorithm 1. By our choice of the threshold parameter  $\tau_{\text{init}}$  as in (16), that is,

$$\tau_{\text{init}} = C \left( \frac{\log^{1/2}(m)}{\sqrt{m}} + \frac{\|B\|_1}{n_{\min}} \right) \sqrt{\frac{n_{\text{ratio}} \left( |B^{-1}|_{0, \text{off}} \vee 1 \right)}{n_{\min}}} + C \sqrt{\log(m)} \| (D^T B^{-1} D)^{-1} \|_2^{1/2},$$

we have a partition  $(\tilde{J}_0, \tilde{J}_1)$  such that  $\tilde{J}_0$  is the set of variables selected for group centering and  $\tilde{J}_1$  is the set of variables selected for global centering. The partition results in a sample covariance matrix  $S(B, \tilde{J}_0, \tilde{J}_1)$  as defined in (S86). Define the event that the Algorithm 1 estimate  $\hat{\gamma}^{\text{init}}$  is close to  $\gamma$  in the sense that

$$\mathcal{E}_{A1} = \left\{ \left\| \hat{\gamma}^{\text{init}} - \gamma \right\|_{\infty} < \tau_{\text{init}} \right\}.$$
(S114)

Note that the event  $\mathcal{E}_{A1}$  implies that the false negatives have small true mean differences. That is, on event  $\mathcal{E}_{A1}$ , by the triangle inequality,

$$\left\|\gamma_{\tilde{J}_{1}}\right\|_{\infty} \leqslant \left\|\gamma_{\tilde{J}_{1}} - \hat{\gamma}_{\tilde{J}_{1}}^{\text{init}}\right\|_{\infty} + \left\|\hat{\gamma}_{\tilde{J}_{1}}^{\text{init}}\right\|_{\infty} \leqslant \tau_{\text{init}} + \tau_{\text{init}} = 2\tau_{\text{init}},\tag{S115}$$

where  $\left\| \widehat{\gamma}_{\widetilde{J}_1}^{\text{init}} \right\|_{\infty} < \tau_{\text{init}}$  by definition of  $\mathcal{E}_{A1}$ , and  $\left\| \gamma_{\widetilde{J}_1} - \widehat{\gamma}_{\widetilde{J}_1}^{\text{init}} \right\|_{\infty} < \tau_{\text{init}}$  by definition of the thresholding set  $\widetilde{J}_1$ .

Under the assumptions of Theorem S21,  $\tau_{\text{init}} \leq \tau_{\text{global}}$  with  $\tau_{\text{global}}$  as defined in (S90), so condition (S90) of Theorem S21 is satisfied. Under the conditions of Theorem S21, event  $\mathcal{E}_6(J_0, J_1)$  as defined in Lemma S22 holds with high probability; that is, the entrywise error in the sample covariance matrix is small.

Let  $\mathcal{E}_B$  denote event (28) in Theorem 4. In view of Theorem S9 and Lemma S10, event

 $\mathcal{E}_B$  holds on  $\mathcal{E}_6(J_0, J_1)$ . Hence

$$P\left(\mathcal{E}_{B}^{c}\right) = P\left(\mathcal{E}_{6}(J_{0}, J_{1})^{c} \mid \mathcal{E}_{A1}\right) P\left(\mathcal{E}_{A1}\right) + P\left(\mathcal{E}_{6}(J_{0}, J_{1})^{c} \mid \mathcal{E}_{A1}^{c}\right) P\left(\mathcal{E}_{A1}^{c}\right)$$
$$\leq P\left(\mathcal{E}_{6}(J_{0}, J_{1})^{c} \mid \mathcal{E}_{A1}\right) + P\left(\mathcal{E}_{A1}^{c}\right)$$
$$\leq 2/m^{d} + 2/m^{d},$$

where the first term is bounded in Lemma S22 and the second in Theorem 3.

Recall the event  $\mathcal{E}_{global}$  as defined in (S112). Event (29) in Theorem 4 holds under the intersection of events  $\mathcal{E}_{global} \cap \mathcal{E}_5(\widetilde{J}_0, \widetilde{J}_1) \cap \mathcal{E}_B \cap \mathcal{E}_{A1}$ . Hence the probability of (29) can be bounded as follows:

$$\begin{split} P(\mathcal{E}_{global}^{c} \cup \mathcal{E}_{5}(\widetilde{J}_{0}, \widetilde{J}_{1})^{c} \cup \mathcal{E}_{B}^{c} \cup \mathcal{E}_{A1}^{c}) \\ &\leq P(\mathcal{E}_{global}^{c}) + P(\mathcal{E}_{B}^{c}) + P(\mathcal{E}_{5}(\widetilde{J}_{0}, \widetilde{J}_{1})^{c} \mid \mathcal{E}_{B})P(\mathcal{E}_{B}) \\ &+ P(\mathcal{E}_{5}(\widetilde{J}_{0}, \widetilde{J}_{1})^{c} \mid \mathcal{E}_{B}^{c})P(\mathcal{E}_{B}^{c}) + P(\mathcal{E}_{A1}^{c}) \\ &\leq P(\mathcal{E}_{global}^{c}) + P(\mathcal{E}_{B}^{c}) + P(\mathcal{E}_{B}^{c}) + P(\mathcal{E}_{5}(\widetilde{J}_{0}, \widetilde{J}_{1})^{c} \mid \mathcal{E}_{B}) + P(\mathcal{E}_{A1}^{c}) \\ &= P(\mathcal{E}_{global}^{c}) + 2P(\mathcal{E}_{B}^{c}) + P(\mathcal{E}_{5}(\widetilde{J}_{0}, \widetilde{J}_{1})^{c} \mid \mathcal{E}_{B}) + P(\mathcal{E}_{A1}^{c}) \,, \end{split}$$

where  $P(\mathcal{E}_{global}^{c})$  and  $P(\mathcal{E}_{5}(\widetilde{J}_{0},\widetilde{J}_{1})^{c} | \mathcal{E}_{B})$  are bounded in Theorem 1,  $P(\mathcal{E}_{B}^{c})$  is bounded above, and  $P(\mathcal{E}_{A1}^{c})$  is bounded in Theorem 3.

# I Proof of Lemmas S22 and S23

We first prove Lemma S22 in Section I.1. The rest of the section contains the proof of Lemma S23, where part I is proved in Sections I.2 and I.3 and part II in Sections I.4 and I.5.

#### I.1 Proof of Lemma S22

The entrywise error in the sample covariance matrix (S86) can be decomposed as

$$\|S(B, J_0, J_1) - B\|_{\infty} \leq \|S(B, J_0, J_1) - B^{\sharp}\|_{\infty} + \|B^{\sharp} - B\|_{\infty}$$
(S116)

$$\leq \left\| S_{\rm IV} - B^{\sharp} \right\|_{\infty} + 2 \left\| S_{\rm III} \right\|_{\infty} + \left\| S_{\rm II} \right\|_{\infty} + \left\| B^{\sharp} - B \right\|_{\infty}.$$
 (S117)

Let  $r_{n,m} = r_1 + r_2 + 2r_3 + r_4$ . By parts I and II of Lemma S23,

$$P(||S(B, J_{0}, J_{1}) - B||_{\infty} \ge r_{n,m})$$

$$\leq P(||S_{IV} - B^{\sharp}||_{\infty} + 2 ||S_{III}||_{\infty} + ||S_{II}||_{\infty} + ||B^{\sharp} - B||_{\infty} \ge r_{n,m}) \quad (by \ (S118))$$

$$\leq P(||S_{IV} - B^{\sharp}||_{\infty} + 2 ||S_{III}||_{\infty} + r_{2} + r_{1} \ge r_{n,m}) \quad (by \ (S110))$$

$$= P(||S_{IV} - B^{\sharp}||_{\infty} + 2 ||S_{III}||_{\infty} \ge r_{4} + 2r_{3})$$

$$\leq P(||S_{IV} - B^{\sharp}||_{\infty} \ge r_{4}) + P(2 ||S_{III}||_{\infty} \ge 2r_{3}) \quad (by \ (S111))$$

$$\leq \frac{2}{m^{d}} + \frac{2}{m^{d}} = \frac{4}{m^{d}}.$$

We show that under the assumptions of Theorem S21, the entrywise error in terms  $S_{\rm II}$ and  $S_{\text{III}}$  is  $O\left(C_A\sqrt{\frac{\log(m)}{m}}\right)$ . Recall that the entrywise rates of convergence of  $S_{\text{II}}$  and  $S_{\text{III}}$  are stated in equations (S107) and (S108), respectively. Let  $s = |\operatorname{supp}(\gamma)|$  denote the sparsity of  $\gamma$ . Let  $m_{01} = |\text{supp}(\gamma_{J_1})|$  denote the number of false negatives.

First, we express the entrywise rate of convergence of  $S_{\rm II}$  in terms of  $\tau_{\rm global}$ . By (S90),  $\|\gamma_{J_1}\|_{\infty} \leq \tau_{\text{global}}$ , which implies that  $\|\gamma_{J_1}\|_2^2 \leq m_{01}\tau_{\text{global}}^2 \leq s\tau_{\text{global}}^2$ , where the last inequality holds because  $m_{01} \leq s$  by definition. Therefore,

$$r_{2} = (4m)^{-1} \|\gamma_{J_{1}}\|_{2}^{2} \leq \frac{s\tau_{\text{global}}^{2}}{4m} \leq C \frac{s\log(m)}{4nm} \|B\|_{2}, \qquad (S118)$$

where the last step holds because  $\tau_{\text{global}} = C\sqrt{\log(m)} \|(D^T B^{-1} D)^{-1}\|_2^{1/2} \approx \sqrt{\frac{\log(m)}{n}} \|B\|_2^{1/2}$  by assumption. Applying (A3) to the right-hand side of (S119) implies that  $r_2 = O\left(C_A\sqrt{\frac{\log(m)}{m}}\right)$ .

Next, consider term  $S_{\rm III}$ . First note that

$$\gamma_{J_1}^T A_{J_1} \gamma_{J_1} \leqslant \|\gamma_{J_1}\|_2^2 \|A_{J_1}\|_2 \leqslant m_{01} \tau_{\text{global}}^2 \|A_{J_1}\|_2, \qquad (S119)$$

where the last inequality holds by (S90). This implies that  $r_3$  is on the order

$$\frac{\log^{1/2}(m)}{m} \left( \check{b}_{\max} \gamma_{J_1}^T A_{J_1} \gamma_{J_1} \right)^{1/2} \leqslant \check{b}_{\max}^{1/2} \|A_{J_1}\|_2^{1/2} \left( \frac{\log^{1/2}(m) m_{01}^{1/2}}{m} \right) \tau_{\text{global}} \\
\leqslant C \frac{\log(m)}{\sqrt{n}} \frac{\sqrt{s}}{m} \|A_{J_1}\|_2^{1/2} \|B\|_2^{1/2} \check{b}_{\max}^{1/2},$$
(S120)

where the last inequality holds because  $m_{01} \leq s \leq m$  and  $\tau_{\text{global}} \approx \sqrt{\frac{\log(m)}{n}} \|B\|_2^{1/2}$ . Under (A2'), the right-hand side of (S121) satisfies

$$\frac{\log(m)}{\sqrt{n}}\frac{\sqrt{s}}{m} \|A_{J_1}\|_2^{1/2} \|B\|_2^{1/2} \check{b}_{\max}^{1/2} \leqslant \sqrt{\log(m)}\frac{\sqrt{s}}{m} C_A \frac{\|A_{J_1}\|_2^{1/2}}{\|A\|_2^{1/2}} \leqslant C_A \sqrt{\frac{\log(m)}{m}}, \qquad (S121)$$

where the last inequality holds because  $s \leq m$ .

#### I.2 Proof of part I of Lemma S23, term I

We bound the entrywise bias,

$$\left\|B^{\sharp} - B\right\|_{\max} = \left\|\frac{\operatorname{tr}\left(A_{J_{0}}\right)}{m}\widetilde{B} + \frac{\operatorname{tr}\left(A_{J_{1}}\right)}{m}\widetilde{B} - B\right\|_{\max}$$
$$\leq \frac{\operatorname{tr}\left(A_{J_{0}}\right)}{m}\left\|\widetilde{B} - B\right\|_{\max} + \frac{\operatorname{tr}\left(A_{J_{1}}\right)}{m}\left\|\widetilde{B} - B\right\|_{\max}.$$
 (S122)

Note that

$$\left\| \check{B} - B \right\|_{\max} = \left\| (I - P_1) B (I - P_1) - B \right\|_{\max} = \left\| P_1 B P_1 - P_1 B - B P_1 \right\|_{\max}$$
  
$$\leq \left\| P_1 B P_1 \right\|_{\max} + \left\| P_1 B \right\|_{\max} + \left\| B P_1 \right\|_{\max}.$$
 (S123)

We bound the first term of (S124) as follows:

$$\left| (P_1 B P_1)_{ij} \right| \leq \left\| p_i^{(1)} \right\|_2 \left\| p_j^{(1)} \right\|_2 \| B \|_2 \leq \frac{\| B \|_2}{n_{\min}}.$$

For the second term of (S124),

$$(P_1B)_{ij} = \left| b_i^T p_j^{(1)} \right| \le \|b_i\|_1 \left\| p_j^{(1)} \right\|_{\infty} \le \|B\|_1 \left\| p_j^{(1)} \right\|_{\infty} \le \frac{\|B\|_1}{n_{\min}},$$

where  $\left\|p_{j}^{(1)}\right\|_{\infty} \leq \frac{1}{n_{\min}}$  by the definition of  $P_{1}$  in (S78). We have shown  $\left\|BP_{1}\right\|_{\max} \leq \frac{\left\|B\right\|_{1}}{n_{\min}}$ . Likewise,  $\left\|BP_{1}\right\|_{\max} \leq \frac{\left\|B\right\|_{1}}{n_{\min}}$ . Therefore,

$$\left\| \check{B} - B \right\|_{\max} \leqslant 3 \frac{\|B\|_1}{n_{\min}}.$$
(S124)

Because the projection matrix  $P_2$  satisfies  $\left\|p_j^{(2)}\right\|_{\infty} \leq \frac{1}{n_{\min}}$ , an analogous proof shows that

$$\left\|\widetilde{B} - B\right\|_{\max} \leqslant \frac{3 \left\|B\right\|_1}{n_{\min}}.$$
(S125)

Substituting (S125) and (S126) into (S123) yields

$$\begin{split} \left\| B^{\sharp} - B \right\|_{\max} &\leq \frac{\operatorname{tr} \left( A_{J_0} \right)}{m} \left\| \widecheck{B} - B \right\|_{\max} + \frac{\operatorname{tr} \left( A_{J_1} \right)}{m} \left\| \widetilde{B} - B \right\|_{\max} \\ &\leq \left( \frac{\operatorname{tr} \left( A_{J_0} \right)}{m} + \frac{\operatorname{tr} \left( A_{J_1} \right)}{m} \right) \frac{3 \left\| B \right\|_1}{n_{\min}} \\ &= \frac{\operatorname{tr} \left( A \right)}{m} \frac{3 \left\| B \right\|_1}{n_{\min}} \\ &= \frac{3 \left\| B \right\|_1}{n_{\min}}. \end{split}$$
(S126)

#### I.3 Proof of part I of Lemma S23, term II

In this section we prove a deterministic entrywise bound on  $S_{\rm II}$ . By (S80), it follows that

$$(I - P_1)M_{J_1}M_{J_1}^T(I - P_1) = (1/4) \|\gamma_{J_1}\|_2^2 \delta_n \delta_n^T,$$

which implies

$$\left\| (I - P_1) M_{J_1} M_{J_1}^T (I - P_1) \right\|_{\infty} = \left\| (1/4) \| \gamma_{J_1} \|_2^2 \,\delta_n \delta_n^T \right\|_{\infty} = (1/4) \| \gamma_{J_1} \|_2^2.$$

Therefore  $S_{\rm II}$  satisfies the maximum entrywise bound

$$\|S_{\mathrm{II}}\|_{\infty} = \|m^{-1}(I-P_1)M_{J_1}M_{J_1}^T(I-P_1)\|_{\infty} = \|(4m)^{-1}\|\gamma_{J_1}\|_2^2 \,\delta_n \delta_n^T\|_{\infty} = (4m)^{-1}\|\gamma_{J_1}\|_2^2,$$

 $\mathbf{SO}$ 

$$\|S_{\mathrm{II}}\|_{\infty} = r_2.$$

Note that if  $J_1$  is chosen so that  $\|\gamma_{J_1}\|_{\infty} \leq \tau$ , then  $\|\gamma_{J_1}\|_2^2 \leq m_{01}\tau^2$ , where  $m_{01}$  is the number of false negatives, so

$$\frac{\|\gamma_1\|_2^2}{4m} \leqslant \frac{m_{01}}{4m} \tau^2 \leqslant \frac{\tau^2}{4}.$$
 (S127)

which implies that the entrywise rate of convergence of  $S_{\rm II}$  is  $O(\tau^2)$ .

## I.4 Proof of part II of Lemma S23, term III

Let  $p_i$  denote the *i*th column of  $P_1^T$ , for i = 1, ..., n. Let  $m_k$  denote the *k*th column of M. Let  $\varepsilon_k$  denote the *k*th column of  $\varepsilon$ . The term  $S_{\text{III}}$  can be expressed as

$$(S_{\text{III}})_{ij} = m^{-1} (e_i - p_i)^T M_{J_1} \varepsilon_{J_1}^T (e_j - p_j)$$
  
=  $m^{-1} \text{tr} \left( \varepsilon_{J_1}^T (e_j - p_j) (e_i - p_i)^T M_{J_1} \right)$   
=  $m^{-1} \sum_{k \in J_1} \varepsilon_k^T (e_j - p_j) (e_i - p_i)^T m_k$   
=  $m^{-1} \text{vec} \{\varepsilon_{J_1}\}^T \left( I_{m_1} \otimes (e_j - p_j) (e_i - p_i)^T \right) \text{vec} \{M_{J_1}\}$   
=  $m^{-1} \text{vec} \{Z\}^T \left( A_{J_1}^{1/2} \otimes B^{1/2} (e_j - p_j) (e_i - p_i)^T \right) \text{vec} \{M_{J_1}\}$   
=  $\text{vec} \{Z\}^T \psi_{ij},$ 

where

$$\psi_{ij} := m^{-1} \left( A_{J_1}^{1/2} \otimes B^{1/2} (e_j - p_j) (e_i - p_i)^T \right) \operatorname{vec} \left\{ M_{J_1} \right\}.$$
(S128)

The squared Euclidean norm of  $\psi_{ij}$  is

$$\begin{aligned} \|\psi_{ij}\|_{2}^{2} &= \operatorname{vec} \left\{ M_{J_{1}} \right\}^{T} \left( A_{J_{1}} \otimes (e_{i} - p_{i})(e_{j} - p_{j})^{T} B(e_{j} - p_{j})(e_{i} - p_{i})^{T} \right) \operatorname{vec} \left\{ M_{J_{1}} \right\} / m^{2} \\ &= \operatorname{vec} \left\{ M_{J_{1}} \right\}^{T} \left( A_{J_{1}} \otimes \check{b}_{jj}(e_{i} - p_{i})(e_{i} - p_{i})^{T} \right) \operatorname{vec} \left\{ M_{J_{1}} \right\} / m^{2} \\ &= \check{b}_{jj} \sum_{k \in J_{1}} \sum_{\ell \in J_{1}} a_{k\ell} m_{k}^{T}(e_{i} - p_{i})(e_{i} - p_{i})^{T} m_{\ell} / m^{2} \\ &= \check{b}_{jj} \sum_{k \in J_{1}} \sum_{\ell \in J_{1}} a_{k\ell} (\delta_{n})_{i} \gamma_{k} (\delta_{n})_{i} \gamma_{\ell} / \left( 4m^{2} \right) \\ &= \check{b}_{jj} \sum_{k \in J_{1}} \sum_{\ell \in J_{1}} a_{k\ell} \gamma_{k} \gamma_{\ell} / \left( 4m^{2} \right) \\ &= \check{b}_{jj} \gamma_{J_{1}}^{T} A_{J_{1}} \gamma_{J_{1}} / \left( 4m^{2} \right). \end{aligned}$$
(S129)

By the Hanson-Wright inequality (Theorem 2.1),

$$\mathbb{P}\left(\left|\operatorname{vec}\left\{Z\right\}^{T}\psi_{ij} - \left\|\psi_{ij}\right\|_{2}\right| > d^{1/2}K^{2}\sqrt{\log(m)}\left\|\psi_{ij}\right\|_{2}\right) \leq 2\exp\left\{-d\log(m)\right\} = 2/m^{d}.$$
 (S130)

Therefore

$$\mathbb{P}\left(\left|(S_{\text{III}})_{ij}\right| > \left(1 + d^{1/2}K^2\sqrt{\log(m)}\right) \|\psi_{ij}\|_2\right) = \mathbb{P}\left(\left|\operatorname{vec}\left\{Z\right\}^T\psi_{ij}\right| > \|\psi_{ij}\|_2 + d^{1/2}K^2\sqrt{\log(m)} \|\psi_{ij}\|_2\right) \\
\leq \mathbb{P}\left(\left|\operatorname{vec}\left\{Z\right\}^T\psi_{ij} - \|\psi_{ij}\|_2\right| > d^{1/2}K^2\sqrt{\log(m)} \|\psi_{ij}\|_2\right) \\
\leq 2/m^d,$$

where the last step follows from (S131). By (S130), it follows that

$$\left(1 + d^{1/2} K^2 \sqrt{\log(m)}\right) \|\psi_{ij}\|_2 \leqslant r_3,$$
(S131)

 $\mathbf{SO}$ 

$$\mathbb{P}\left(|(S_{\text{III}})_{ij}| > r_3\right) \le \mathbb{P}\left(|(S_{\text{III}})_{ij}| > \left(1 + d^{1/2} K^2 \sqrt{\log(m)}\right) \|\psi_{ij}\|_2\right) \le 2/m^d,$$
(S132)

by (S132). By the union bound,

$$\mathbb{P}(\|S_{\text{III}}\|_{\infty} > r_3) \leq \sum_{i=1}^{m} \sum_{j=1}^{m} \mathbb{P}(|(S_{\text{III}})_{ij}| > r_3) \leq 2/m^{d-2}.$$

#### I.5 Proof of part II of Lemma S23, term IV

We now analyze term  $S_{IV}$ . To do so, we express  $S_{IV}$  as a quadratic form in order to apply the Hanson-Wright inequality.

Let  $p_i^{(1)}$  denote the *i*th column of  $P_1^T$ . Let  $p_i^{(2)}$  denote the *i*th column of  $P_2^T$ . Define

$$H_{\text{group}}^{ij} = I_{m_0} \otimes \left(e_j - p_j^{(2)}\right) \left(e_j - p_j^{(2)}\right)^T \quad \text{and} \quad H_{\text{global}}^{ij} = I_{m_1} \otimes \left(e_j - p_j^{(1)}\right) \left(e_j - p_j^{(1)}\right)^T,$$
(S133)

and let

$$H^{ij}(J_0, J_1) = \begin{bmatrix} H^{ij}_{\text{group}} \\ & H^{ij}_{\text{global}} \end{bmatrix}, \qquad (S134)$$

where  $H_{\text{group}}^{ij} \in \mathbb{R}^{m_0 n \times m_0 n}$ ,  $H_{\text{global}}^{ij} \in \mathbb{R}^{m_1 n \times m_1 n}$ , and  $H^{ij}(J_0, J_1) \in \mathbb{R}^{m n \times m n}$ . Recall that

$$S_{\rm IV} = m^{-1}(I - P_2)\varepsilon_{J_0}\varepsilon_{J_0}^T(I - P_2) + m^{-1}(I - P_1)\varepsilon_{J_1}\varepsilon_{J_1}^T(I - P_1).$$

The second term of  $S_{\text{IV}}$  can be expressed as a quadratic form, as follows (where  $\varepsilon_k$  denotes the *k*th column of  $\varepsilon \in \mathbb{R}^{n \times m}$ ):

$$m^{-1}(I - P_{1})\varepsilon_{J_{1}}\varepsilon_{J_{1}}^{T}(I - P_{1}) = m^{-1}\sum_{k\in J_{1}} \left(e_{i} - p_{i}^{(1)}\right)^{T}\varepsilon_{k}\varepsilon_{k}^{T}\left(e_{j} - p_{j}^{(1)}\right)$$

$$= m^{-1}\sum_{k\in J_{1}}\operatorname{tr}\left(\left(e_{i} - p_{i}^{(1)}\right)^{T}\varepsilon_{k}\varepsilon_{k}^{T}\left(e_{j} - p_{j}^{(1)}\right)\right)$$

$$= m^{-1}\sum_{k\in J_{1}}\varepsilon_{k}^{T}\left(e_{j} - p_{j}^{(1)}\right)\left(e_{i} - p_{i}^{(1)}\right)^{T}\varepsilon_{k}$$

$$= m^{-1}\operatorname{vec}\left\{\varepsilon_{J_{1}}\right\}^{T}\left(I_{m_{1}}\otimes\left(e_{j} - p_{j}^{(1)}\right)\left(e_{i} - p_{i}^{(1)}\right)^{T}\right)\operatorname{vec}\left\{\varepsilon_{J_{1}}\right\}^{T}$$

$$= m^{-1}\operatorname{vec}\left\{\varepsilon_{J_{1}}\right\}^{T}H_{\text{global}}^{ij}\operatorname{vec}\left\{\varepsilon_{J_{1}}\right\}^{T}.$$
(S135)

Analogously, the first term of  $S_{\rm IV}$  can be expressed as a quadratic form:

$$m^{-1}(I - P_2)\varepsilon_{J_0}\varepsilon_{J_0}^T(I - P_2) = m^{-1}\sum_{k\in J_0} \left(e_i - p_i^{(2)}\right)^T \varepsilon_k \varepsilon_k^T \left(e_j - p_j^{(2)}\right)$$
$$= m^{-1}\operatorname{vec}\left\{\varepsilon_{J_0}\right\}^T H_{\text{group}}^{ij}\operatorname{vec}\left\{\varepsilon_{J_0}\right\}^T.$$
(S136)

We now express  $S_{IV}$  as a quadratic form. Let  $\pi(X)$  denote the matrix X with reordered columns:

$$\pi(X) = \begin{bmatrix} X_{J_0} & X_{J_1} \end{bmatrix} \quad \text{and} \quad \pi(A) = \operatorname{Cov}\left(\operatorname{vec}\left\{\pi(X)\right\}\right). \tag{S137}$$

Then by (S136) and (S137),

$$(S_{\rm IV})_{ij} = m^{-1} \operatorname{vec} \{\varepsilon_{J_0}\}^T H^{ij}_{\rm group} \operatorname{vec} \{\varepsilon_{J_0}\}^T + m^{-1} \operatorname{vec} \{\varepsilon_{J_1}\}^T H^{ij}_{\rm global} \operatorname{vec} \{\varepsilon_{J_1}\}^T = m^{-1} \operatorname{vec} \{\pi(\varepsilon)\}^T H^{ij}(J_0, J_1) \operatorname{vec} \{\pi(\varepsilon)\} = m^{-1} \operatorname{vec} \{Z\}^T \left( \left(\pi(A)^{1/2} \otimes B^{1/2}\right) H^{ij}(J_0, J_1) \left(\pi(A)^{1/2} \otimes B^{1/2}\right) \right) \operatorname{vec} \{Z\},$$

where the last step holds by decorrelation, with  $Z \in \mathbb{R}^{n \times m}$  as a random matrix with independent subgaussian entries.

Note that the (i, j)th entry of  $S_{\text{IV}}$  can be expressed as

$$(S_{\rm IV})_{ij} = \operatorname{vec} \{Z\}^T \Phi_{i,j} \operatorname{vec} \{Z\}, \qquad (S138)$$

with

$$\Phi_{i,j} = m^{-1} \left( \pi(A)^{1/2} \otimes B^{1/2} \right) H^{ij}(J_0, J_1) \left( \pi(A)^{1/2} \otimes B^{1/2} \right).$$
(S139)

Having expressed  $(S_{IV})_{ij}$  as a quadratic form in (S139), we find the trace and Frobenius norm of  $\Phi_{i,j}$ , then apply the Hanson-Wright inequality. First we find the trace of  $\Phi_{i,j}$ . Let

$$\mathcal{I}_{0} = \begin{bmatrix} I_{m_{0} \times m_{0}} & 0_{m_{0} \times m_{1}} \\ 0_{m_{1} \times m_{0}} & 0_{m_{1} \times m_{1}} \end{bmatrix} \quad \text{and} \quad \mathcal{I}_{1} = \begin{bmatrix} 0_{m_{0} \times m_{0}} & 0_{m_{0} \times m_{1}} \\ 0_{m_{1} \times m_{0}} & I_{m_{1} \times m_{1}} \end{bmatrix}.$$
 (S140)

Note that  $H^{ij}(J_0, J_1)$  can be written as a sum of Kronecker products,

$$H^{ij}(J_0, J_1) = \mathcal{I}_0 \otimes \left(e_j - p_j^{(2)}\right) \left(e_i - p_i^{(2)}\right)^T + \mathcal{I}_1 \otimes \left(e_j - p_j^{(1)}\right) \left(e_i - p_i^{(1)}\right)^T, \quad (S141)$$

hence (S140) can be expressed as

$$m^{-1}\left(\pi(A)^{1/2} \otimes B^{1/2}\right) \left(\mathcal{I}_0 \otimes \left(e_j - p_j^{(2)}\right) \left(e_i - p_i^{(2)}\right)^T\right) \left(\pi(A)^{1/2} \otimes B^{1/2}\right) \tag{S142}$$

$$+ m^{-1} \left( \pi(A)^{1/2} \otimes B^{1/2} \right) \left( \mathcal{I}_1 \otimes \left( e_j - p_j^{(1)} \right) \left( e_i - p_i^{(1)} \right)^T \right) \left( \pi(A)^{1/2} \otimes B^{1/2} \right).$$
(S143)

The trace of the term (S143) is

$$m^{-1} \operatorname{tr} \left( \left( \pi(A)^{1/2} \otimes B^{1/2} \right) \left( \mathcal{I}_0 \otimes \left( e_j - p_j^{(2)} \right) \left( e_i - p_i^{(2)} \right)^T \right) \left( \pi(A)^{1/2} \otimes B^{1/2} \right) \right)$$
  

$$= m^{-1} \operatorname{tr} \left( \pi(A)^{1/2} \mathcal{I}_0 \pi(A)^{1/2} \otimes B^{1/2} \left( e_j - p_j^{(2)} \right) \left( e_i - p_i^{(2)} \right)^T B^{1/2} \right)$$
  

$$= m^{-1} \operatorname{tr} \left( \pi(A)^{1/2} \mathcal{I}_0 \pi(A)^{1/2} \right) \operatorname{tr} \left( B^{1/2} \left( e_j - p_j^{(2)} \right) \left( e_i - p_i^{(2)} \right)^T B^{1/2} \right)$$
  

$$= m^{-1} \operatorname{tr} \left( \mathcal{I}_0 \pi(A) \right) \left( \left( e_i - p_i^{(2)} \right)^T B \left( e_j - p_j^{(2)} \right) \right)$$
  

$$= m^{-1} \operatorname{tr} \left( A_{J_0} \right) \left[ (I - P_2) B(I - P_2) \right]_{ij}$$
  

$$= m^{-1} \operatorname{tr} \left( A_{J_0} \right) \widetilde{b}_{ij}.$$

Analogously, the trace of the term (S144) is

$$m^{-1} \operatorname{tr} \left( \left( \pi(A)^{1/2} \otimes B^{1/2} \right) \left( \mathcal{I}_1 \otimes \left( e_j - p_j^{(1)} \right) \left( e_i - p_i^{(1)} \right)^T \right) \left( \pi(A)^{1/2} \otimes B^{1/2} \right) \right)$$
  
=  $m^{-1} \operatorname{tr} \left( A_{J_1} \right) \left[ (I - P_1) B(I - P_1) \right]_{ij}$   
=  $m^{-1} \operatorname{tr} \left( A_{J_1} \right) \check{b}_{ij}.$ 

Let  $b_{ij}^{\sharp}$  denote the (i, j)th entry of  $B^{\sharp}$  defined in (S87). We have shown that the trace of  $\Phi_{i,j}$  (as defined in (S140)) is

$$\operatorname{tr}(\Phi_{i,j}) = m^{-1} \operatorname{tr}(A_{J_0}) \widetilde{b}_{ij} + m^{-1} \operatorname{tr}(A_{J_1}) \widecheck{b}_{ij} = b_{ij}^{\sharp}.$$
 (S144)

Next, we find the Frobenius norm of  $\Phi_{i,j}$ . For convenience, define

$$\mathcal{A}_{0} = \pi(A)^{1/2} \mathcal{I}_{0} \pi(A)^{1/2} \quad \text{and} \quad \mathcal{A}_{1} = \pi(A)^{1/2} \mathcal{I}_{1} \pi(A)^{1/2}$$
(S145)  
$$\mathcal{B}_{2,ij} = B^{1/2} \left( e_{j} - p_{j}^{(2)} \right) \left( e_{i} - p_{i}^{(2)} \right)^{T} B^{1/2} \quad \text{and} \quad \mathcal{B}_{1,ij} = B^{1/2} \left( e_{j} - p_{j}^{(1)} \right) \left( e_{i} - p_{i}^{(1)} \right)^{T} B^{1/2}.$$
(S146)

Then

$$\begin{split} \|\Phi_{i,j}\|_{F}^{2} &= \left\|m^{-1}\left(\pi(A)^{1/2}\otimes B^{1/2}\right)H^{ij}(J_{0},J_{1})\left(\pi(A)^{1/2}\otimes B^{1/2}\right)\right\|_{F}^{2}\\ &= m^{-2}\left\|\mathcal{A}_{0}\otimes\mathcal{B}_{2,ij}+\mathcal{A}_{1}\otimes\mathcal{B}_{1,ij}\right\|_{F}^{2}\\ &= m^{-2}\operatorname{tr}\left(\left(\mathcal{A}_{0}\otimes\mathcal{B}_{2,ij}+\mathcal{A}_{1}\mathcal{B}_{1,ij}\right)^{T}\left(\mathcal{A}_{0}\otimes\mathcal{B}_{2,ij}+\mathcal{A}_{1}\otimes\mathcal{B}_{1,ij}\right)\right)\\ &= m^{-2}\operatorname{tr}\left(\mathcal{A}_{0}^{T}\mathcal{A}_{0}\otimes\mathcal{B}_{2,ij}^{T}\mathcal{B}_{2,ij}\right)+m^{-2}\operatorname{tr}\left(\mathcal{A}_{1}^{T}\mathcal{A}_{1}\otimes\mathcal{B}_{1,ij}^{T}\mathcal{B}_{1,ij}\right)\\ &+m^{-2}\operatorname{tr}\left(\mathcal{A}_{0}^{T}\mathcal{A}_{1}\otimes\mathcal{B}_{2,ij}^{T}\mathcal{B}_{1,ij}\right)+m^{-2}\operatorname{tr}\left(\mathcal{A}_{1}^{T}\mathcal{A}_{0}\otimes\mathcal{B}_{1,ij}^{T}\mathcal{B}_{2,ij}\right). \end{split}$$
(S147)

We now find the traces of each of the terms in (S148). First, note that

$$\operatorname{tr}\left(\mathcal{A}_{0}^{T}\mathcal{A}_{0}\right) = \operatorname{tr}\left(\mathcal{I}_{0}\pi(A)\mathcal{I}_{0}\pi(A)\right) = \operatorname{tr}\left(A_{J_{0}}^{2}\right) = \left\|A_{J_{0}}\right\|_{F}^{2}.$$
(S148)

Analogously,

$$\operatorname{tr}\left(\mathcal{A}_{1}^{T}\mathcal{A}_{1}\right) = \left\|A_{J_{1}}\right\|_{F}^{2}.$$
(S149)

For the cross-term, let  $A_{J_0J_1}$  denote the  $m_0 \times m_1$  submatrix of  $\pi(A)$  given by columns of A in  $J_0$  and rows of A in  $J_1$ . Then

$$\operatorname{tr} \left( \mathcal{A}_{0}^{T} \mathcal{A}_{1} \right) = \operatorname{tr} \left( \mathcal{I}_{0} \pi(A) \mathcal{I}_{1} \pi(A) \right)$$
$$= \operatorname{tr} \left( \begin{bmatrix} 0_{m_{0} \times m_{0}} & A_{J_{0}J_{1}} \\ 0_{m_{1} \times m_{0}} & 0_{m_{1} \times m_{1}} \end{bmatrix} \pi(A) \right)$$
$$= \operatorname{tr} \left( A_{J_{0}J_{1}}^{T} A_{J_{0}J_{1}} \right)$$
$$= \| A_{J_{0}J_{1}} \|_{F}^{2}.$$
(S150)

Next,

$$\operatorname{tr}\left(\mathcal{B}_{1,ij}^{T}\mathcal{B}_{1,ij}\right) = \operatorname{tr}\left(B^{1/2}\left(e_{i}-p_{i}^{(1)}\right)\left(e_{j}-p_{j}^{(1)}\right)^{T}B\left(e_{j}-p_{j}^{(1)}\right)\left(e_{i}-p_{i}^{(1)}\right)^{T}B^{1/2}\right)$$
$$= \left(\left(e_{j}-p_{j}^{(1)}\right)^{T}B\left(e_{j}-p_{j}^{(1)}\right)\right)\left(\left(e_{i}-p_{i}^{(1)}\right)^{T}B\left(e_{i}-p_{i}^{(1)}\right)\right)$$
$$= \check{b}_{jj}\check{b}_{ii}.$$
(S151)

Analogously,

$$\operatorname{tr}\left(\mathcal{B}_{2,ij}^{T}\mathcal{B}_{2,ij}\right) = \left(\left(e_{j} - p_{j}^{(2)}\right)^{T}B\left(e_{j} - p_{j}^{(2)}\right)\right)\left(\left(e_{i} - p_{i}^{(2)}\right)^{T}B\left(e_{i} - p_{i}^{(2)}\right)\right)$$
$$= \widetilde{b}_{jj}\widetilde{b}_{ii}.$$
(S152)

The cross-terms yield

$$\operatorname{tr}\left(\mathcal{B}_{1,ij}^{T}\mathcal{B}_{2,ij}\right) = \left(\left(e_{j} - p_{j}^{(1)}\right)^{T} B\left(e_{j} - p_{j}^{(2)}\right)\right) \left(\left(e_{i} - p_{i}^{(2)}\right)^{T} B\left(e_{i} - p_{i}^{(1)}\right)\right) = \breve{b}_{ii}\breve{b}_{jj}.$$
 (S153)

The squared Frobenius norm of  $\Phi_{i,j}$  is

$$\begin{split} \|\Phi_{i,j}\|_{F}^{2} &= \frac{1}{m^{2}} \left( \|A_{J_{0}}\|_{F}^{2} \breve{b}_{ii}\breve{b}_{jj} + \|A_{J_{1}}\|_{F}^{2} \widetilde{b}_{ii}\widetilde{b}_{jj} + 2 \|A_{J_{0},J_{1}}\|_{F}^{2} \breve{b}_{ii}\breve{b}_{jj} \right) \\ &\leqslant \frac{1}{m^{2}} C \left( \|A_{J_{0}}\|_{F}^{2} + \|A_{J_{1}}\|_{F}^{2} + 2 \|A_{J_{0}J_{1}}\|_{F}^{2} \right) \|B\|_{2}^{2} \\ &= C \frac{1}{m^{2}} \|A\|_{F}^{2} \|B\|_{2}^{2}. \end{split}$$

We now apply the Hanson-Wright inequality,

$$\mathbb{P}\left(\left|(S_{\mathrm{I}})_{ij} - b_{ij}^{\sharp}\right| > r_{4}\right) = \mathbb{P}\left(\left|\operatorname{vec}\left\{Z\right\}^{T} \Phi_{i,j} \operatorname{vec}\left\{Z\right\} - \operatorname{tr}\left(\Phi_{i,j}\right)\right| > r_{4}\right)$$

$$\leq 2 \exp\left(-c \min\left\{d \log(m), d^{1/2} \sqrt{\log(m)} \frac{\|\Phi_{i,j}\|_{F}}{\|\Phi_{i,j}\|_{2}}\right\}\right)$$

$$\leq 2 \max\left(m^{-d}, \exp\left(d^{1/2} \sqrt{\log(m)} r^{1/2}(\Phi_{i,j})\right)\right).$$

The first step holds by (S139) and (S145).

## J Comparisons to related methods

The most similar existing method to ours is the sphering approach from Allen and Tibshirani (2012). Both methods use a preliminary demeaned version of the data to generate covariance estimates, then use these estimates to adjust the gene-wise t-tests. The largest difference between the procedures lies in this last step. The sphering approach produces an adjusted data set based on decorrelating residuals from a preliminary mean estimate and performs testing and mean estimation on this adjusted data using traditional OLS techniques. Though their approach is well-motivated at the population level, they do not provide theoretical support for their plug-in procedure, and in particular do not explore how noise in the initial mean estimate may complicate their decorrelation procedure. In contrast, our approach uses a generalized least squares approach motivated by classical statistical results including the Gauss Markov theorem.

The sphering approach also involves decorrelating a data matrix along both axes. Our work, including the theoretical analysis in Zhou (2014a), suggests that when the data matrix is non-square, attempting to decorrelate along the longer axis generally degrades performance. For genetics applications, where there are usually many more genes than samples, this suggests that decorrelating along the genes may hurt the performance of the sphering method. Fortunately, for gene-level analyses it is not necessary to decorrelate along the gene axis, since inference methods like false discovery rate are robust to a certain level of dependence among the variables (genes) (Benjamini and Yekutieli, 2001). Therefore, we also consider a modification of the sphering algorithm that only decorrelates along the samples. Confounder adjustment is another related line of work that deals with similar issues when attempting to discover mean differences. In particular, a part of that literature posits models where row-wise connections arise from the additive effects of potential latent variables. Sun et al. (2012) and Wang et al. (2015) use models of the form

$$X_{n \times m} = D_{n \times 1} \beta_{m \times 1}^T + Z_{n \times r} \Gamma_{m \times r}^T + E_{n \times m}$$
$$Z_{n \times r} = D_{n \times 1} \alpha_{r \times 1}^T + W_{n \times r}$$

where Z is an unobserved matrix of r latent factors. Rewriting these equations into the following form lets us better contrast the confounder model to our matrix-variate setup:

$$X = D(\beta + \Gamma \alpha)^T + W\Gamma^T + E.$$
(S154)

These models are generally estimated by using some form of factor analysis to estimate  $\Gamma$ and then using regression methods with additive outlier detection to identify  $\beta$ , methodology that is quite different from our GLS-based methods.

For the two-group model, in the case of a globally centered data matrix X, the design matrix D in (S155) takes the form

$$D_{n\times 1}^{T} = \begin{bmatrix} -1 & \cdots & -1 & 1 \cdots & 1 \end{bmatrix} = \begin{bmatrix} -1_{n_{1}}^{T} & 1_{n_{2}}^{T} \end{bmatrix},$$
 (S155)

and  $2\beta$  represents the vector of true mean differences between the groups. The vector  $\beta$  is estimated via OLS, yielding  $\hat{\beta}_{OLS}$ , and CATE considers whether the residual  $X - D_{n\times 1}\hat{\beta}_{OLS}$ has a low-rank covariance structure plus noise. If so,  $\hat{\Gamma}\hat{\alpha}$  aims to take out the residual lowrank structure through  $D(\widehat{\Gamma\alpha})^T$ . As illustrated in simulation and data analysis, this improves upon inference based only on  $\hat{\beta}_{OLS}$ . When applying the CATE and related methods to data originated from the generative model as described in the present paper, CATE (and in particular, the related LEAPP) method essentially seeks a sparse approximation of  $\hat{\beta}_{OLS}$ ; Moreover in LEAPP, this is essentially achieved via hard thresholding of coefficients of  $\hat{\beta}_{OLS}$ , leading to improvements in performance in variable selection and its subsequence inference when the vector of true mean differences is presumed to be sparse. In our setting, we improve upon OLS using GLS.

#### J.1 Simulation results

Figure S3 compares the performance of Algorithm 2 to the sphering method of Allen and Tibshirani (2012) and the robust regression confounder adjustment method of Wang et al. (2015) on simulated matrix variate data motivated by the ulcerative colitis dataset described in Section 5. Note that this robust regression confounder adjustment is a minor modification of the LEAPP algorithm introduced in Sun et al. (2012). As discussed above, we also consider a modification of Allen and Tibshirani (2012) that only decorrelates along the rows.

We can see that across a range of dataset sizes our method consistently outperforms sphering in terms of sensitivity and specificity for identifying mean differences. In some settings, CATE improves on Tsphere and t-statistics despite being applied on misspecified models, because CATE takes out the additional rank two structure from the mean after OLS regression and does some approximate thresholding on the coefficients. Our method using GLS performs significantly better than CATE in the setting of non-identity B, with edges present both within and between groups.

Figure S5 fixes the sample size and repeats these comparisons on different sample correlation structures (which are described in Section 4). Figure S6 is analogous to Figure S5, but with A as the identity matrix. Algorithm 2 is competitive or superior to the competing methods across a range of topologies.

#### J.2 Comparison on UC data

We apply both Algorithm 2 and CATE on the ulcerative colitis data to compare their respective findings on real data. Figure S7 presents the test statistics from these algorithms. The test statistics have a correlation of 0.75. As expected, both methods find that the bulk of genes have small test statistics. Note that the regression line of the CATE test statistics on Algorithm 2's test statistics has a slope less than 1. This implies that Algorithm 2 generates more dispersed test statistics than CATE, and, given that we have shown in Figures 5 and 8 that Algorithm 2 provides well-calibrated test statistics, that it also has more power in this situation.

Using a threshold of FDR adjusted p-values smaller than 0.1, both methods find four genes with significant mean differences. However, there is only one gene (DPP10-AS1) that



Figure S3: Performance of Algorithm 2 (GLS) relative to sphering and confounder adjustment methods, labeled as **tsphere** and **cate**, respectively. These are ROC curves for identifying true mean differences. An implementation of the sphering algorithm that does not adjust for A is also included, labeled as **tsphere\_noA**. Each panel shows the average ROC curves over 200 simulations. We simulate matrix variate data with gene correlations from an AR1(0.8) model and let s = 10 genes have true mean differences of 0.8, 0.6, and 0.4 for the first, second and third rows, respectively. For all of these the true B is set to  $\hat{B}$  from the ulcerative colitis data (using a repeated block structure for larger n values), described in Section 5 and evenly-sized groups are assigned randomly.



Figure S4: Performance of Algorithm 2 (GLS) relative to sphering and confounder adjustment, labeled as tsphere and cate, respectively. These are ROC curves for identifying true mean differences. An implementation of the sphering algorithm that does not adjust for A is also included, labeled as tsphere\_noA. Each panel shows the average ROC curves over 200 simulations. We simulate matrix variate data with no gene-wise correlations (A = I) and let s = 10 genes have true mean differences of 0.8, 0.6, and 0.4 for the first, second and third rows, respectively. For all of these the true B is set to  $\hat{B}$  from the ulcerative colitis data (using a repeated block structure for larger n values), described in Section 5 and evenly-sized groups are assigned randomly.

both methods identify. So, although there is significant correlation between the test statistics, the methods do not necessarily identify the same genes.



Figure S5: Performance of Algorithm 2 (GLS) relative to sphering and confounder adjustment, labeled as tsphere and cate, respectively. These are ROC curves for identifying true mean differences. An implementation of the sphering algorithm that does not adjust for A is also included, labeled as tsphere\_noA. Each panel shows the average ROC curves over 200 simulations. We simulate matrix variate data with an AR1(0.8) model for A and let s = 10genes have true mean differences of 0.8. B is constructed according to a Star-Block model with blocks of size 4, an AR1(0.8), and an Erdős-Rényi random graph with  $d = n \log n$  edges. All of these use n = 20 and randomly assign 10 observations to each group.



Figure S6: Performance of Algorithm 2 (GLS) relative to sphering and confounder adjustment, labeled as **tsphere** and **cate**, respectively. These are ROC curves for identifying true mean differences. An implementation of the sphering algorithm that does not adjust for A is also included, labeled as **tsphere\_noA**. Each panel shows the average ROC curves over 200 simulations. We simulate matrix variate data with no gene-wise correlations (A = I) and let s = 10 genes have true mean differences of 0.6. B is constructed according to a Star-Block model with blocks of size 4, an AR1(0.8), and an Erdős-Rényi random graph with  $d = n \log n$ edges. All of these use n = 40 and randomly assign 20 observations to each group.



Figure S7: Scatterplot of t-statistics for CATE and Algorithm 2 applied on the ulcerative colitis data. The 45-degree line is included in black while red dashed line is the linear fit.

## References

- ALLEN, G. I. and TIBSHIRANI, R. (2012). Inference with transposable data: modelling the effects of row and column correlations. *Journal of the Royal Statistical Society: Series B* (Statistical Methodology) **74** 721–743.
- BENJAMINI, Y. and YEKUTIELI, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of statistics* 1165–1188.
- ROTHMAN, A., BICKEL, P., LEVINA, E. and ZHU, J. (2008). Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics* **2** 494–515.
- SUN, Y., ZHANG, N. R. and OWEN, A. B. (2012). Multiple hypothesis testing adjusted for latent variables, with an application to the agemap gene expression data. *The Annals* of Applied Statistics 1664–1688.
- WANG, J., ZHAO, Q., HASTIE, T. and OWEN, A. B. (2015). Confounder adjustment in multiple hypothesis testing. arXiv preprint arXiv:1508.04178.
- ZHOU, S. (2014a). Gemini: Graph estimation with matrix variate normal instances. Annals of Statistics 42 532–562.
- ZHOU, S. (2014b). Supplement to "Gemini: Graph estimation with matrix variate normal instances". Annals of Statistics DOI:10.1214/13-AOS1187SUPP.
- ZHOU, S., RÜTIMANN, P., XU, M. and BÜHLMANN, P. (2011). High-dimensional covariance estimation based on Gaussian graphical models. *Journal of Machine Learning Research* 12 2975–3026.