# Supplement to "User-Friendly Covariance Estimation for Heavy-Tailed Distributions"

Yuan Ke,[*] Stanislav Minsker,[†] Zhao Ren,[‡] Qiang Sun[§] and Wen-Xin Zhou[¶]

In Sections A–C, we provide proofs of all the theoretical results in the main text. In addition, we investigate robust covariance estimation and inference under factor models in Section D, which might be of independent interest.

## A  Proof of Proposition 2.1

Without loss of generality we assume $\boldsymbol{\mu} = \boldsymbol{0}$. We construct a random vector $\boldsymbol{X} \in \mathbb{R}^d$ that follows the distribution below:

$$\mathbb{P}\Big\{\boldsymbol{X} = (0,\ldots,0,\underbrace{n\eta}_{j\text{th}},0,\ldots,0)^{\mathsf{T}}\Big\} = \mathbb{P}\Big\{\boldsymbol{X} = (0,\ldots,0,\underbrace{-n\eta}_{j\text{th}},0,\ldots,0)^{\mathsf{T}}\Big\} = \frac{\sigma^2}{2n^2\eta^2}$$

for each $j = 1,\ldots,d$, and

$$\mathbb{P}(\boldsymbol{X} = \boldsymbol{0}) = 1 - \frac{d\sigma^2}{n^2\eta^2}.$$

Here we assume $\eta^2 > d\sigma^2/n^2$ so that $\mathbb{P}(\boldsymbol{X} = \boldsymbol{0}) > 0$. In other words, the number of non-zero elements of $\boldsymbol{X}$ is at most 1. It is easy to see that the mean and covariance matrix of $\boldsymbol{X}$ are $\boldsymbol{0}$ and $\sigma^2 \mathbf{I}_d$, respectively.

Consider the empirical mean $\bar{\boldsymbol{X}} = (1/n)\sum_{i=1}^n \boldsymbol{X}_i$, where $\boldsymbol{X}_1,\ldots,\boldsymbol{X}_n$ are i.i.d. from $\boldsymbol{X}$. It follows that

$$\mathbb{P}(\|\bar{\boldsymbol{X}}\|_\infty \geq \eta) \geq \mathbb{P}\Big(\text{exactly one of the } n \text{ samples is not equal to } \boldsymbol{0}\Big)$$

$$= \frac{d\sigma^2}{n\eta^2}\Big(1 - \frac{d\sigma^2}{n^2\eta^2}\Big)^{n-1}.$$

Therefore, as long as $\delta < (2e)^{-1}$, the following bound

$$\|\bar{\boldsymbol{X}}\|_\infty \geq \sigma\sqrt{\frac{d}{n\delta}}\Big(1 - \frac{2e\delta}{n}\Big)^{(n-1)/2}$$

holds with probability at least $\delta$, as claimed. $\qquad\square$

---
[*]Department of Statistics, University of Georgia, Athens, GA 30602, USA. E-mail: yuan.ke@uga.edu.

[†]Department of Mathematics, University of Southern California, Los Angeles, CA 90089, USA. E-mail: minsker@usc.edu.

[‡]Department of Statistics, University of Pittsburgh, Pittsburgh, PA 15260, USA. E-mail: zren@pitt.edu.

[§]Department of Statistical Sciences, University of Toronto, Toronto, ON M5S 3G3, Canada. E-mail: qsun@utstat.toronto.edu.

[¶]Department of Mathematics, University of California, San Diego, La Jolla, CA 92093, USA. E-mail: wez243@ucsd.edu.

# B Proofs for Section 3

## B.1 Proof of Theorem 3.1

For each $1 \le k \le \ell \le d$, note that $\widehat{\sigma}_{1,k\ell}^{\mathcal{T}}$ is a $U$-statistic with a bounded kernel of order two, say $\widehat{\sigma}_{1,k\ell}^{\mathcal{T}} = \binom{n}{2}^{-1} \sum_{1 \le i < j \le n} h_{k\ell}(X_i, X_j)$. According to Hoeffding (1963), $\widehat{\sigma}_{1,k\ell}^{\mathcal{T}}$ can be represented as an average of (dependent) sums of independent random variables. Specifically, define

$$W(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n) = \frac{h_{k\ell}(\boldsymbol{x}_1, \boldsymbol{x}_2) + h_{k\ell}(\boldsymbol{x}_3, \boldsymbol{x}_4) + \cdots + h_{k\ell}(\boldsymbol{x}_{2m-1}, \boldsymbol{x}_{2m})}{m}$$

for $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n \in \mathbb{R}^d$. Let $\sum_{\mathcal{P}}$ denote the summation over all $n!$ permutations $(i_1, \ldots, i_n)$ of $[n] := \{1, \ldots, n\}$ and $\sum_C$ denote the summation over all $\binom{n}{2}$ pairs $(i_1, i_2)$ $(i_1 < i_2)$ from $[n]$. Then we have $m \sum_{\mathcal{P}} W(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n) = m2!(n-2)! \sum_C h_{k\ell}(\boldsymbol{x}_{i_1}, \boldsymbol{x}_{i_2})$ and hence

$$\widehat{\sigma}_{1,k\ell}^{\mathcal{T}} = \frac{1}{n!} \sum_{\mathcal{P}} W(X_{i_1}, \ldots, X_{i_n}). \tag{B.1}$$

Write $\tau = \tau_{k\ell}$ and $v = v_{k\ell}$ for simplicity. For any $y > 0$, by Markov's inequality, (B.1), convexity and independence, we derive that

$$\mathbb{P}(\widehat{\sigma}_{1,k\ell}^{\mathcal{T}} - \sigma_{k\ell} \ge y) \le e^{-(m/\tau)(y+\sigma_{k\ell})} \mathbb{E} e^{(m/\tau)\widehat{\sigma}_{1,k\ell}^{\mathcal{T}}}$$

$$\le e^{-(m/\tau)(y+\sigma_{k\ell})} \frac{1}{n!} \sum_{\mathcal{P}} \mathbb{E} e^{(1/\tau) \sum_{j=1}^{m} h_{k\ell}(X_{i_{2j-1}}, X_{i_{2j}})}$$

$$= e^{-(m/\tau)(y+\sigma_{k\ell})} \frac{1}{n!} \sum_{\mathcal{P}} \prod_{j=1}^{m} \mathbb{E} e^{(1/\tau) h_{k\ell}(X_{i_{2j-1}}, X_{i_{2j}})}.$$

Note that $h_{k\ell}(X_{i_{2j-1}}, X_{i_{2j}}) = \psi_{\tau}(Y_{\pi k} Y_{\pi \ell}/2) = \tau \psi_1(Y_{\pi k} Y_{\pi \ell}/(2\tau))$ for some $1 \le \pi \le N$. In addition, it is easy to verify the inequality that

$$-\log(1 - x + x^2) \le \psi_1(x) \le \log(1 + x + x^2) \text{ for all } x \in \mathbb{R}. \tag{B.2}$$

Therefore, we have

$$\mathbb{E} e^{(1/\tau) h_{k\ell}(X_{i_{2j-1}}, X_{i_{2j}})} \le \mathbb{E}\{1 + Y_{\pi k} Y_{\pi \ell}/(2\tau) + (Y_{\pi k} Y_{\pi \ell})^2/(2\tau)^2\}$$

$$= 1 + \sigma_{k\ell}/\tau + (1/\tau)^2 \mathbb{E}(Y_{\pi k} Y_{\pi \ell}/2)^2 \le e^{\sigma_{k\ell}/\tau + (v/\tau)^2}.$$

Combining the above calculations gives

$$\mathbb{P}(\widehat{\sigma}_{1,k\ell}^{\mathcal{T}} - \sigma_{k\ell} \ge y) \le e^{-(m/\tau)y + m(v/\tau)^2} = e^{-my^2/(4v^2)},$$

where the equality holds by taking $\tau = 2v^2/y$. Similarly, it can be shown that $\mathbb{P}(\widehat{\sigma}_{1,k\ell}^{\mathcal{T}} - \sigma_{k\ell} \le -y) \le e^{-my^2/(4v^2)}$. Consequently, for $\delta \in (0, 1)$, taking $y = 2v\sqrt{(2\log d + \log \delta^{-1})/m}$, or equivalently, $\tau = v\sqrt{m/(2\log d + \log \delta^{-1})}$, we arrive at

$$\mathbb{P}\left(|\widehat{\sigma}_{1,k\ell}^{\mathcal{T}} - \sigma_{k\ell}| \ge 2v\sqrt{\frac{2\log d + \log \delta^{-1}}{m}}\right) \le \frac{2\delta}{d^2}.$$

From the union bound it follows that

$$\mathbb{P}\left(\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{T}} - \boldsymbol{\Sigma}\|_{\max} > 2 \max_{1 \le k, \ell \le d} v_{k\ell} \sqrt{\frac{2\log d + \log \delta^{-1}}{n}}\right) \le (1 + d^{-1})\delta.$$

This proves (3.5). $\qquad\square$

2

## B.2 Proof of Theorem 3.2

To begin with, note that $\widehat{\Sigma}_2^{\mathcal{T}}$ can be written as a $U$-statistic of order 2. Define the index set $\mathcal{I} = \{(i,j) : 1 \le i < j \le n\}$ with cardinality $\binom{n}{2}$. Let $h(X_i, X_j) = (X_i - X_j)(X_i - X_j)^\intercal/2$ and $\mathbf{Z}_{i,j} = \tau^{-1}\psi_\tau(h(X_i, X_j)) = \psi_1(\tau^{-1}h(X_i, X_j))$, such that

$$\widetilde{\Sigma} := \frac{1}{\tau}\widehat{\Sigma}_2^{\mathcal{T}} = \frac{1}{\binom{n}{2}}\sum_{(i,j)\in\mathcal{I}} \mathbf{Z}_{i,j}.$$

We now rewrite the $U$-statistic $\widetilde{\Sigma}$ as a convex combination of sums of independent random matrices. As in the proof of Theorem 3.1, we define

$$\mathbf{W}_{(1,\dots,n)} = m^{-1}(\mathbf{Z}_{1,2} + \mathbf{Z}_{3,4} + \dots + \mathbf{Z}_{2m-1,2m}).$$

For every permutation $\pi = (i_1,\dots,i_n)$, we adopt the notation $\mathbf{W}_\pi = \mathbf{W}_{(i_1,\dots,i_n)}$ such that $\widetilde{\Sigma}^\tau = (n!)^{-1}\sum_{\pi\in\mathcal{P}} \mathbf{W}_\pi$. Using the convexity of the mappings $\mathbf{A} \mapsto \lambda_{\max}(\mathbf{A})$ and $x \mapsto e^x$, we obtain that

$$\exp\{\lambda_{\max}(\widetilde{\Sigma} - \Sigma^\tau)\} \le \frac{1}{n!}\sum_{\pi\in\mathcal{P}} \exp\{\lambda_{\max}(\mathbf{W}_\pi - \Sigma^\tau)\},$$

where $\Sigma^\tau := \tau^{-1}\Sigma$. Combined with Markov's inequality and the inequality $e^{\lambda_{\max}(\mathbf{A})} \le \operatorname{tr} e^{\mathbf{A}}$, this further implies

$$\mathbb{P}\{\sqrt{m}\,\lambda_{\max}(\widehat{\Sigma}_2^{\mathcal{T}} - \Sigma) \ge y\} = \mathbb{P}\left\{e^{\lambda_{\max}(m\widetilde{\Sigma} - m\Sigma^\tau)} \ge e^{y\sqrt{m}/\tau}\right\}$$

$$\le e^{-y\sqrt{m}/\tau}\frac{1}{n!}\sum_{\pi\in\mathcal{P}} \mathbb{E}\exp\{\lambda_{\max}(m\mathbf{W}_\pi - m\Sigma^\tau)\}$$

$$\le e^{-y\sqrt{m}/\tau}\frac{1}{n!}\sum_{\pi\in\mathcal{P}} \mathbb{E}\operatorname{tr}\exp(m\mathbf{W}_\pi - m\Sigma^\tau).$$

For every $\pi = (i_1,\dots,i_n) \in \mathcal{P}$, define $\mathbf{Z}_{\pi,j} = \mathbf{Z}_{i_{2j-1},i_{2j}}$ and $\mathbf{H}_{\pi,j} = h(X_{i_{2j-1}}, X_{i_{2j}})$, such that $\mathbf{Z}_{\pi,1},\dots,\mathbf{Z}_{\pi,m}$ are independent and $\mathbb{E}\mathbf{H}_{\pi,j} = \Sigma$. Then $\mathbf{W}_\pi$ can be written as $\mathbf{W}_\pi = m^{-1}(\mathbf{Z}_{\pi,1} + \dots + \mathbf{Z}_{\pi,m})$. Recall that $\psi_\tau(x) = \tau\psi_1(x/\tau)$. In view of (B.2), we have the matrix inequality

$$-\log(\mathbf{I} - \tau^{-1}\mathbf{H}_{\pi,j} + \tau^{-2}\mathbf{H}_{\pi,j}^2) \preceq \mathbf{Z}_{\pi,j} \preceq \log(\mathbf{I} + \tau^{-1}\mathbf{H}_{\pi,j} + \tau^{-2}\mathbf{H}_{\pi,j}^2).$$

Then we can bound $\mathbb{E}\exp\operatorname{tr}(m\mathbf{W}_\pi - m\Sigma^\tau)$ by

$$\mathbb{E}_{[m-1]}\mathbb{E}_m \operatorname{tr}\exp\left(\sum_{j=1}^{m-1}\mathbf{Z}_{\pi,j} - m\Sigma^\tau + \mathbf{Z}_{\pi,m}\right)$$

$$\le \mathbb{E}_{[m-1]}\mathbb{E}_m \operatorname{tr}\exp\left\{\sum_{j=1}^{m-1}\mathbf{Z}_{\pi,j} - m\Sigma^\tau + \log(\mathbf{I} + \tau^{-1}\mathbf{H}_{\pi,m} + \tau^{-2}\mathbf{H}_{\pi,m}^2)\right\}, \tag{B.3}$$

where the expectation $\mathbb{E}_m$ is taken with respect to $\{X_{i_{2m-1}}, X_{i_{2m}}\}$ and the expectation $\mathbb{E}_{[m-1]}$ is taken with respect to $\{X_{i_1}, \dots, X_{i_{2m-2}}\}$. To bound the right-hand side of (B.3), we follow a similar argument as in Minsker (2018). By Lieb's concavity theorem (see, e.g. Fact 2.5 in Minsker (2018)) and

Jensen's inequality, we arrive at

$$\mathbb{E}\operatorname{tr}\exp(m\mathbf{W}_\pi - m\boldsymbol{\Sigma}^\tau)$$

$$\leq \mathbb{E}\operatorname{tr}\exp\Big\{\sum_{j=1}^{m-1}\mathbf{Z}_{\pi,j} - m\boldsymbol{\Sigma}^\tau + \log(\mathbf{I} + \tau^{-1}\mathbb{E}\mathbf{H}_{\pi,m} + \tau^{-2}\mathbb{E}\mathbf{H}_{\pi,m}^2)\Big\}$$

$$\leq \operatorname{tr}\exp\Big\{\sum_{j=1}^{m}\log(\mathbf{I} + \tau^{-1}\mathbb{E}\mathbf{H}_{\pi,j} + \tau^{-2}\mathbb{E}\mathbf{H}_{\pi,j}^2) - m\boldsymbol{\Sigma}^\tau\Big\}$$

$$\leq \operatorname{tr}\exp\Big(\frac{1}{\tau^2}\sum_{j=1}^{m}\mathbb{E}\mathbf{H}_{\pi,j}^2\Big)$$

$$\leq d\exp(m\tau^{-2}\|\mathbb{E}\mathbf{H}_{\pi,1}^2\|_2)$$

$$= d\exp(m\tau^{-2}v^2),$$

where we used the bound $\operatorname{tr}e^{\mathbf{A}} \leq de^{\|\mathbf{A}\|}$ in the last inequality and the definition $v^2$ from (3.7) in the last equality.

Letting $\tau = 2v^2\sqrt{m}/y$, we get

$$\mathbb{P}\{\sqrt{m}\,\lambda_{\max}(\widehat{\boldsymbol{\Sigma}}_2^{\mathcal{T}} - \boldsymbol{\Sigma}) \geq y\} \leq d\exp\Big(-\frac{y\sqrt{m}}{\tau} + \frac{mv^2}{\tau^2}\Big) \leq de^{-y^2/(4v^2)}.$$

Similarly, it can be shown that

$$\mathbb{P}\{\sqrt{m}\,\lambda_{\min}(\widehat{\boldsymbol{\Sigma}}_2^{\mathcal{T}} - \boldsymbol{\Sigma}) \leq -y\} \leq de^{-y^2/(4v^2)}.$$

Finally, taking $y = 2v\sqrt{\log(2d) + \log\delta^{-1}}$ in the last two displays proves (3.9). □

## B.3  Proof of Theorem 3.3

Let $v_{\max} = \max_{1\leq k,\ell\leq d} v_{k\ell}$. By the union bound, for any $y > 0$ it holds

$$\mathbb{P}(\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\max} \geq v_{\max}\,y)$$

$$\leq \sum_{1\leq k\leq\ell\leq d}\mathbb{P}(|\widehat{\sigma}_{1,k\ell}^{\mathcal{H}} - \sigma_{k\ell}| \geq v_{k\ell}\,y) \leq \frac{d(d+1)}{2}\max_{1\leq k\leq\ell\leq d}\mathbb{P}(|\widehat{\sigma}_{1,k\ell}^{\mathcal{H}} - \sigma_{k\ell}| \geq v_{k\ell}\,y). \qquad (B.4)$$

In the rest of the proof, we fix $(k,\ell) \in [d]\times[d]$ and write $\tau = \tau_{k\ell}$ and $v = v_{k\ell}$ for simplicity. Moreover, define the index set $\mathcal{I} = \{(i,j) : 1 \leq i < j \leq n\}$, the collection $\{Z_{i,j} = (X_{ik} - X_{jk})(X_{i\ell} - X_{j\ell})/2 : (i,j) \in \mathcal{I}\}$ of random variables indexed by $\mathcal{I}$ and the loss function $\mathcal{L}(\theta) = \sum_{(i,j)\in\mathcal{I}}\ell_\tau(Z_{i,j} - \theta)$. With this notation, we have

$$\widehat{\sigma}_{1,k\ell}^{\mathcal{H}} = \widehat{\theta} := \operatorname*{argmin}_{\theta\in\mathbb{R}}\mathcal{L}(\theta).$$

Without loss of generality, we assume $\boldsymbol{\mu} = (\mu_1,\dots,\mu_d)^\top = \mathbf{0}$; otherwise, we can simply replace $X_{ik}$ by $X_{ik} - \mu_k$ for all $i \in [n]$ and $k \in [d]$.

Note that $\widehat{\theta}$ is the unique solution of the equation

$$\Psi(\theta) := \frac{1}{\binom{n}{2}}\sum_{(i,j)\in\mathcal{I}}\psi_\tau(Z_{i,j} - \theta) = 0,$$

where $\psi_\tau(\cdot)$ is defined in (3.1). Similarly to the proof of Theorem 3.1, we define

$$w_{(1,\dots,n)}(\theta) = \frac{1}{m\tau}\{\psi_\tau(Z_{1,2} - \theta) + \psi_\tau(Z_{3,4} - \theta) + \dots + \psi_\tau(Z_{2m-1,2m} - \theta)\}.$$

Denote by $\mathcal{P}$ the class of all $n!$ permutations on $[n]$ and let $\pi = (i_1, \dots, i_n)$ be a permutation, i.e., $\pi(j) = i_j$ for $j = 1, \dots, n$. Put $w_\pi(\theta) = w_{(i_1,\dots,i_n)}(\theta)$ for $\pi \in \mathcal{P}$, such that $\tau^{-1} m\Psi(\theta) = (n!)^{-1} \sum_{\pi \in \mathcal{P}} m w_\pi(\theta)$. By convexity, we have

$$\mathbb{E}\{e^{\tau^{-1} m\Psi(\theta)}\} \le \frac{1}{n!} \sum_{\pi \in \mathcal{P}} \mathbb{E}\{e^{m w_\pi(\theta)}\}.$$

Recall that $\mathbb{E}Z_{i,j} = \sigma_{k\ell}$ for any $(i, j) \in \mathcal{I}$. By (3.14),

$$v^2 = \mathrm{var}(Z_{1,2}) = \frac{1}{2}\{\mathbb{E}((X_k - \mu_k)^2(X_\ell - \mu_\ell)^2) + \sigma_{kk}\sigma_{\ell\ell}\}.$$

For $\pi = (1, \dots, n)$, by (B.2) and the fact that $\tau^{-1}\psi_\tau(x) = \psi_1(x/\tau)$, we have

$$
\begin{aligned}
\mathbb{E}\{e^{m w_\pi(\theta)}\} &= \prod_{j=1}^m \mathbb{E}\exp\{\psi_1((Z_{2j-1,2j} - \theta)/\tau)\} \\
&\le \prod_{j=1}^m \mathbb{E}\{1 + \tau^{-1}(Z_{2j-1,2j} - \theta) + \tau^{-2}(Z_{2j-1,2j} - \theta)^2\} \\
&\le \prod_{j=1}^m [1 + \tau^{-1}(\sigma_{k\ell} - \theta) + \tau^{-2}\{v^2 + (\sigma_{k\ell} - \theta)^2\}] \\
&\le \exp[m\tau^{-1}(\sigma_{k\ell} - \theta) + m\tau^{-2}\{v^2 + (\sigma_{k\ell} - \theta)^2\}]. \quad\quad\quad\text{(B.5)}
\end{aligned}
$$

Similarly, it can be shown that

$$\mathbb{E}\{-e^{m w_\pi(\theta)}\} \le \exp[-m\tau^{-1}(\sigma_{k\ell} - \theta) + m\tau^{-2}\{v^2 + (\sigma_{k\ell} - \theta)^2\}]. \quad\quad\quad\text{(B.6)}$$

Inequalities (B.5) and (B.6) hold for every permutation $\pi \in \mathcal{P}$. For $\eta \in (0, 1)$, define

$$
\begin{aligned}
B_+(\theta) &= \sigma_{k\ell} - \theta + \frac{v^2 + (\sigma_{k\ell} - \theta)^2}{\tau} + \frac{\tau \log \eta^{-1}}{m}, \\
B_-(\theta) &= \sigma_{k\ell} - \theta - \frac{v^2 + (\sigma_{k\ell} - \theta)^2}{\tau} - \frac{\tau \log \eta^{-1}}{m}.
\end{aligned}
$$

Together, (B.5), (B.6) and Markov's inequality imply

$$\mathbb{P}\{\Psi(\theta) > B_+(\theta)\} \le e^{-\tau^{-1} m B_+(\theta)} \mathbb{E}\{e^{\tau^{-1} m\Psi(\theta)}\} \le \eta,$$

$$\text{and } \mathbb{P}\{\Psi(\theta) < B_-(\theta)\} \le e^{-\tau^{-1} m B_-(\theta)} \mathbb{E}\{-e^{\tau^{-1} m\Psi(\theta)}\} \le \eta.$$

Recall that $\Psi(\widehat{\theta}) = 0$. Let $\theta_+$ be the smallest solution of the quadratic equation $B_+(\theta_+) = 0$, and $\theta_-$ be the largest solution of the equation $B_-(\theta_-) = 0$. Noting that $\Psi(\cdot)$ is decreasing, it follows from the last display that

$$\mathbb{P}(\theta_- \le \widehat{\theta} \le \theta_+) \ge 1 - 2\eta.$$

Similarly to the proof of Proposition 2.4 in Catoni (2012), it can be shown that with $\tau = v\sqrt{m/\log\eta^{-1}}$,

$$\theta_+ \leq \sigma_{k\ell} + 2\left(\frac{v^2}{\tau} + \frac{\tau\log\eta^{-1}}{m}\right) \text{ and } \theta_- \geq \sigma_{k\ell} - 2\left(\frac{v^2}{\tau} + \frac{\tau\log\eta^{-1}}{m}\right)$$

as long as $m \geq 8\log\eta^{-1}$. Consequently, we obtain that with probability at least $1 - 2\eta$, $|\widehat{\sigma}_{1,k\ell}^{\mathcal{H}} - \sigma_{k\ell}| \leq 4v\sqrt{(\log\eta^{-1})/m}$.

Taking $y = 4\sqrt{(\log\eta^{-1})/m}$ in (B.4) yields $\|\widehat{\Sigma}_1^{\mathcal{H}} - \Sigma\|_{\max} \leq 4v_{\max}\sqrt{(\log\eta^{-1})/m}$ with probability at least $1 - d(d+1)\eta$. Finally, taking $\delta = d^2\eta$ proves (3.16). $\qquad\square$

## B.4  Proof of Theorem 3.5

We will use Theorem 1 and Lemma 1 in Minsker and Strawn (2017) that connect the performance of $\widehat{\sigma}_{\ell m}^{\mathrm{MOM}}$ to the rate of convergence of $\widehat{\sigma}_{\ell m}^{(1)}$ to the normal law. It is well known that, whenever the 4th moments of the entries of $X$ are finite, $\sqrt{|G_1|}\frac{\widehat{\sigma}_{\ell m}^{(1)} - \sigma_{\ell m}}{\Delta_{\ell m}}$ converges in distribution to the standard normal distribution. The rate of this convergence can be obtained via an analogue of the Berry-Esseen theorem for the sample covariance. Specifically, for any $1 \leq \ell, m \leq d$, we seek an upper bound on

$$\sup_{t\in\mathbb{R}}\left|\mathbb{P}\left(\sqrt{|G_1|}\frac{\widehat{\sigma}_{\ell m}^{(1)} - \sigma_{\ell m}}{\Delta_{\ell m}} \leq t\right) - \mathbb{P}(Z \leq t)\right|,$$

where $Z \sim \mathcal{N}(0,1)$. To this end, we will use Theorem 2.9 in Pinelis and Molzon (2016). Using the notation therein, we take $V = (X_\ell - \mathbb{E}X_\ell, X_m - \mathbb{E}X_m, X_\ell X_m - \mathbb{E}(X_\ell X_m))^\intercal$, $f(x_1, x_2, x_3) = x_3 - x_1 \cdot x_2$, and deduce that

$$\sup_{t\in\mathbb{R}}\left|\mathbb{P}\left(\sqrt{|G_1|}\frac{\widehat{\sigma}_{\ell m}^{(1)} - \sigma_{\ell m}}{\Delta_{\ell m}} \leq t\right) - \mathbb{P}(Z \leq t)\right| \leq \frac{\mathfrak{C}_{\ell m}}{\sqrt{|G_1|}}, \tag{B.7}$$

where $\mathfrak{C}_{\ell m} > 0$ is a constant depending on $\Delta_{\ell m}$ and $\mathbb{E}|(X_\ell - \mathbb{E}X_\ell)(X_m - \mathbb{E}X_m)|^3$. Together with Theorem 1 and Lemma 1 of Minsker and Strawn (2017), (B.7) implies that

$$\left|\widehat{\sigma}_{\ell m}^{\mathrm{MOM}} - \sigma_{\ell m}\right| \leq 3\Delta_{\ell m}\sqrt{\frac{k}{n}}\left(\sqrt{\frac{s}{k}} + \mathfrak{C}_{\ell m}\sqrt{\frac{k}{n}}\right)$$

with probability at least $1 - 4e^{-2s}$ for all $s > 0$ satisfying

$$\sqrt{\frac{s}{k}} + \mathfrak{C}_{\ell m}\sqrt{\frac{k}{n}} \leq 0.33. \tag{B.8}$$

Taking the union bound over all $\ell, m$, we obtain that with probability at least $1 - 2d(d+1)e^{-2s}$,

$$\|\widehat{\Sigma}^{\mathrm{MOM}} - \Sigma\|_{\max} \leq 3\max_{\ell,m}\Delta_{\ell m}\left(\sqrt{\frac{s}{n}} + \max_{\ell,m}\mathfrak{C}_{\ell m}\frac{k}{n}\right)$$

for all $s > 0$ satisfying (B.8). The latter is equivalent to the statement of the theorem. $\qquad\square$

## B.5  Proof of Corollary 3.1

From the proof of Theorem 3.2, we find that

$$\|\mathbb{E}\{(X_1 - X_2)(X_1 - X_2)^\top\}^2\|_2 = 2\|\mathbb{E}\{(X - \mu)(X - \mu)^\top\}^2 + \mathrm{Tr}(\Sigma)\Sigma + 2\Sigma^2\|_2.$$

Under the bounded kurtosis condition that $K = \sup_{u \in \mathbb{S}^{d-1}} \mathrm{kurt}(u^\top X) < \infty$, it follows from Lemma 4.1 in Minsker and Wei (2018) that

$$\|\mathbb{E}\{(X - \mu)(X - \mu)^\top\}^2\|_2 \le K\mathrm{Tr}(\Sigma)\|\Sigma\|_2.$$

Together, the last two displays imply

$$\|\mathbb{E}\{(X_1 - X_2)(X_1 - X_2)^\top\}^2\|_2 \le 2\|\Sigma\|_2\{(K + 1)\mathrm{Tr}(\Sigma) + 2\|\Sigma\|_2\}.$$

Taking $v = \|\mathbb{E}\{(X_1 - X_2)(X_1 - X_2)^\top\}^2\|_2^{1/2}/2$ that scales with $\mathrm{Tr}(\Sigma)^{1/2}\|\Sigma\|_2^{1/2} = \mathrm{r}(\Sigma)^{1/2}\|\Sigma\|_2$, the resulting estimator satisfies

$$\|\widehat{\Sigma}_2^{\mathcal{T}} - \Sigma\|_2 \lesssim K^{1/2}\|\Sigma\|_2 \sqrt{\frac{\mathrm{r}(\Sigma)(\log d + t)}{n}} \tag{B.9}$$

with probability at least $1 - e^{-t}$. $\qquad\square$

# C  Proofs for Section 5

## C.1  Proof of Theorem 5.1

Define each principal submatrix of $\Sigma$ as $\Sigma^{(p,q)} = \mathbb{E}Z_1^{(p,q)}Z_1^{(p,q)\top}/2$, which is estimated by $\widehat{\Sigma}_2^{(p,q),\mathcal{T}}$. As a result, we expect the final estimator $\widehat{\Sigma}_q$ to be close to

$$\Sigma_q = \sum_{j=-1}^{\lceil (d-1)/q \rceil} \mathbf{E}_{jq+1}^d(\Sigma^{(jq+1,2q)}) - \sum_{j=0}^{\lceil (d-1)/q \rceil} \mathbf{E}_{jq+1}^d(\Sigma^{(jq+1,q)}).$$

By the triangle inequality, we have $\|\widehat{\Sigma}_q - \Sigma\|_2 \le \|\widehat{\Sigma}_q - \Sigma_q\|_2 + \|\Sigma_q - \Sigma\|_2$. We first establish an upper bound for the bias term $\|\Sigma_q - \Sigma\|_2$. According to the decomposition illustrated by Figure 2, $\Sigma_q$ is a banded version of the population covariance with bandwidth between $q$ and $2q$. Therefore, we bound the spectral norm of $\Sigma_q - \Sigma$ by the $\|\cdot\|_{1,1}$ norm as follows:

$$\|\Sigma_q - \Sigma\|_2 \le \max_{1 \le \ell \le d} \sum_{k:|k-\ell|>q} |\sigma_{k\ell}| \le \frac{M}{q^\alpha}.$$

It remains to control the estimation error $\|\widehat{\Sigma}_q - \Sigma_q\|_2$. Define $\mathbf{D}^{(p,q)} = \widehat{\Sigma}_2^{(p,q),\mathcal{T}} - \Sigma^{(p,q)}$,

$$\mathbf{S}_1 = \sum_{j=-1:j \text{ is odd}}^{\lceil (d-1)/q \rceil} \mathbf{E}_{jq+1}^d\{\mathbf{D}^{(jq+1,2q)}\}, \quad \mathbf{S}_2 = \sum_{j=0:j \text{ is even}}^{\lceil (d-1)/q \rceil} \mathbf{E}_{jq+1}^d\{\mathbf{D}^{(jq+1,2q)}\},$$

and $\mathbf{S}_3 = \sum_{j=0}^{\lceil (d-1)/q \rceil} \mathbf{E}_{jq+1}^d\{\mathbf{D}^{(jq+1,q)}\}$. Note that each $\mathbf{S}_i$ above is a sum of disjoint block diagonal matrices. Therefore,

$$\|\widehat{\Sigma}_q - \Sigma_q\|_2 \le \|\mathbf{S}_1\|_2 + \|\mathbf{S}_3\|_2 + \|\mathbf{S}_3\|_2$$
$$\le 3 \max_{j=-1}^{\lceil (d-1)/q \rceil} \{\|\mathbf{D}^{(jq+1,2q)}\|_2, \|\mathbf{D}^{(jq+1,q)}\|_2\}. \tag{C.1}$$

Applying Theorem 3.2 to each principal submatrix with the choice of $\delta = (n^{c_0}d)^{-1}$ in $\tau$, and by the union bound, we obtain that with probability at least $1 - 2d\delta = 1 - 2n^{-c_0}$,

$$\max_{j=-1}^{\lceil (d-1)/q \rceil} \{\|\mathbf{D}^{(jq+1,2q)}\|_2, \|\mathbf{D}^{(jq+1,q)}\|_2\}$$

$$\leq 2\|\mathbf{\Sigma}\|_2^{1/2} \{(M_1 + 1)q\|\mathbf{\Sigma}\|_2 + \|\mathbf{\Sigma}\|_2\}^{1/2} \sqrt{\frac{\log(4q) + \log \delta^{-1}}{m}}$$

$$\leq 2M_0 \sqrt{1 + (M_1 + 1)q} \sqrt{\frac{\log(4d) + c_0 \log(nd)}{n}},$$

where we used the inequalities $\text{tr}(\mathbf{D}^{(jq+1,2q)}) \leq 2q\|\mathbf{\Sigma}\|_2$ and $\|\mathbf{\Sigma}\|_2 \leq M_0$. Plugging this into (C.1), we obtain that with probability at least $1 - 2n^{-c_0}$,

$$\|\widehat{\mathbf{\Sigma}}_q - \mathbf{\Sigma}_q\|_2 \leq 6M_0 \sqrt{1 + (M_1 + 1)q} \sqrt{\frac{\log(4d) + c_0 \log(nd)}{n}}.$$

In view of the upper bounds on $\|\widehat{\mathbf{\Sigma}}_q - \mathbf{\Sigma}_q\|_2$ and $\|\mathbf{\Sigma}_q - \mathbf{\Sigma}\|_2$, the optimal bandwidth $q$ is of order $\{n/\log(nd)\}^{1/(2\alpha+1)} \wedge d$, which leads to the desired result. $\quad\square$

## C.2   Proof of Theorem 5.3

Define the symmetrized Bregman divergence for the loss function $\mathcal{L}(\mathbf{\Theta}) = \langle \mathbf{\Theta}^2, \widehat{\mathbf{\Sigma}}_1^T \rangle - \text{tr}(\mathbf{\Theta})$ as $D_{\mathcal{L}}^s(\mathbf{\Theta}_1, \mathbf{\Theta}_2) = \langle \nabla \mathcal{L}(\mathbf{\Theta}_1) - \nabla \mathcal{L}(\mathbf{\Theta}_2), \mathbf{\Theta}_1 - \mathbf{\Theta}_2 \rangle$. We first need the following two lemmas.

**Lemma C.1.** Provided $\lambda \geq 2\|\nabla \mathcal{L}(\mathbf{\Theta}^*)\|_{\max}$, $\widehat{\mathbf{\Theta}}$ falls in the $\ell_1$-cone

$$\|\widehat{\mathbf{\Theta}}_{\mathcal{S}^c} - \mathbf{\Theta}_{\mathcal{S}^c}^*\|_{\ell_1} \leq 3\|\widehat{\mathbf{\Theta}}_{\mathcal{S}} - \mathbf{\Theta}_{\mathcal{S}}^*\|_{\ell_1}.$$

*Proof of Lemma C.1.* Set $\widehat{\mathbf{\Gamma}} = (\widehat{\Gamma}_{k\ell})_{1 \leq k,\ell \leq d} \in \mathbb{R}^{d \times d}$, where $\widehat{\Gamma}_{k\ell} = \partial|\widehat{\Theta}_{k\ell}| \in [1, 1]$ whenever $k \neq \ell$, and $\widehat{\Gamma}_{k\ell} = 0$ whenever $k = \ell$. Here $\partial f(x_0)$ denotes the subdifferential of $f$ at $x_0$. By the convexity of the loss function and the optimality condition, we have

$$0 \leq \langle \nabla \mathcal{L}(\widehat{\mathbf{\Theta}}) - \nabla \mathcal{L}(\mathbf{\Theta}^*), \widehat{\mathbf{\Theta}} - \mathbf{\Theta}^* \rangle$$

$$= \langle -\lambda \widehat{\mathbf{\Gamma}} - \nabla \mathcal{L}(\mathbf{\Theta}^*), \widehat{\mathbf{\Theta}} - \mathbf{\Theta}^* \rangle$$

$$= -\langle \lambda \widehat{\mathbf{\Gamma}}, \widehat{\mathbf{\Theta}} - \mathbf{\Theta}^* \rangle - \langle \nabla \mathcal{L}(\mathbf{\Theta}^*), \widehat{\mathbf{\Theta}} - \mathbf{\Theta}^* \rangle$$

$$\leq -\lambda \|\widehat{\mathbf{\Theta}}_{\mathcal{S}^c} - \mathbf{\Theta}_{\mathcal{S}^c}^*\|_{\ell_1} + \lambda \|\widehat{\mathbf{\Theta}}_{\mathcal{S}} - \mathbf{\Theta}_{\mathcal{S}}^*\|_{\ell_1} + \frac{\lambda}{2}\|\widehat{\mathbf{\Theta}} - \mathbf{\Theta}^*\|_{\ell_1}.$$

Rearranging terms proves the stated result. $\quad\square$

**Lemma C.2.** Under the restricted eigenvalue condition, it holds

$$D_{\mathcal{L}}^s(\widehat{\mathbf{\Theta}}, \mathbf{\Theta}^*) \geq \kappa_- \|\widehat{\mathbf{\Theta}} - \mathbf{\Theta}^*\|_F^2.$$

*Proof.* We use $\text{vec}(\mathbf{A})$ to denote the vectorized form of matrix $\mathbf{A}$. Let $\mathbf{\Delta} = \widehat{\mathbf{\Theta}} - \mathbf{\Theta}^*$. Then by the mean value theorem, there exists a $\gamma \in [0, 1]$ such that

$$D_{\mathcal{L}}^s(\widehat{\mathbf{\Theta}}, \mathbf{\Theta}^*) = \langle \nabla \mathcal{L}(\widehat{\mathbf{\Theta}}) - \nabla \mathcal{L}(\mathbf{\Theta}^*), \widehat{\mathbf{\Theta}} - \mathbf{\Theta}^* \rangle$$

$$= \text{vec}(\widehat{\mathbf{\Theta}} - \mathbf{\Theta}^*)^\top \nabla^2 \mathcal{L}(\widehat{\mathbf{\Theta}} + \gamma \mathbf{\Delta}) \text{vec}(\widehat{\mathbf{\Theta}} - \mathbf{\Theta}^*)$$

$$\geq \kappa_- \|\mathbf{\Delta}\|_F^2,$$

where the last step is due to the restricted eigenvalue condition and Lemma C.1. This completes the proof. $\qquad\square$

Applying Lemma C.2 gives

$$\kappa_-\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_F^2 \leq \langle \nabla\mathcal{L}(\widehat{\boldsymbol{\Theta}}) - \nabla\mathcal{L}(\boldsymbol{\Theta}^*), \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^* \rangle. \tag{C.2}$$

Next, note that the sub-differential of the norm $\|\cdot\|_{\ell_1}$ evaluated at $\boldsymbol{\Psi} = (\Psi_{k\ell})_{1\leq k,\ell\leq d}$ consists the set of all symmetric matrices $\boldsymbol{\Gamma} = (\Gamma_{k\ell})_{1\leq k,\ell\leq d}$ such that $\Gamma_{k\ell} = 0$ if $k = \ell$, $\Gamma_{k\ell} = \text{sign}(\Psi_{k\ell})$ if $k \neq \ell$ and $\Psi_{k\ell} \neq 0$, $\Gamma_{k\ell} \in [-1, +1]$ if $k \neq \ell$ and $\Psi_{k\ell} = 0$. Then by the Karush-Kuhn-Tucker conditions, there exists some $\widehat{\boldsymbol{\Gamma}} \in \partial\|\widehat{\boldsymbol{\Theta}}\|_{\ell_1}$ such that

$$\nabla\mathcal{L}(\widehat{\boldsymbol{\Theta}}) + \lambda\widehat{\boldsymbol{\Gamma}} = \mathbf{0}.$$

Plugging the above equality into (C.2) and rearranging terms, we obtain

$$\kappa_-\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_F^2 + \underbrace{\langle \nabla\mathcal{L}(\boldsymbol{\Theta}^*), \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^* \rangle}_{\text{I}} + \underbrace{\langle \lambda\widehat{\boldsymbol{\Gamma}}, \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^* \rangle}_{\text{II}} \leq 0. \tag{C.3}$$

We bound terms I and II separately, starting with I. Our first observation is

$$\nabla\mathcal{L}(\boldsymbol{\Theta}^*) = (\boldsymbol{\Theta}^*\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{T}} - \mathbf{I})/2 + (\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{T}}\boldsymbol{\Theta}^* - \mathbf{I})/2.$$

By Theorem 3.1, we obtain that with probability at least $1 - 2\delta$,

$$\|\nabla\mathcal{L}(\boldsymbol{\Theta}^*)\|_{\max} \leq \|\boldsymbol{\Theta}^*\|_{1,1}\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{T}} - \boldsymbol{\Sigma}\|_{\max} \leq 2M\|\mathbf{V}\|_{\max}\sqrt{\frac{2\log d + \log\delta^{-1}}{\lfloor n/2\rfloor}} \leq \lambda/2.$$

Let $\mathcal{S}$ be the support of nonzero elements of $\boldsymbol{\Theta}^*$ and $\mathcal{S}^c$ be its complement with respect to the full index set $\{(k,\ell) : 1 \leq k, \ell \leq d\}$. For term I, separating the support of $\nabla\mathcal{L}(\boldsymbol{\Theta}^*)$ and $\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*$ to $\mathcal{S}$ and $\mathcal{S}^c$ and applying the matrix Hölder inequality, we obtain

$$\begin{aligned}
\langle \nabla\mathcal{L}(\boldsymbol{\Theta}^*), \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^* \rangle &= \langle (\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}}, (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}} \rangle + \langle (\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}^c}, (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}^c} \rangle \\
&\geq -\|(\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}}\|_F\|(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}}\|_F - \|(\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}^c}\|_F\|(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}^c}\|_F.
\end{aligned}$$

For term II, separating the support of $\lambda\widehat{\boldsymbol{\Gamma}}$ and $\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*$ to $\mathcal{S}$ and $\mathcal{S}^c$, we have

$$\langle \lambda\widehat{\boldsymbol{\Gamma}}, \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^* \rangle = \langle \lambda\widehat{\boldsymbol{\Gamma}}_{\mathcal{S}}, (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}} \rangle + \langle \lambda\widehat{\boldsymbol{\Gamma}}_{\mathcal{S}^c}, (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}^c} \rangle. \tag{C.4}$$

Let $1_{\mathcal{A}} \in \mathbb{R}^{d\times d}$ be a $d$-by-$d$ matrix such that $1_{k\ell} = 1$ if $(k,\ell) \in \mathcal{A}$, $1_{k\ell} = 0$ otherwise. For the last term in the above equality, we have

$$\langle \lambda\widehat{\boldsymbol{\Gamma}}_{\mathcal{S}^c}, (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta})_{\mathcal{S}^c} \rangle = \langle \lambda \cdot 1_{\mathcal{S}^c}, |\widehat{\boldsymbol{\Theta}}_{\mathcal{S}^c}| \rangle = \langle \lambda \cdot 1_{\mathcal{S}^c}, |(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta})_{\mathcal{S}^c}| \rangle. \tag{C.5}$$

Plugging (C.5) into (C.4) and applying the matrix Hölder inequality yields

$$\begin{aligned}
\langle \lambda\widehat{\boldsymbol{\Gamma}}, \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^* \rangle &= \langle \lambda\widehat{\boldsymbol{\Gamma}}_{\mathcal{S}}, (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}} \rangle + \langle \lambda \cdot 1_{\mathcal{S}^c}, |(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}^c}| \rangle \\
&= \langle \lambda\widehat{\boldsymbol{\Gamma}}_{\mathcal{S}}, (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}} \rangle + \|\lambda \cdot 1_{\mathcal{S}^c}\|_F\|(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}^c}\|_F \\
&\geq -\|\lambda \cdot 1_{\mathcal{S}}\|_F\|(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}}\|_F + \lambda\sqrt{s}\|(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}^c}\|_F.
\end{aligned}$$

9

Plugging the bounds for I and II back into (C.3), we find

$$\kappa_-\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_{\mathrm{F}}^2 + (\|\lambda \cdot 1_{\mathcal{S}^c}\|_{\mathrm{F}} - \|(\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}^c}\|_{\mathrm{F}})\|(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*)_{\mathcal{S}^c}\|_{\mathrm{F}}$$
$$\leq (\|(\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}}\|_{\mathrm{F}} + \|\lambda \cdot 1_{\mathcal{S}}\|_{\mathrm{F}})\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_{\mathrm{F}}.$$

Since $\|(\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}^c}\|_{\mathrm{F}} \leq |\mathcal{S}^c|^{1/2}\|\nabla\mathcal{L}(\boldsymbol{\Theta}^*)\|_{\max} \leq |\mathcal{S}^c|^{1/2}\lambda = \|\lambda \cdot 1_{\mathcal{S}^c}\|_{\mathrm{F}}$, it follows that

$$\kappa_-\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_{\mathrm{F}}^2 \leq (\|(\nabla\mathcal{L}(\boldsymbol{\Theta}^*))_{\mathcal{S}}\|_{\mathrm{F}} + \|\lambda \cdot 1_{\mathcal{S}}\|_{\mathrm{F}})\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_{\mathrm{F}}.$$

Canceling $\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_{\mathrm{F}}$ on both sides yields

$$\kappa_-\|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_{\mathrm{F}} \leq \|\lambda \cdot 1_S\|_{\mathrm{F}} + \|\nabla\mathcal{L}(\boldsymbol{\Theta}^*)_{\mathcal{S}}\|_{\mathrm{F}} \leq 3\lambda\sqrt{s}/2$$

under the scaling $\lambda \geq 2\|\nabla\mathcal{L}(\boldsymbol{\Theta}^*)\|_{\max}$. Plugging $\lambda$ completes the proof. □

# D Robust estimation and inference under factor models

As a complement to the three examples considered in the main text, in this section we discuss robust covariance estimation (Section D.1) and inference (Section D.2) under factor models, which might be of independent interest. In Section D.2, we provide a self-contained analysis to prove the consistency of estimating the false discovery proportion, while there is no such a theoretical guarantee in Fan et al. (2019) without using sample splitting.

## D.1 Covariance estimation through factor models

Consider the approximate factor model of the form $X = (X_1, \ldots, X_d)^\mathsf{T} = \boldsymbol{\mu} + \mathbf{B}\boldsymbol{f} + \boldsymbol{\varepsilon}$, from which we observe

$$X_i = (X_{i1}, \ldots, X_{id})^\mathsf{T} = \boldsymbol{\mu} + \mathbf{B}\boldsymbol{f}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \ldots, n, \tag{D.1}$$

where $\boldsymbol{\mu}$ is a $d$-dimensional unknown mean vector, $\mathbf{B} = (\boldsymbol{b}_1, \ldots, \boldsymbol{b}_d)^\mathsf{T} \in \mathbb{R}^{d \times r}$ is the factor loading matrix, $\boldsymbol{f}_i \in \mathbb{R}^r$ is a vector of common factors to the $i$th observation and is independent of the idiosyncratic noise $\boldsymbol{\varepsilon}_i$. For more details about factor analysis, we refer the readers to Anderson and Rubin (1956), Chamberlain and Rothschild (1983), Bai and Li (2012) and Fan and Han (2017), among others. Factor pricing model has been widely used in financial economics, where $X_{ik}$ is the excess return of fund/asset $k$ at time $i$, $\boldsymbol{f}_i$'s are the systematic risk factors related to some specific linear pricing model, such as the capital asset pricing model (CAPM) (Sharpe, 1964), and the Fama-French three-factor model (Fama and French, 1993).

Under model (D.1), the covariance matrix of $X$ can be written as

$$\boldsymbol{\Sigma} = (\sigma_{k\ell})_{1 \leq k, \ell \leq d} = \mathbf{B}\mathrm{cov}(\boldsymbol{f})\mathbf{B}^\mathsf{T} + \boldsymbol{\Sigma}_{\varepsilon}, \tag{D.2}$$

where $\boldsymbol{\Sigma}_{\varepsilon} = (\sigma_{\varepsilon,k\ell})_{1 \leq k, \ell \leq d}$ denotes the covariance matrix of $\boldsymbol{\varepsilon} = (\varepsilon_1, \ldots \varepsilon_d)^\mathsf{T}$, which is typically assumed to be sparse. When $\boldsymbol{\Sigma}_{\varepsilon} = \mathbf{I}_d$, model (D.1) is known as the strict factor model. To make the model identifiable, following Bai and Li (2012) we assume that $\mathrm{cov}(\boldsymbol{f}) = \mathbf{I}_r$ and that the columns of $\mathbf{B}$ are orthogonal.

We consider the robust estimation of $\mathbf{\Sigma}$ based on independent observations $X_1, \ldots, X_n$ from model (D.1). By (D.2) and the identifiability condition, $\mathbf{\Sigma}$ is comprised of two components: the low-rank component $\mathbf{B}\mathbf{B}^\mathsf{T}$ and the sparse component $\mathbf{\Sigma}_\varepsilon$. Using a pilot robust covariance estimator $\widehat{\mathbf{\Sigma}}_1^{\mathcal{T}}$ given in (3.3) or $\widehat{\mathbf{\Sigma}}_1^{\mathcal{H}}$ given in (3.13), we propose the following robust version of the principal orthogonal complement thresholding (POET) procedure (Fan, Liao and Mincheva, 2013):

(i) Let $\widehat{\lambda}_1 \geq \widehat{\lambda}_2 \geq \cdots \geq \widehat{\lambda}_r$ be the top $r$ eigenvalues of $\widehat{\mathbf{\Sigma}}_1^{\mathcal{H}}$ (or $\widehat{\mathbf{\Sigma}}_1^{\mathcal{T}}$) with corresponding eigenvectors $\widehat{\mathbf{v}}_1, \widehat{\mathbf{v}}_2, \ldots, \widehat{\mathbf{v}}_r$. Compute the principal orthogonal complement

$$\widehat{\mathbf{\Sigma}}_\varepsilon = (\widehat{\sigma}_{\varepsilon,k\ell})_{1 \leq k,\ell \leq d} = \widehat{\mathbf{\Sigma}}_1^{\mathcal{H}} - \widehat{\mathbf{V}}\widehat{\mathbf{\Lambda}}\widehat{\mathbf{V}}^\mathsf{T}, \tag{D.3}$$

where $\widehat{\mathbf{V}} = (\widehat{\mathbf{v}}_1, \ldots, \widehat{\mathbf{v}}_r)$ and $\widehat{\mathbf{\Lambda}} = \mathrm{diag}(\widehat{\lambda}_1, \ldots, \widehat{\lambda}_r)$.

(ii) To achieve sparsity, apply the adaptive thresholding method (Rothman, Levina and Zhu, 2009; Cai and Liu, 2011) to $\widehat{\mathbf{\Sigma}}_\varepsilon$ and obtain $\widehat{\mathbf{\Sigma}}_\varepsilon^{\mathcal{T}} = (\widehat{\sigma}_{\varepsilon,k\ell}^{\mathcal{T}})_{1 \leq k,\ell \leq d}$ such that

$$\widehat{\sigma}_{\varepsilon,k\ell}^{\mathcal{T}} = \begin{cases} \widehat{\sigma}_{\varepsilon,k\ell} & \text{if } k = \ell, \\ s_{k\ell}(\widehat{\sigma}_{\varepsilon,k\ell}) & \text{if } k \neq \ell, \end{cases} \tag{D.4}$$

where $s_{k\ell}(z) = \mathrm{sign}(z)(|z| - \lambda_{k\ell})$, $z \in \mathbb{R}$ is the soft thresholding function with $\lambda_{k\ell} = \lambda(\widehat{\sigma}_{\varepsilon,kk}\,\widehat{\sigma}_{\varepsilon,\ell\ell})^{1/2}$ and $\lambda > 0$ being a regularization parameter.

(iii) Obtain the final estimator of $\mathbf{\Sigma}$ as $\widehat{\mathbf{\Sigma}} = \widehat{\mathbf{V}}\widehat{\mathbf{\Lambda}}\widehat{\mathbf{V}}^\mathsf{T} + \widehat{\mathbf{\Sigma}}_\varepsilon^{\mathcal{T}}$.

**Remark 1.** The POET method (Fan, Liao and Mincheva, 2013) employs the sample covariance matrix as an initial estimator and has desirable properties for sub-Gaussian data. For elliptical distributions, Fan, Liu and Wang (2018) proposed to use the marginal Kendall's tau to estimate $\mathbf{\Sigma}$, and to use its top $r$ eigenvalues and the spatial Kendall's tau to estimate the corresponding leading eigenvectors. In the above robust POET procedure, we only need to compute one initial estimator of $\mathbf{\Sigma}$ and moreover, optimal convergence rates can be achieved in high dimensions under finite fourth moment conditions; see Theorem D.1.

**Condition D.1.** Under model (D.1), the latent factor $\mathbf{f} \in \mathbb{R}^r$ and the idiosyncratic noise $\boldsymbol{\varepsilon} \in \mathbb{R}^d$ are independent. Moreover,

(i) (Identifiability) $\mathrm{cov}(\mathbf{f}) = \mathbf{I}_r$ and the columns of $\mathbf{B}$ are orthogonal;

(ii) (Pervasiveness) there exist positive constants $c_l, c_u$ and $C_1$ such that

$$c_l \leq \min_{1 \leq \ell \leq r}\{\lambda_\ell(\mathbf{B}^\mathsf{T}\mathbf{B}/d) - \lambda_{\ell+1}(\mathbf{B}^\mathsf{T}\mathbf{B}/d)\} \leq c_u \text{ with } \lambda_{r+1}(\mathbf{B}^\mathsf{T}\mathbf{B}/d) = 0,$$

and $\max\{\|\mathbf{B}\|_{\max}, \|\mathbf{\Sigma}_\varepsilon\|_2\} \leq C_1$;

(iii) (Moment condition) $\max_{1 \leq \ell \leq d} \mathrm{kurt}(X_\ell) \leq C_2$ for some constant $C_2 > 0$;

(iv) (Sparsity) $\mathbf{\Sigma}_\varepsilon$ is sparse in the sense that $s := \max_{1 \leq k \leq d} \sum_{\ell=1}^d I(\sigma_{\varepsilon,k\ell} \neq 0)$ satisfies

$$s^2 \log d = o(n) \text{ and } s^2 = o(d) \text{ as } n, d \to \infty.$$

11

**Theorem D.1.** Under Condition D.1, the robust POET estimator with

$$\tau_{k\ell} \asymp \sqrt{n/(\log d)}, \ 1 \le k, \ell \le d, \ \text{ and } \ \lambda \asymp w_{n,d} := \sqrt{\log(d)/n} + d^{-1/2}$$

satisfies

$$\|\widehat{\Sigma}_{\varepsilon}^{\mathcal{T}} - \Sigma_{\varepsilon}\|_{\max} = O_{\mathbb{P}}(w_{n,d}), \quad \|\widehat{\Sigma}_{\varepsilon}^{\mathcal{T}} - \Sigma_{\varepsilon}\|_2 = O_{\mathbb{P}}(s w_{n,d}), \tag{D.5}$$

$$\|\widehat{\Sigma} - \Sigma\|_{\max} = O_{\mathbb{P}}(w_{n,d}) \ \text{ and } \ \|\widehat{\Sigma} - \Sigma\|_2 = O_{\mathbb{P}}(d w_{n,d}) \tag{D.6}$$

as $n, d \to \infty$.

## D.2 Factor-adjusted multiple testing

Here we consider the problem of simultaneously testing the hypotheses

$$H_{0k} : \mu_k = 0 \ \text{ versus } \ H_{1k} : \mu_k \ne 0, \ \text{ for } k = 1, \ldots, d, \tag{D.7}$$

under model (D.1). Although the key implication from the multi-factor pricing theory is that the intercept $\mu_k$ should be zero, known as the "mean-variance efficiency" pricing, for any asset $k$, an important question is whether such a pricing theory can be validated by empirical data. In fact, a very small proportion of $\mu_k$'s might be nonzero according to the Berk and Green equilibrium (Berk and Green, 2004). Various statistical methods have been proposed to identify those positive $\mu_k$'s (Barras, Scaillet and Wermer, 2010; Fan and Han, 2017; Lan and Du, 2019). These works assume that both the factor and idiosyncratic noise follow multivariate normal distributions. To accommodate the heavy-tailed character of empirical data, we develop a robust multiple testing procedure that controls the overall false discovery rate or false discovery proportion.

For each $1 \le k \le d$, let $T_k$ be a generic test statistic for testing the individual hypothesis $H_{0k} : \mu_k = 0$. For any threshold level $z > 0$, we reject the $j$th hypothesis whenever $|T_j| \ge z$. The numbers of total discoveries $R(z)$ and false discoveries $V(z)$ are defined by

$$R(z) = \sum_{k=1}^{d} I(|T_k| \ge z) \ \text{ and } \ V(z) = \sum_{k \in \mathcal{H}_0} I(|T_k| \ge z), \tag{D.8}$$

respectively, where $\mathcal{H}_0 = \{1 \le k \le d : \mu_k = 0\}$. The main object of interest is the false discovery proportion (FDP), given by

$$\text{FDP}(z) = V(z)/R(z).$$

Throughout we use the convention $0/0 = 0$. Note that $R(z)$ is observable given all the test statistics, while $V(z)$ is an unobservable random variable that needs to be estimated. For testing individual hypotheses $H_{0k}$, the standardized means $Z_k$, where $Z_k = n^{-1/2} \sum_{i=1}^{n} X_{ik}$, are sensitive to the tails of the sampling distributions. In particular, when the number of features $d$ is large, stochastic outliers from the test statistics $Z_k$ can be so large that they are mistakenly regarded as discoveries. Motivated by recent advances on robust estimation and inference (Catoni, 2012; Zhou et al., 2018), we consider the following robust $M$-estimator of $\mu_k$:

$$\widehat{\mu_k} = \operatorname*{argmin}_{\theta \in \mathbb{R}} \sum_{i=1}^{n} \ell_{\tau_k}(X_{ik} - \theta) \ \text{ for some } \ \tau_k > 0. \tag{D.9}$$

The corresponding test statistic is then given by $T_k = \sqrt{n}\,\widehat{\mu}_k$ for $k = 1, \ldots, d$.

Based on the law of large numbers, we define the approximate FDP by

$$\mathrm{FDP_A}(z) = \frac{1}{R(z)} \sum_{k=1}^{d} \left\{ \Phi\left( \frac{-z + \sqrt{n}\,\boldsymbol{b}_k^\top \overline{\boldsymbol{f}}}{\sqrt{\sigma_{kk} - \|\boldsymbol{b}_k\|_2^2}} \right) + \Phi\left( \frac{-z - \sqrt{n}\,\boldsymbol{b}_k^\top \overline{\boldsymbol{f}}}{\sqrt{\sigma_{kk} - \|\boldsymbol{b}_k\|_2^2}} \right) \right\}, \qquad (\mathrm{D}.10)$$

where $\overline{\boldsymbol{f}} = (1/n) \sum_{i=1}^{n} \boldsymbol{f}_i$. It is shown in the appendix that the approximate FDP in (D.10) serves as a conservative surrogate for the true FDP.

Note that the approximate FDP defined in (D.10) depends on a number of unknown parameters, say $\{\boldsymbol{b}_k, \sigma_{kk}\}_{k=1}^{d}$ and $\overline{\boldsymbol{f}}$. In this section, we describe robust procedures to estimate these quantities using the only observations $\boldsymbol{X}_1, \ldots, \boldsymbol{X}_n$.

(a) Compute the Huber-type covariance estimator $\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} = (\widehat{\sigma}_{1,k\ell}^{\mathcal{H}})_{1 \le k, \ell \le d}$ (or the truncated estimator $\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{T}}$), and let $\widehat{\lambda}_1 \ge \cdots \ge \widehat{\lambda}_r$ and $\widehat{\boldsymbol{v}}_1, \ldots, \widehat{\boldsymbol{v}}_r$ be its top $r$ eigenvalues and the corresponding eigenvectors, respectively.

(b) Compute $\widehat{\mathbf{B}} = (\widehat{\lambda}_1^{1/2}\widehat{\boldsymbol{v}}_1, \ldots, \widehat{\lambda}_r^{1/2}\widehat{\boldsymbol{v}}_r) \in \mathbb{R}^{d \times r}$ and $\widehat{\boldsymbol{u}} = \sqrt{n}\,(\widehat{\mathbf{B}}^\top\widehat{\mathbf{B}})^{-1}\widehat{\mathbf{B}}^\top\overline{\boldsymbol{X}} \in \mathbb{R}^r$, which serve as estimators of $\mathbf{B}$ and $\sqrt{n}\overline{\boldsymbol{f}}$, respectively. Here $\overline{\boldsymbol{X}} = (1/n) \sum_{i=1}^{n} \boldsymbol{X}_i$.

(c) Denote by $\widehat{\boldsymbol{b}}_1, \ldots, \widehat{\boldsymbol{b}}_d$ the $d$ rows of $\widehat{\mathbf{B}}$. For any $z \ge 0$, we estimate the approximate FDP $\mathrm{FDP_A}(z)$ by

$$\widehat{\mathrm{FDP}}_\mathrm{A}(z) = \frac{1}{R(z)} \sum_{k=1}^{d} \left\{ \Phi\left( \frac{-z + \widehat{\boldsymbol{b}}_k^\top\widehat{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}} \right) + \Phi\left( \frac{-z - \widehat{\boldsymbol{b}}_k^\top\widehat{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}} \right) \right\}, \qquad (\mathrm{D}.11)$$

where $\widehat{\sigma}_{\varepsilon,kk} = \widehat{\sigma}_{1,kk}^{\mathcal{H}} - \|\widehat{\boldsymbol{b}}_k\|_2^2$ for $k = 1, \ldots, d$.

The construction of $\widehat{\mathbf{B}}$ is based on the observation that principal component analysis and factor analysis are approximately equivalent under the pervasive assumption in high dimensions (Fan, Liao and Mincheva, 2013). To estimate $\overline{\boldsymbol{f}}$, note from model (D.1) that $\overline{\boldsymbol{X}} = \boldsymbol{\mu} + \mathbf{B}\overline{\boldsymbol{f}} + \overline{\boldsymbol{\varepsilon}}$, where $\boldsymbol{\mu}$ is assumed to be sparse and therefore is ignored for simplicity.

**Theorem D.2.** Under model (D.1), assume that $\boldsymbol{f}$ and $\boldsymbol{\varepsilon}$ are independent zero-mean random vectors and let $s_1 = \|\boldsymbol{\mu}\|_0$. Assume (i)–(iii) of Condition D.1 hold, and that $(n, d, s_1)$ satisfies $\log d = o(n)$ and $n s_1 = o(d)$ as $n, d \to \infty$. Then for any $z \ge 0$,

$$\widehat{\mathrm{FDP}}_\mathrm{A}(z)/\mathrm{FDP_A}(z) \xrightarrow{\mathbb{P}} 1 \quad \text{as} \quad n, d \to \infty. \qquad (\mathrm{D}.12)$$

## D.3   Proof of Theorem D.1

The proof is based on Theorem 2.1 and (A.1) in Fan, Liu and Wang (2018), which provides high level results for the generic POET procedure. To that end, it suffices to show that with properly chosen $\tau_{k\ell}$,

$$\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\max} = O_{\mathbb{P}}\{n^{-1/2}(\log d)^{1/2}\}, \quad \max_{1 \le \ell \le r} |\widehat{\lambda}_\ell/\lambda_\ell - 1| = O_{\mathbb{P}}\{n^{-1/2}(\log d)^{1/2}\} \qquad (\mathrm{D}.13)$$

$$\text{and} \quad \max_{1 \le \ell \le r} \|\widehat{\boldsymbol{v}}_\ell - \boldsymbol{v}_\ell\|_\infty = O_{\mathbb{P}}\{(nd)^{-1/2}(\log d)^{1/2}\}, \qquad (\mathrm{D}.14)$$

where $\lambda_1 \geq \cdots \geq \lambda_r$ are the top $r$ eigenvalues of $\boldsymbol{\Sigma}$ and $\boldsymbol{v}_1, \ldots, \boldsymbol{v}_r$ are the corresponding eigenvectors.

First, applying Theorem 3.3 with $\tau_{k\ell} \asymp \sqrt{n/(\log d)}$ implies that $\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\max} \lesssim n^{-1/2}(\log d)^{1/2}$ with probability at least $1 - d^{-1}$. This verifies the first criterion in (D.13). Next, by Weyl's inequality and the inequality $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_{1,1}$ for symmetric matrices, we have

$$\max_{1 \leq \ell \leq r} |\widehat{\lambda}_\ell - \lambda_\ell| \leq \|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_2 \leq \|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{1,1} \leq d\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\max}.$$

Let $\overline{\lambda}_1 \geq \cdots \geq \overline{\lambda}_r$ be the top $r$ eigenvalues of $\mathbf{B}\mathbf{B}^\intercal$, and therefore of $\mathbf{B}^\intercal\mathbf{B}$. Note that, by Weyl's inequality, $\max_{1 \leq \ell \leq r} |\lambda_\ell - \overline{\lambda}_\ell| \leq \|\boldsymbol{\Sigma}_\varepsilon\|_2$. It thus follows from Condition D.1 that

$$\min_{1 \leq \ell \leq r-1} |\lambda_\ell - \lambda_{\ell+1}| \asymp d \text{ and } \lambda_r \asymp d \text{ as } d \to \infty.$$

Together, the last two displays imply $\max_{1 \leq \ell \leq r} |\widehat{\lambda}_\ell/\lambda_\ell - 1| \lesssim n^{-1/2}(\log d)^{1/2}$ with probability at least $1 - d^{-1}$. Therefore, the second criterion in (D.13) is fulfilled.

For (D.14), applying Theorem 3 and Proposition 3 in Fan, Wang and Zhong (2018) we arrive at

$$\max_{1 \leq \ell \leq r} \|\widehat{\boldsymbol{v}}_\ell - \boldsymbol{v}_\ell\|_\infty$$
$$\lesssim d^{-3/2}(r^4\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\infty,\infty} + r^{3/2}\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_2) \lesssim r^4 d^{-1/2}\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\max}.$$

This validates (D.14).

In summary, (D.5) and the first bound in (D.6) follow from Theorem 2.1 and (A.1) in Fan, Liu and Wang (2018), and the second bound in (D.6) follows directly from the fact that $\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_2 \leq \|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{1,1} \leq d\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\max}$. $\qquad\square$

## D.4 Asymptotic property of FDP

In this section, we show that the approximate FDP in (D.10) serves as a conservative surrogate for the true FDP.

**Condition D.2.** Under model (D.1), $\boldsymbol{f}$ are $\boldsymbol{\varepsilon}$ are independent zero-mean random vectors. (i) $\mathrm{cov}(\boldsymbol{f}) = \mathbf{I}_r$ and $\|\boldsymbol{f}\|_{\psi_2} = \sup_{\boldsymbol{u} \in \mathbb{S}^{r-1}} \|\boldsymbol{u}^\intercal\boldsymbol{f}\|_{\psi_2} \leq C_f$ for some constant $C_f > 0$; (ii) the correlation matrix $\mathbf{R}_\varepsilon = (\varrho_{\varepsilon,k\ell})_{1 \leq k,\ell \leq d}$ of $\boldsymbol{\varepsilon}$ satisfies $d^{-2} \sum_{1 \leq k \neq \ell \leq d} \varrho_{\varepsilon,k\ell} \leq C_0 d^{-\delta_0}$ for some constants $C_0, \delta_0 > 0$; (iii) $d = d(n) \to \infty$ and $\log d = o(n^{1/2})$ as $n \to \infty$, and $\liminf_{n \to \infty} \frac{d_0}{d} > 0$, where $d_0 = \sum_{k=1}^d I(\mu_k = 0)$; (iv) $C_l \leq \sigma_{\varepsilon,kk} \leq v_k^{1/2} \vee \sigma_{kk} \leq C_u$ for all $1 \leq k \leq d$, where $v_k = \mathbb{E}(\varepsilon_k^4)$ and $C_u, C_l$ are positive constants.

**Theorem D.3.** Assume that Condition D.2 holds. In (D.9), let $\tau_k = a_k\sqrt{n/\log(nd)}$ with $a_k \geq \sigma_{kk}^{1/2}$ for $k = 1, \ldots, d$. Then, as $n, d \to \infty$,

$$\frac{V(z)}{d_0} = \frac{1}{d_0} \sum_{k \in \mathcal{H}_0} \left\{ \Phi\left(\frac{-z + \sqrt{n}\,\boldsymbol{b}_k^\intercal\overline{\boldsymbol{f}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) + \Phi\left(\frac{z - \sqrt{n}\,\boldsymbol{b}_k^\intercal\overline{\boldsymbol{f}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right\}$$
$$+ O_{\mathbb{P}}\left[ \frac{1}{d^{(\delta_0 \wedge 1)/2}} + \frac{\log(nd)}{\sqrt{n}} + \left\{\frac{\log(nd)}{n}\right\}^{1/4} \right] \tag{D.15}$$

and

$$\frac{R(z)}{d} = \frac{1}{d} \sum_{k=1}^{d} \left\{ \Phi\left( \frac{-z + \sqrt{n}(\mu_k + \boldsymbol{b}_k^\mathsf{T}\overline{\boldsymbol{f}})}{\sqrt{\sigma_{\varepsilon,kk}}} \right) + \Phi\left( \frac{-z - \sqrt{n}(\mu_k + \boldsymbol{b}_k^\mathsf{T}\overline{\boldsymbol{f}})}{\sqrt{\sigma_{\varepsilon,kk}}} \right) \right\}$$
$$+ O_\mathbb{P}\left[ \frac{1}{d^{(\delta_0 \wedge 1)/2}} + \frac{\log(nd)}{\sqrt{n}} + \left\{ \frac{\log(nd)}{n} \right\}^{1/4} \right] \tag{D.16}$$

uniformly over $z \geq 0$. In addition, for any $z \geq 0$ it holds

$$\mathrm{FDP}(z) = \mathrm{FDP}_{\mathrm{orc}}(z) + o_\mathbb{P}(1) \quad \text{as } n, d \to \infty, \tag{D.17}$$

where

$$\mathrm{FDP}_{\mathrm{orc}}(z) := \frac{1}{R(z)} \sum_{k \in \mathcal{H}_0} \left\{ \Phi\left( \frac{-z + \sqrt{n}\,\boldsymbol{b}_k^\mathsf{T}\overline{\boldsymbol{f}}}{\sqrt{\sigma_{kk} - \|\boldsymbol{b}_k\|_2^2}} \right) + \Phi\left( \frac{-z - \sqrt{n}\,\boldsymbol{b}_k^\mathsf{T}\overline{\boldsymbol{f}}}{\sqrt{\sigma_{kk} - \|\boldsymbol{b}_k\|_2^2}} \right) \right\}.$$

### D.4.1 Preliminaries

To prove Theorem D.3, we need the following results on the robust estimators $\mu_k$'s given in (D.9). Define $u_k = X_k - \mu_k = \boldsymbol{b}_k^\mathsf{T}\boldsymbol{f} + \varepsilon_k$ for $k = 1, \ldots, d$. Assume that $\mathbb{E}(\boldsymbol{f}) = \boldsymbol{0}$, $\mathbb{E}(\varepsilon_k) = 0$ and $\boldsymbol{f}$ are $\varepsilon_k$ are independent. Then we have $\mathbb{E}(u_k) = 0$ and $\mathbb{E}(u_k^2) = \sigma_{kk} = \|\boldsymbol{b}_k\|_2^2 + \sigma_{\varepsilon,kk}$.

The first lemma is Theorem 5 in Fan, Li and Wang (2017) regarding the concentration of the robust mean estimator.

**Lemma D.1.** For every $1 \leq k \leq d$ and $t > 0$, the estimator $\widehat{\mu}_k$ in (D.9) with $\tau_k = a_k(n/t)^{1/2}$ for $a_k \geq \sigma_{kk}^{1/2}$ satisfies $|\widehat{\mu}_k - \mu_k| \leq 4a_k(t/n)^{1/2}$ with probability at least $1 - 2e^{-t}$ provided $n \geq 8t$.

The next result provides a nonasymptotic Bahadur representation for $\widehat{\mu}_k$, which follows directly from Lemma D.1 and Theorem 2.1 in Zhou et al. (2018). Let $u_{ik} = \boldsymbol{b}_k^\mathsf{T}\boldsymbol{f}_i + \varepsilon_{ik}$ for $i = 1, \ldots, n$ and $k = 1, \ldots, d$.

**Lemma D.2.** Under the conditions of Lemma D.1, it holds for every $1 \leq k \leq d$ that

$$\left| \sqrt{n}\,(\widehat{\mu}_k - \mu_k) - \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \psi_{\tau_k}(u_{ik}) \right| \leq C \frac{a_k t}{\sqrt{n}} \tag{D.18}$$

with probability at least $1 - 3e^{-t}$ as long as $n \geq 8t$, where $C > 0$ is an absolute constant and $\psi_\tau(\cdot)$ is given in (3.1).

Let $\tau_k$ be as in Lemma D.1 and write

$$v_k = \mathbb{E}(\varepsilon_k^4), \quad \xi_k = \psi_k(u_k) \text{ for } k = 1, \ldots, d. \tag{D.19}$$

Here $\xi_k$ are truncated versions of $u_k$. The next result shows that the differences between the first two (conditional) moments of $u_k$ and $\xi_k$ given $\boldsymbol{f}$ decay as $\tau_k$ grows.

**Lemma D.3.** Assume that $v_k < \infty$ for $k = 1, \ldots, d$.

15

1. On the event $\mathcal{G}_k := \{|\boldsymbol{b}_k^\top \boldsymbol{f}| \leq \tau_k/2\}$, the following inequalities hold almost surely:

$$|\mathbb{E}_f(\xi_k) - \boldsymbol{b}_k^\top \boldsymbol{f}| \leq \min(2\tau_k^{-1}\sigma_{\varepsilon,kk}, 8\tau_k^{-3}v_k) \tag{D.20}$$

$$\text{and } \sigma_{\varepsilon,kk} - 4\tau_k^{-2}(v_k + \sigma_{\varepsilon,kk}^2) \leq \text{var}_f(\xi_k) \leq \sigma_{\varepsilon,kk}, \tag{D.21}$$

where $\mathbb{E}_f(\cdot)$ and $\text{var}_f(\cdot)$ denote the conditional mean and variance, separately.

2. On the event $\mathcal{G}_k \cap \mathcal{G}_\ell$, the following holds almost surely:

$$|\text{cov}_f(\xi_k, \xi_\ell) - \sigma_{\varepsilon,k\ell}| \leq C \frac{v_k \vee v_\ell}{(\tau_k \wedge \tau_\ell)^2}, \tag{D.22}$$

where $C > 0$ is an absolute constant.

*Proof of Lemma D.3.* First we prove (D.20) and (D.21). Fix $k$ and let $\tau = \tau_k$ for simplicity. Since $\varepsilon_k$ and $\boldsymbol{f}$ are independent, we have

$$\mathbb{E}_f \xi_k - \boldsymbol{b}_k^\top \boldsymbol{f}$$
$$= -\mathbb{E}_f(\varepsilon_k + \boldsymbol{b}_k^\top \boldsymbol{f} - \tau)I(\varepsilon_k > \tau - \boldsymbol{b}_k^\top \boldsymbol{f}) + \mathbb{E}_f(-\varepsilon_k - \boldsymbol{b}_k^\top \boldsymbol{f} - \tau)I(\varepsilon_k < -\tau - \boldsymbol{b}_k^\top \boldsymbol{f}).$$

Therefore, on the event $\mathcal{G}_k$, it holds for any $2 \leq q \leq 4$ that

$$|\mathbb{E}_f \xi_k - \boldsymbol{b}_k^\top \boldsymbol{f}| \leq \mathbb{E}_f\{|\varepsilon_k|I(|\varepsilon_k| > \tau - |\boldsymbol{b}_k^\top \boldsymbol{f}|)\} \leq (\tau - |\boldsymbol{b}_k^\top \boldsymbol{f}|)^{1-q}\mathbb{E}(|\varepsilon_k|^q)$$

almost surely. This proves (D.20) by taking $q$ to be 2 and 4. For the conditional variance, by (D.20) and the decomposition $\mathbb{E}_f(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})^2 = \text{var}_f(\xi_k) + (\mathbb{E}_f \xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})^2$, we have

$$\mathbb{E}_f(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})^2 - \frac{\sigma_{\varepsilon,kk}^2}{(\tau - |\boldsymbol{b}_k^\top \boldsymbol{f}|)^2} \leq \text{var}_f(\xi_k) \leq \mathbb{E}_f(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})^2. \tag{D.23}$$

Note that $\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f}$ can be written as

$$\varepsilon_k I(|\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k| \leq \tau) + (\tau - \boldsymbol{b}_k^\top \boldsymbol{f})I(\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k > \tau) - (\tau + \boldsymbol{b}_k^\top \boldsymbol{f})I(\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k < -\tau).$$

It follows that

$$(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})^2$$
$$= \varepsilon_k^2 I(|\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k| \leq \tau) + (\tau - \boldsymbol{b}_k^\top \boldsymbol{f})^2 I(\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k > \tau) + (\tau + \boldsymbol{b}_k^\top \boldsymbol{f})^2 I(\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k < -\tau).$$

Taking conditional expectations on both sides gives

$$\mathbb{E}_f(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})^2 = \mathbb{E}(\varepsilon_k^2) - \mathbb{E}_f\{\varepsilon_k^2 I(|\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k| > \tau)\}$$
$$+ (\tau - \boldsymbol{b}_k^\top \boldsymbol{f})^2 \mathbb{P}_f(\varepsilon_k > \tau - \boldsymbol{b}_k^\top \boldsymbol{f}) + (\tau + \boldsymbol{b}_k^\top \boldsymbol{f})^2 \mathbb{P}_f(\varepsilon_k < -\tau - \boldsymbol{b}_k^\top \boldsymbol{f}).$$

Using the identity that $u^2 = 2\int_0^u t\, dt$ for any $u > 0$, we have

$$\mathbb{E}_f\{\varepsilon_k^2 I(\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k > \tau)\}$$
$$= 2\mathbb{E}_f \int_0^\infty I(\varepsilon_k > t)I(\varepsilon_k > \tau - \boldsymbol{b}_k^\top \boldsymbol{f})t\, dt$$
$$= 2\mathbb{E}_f \int_0^{\tau - \boldsymbol{b}_k^\top \boldsymbol{f}} I(\varepsilon_k > \tau - \boldsymbol{b}_k^\top \boldsymbol{f})t\, dt + 2\mathbb{E}_f \int_{\tau - \boldsymbol{b}_k^\top \boldsymbol{f}}^\infty I(\varepsilon_k > t)t\, dt$$
$$= (\tau - \boldsymbol{b}_k^\top \boldsymbol{f})^2 \mathbb{P}_f(\varepsilon_k > \tau - \boldsymbol{b}_k^\top \boldsymbol{f}) + 2\int_{\tau - \boldsymbol{b}_k^\top \boldsymbol{f}}^\infty \mathbb{P}(\varepsilon_k > t)t\, dt.$$

16

It can be similarly shown that

$$\mathbb{E}_f\{\varepsilon_k^2 I(\boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k < -\tau)\} = (\tau + \boldsymbol{b}_k^\top \boldsymbol{f})^2 \mathbb{P}_f(\varepsilon_k < -\tau - \boldsymbol{b}_k^\top \boldsymbol{f}) + 2 \int_{\tau + \boldsymbol{b}_k^\top \boldsymbol{f}}^\infty \mathbb{P}(-\varepsilon_k > t)t\, dt.$$

Together, the last three displays imply

$$\begin{aligned}
0 &\geq \mathbb{E}_f(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})^2 - \mathbb{E}(\varepsilon_k^2) \\
&\geq -2 \int_{\tau - |\boldsymbol{b}_k^\top \boldsymbol{f}|}^\infty \mathbb{P}(|\varepsilon_k| > t)t\, dt \geq -2v_k \int_{\tau - |\boldsymbol{b}_k^\top \boldsymbol{f}|}^\infty \frac{dt}{t^3} = -\frac{v_k}{(\tau - |\boldsymbol{b}_k^\top \boldsymbol{f}|)^2}.
\end{aligned}$$

Combining this with (D.23) and (D.20) proves (D.21).

Next, we study the covariance $\mathrm{cov}_f(\xi_k, \xi_\ell)$ for $k \neq \ell$, which can be written as

$$\begin{aligned}
\mathrm{cov}_f(\xi_k, \xi_\ell) &= \mathbb{E}_f(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f} + \boldsymbol{b}_k^\top \boldsymbol{f} - \mathbb{E}_f \xi_k)(\xi_\ell - \boldsymbol{b}_\ell^\top \boldsymbol{f} + \boldsymbol{b}_\ell^\top \boldsymbol{f} - \mathbb{E}_f \xi_\ell) \\
&= \underbrace{\mathbb{E}_f(\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})(\xi_\ell - \boldsymbol{b}_\ell^\top \boldsymbol{f})}_{\Pi_1} - \underbrace{(\mathbb{E}_f \xi_k - \boldsymbol{b}_k^\top \boldsymbol{f})(\mathbb{E}_f \xi_\ell - \boldsymbol{b}_\ell^\top \boldsymbol{f})}_{\Pi_2}.
\end{aligned}$$

For $\Pi_2$, it follows immediately from (D.20) that $|\Pi_2| \lesssim (\tau_k \tau_\ell)^{-1} \sigma_{\varepsilon,kk} \sigma_{\varepsilon,\ell\ell}$ almost surely on the event $\mathcal{G}_{k\ell} := \{|\boldsymbol{b}_k^\top \boldsymbol{f}| \leq \tau_k/2\} \cap \{|\boldsymbol{b}_\ell^\top \boldsymbol{f}| \leq \tau_\ell/2\}$. It remains to consider $\Pi_1$. Recall that $\xi_k - \boldsymbol{b}_k^\top \boldsymbol{f} = \varepsilon_k I(|u_k| \leq \tau_k) + (\tau_k - \boldsymbol{b}_k^\top \boldsymbol{f})I(u_k > \tau_k) - (\tau_k + \boldsymbol{b}_k^\top \boldsymbol{f})I(u_k < -\tau_k)$, where $u_k = \boldsymbol{b}_k^\top \boldsymbol{f} + \varepsilon_k$. Then, $\Pi_1$ can be written as

$$\begin{aligned}
&\mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_k| \leq \tau_k, |u_\ell| \leq \tau_\ell) + (\tau_\ell - \boldsymbol{b}_\ell^\top \boldsymbol{f})\mathbb{E}_f \varepsilon_k I(|u_k| \leq \tau, u_\ell > \tau) \\
&\quad - (\tau_\ell + \boldsymbol{b}_\ell^\top \boldsymbol{f})\mathbb{E}_f \varepsilon_k I(|u_k| \leq \tau_k, u_\ell < -\tau_\ell) + (\tau_k - \boldsymbol{b}_k^\top \boldsymbol{f})\mathbb{E}_f \varepsilon_\ell I(u_k > \tau_k, |u_\ell| \leq \tau_\ell) \\
&\quad + (\tau_k - \boldsymbol{b}_k^\top \boldsymbol{f})(\tau_\ell - \boldsymbol{b}_\ell^\top \boldsymbol{f})\mathbb{E}_f I(u_k > \tau_k, u_\ell > \tau_\ell) - (\tau_k - \boldsymbol{b}_k^\top \boldsymbol{f})(\tau_\ell + \boldsymbol{b}_\ell^\top \boldsymbol{f})\mathbb{E}_f I(u_k > \tau_k, u_\ell < -\tau_\ell) \\
&\quad - (\tau_k + \boldsymbol{b}_k^\top \boldsymbol{f})\mathbb{E}_f \varepsilon_\ell I(u_k < -\tau_k, |u_\ell| \leq \tau_\ell) - (\tau_k + \boldsymbol{b}_k^\top \boldsymbol{f})(\tau_\ell - \boldsymbol{b}_\ell^\top \boldsymbol{f})\mathbb{E}_f I(u_k < -\tau_k, u_\ell > \tau_\ell) \\
&\quad + (\tau_k + \boldsymbol{b}_k^\top \boldsymbol{f})(\tau_\ell + \boldsymbol{b}_\ell^\top \boldsymbol{f})\mathbb{E}_f I(u_k < -\tau_k, u_\ell < -\tau_\ell). \tag{D.24}
\end{aligned}$$

For the first term in (D.24), note that

$$\begin{aligned}
&\mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_k| \leq \tau_k, |u_\ell| \leq \tau_\ell) \\
&= \mathrm{cov}(\varepsilon_k, \varepsilon_\ell) - \mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_k| > \tau_k) - \mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_\ell| > \tau_\ell) + \mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_k| > \tau_k, |u_\ell| > \tau_\ell),
\end{aligned}$$

where $|\mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_k| > \tau_k)| \leq (\tau_k - |\boldsymbol{b}_k^\top \boldsymbol{f}|)^{-2}\mathbb{E}(|\varepsilon_k|^3 |\varepsilon_\ell|) \leq 4\tau_k^{-2} v_k^{3/4} v_\ell^{1/4}$ and

$$|\mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_k| > \tau_k, |u_\ell| > \tau_\ell)| \leq (\tau_k - |\boldsymbol{b}_k^\top \boldsymbol{f}|)^{-1}(\tau_\ell - |\boldsymbol{b}_\ell^\top \boldsymbol{f}|)^{-1}\mathbb{E}(\varepsilon_k^2 \varepsilon_\ell^2) \leq 4\tau_k^{-1}\tau_\ell^{-1} v_k^{1/2} v_\ell^{1/2}$$

almost surely on $\mathcal{G}_{k\ell}$. Hence,

$$|\mathbb{E}_f \varepsilon_k \varepsilon_\ell I(|u_k| \leq \tau_k, |u_\ell| \leq \tau_\ell) - \mathrm{cov}(\varepsilon_k, \varepsilon_\ell)| \lesssim (\tau_k \wedge \tau_\ell)^{-2}$$

holds almost surely on the same event. For the remaining terms in (D.24), it can be similarly derived that, almost surely on the same event,

$$\begin{aligned}
|\mathbb{E}_f \varepsilon_k I(|u_k| \leq \tau_k, u_\ell > \tau_\ell)| &\leq |\tau_\ell - \boldsymbol{b}_\ell^\top \boldsymbol{f}|^{-3}\mathbb{E}(|\varepsilon_k||\varepsilon_\ell|^3), \\
|\mathbb{E}_f \varepsilon_k I(|u_k| \leq \tau_k, u_\ell < -\tau_\ell)| &\leq |\tau_\ell + \boldsymbol{b}_\ell^\top \boldsymbol{f}|^{-3}\mathbb{E}(|\varepsilon_k||\varepsilon_\ell|^3), \\
\text{and } \mathbb{E}_f I(u_k > \tau_k, \xi_\ell < -\tau_\ell) &\leq |\tau_k - \boldsymbol{b}_k^\top \boldsymbol{f}|^{-2}|\tau_\ell + \boldsymbol{b}_\ell^\top \boldsymbol{f}|^{-2}\mathbb{E}(\varepsilon_k^2 \varepsilon_\ell^2).
\end{aligned}$$

Putting the pieces together, we arrive at $|\Pi_1 - \mathrm{cov}(\varepsilon_k, \varepsilon_\ell)| \lesssim (\tau_k \wedge \tau_\ell)^{-2}(v_k \vee v_\ell)$ almost surely on $\mathcal{G}_{k\ell}$. This proves the stated result (D.22). $\qquad\square$

The next lemma provides several concentration results regarding the factors $f_i$'s and their functionals.

**Lemma D.4.** Assume that (i) of Condition D.2 holds. Then, for any $t > 0$,

$$\mathbb{P}\left\{\max_{1\leq i\leq n} \|f_i\|_2 > C_1 C_f (r + \log n + t)^{1/2}\right\} \leq e^{-t}, \tag{D.25}$$

$$\mathbb{P}\{\|\sqrt{n}\,\overline{f}\|_2 > C_2 C_f (r + t)^{1/2}\} \leq e^{-t}, \tag{D.26}$$

$$\text{and}\quad \mathbb{P}[\|\widehat{\Sigma}_f - \mathbf{I}_r\|_2 > \max\{C_3 C_f^2 n^{-1/2}(r + t)^{1/2}, C_3^2 C_f^4 n^{-1}(r + t)\}] \leq 2e^{-t}, \tag{D.27}$$

where $\overline{f} = (1/n)\sum_{i=1}^n f_i$, $\widehat{\Sigma}_f = (1/n)\sum_{i=1}^n f_i f_i^{\mathsf{T}}$ and $C_1$–$C_3$ are absolute constants.

*Proof of Lemma D.4.* Under (i) of Condition D.2, it holds for any $u \in \mathbb{R}^r$ that

$$\mathbb{E}e^{u^{\mathsf{T}} f_i} \leq e^{CC_f^2\|u\|_2^2}\ \text{ for } i = 1,\ldots,n,$$

$$\text{and}\quad \mathbb{E}e^{\sqrt{n}\,u^{\mathsf{T}}\overline{f}} = \prod_{i=1}^n \mathbb{E}e^{n^{-1/2}u^{\mathsf{T}} f_i} \leq \prod_{i=1}^n e^{CC_f^2 n^{-1}\|u\|_2^2} \leq e^{CC_f^2\|u\|_2^2},$$

where $C > 0$ is an absolute constant. Using Theorem 2.1 in Hsu, Kakade and Zhang (2012), we derive that for any $t > 0$,

$$\mathbb{P}\{\|f_i\|_2^2 > 2CC_f^2(r + 2\sqrt{rt} + 2t)\} \leq e^{-t}\ \text{ and }\ \mathbb{P}\{\|\sqrt{n}\,\overline{f}\|_2^2 > 2CC_f^2(r + 2\sqrt{rt} + 2t)\} \leq e^{-t}.$$

Then, (D.25) follows from the first inequality and the union bound, and the second inequality leads to (D.26) immediately. Finally, applying Theorem 5.39 in Vershynin (2012) gives (D.27). □

### D.4.2 Proof of Theorem D.3

First we introduce the following notations:

$$v_k = \mathbb{E}(\varepsilon_k^4), \quad \kappa_{\varepsilon,k} = v_k/\sigma_{\varepsilon,kk}^2, \quad u_{ik} = b_k^{\mathsf{T}} f_i + \varepsilon_{ik}, \quad k = 1,\ldots,d, \ i = 1,\ldots,n.$$

Let $t \geq 1$ and set $\tau_k = a_k(n/t)^{1/2}$ with $a_k \geq \sigma_{kk}^{1/2}$ for $k = 1,\ldots,d$. In view of Lemma D.2, define the event

$$\mathcal{E}_1(t) = \bigcap_{k=1}^d \left\{\left|\sqrt{n}\,(\widehat{\mu}_k - \mu_k) - \frac{1}{\sqrt{n}}\sum_{i=1}^n \psi_{\tau_k}(u_{ik})\right| \leq C\frac{a_k t}{\sqrt{n}}\right\}, \tag{D.28}$$

such that $\mathbb{P}\{\mathcal{E}_1(t)^c\} \leq 3de^{-t}$. Moreover, by Lemma D.4, let $\mathcal{E}_2(t)$ be the event that the following hold:

$$\max_{1\leq i\leq n}\|f_i\|_2 \leq C_1 C_f (r + \log n + t)^{1/2}, \quad \|\sqrt{n}\,\overline{f}\|_2 \leq C_2 C_f (r + t)^{1/2},$$

$$\text{and}\quad \|\widehat{\Sigma}_f - \mathbf{I}_r\|_2 \leq \max\{C_3 C_f^2 n^{-1/2}(r + t)^{1/2}, C_3^2 C_f^4 n^{-1}(r + t)\}. \tag{D.29}$$

By the union bound, $\mathbb{P}\{\mathcal{E}_2(t)^c\} \leq 4e^{-t}$.

Now we are ready to prove (D.15). The proof of (D.16) follows the same argument, and thus is omitted. For $k = 1,\ldots,d$, define

$$B_k = \sqrt{n}\,b_k^{\mathsf{T}}\overline{f}, \quad V_k = \frac{1}{\sqrt{n}}\sum_{i=1}^n V_{ik} := \frac{1}{\sqrt{n}}\sum_{i=1}^n \{\psi_{\tau_k}(u_{ik}) - \mathbb{E}_{f_i}\psi_{\tau_k}(u_{ik})\}, \tag{D.30}$$

and $R_k = n^{-1/2} \sum_{i=1}^n \{\mathbb{E}_{f_i} \psi_{\tau_k}(u_{ik}) - \boldsymbol{b}_k^\mathsf{T} \boldsymbol{f}_i\}$, where $\mathbb{E}_{f_i}(\cdot) := \mathbb{E}(\cdot | \boldsymbol{f}_i)$. On the event $\mathcal{E}_1(t)$,

$$|T_k - (V_k + B_k + R_k)| \le C a_k n^{-1/2} t \quad \text{for all } 1 \le k \le d. \tag{D.31}$$

On $\mathcal{E}_2(t)$, it holds $\max_{1 \le i \le n} |\boldsymbol{b}_k^\mathsf{T} \boldsymbol{f}_i| \le C_1 C_f \|\boldsymbol{b}_k\|_2 (r + \log n + t)^{1/2} \le C_1 C_f \sigma_{kk}^{1/2} (r + \log n + t)^{1/2}$, which further implies

$$\max_{1 \le i \le n} |\boldsymbol{b}_k^\mathsf{T} \boldsymbol{f}_i| \le \tau_k / 2 \quad \text{for all } 1 \le k \le d \tag{D.32}$$

as long as $n \ge 4(C_1 C_f)^2 (r + \log n + t) t$. Then, it follows from Lemma D.3 that

$$|R_k| \le \sqrt{n} \max_{1 \le i \le n} |\mathbb{E}_{f_i} \psi_{\tau_k}(u_{ik}) - \boldsymbol{b}_k^\mathsf{T} \boldsymbol{f}_i| \le 8\sqrt{n} \, \tau_k^{-3} v_k \le 8\sigma_{kk}^{-3/2} v_k n^{-1} t^{3/2} \tag{D.33}$$

holds almost surely on $\mathcal{E}_2(t)$ for all $1 \le k \le d$. Together, (D.31) and (D.33) imply that for any $z \ge 0$,

$$\sum_{k \in \mathcal{H}_0} I\left(|V_k + B_k| \ge z + C a_k n^{-1/2} t + 8\kappa_{\varepsilon,k} \sigma_{\varepsilon,kk}^{1/2} n^{-1} t^{3/2}\right)$$

$$\le V(z) \le \sum_{k \in \mathcal{H}_0} I\left(|V_k + B_k| \ge z - C a_k n^{-1/2} t - 8\kappa_{\varepsilon,k} \sigma_{\varepsilon,kk}^{1/2} n^{-1} t^{3/2}\right) \tag{D.34}$$

holds almost surely on $\mathcal{E}_1(t) \cap \mathcal{E}_2(t)$. In view of (D.34), we will instead deal with $\widetilde{V}_+(z)$ and $\widetilde{V}_-(z)$, where

$$\widetilde{V}_+(z) := \sum_{k \in \mathcal{H}_0} I(V_k \ge z - B_k) \quad \text{and} \quad \widetilde{V}_-(z) := \sum_{k \in \mathcal{H}_0} I(V_k \le -z - B_k)$$

are such that $\widetilde{V}_+(z) + \widetilde{V}_-(z) = \sum_{k \in \mathcal{H}_0} I(|V_k + B_k| \ge z)$.

In the following, we will focus on $\widetilde{V}_+(z)$ ($\widetilde{V}_-(z)$ can be dealt with in the same way). Observe that, conditional on $\mathcal{F}_n := \{\boldsymbol{f}_i\}_{i=1}^n$, $I(V_1 \ge z - B_1), \ldots, I(V_d \ge z - B_d)$ are weakly correlated random variables. Define $Y_k = I(V_k \ge z - B_k)$ and $P_k = \mathbb{E}(Y_k | \mathcal{F}_n)$ for $k = 1, \ldots, d$. To prove the consistency of $\widetilde{V}_+(z)$, we calculate its variance:

$$\mathrm{var}\left(\frac{1}{d_0} \sum_{k \in \mathcal{H}_0} Y_k \Big| \mathcal{F}_n\right) = \mathbb{E}\left[\left\{\frac{1}{d_0} \widetilde{V}_+(z) - \frac{1}{d_0} \sum_{k \in \mathcal{H}_0} P_k\right\}^2 \Big| \mathcal{F}_n\right]$$

$$= \frac{1}{d_0^2} \sum_{k \in \mathcal{H}_0} \mathrm{var}(Y_k | \mathcal{F}_n) + \frac{1}{d_0^2} \sum_{k,\ell \in \mathcal{H}_0 : k \ne \ell} \mathrm{cov}(Y_k, Y_\ell | \mathcal{F}_n)$$

$$\le \frac{1}{4d_0} + \frac{1}{d_0^2} \sum_{k,\ell \in \mathcal{H}_0 : k \ne \ell} \{\mathbb{E}(Y_k Y_\ell | \mathcal{F}_n) - P_k P_\ell\} \tag{D.35}$$

almost surely. In what follows, we study $P_k$ and $\mathbb{E}(Y_k Y_\ell | \mathcal{F}_n)$ separately, starting with the former. For each $k$, $V_k$ is a sum of conditionally independent zero-mean random variables given $\mathcal{F}_n$. Define

$$v_k^2 = \mathrm{var}(V_k | \mathcal{F}_n) = \frac{1}{n} \sum_{i=1}^n v_{ik}^2 \quad \text{with} \quad v_{ik}^2 = \mathrm{var}(V_{ik} | \mathcal{F}_n).$$

Then, it follows from the Berry-Esseen theorem that

$$\sup_{z \in \mathbb{R}} |\mathbb{P}(V_k \le v_k x | \mathcal{F}_n) - \Phi(x)|$$

$$\lesssim \frac{1}{v_k^3 n^{3/2}} \sum_{i=1}^n \mathbb{E}\{|\psi_{\tau_k}(u_{ik})|^3 | \mathcal{F}_n\} \lesssim \frac{1}{v_k^3 n^{3/2}} \sum_{i=1}^n (|\boldsymbol{b}_k^\mathsf{T} \boldsymbol{f}_i|^3 + \mathbb{E}|\varepsilon_{ik}|^3) \tag{D.36}$$

19

almost surely. By (D.32) and (D.21), it holds

$$1 - 4(1 + \kappa_{\varepsilon,k})n^{-1}t \le \sigma_{\varepsilon,kk}^{-1}v_k^2 \le 1 \tag{D.37}$$

almost surely on $\mathcal{E}_2(t)$ for all $1 \le k \le d$. Combining this with (D.29) and (D.36) gives

$$\left| P_k - \Phi\left(\frac{-z + B_k}{v_k}\right) \right| \lesssim \kappa_{\varepsilon,k}^{3/4} \frac{1}{\sqrt{n}} + \sigma_{\varepsilon,kk}^{-3/2} \|\boldsymbol{b}_k\|_2^3 \sqrt{\frac{r + \log n + t}{n}}$$

almost surely on $\mathcal{E}_2(t)$ for all $1 \le k \le d$ as long as $n \ge 2C_3^2 C_f^4(r + t) \vee 8(1 + \kappa_{\varepsilon,\max})t$, where $\kappa_{\varepsilon,\max} = \max_{1 \le \ell \le d} \kappa_{\varepsilon,\ell}$. By the mean value theorem, there exists some $\eta_k \in [\sigma_{\varepsilon,kk}^{-1/2}, v_k^{-1}]$ such that

$$\left| \Phi\left(\frac{-z + B_k}{v_k}\right) - \Phi\left(\frac{-z + B_k}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right| = \phi(\eta_k|z - B_k|)\frac{\eta_k|z - B_k|}{\eta_k}\left| \frac{1}{v_k} - \frac{1}{\sqrt{\sigma_{\varepsilon,kk}}} \right| \lesssim \frac{\kappa_{\varepsilon,k}t}{n}.$$

Together, the last two displays imply that almost surely on $\mathcal{E}_2(t)$,

$$\left| P_k - \Phi\left(\frac{-z + B_k}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right| \lesssim \kappa_{\varepsilon,k}\left(\frac{1}{\sqrt{n}} + \frac{t}{n}\right) + \sigma_{\varepsilon,kk}^{-3/2}\|\boldsymbol{b}_k\|_2^3 \sqrt{\frac{r + \log n + t}{n}} \tag{D.38}$$

uniformly for all $1 \le k \le d$ and $z \ge 0$.

Next we consider $\mathbb{E}(Y_k Y_\ell | \mathcal{F}_n) = \mathbb{P}(V_k \ge z - B_k, V_\ell \ge z - B_\ell | \mathcal{F}_n)$. Define bivariate random vectors $\boldsymbol{V}_i = (v_k^{-1} V_{ik}, v_\ell^{-1} V_{i\ell})^\intercal$ for $i = 1, \ldots, n$, where $V_{ik}, V_{i\ell}$ are as in (D.30). Observe that $\boldsymbol{V}_1, \ldots, \boldsymbol{V}_n$ are conditionally independent random vectors given $\mathcal{F}_n$. Denote by $\boldsymbol{\Theta} = (\theta_{uv})_{1 \le u, v \le 2}$ the conditional covariance matrix of $n^{-1/2} \sum_{i=1}^n \boldsymbol{V}_i = (v_k^{-1} V_k, v_\ell^{-1} V_\ell)^\intercal$ given $\mathcal{F}_n$, such that $\theta_{11} = \theta_{22} = 1$ and $\theta_{12} = \theta_{21} = (v_k v_\ell n)^{-1} \sum_{i=1}^n \operatorname{cov}_{f_i}(V_{ik}, V_{i\ell})$. By (D.22), (D.32) and (D.37),

$$|\theta_{12} - r_{\varepsilon,k\ell}| \lesssim (\kappa_{\varepsilon,k} \vee \kappa_{\varepsilon,\ell})n^{-1}t \tag{D.39}$$

holds almost surely on $\mathcal{E}_2(t)$ for all $1 \le k \ne \ell \le d$ and sufficient large $n$, say $n \gtrsim \kappa_{\varepsilon,\max}t$. Let $\boldsymbol{G} = (G_1, G_2)^\intercal$ be a Gaussian random vector with $\mathbb{E}(\boldsymbol{G}) = \boldsymbol{0}$ and $\operatorname{cov}(\boldsymbol{G}) = \boldsymbol{\Theta}$. Applying Theorem 1.1 in Bentkus (2005) and (D.29), we have

$$\sup_{x,y \in \mathbb{R}} |\mathbb{P}(V_k \ge v_k x, V_\ell \ge v_\ell y | \mathcal{F}_n) - \mathbb{P}(G_1 \ge x, G_2 \ge y)|$$

$$\lesssim \frac{1}{n^{3/2}} \sum_{i=1}^n \mathbb{E}\|\boldsymbol{\Theta}^{-1/2} \boldsymbol{V}_i\|_2^3$$

$$\lesssim \frac{1}{(\sigma_{\varepsilon,kk}n)^{3/2}} \sum_{i=1}^n (\mathbb{E}|\varepsilon_{ik}|^3 + |\boldsymbol{b}_k^\intercal \boldsymbol{f}_i|^3) + \frac{1}{(\sigma_{\varepsilon,\ell\ell}n)^{3/2}} \sum_{i=1}^n (\mathbb{E}|\varepsilon_{i\ell}|^3 + |\boldsymbol{b}_\ell^\intercal \boldsymbol{f}_i|^3)$$

$$\lesssim \frac{\kappa_{\varepsilon,k} + \kappa_{\varepsilon,\ell}}{\sqrt{n}} + (\sigma_{\varepsilon,kk}^{-3/2}\|\boldsymbol{b}_k\|_2^3 + \sigma_{\varepsilon,\ell\ell}^{-3/2}\|\boldsymbol{b}_\ell\|_2^3) \sqrt{\frac{r + \log n + t}{n}}$$

almost surely on $\mathcal{E}_2(t)$ for all $1 \le k \ne \ell \le d$. In particular, taking $x = v_k^{-1}(z - B_k)$ and $y = v_\ell^{-1}(z - B_\ell)$ gives

$$|\mathbb{E}(Y_k Y_\ell | \mathcal{F}_n) - \mathbb{P}\{G_1 \ge v_k^{-1}(z - B_k), G_2 \ge v_\ell^{-1}(z - B_\ell)\}|$$

$$\lesssim \frac{\kappa_{\varepsilon,k} + \kappa_{\varepsilon,\ell}}{\sqrt{n}} + (\sigma_{\varepsilon,kk}^{-3/2}\|\boldsymbol{b}_k\|_2^3 + \sigma_{\varepsilon,\ell\ell}^{-3/2}\|\boldsymbol{b}_\ell\|_2^3) \sqrt{\frac{r + \log n + t}{n}} \tag{D.40}$$

almost surely on $\mathcal{E}_2(t)$. In addition, it follows from Corollary 2.1 in Li and Shao (2002) that

$$|\mathbb{P}(G_1 \geq x, G_2 \geq y) - \{1 - \Phi(x)\}\{1 - \Phi(y)\}| \leq \frac{|\theta_{12}|}{4} e^{-(x^2+y^2)/(2+2|\theta_{12}|)} \leq \frac{|\theta_{12}|}{4} \qquad (\text{D.41})$$

for all $x, y \in \mathbb{R}$.

Substituting the bounds (D.38), (D.39), (D.40) and (D.41) into (D.35), we obtain

$$\mathbb{E}\left[\left\{\frac{1}{d_0}\widetilde{V}_+(z) - \frac{1}{d_0}\sum_{k \in \mathcal{H}_0} P_k\right\}^2 \middle| \mathcal{F}_n\right]$$

$$\lesssim \frac{1}{d_0^2}\sum_{k,\ell \in \mathcal{H}_0 : k \neq \ell} \varrho_{\varepsilon,k\ell} + \frac{1}{d_0} + \frac{\kappa_{\varepsilon,\max}}{\sqrt{n}} + \frac{\kappa_{\varepsilon,\max}t}{n} + \max_{1 \leq k \leq d} \sigma_{\varepsilon,kk}^{-3/2}\|\boldsymbol{b}_k\|_2^3 \sqrt{\frac{r + \log n + t}{n}}$$

and

$$\left|\frac{1}{d_0}\sum_{k \in \mathcal{H}_0} P_k - \frac{1}{d_0}\sum_{k \in \mathcal{H}_0} \Phi\left(\frac{-z + B_k}{\sqrt{\sigma_{\varepsilon,kk}}}\right)\right|$$

$$\lesssim \frac{\kappa_{\varepsilon,\max}}{\sqrt{n}} + \frac{\kappa_{\varepsilon,\max}t}{n} + \max_{1 \leq k \leq d} \sigma_{\varepsilon,kk}^{-3/2}\|\boldsymbol{b}_k\|_2^3 \sqrt{\frac{r + \log n + t}{n}}$$

almost surely on $\mathcal{E}_2(t)$. Similar bounds can be derived for

$$\text{var}\left(\frac{1}{d_0}\widetilde{V}_-(z)\middle|\mathcal{F}_n\right) \quad \text{and} \quad \frac{1}{d_0}\mathbb{E}\{\widetilde{V}_-(z)|\mathcal{F}_n\}.$$

Taking $t = \log(nd)$ so that $\mathbb{P}\{\mathcal{E}_1(t)^c\} \leq 3n^{-1}$ and $\mathbb{P}\{\mathcal{E}_2(t)^c\} \leq 4(nd)^{-1}$. Under Condition D.2, it follows that

$$\frac{1}{d_0}\widetilde{V}_+(z) = \frac{1}{d_0}\sum_{k \in \mathcal{H}_0} \Phi\left(\frac{-z + B_k}{\sqrt{\sigma_{\varepsilon,kk}}}\right) + O_{\mathbb{P}}[d^{-(1 \wedge \delta_0)/2} + n^{-1/4}\{\log(nd)\}^{1/4}],$$

$$\frac{1}{d_0}\widetilde{V}_-(z) = \frac{1}{d_0}\sum_{k \in \mathcal{H}_0} \Phi\left(\frac{-z - B_k}{\sqrt{\sigma_{\varepsilon,kk}}}\right) + O_{\mathbb{P}}[d^{-(1 \wedge \delta_0)/2} + n^{-1/4}\{\log(nd)\}^{1/4}]$$

uniformly over all $z \geq 0$. This, together with (D.34) and the fact that $|\Phi(z_1) - \Phi(z_2)| \leq (2\pi)^{-1/2}|z_1 - z_2|$, proves (D.15). $\qquad\square$

## D.5   Proof of Theorem D.2

Let $\boldsymbol{b}^{(1)}, \ldots, \boldsymbol{b}^{(r)} \in \mathbb{R}^d$ be the columns of $\mathbf{B}$. Without loss of generality, assume that $\|\boldsymbol{b}^{(1)}\|_2 \geq \cdots \geq \|\boldsymbol{b}^{(r)}\|_2$. Under (i) of Condition D.1, $\mathbf{B}\mathbf{B}^\intercal$ has non-vanishing eigenvalues $\{\overline{\lambda}_\ell := \|\boldsymbol{b}^{(\ell)}\|_2^2\}_{\ell=1}^r$ with eigenvectors $\{\overline{\boldsymbol{v}}_\ell := \boldsymbol{b}^{(\ell)}/\|\boldsymbol{b}^{(\ell)}\|_2\}_{\ell=1}^r$, and $\mathbf{B} = (\boldsymbol{b}_1, \ldots, \boldsymbol{b}_d)^\intercal = (\overline{\lambda}_1^{1/2}\overline{\boldsymbol{v}}_1, \ldots, \overline{\lambda}_r^{1/2}\overline{\boldsymbol{v}}_r)$. Moreover, write $\widehat{\mathbf{B}} = (\widehat{\boldsymbol{b}}_1, \ldots, \widehat{\boldsymbol{b}}_d)^\intercal = (\widehat{\lambda}_1^{1/2}\widehat{\boldsymbol{v}}_1, \ldots, \widehat{\lambda}_r^{1/2}\widehat{\boldsymbol{v}}_r)$ with $\widehat{\boldsymbol{v}}_\ell = (\widehat{v}_{\ell 1}, \ldots, \widehat{v}_{\ell d})^\intercal$ for $\ell = 1, \ldots, r$ and $\widehat{\boldsymbol{b}}_k = (\widehat{\lambda}_1^{1/2}\widehat{v}_{1k}, \ldots, \widehat{\lambda}_r^{1/2}\widehat{v}_{rk})^\intercal$ for $k = 1, \ldots, d$.

A key step in proving (D.12) is the derivation of an upper bound on the estimation error $\Delta_d := \max_{1 \leq k \leq d} \|\widehat{\boldsymbol{b}}_k\|_2 - \|\boldsymbol{b}_k\|_2$. By Weyl's inequality and the decomposition that $\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} = \mathbf{B}\mathbf{B}^\intercal + (\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}) + \boldsymbol{\Sigma}_\varepsilon$,

$$\max_{1 \leq \ell \leq r} |\widehat{\lambda}_\ell - \overline{\lambda}_\ell| \leq \|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_2 + \|\boldsymbol{\Sigma}_\varepsilon\|_2. \qquad (\text{D.42})$$

Applying Corollary 1 in Yu, Wang and Samworth (2015) yields that, for every $1 \le \ell \le r$,

$$\|\widehat{\boldsymbol{v}}_\ell - \overline{\boldsymbol{v}}_\ell\|_2 \le \frac{2^{3/2}(\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_2 + \|\boldsymbol{\Sigma}_\varepsilon\|_2)}{\min(\overline{\lambda}_{\ell-1} - \overline{\lambda}_\ell, \overline{\lambda}_\ell - \overline{\lambda}_{\ell+1})},$$

where we set $\overline{\lambda}_0 = \infty$ and $\overline{\lambda}_{r+1} = 0$. Under (ii) of Condition D.1, it follows that

$$\max_{1 \le \ell \le r} \|\widehat{\boldsymbol{v}}_\ell - \overline{\boldsymbol{v}}_\ell\|_2 \lesssim d^{-1}(\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_2 + \|\boldsymbol{\Sigma}_\varepsilon\|_2). \tag{D.43}$$

Moreover, apply Theorem 3 and Proposition 3 in Fan, Wang and Zhong (2018) to reach

$$\max_{1 \le \ell \le r} \|\widehat{\boldsymbol{v}}_\ell - \overline{\boldsymbol{v}}_\ell\|_\infty \lesssim r^4(d^{-1/2}\|\widehat{\boldsymbol{\Sigma}}_1^{\mathcal{H}} - \boldsymbol{\Sigma}\|_{\max} + d^{-1}\|\boldsymbol{\Sigma}_\varepsilon\|_2). \tag{D.44}$$

Note that, under (ii) of Condition D.1,

$$\|\overline{\boldsymbol{v}}_\ell\|_\infty = \|\boldsymbol{b}^{(\ell)}\|_\infty/\|\boldsymbol{b}^{(\ell)}\|_2 \le \|\mathbf{B}\|_{\max}/\|\boldsymbol{b}^{(\ell)}\|_2 \lesssim d^{-1/2} \text{ for all } \ell = 1, \dots, r. \tag{D.45}$$

Define $\widetilde{\boldsymbol{b}}_k = (\overline{\lambda}_1^{1/2}\widehat{v}_{1k}, \dots, \overline{\lambda}_r^{1/2}\widehat{v}_{rk})^\intercal$. By the triangular inequality,

$$\begin{aligned}
\|\widehat{\boldsymbol{b}}_k - \boldsymbol{b}_k\|_2 &\le \|\widehat{\boldsymbol{b}}_k - \widetilde{\boldsymbol{b}}_k\|_2 + \|\widetilde{\boldsymbol{b}}_k - \boldsymbol{b}_k\|_2 \\
&= \left\{ \sum_{\ell=1}^r (\widehat{\lambda}_\ell^{1/2} - \overline{\lambda}_\ell^{1/2})^2 \widehat{v}_{\ell k}^2 \right\}^{1/2} + \left\{ \sum_{\ell=1}^r \overline{\lambda}_\ell (\widehat{v}_{\ell k} - \overline{v}_{\ell k})^2 \right\}^{1/2} \\
&\le r^{1/2}\left( \max_{1 \le \ell \le r} |\widehat{\lambda}_\ell^{1/2} - \overline{\lambda}_\ell^{1/2}| \|\widehat{\boldsymbol{v}}_\ell\|_\infty + \max_{1 \le \ell \le r} \overline{\lambda}_\ell^{1/2} \|\widehat{\boldsymbol{v}}_\ell - \overline{\boldsymbol{v}}_\ell\|_\infty \right).
\end{aligned}$$

This, together with (D.42)–(D.45) and Theorem 3.3, implies

$$\Delta_d \le \max_{1 \le k \le d} \|\widehat{\boldsymbol{b}}_k - \boldsymbol{b}_k\|_2 = O_{\mathbb{P}}(w_{n,d}). \tag{D.46}$$

With the above preparations, now we are ready to prove (D.12). To that end, define $\widetilde{\boldsymbol{u}} = \sqrt{n}(\mathbf{B}^\intercal\mathbf{B})^{-1}\mathbf{B}^\intercal\overline{X}$, so that for every $1 \le k \le d$,

$$\boldsymbol{b}_k^\intercal\widetilde{\boldsymbol{u}} = \sqrt{n}\,\boldsymbol{b}_k^\intercal\overline{\boldsymbol{f}} + \sqrt{n}\,\boldsymbol{b}_k^\intercal(\mathbf{B}^\intercal\mathbf{B})^{-1}\mathbf{B}^\intercal\boldsymbol{\mu} + \sqrt{n}\,\boldsymbol{b}_k^\intercal(\mathbf{B}^\intercal\mathbf{B})^{-1}\mathbf{B}^\intercal\overline{\boldsymbol{\varepsilon}}.$$

Consider the decomposition

$$\begin{aligned}
&\left| \Phi\left(\frac{-z + \widehat{\boldsymbol{b}}_k^\intercal\widehat{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\right) - \Phi\left(\frac{-z + \sqrt{n}\,\boldsymbol{b}_k^\intercal\overline{\boldsymbol{f}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right| \\
&\le \left| \Phi\left(\frac{-z + \widehat{\boldsymbol{b}}_k^\intercal\widehat{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\right) - \Phi\left(\frac{-z + \boldsymbol{b}_k^\intercal\widetilde{\boldsymbol{u}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right| + \left| \Phi\left(\frac{-z + \boldsymbol{b}_k^\intercal\widetilde{\boldsymbol{u}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) - \Phi\left(\frac{-z + \sqrt{n}\,\boldsymbol{b}_k^\intercal\overline{\boldsymbol{f}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right| \\
&\le \left| \Phi\left(\frac{-z + \boldsymbol{b}_k^\intercal\widetilde{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\right) - \Phi\left(\frac{-z + \boldsymbol{b}_k^\intercal\widetilde{\boldsymbol{u}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right| \\
&\quad + \left| \Phi\left(\frac{-z + \widehat{\boldsymbol{b}}_k^\intercal\widehat{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\right) - \Phi\left(\frac{-z + \boldsymbol{b}_k^\intercal\widetilde{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\right) \right| + \left| \Phi\left(\frac{-z + \boldsymbol{b}_k^\intercal\widetilde{\boldsymbol{u}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) - \Phi\left(\frac{-z + \sqrt{n}\,\boldsymbol{b}_k^\intercal\overline{\boldsymbol{f}}}{\sqrt{\sigma_{\varepsilon,kk}}}\right) \right| \\
&:= \Delta_{k1} + \Delta_{k2} + \Delta_{k3}.
\end{aligned} \tag{D.47}$$

In the following, we deal with $\Delta_{k1}$, $\Delta_{k2}$ and $\Delta_{k3}$ separately.

By the mean value theorem, there exists some $\xi_k$ between $\widehat{\sigma}_{\varepsilon,kk}^{-1/2}$ and $\sigma_{\varepsilon,kk}^{-1/2}$ such that

$$\Delta_{k1} = \phi(\xi_k|z - \boldsymbol{b}_k^\top \widetilde{\boldsymbol{u}}|)|z - \boldsymbol{b}_k^\top \widetilde{\boldsymbol{u}}|\left|\frac{1}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}} - \frac{1}{\sqrt{\sigma_{\varepsilon,kk}}}\right|,$$

where $\phi(\cdot) = \Phi'(\cdot)$. With $\tau_{kk} \asymp \sqrt{n/\log(nd)}$ for $k = 1, \ldots, d$, it follows that the event

$$\mathcal{E}_0 = \left\{\max_{1 \le k \le d}|\widehat{\sigma}_{kk} - \sigma_{kk}| \lesssim \sqrt{\log(d)/n}\right\}$$

satisfies $\mathbb{P}(\mathcal{E}_0^c) \lesssim n^{-1}$. On $\mathcal{E}_0$, it holds $\widehat{\sigma}_{\varepsilon,kk}^{-1} \ge (2\sigma_{kk})^{-1}$, $\sigma_{\varepsilon,kk}^{-1} \ge \sigma_{kk}^{-1}$ and therefore $\xi_k \ge (2\sigma_{kk})^{-1/2}$ uniformly for all $1 \le k \le d$ as long as $n \gtrsim \log d$. This further implies $\max_{1 \le k \le d}\max_{z \ge 0}\phi(\xi_k|z - \boldsymbol{b}_k^\top \widetilde{\boldsymbol{u}}|)|z - \boldsymbol{b}_k^\top \widetilde{\boldsymbol{u}}| = O_{\mathbb{P}}(1)$. By (D.46),

$$\left|\frac{1}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}} - \frac{1}{\sqrt{\sigma_{\varepsilon,kk}}}\right| = O_{\mathbb{P}}(|\widehat{\sigma}_{kk} - \sigma_{kk}| + \|\widehat{\boldsymbol{b}}_k - \boldsymbol{b}_k\|_2) = O_{\mathbb{P}}(w_{n,d}) \tag{D.48}$$

uniformly over $k = 1, \ldots, d$, where $w_{n,d} = \sqrt{\log(d)/n} + d^{-1/2}$. Putting the above calculations together, we arrive at

$$\frac{1}{d}\sum_{k=1}^{d}\Delta_{k1} = O_{\mathbb{P}}(w_{n,d}). \tag{D.49}$$

Turning to $\Delta_{k2}$, again by the mean value theorem, there exists some $\eta_k$ between $\widehat{\boldsymbol{b}}_k^\top \widetilde{\boldsymbol{u}}$ and $\boldsymbol{b}_k^\top \widetilde{\boldsymbol{u}}$ such that

$$\Delta_{k2} = \phi\left(\frac{-z + \eta_k}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\right)\frac{\widehat{\boldsymbol{b}}_k^\top \widetilde{\boldsymbol{u}} - \boldsymbol{b}_k^\top \widetilde{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}.$$

In view of (D.48),

$$\max_{1 \le k \le d}\phi\left(\frac{-z + \eta_k}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\right)\frac{1}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}} = O_{\mathbb{P}}(1).$$

Observe that $\widehat{\mathbf{B}}\widehat{\boldsymbol{u}} = (\sum_{\ell=1}^{r}\widehat{\boldsymbol{v}}_\ell\widehat{\boldsymbol{v}}_\ell^\top)\boldsymbol{Z}$ and $\mathbf{B}\widetilde{\boldsymbol{u}} = (\sum_{\ell=1}^{r}\overline{\boldsymbol{v}}_\ell\overline{\boldsymbol{v}}_\ell^\top)\boldsymbol{Z}$, where $\boldsymbol{Z} = \sqrt{n}\overline{\boldsymbol{X}}$. By the Cauchy-Schwarz inequality,

$$\sum_{k=1}^{d}|\widehat{\boldsymbol{b}}_k^\top \widehat{\boldsymbol{u}} - \boldsymbol{b}_k^\top \widetilde{\boldsymbol{u}}| \le d^{1/2}\|\widehat{\mathbf{B}}\widehat{\boldsymbol{u}} - \mathbf{B}\widetilde{\boldsymbol{u}}\|_2$$

$$\le d^{1/2}\left\|\sum_{\ell=1}^{r}(\widehat{\boldsymbol{v}}_\ell\widehat{\boldsymbol{v}}_\ell^\top - \overline{\boldsymbol{v}}_\ell\overline{\boldsymbol{v}}_\ell^\top)\right\|_2\|\boldsymbol{Z}\|_2 \le 2rd^{1/2}\max_{1 \le \ell \le r}\|\widehat{\boldsymbol{v}}_\ell - \overline{\boldsymbol{v}}_\ell\|_2\|\boldsymbol{Z}\|_2. \tag{D.50}$$

For $\|\boldsymbol{Z}\|_2$, we calculate $\mathbb{E}\|\boldsymbol{Z}\|_2^2 = n\|\boldsymbol{\mu}\|_2^2 + \sum_{k=1}^{d}\sigma_{kk}$, indicating that $\|\boldsymbol{Z}\|_2 = O_{\mathbb{P}}(\sqrt{n}\|\boldsymbol{\mu}\|_2 + d^{1/2})$. Combining this with (D.43), (D.50) and Theorem 3.3, we conclude that

$$\frac{1}{d}\sum_{k=1}^{d}\Delta_{k2} = O_{\mathbb{P}}\{(d^{-1/2}\sqrt{n}\|\boldsymbol{\mu}\|_2 + 1)w_{n,d}\}. \tag{D.51}$$

For $\Delta_{k3}$, following the same arguments as above, it suffices to consider

$$\frac{\sqrt{n}}{d}\sum_{k=1}^{d}|\boldsymbol{b}_k^\top(\mathbf{B}^\top\mathbf{B})^{-1}\mathbf{B}^\top\boldsymbol{\mu} + \boldsymbol{b}_k^\top(\mathbf{B}^\top\mathbf{B})^{-1}\mathbf{B}^\top\overline{\boldsymbol{\varepsilon}}|,$$

which, by the Cauchy-Schwarz inequality, is bounded by

$$\sqrt{\frac{n}{d}}\Big\|\sum_{\ell=1}^{r}\overline{\boldsymbol{v}}_\ell\overline{\boldsymbol{v}}_\ell^\top\Big\|_2\|\boldsymbol{\mu}\|_2 + \max_{1\le k\le d}\|\boldsymbol{b}_k\|_2\|\boldsymbol{u}\|_2 \le \sqrt{\frac{n}{d}}\|\boldsymbol{\mu}\|_2 + \max_{1\le k\le d}\sqrt{\sigma_{kk}}\|\boldsymbol{u}\|_2,$$

where $\boldsymbol{u} = \sqrt{n}\,(\mathbf{B}^\top\mathbf{B})^{-1}\mathbf{B}^\top\overline{\boldsymbol{\varepsilon}} \in \mathbb{R}^r$ is a zero-mean random vector with covariance matrix $\boldsymbol{\Sigma}_u = (\mathbf{B}^\top\mathbf{B})^{-1}\mathbf{B}^\top\boldsymbol{\Sigma}_\varepsilon\mathbf{B}(\mathbf{B}^\top\mathbf{B})^{-1}$. Recall that $\mathbf{B}^\top\mathbf{B} \in \mathbb{R}^{r\times r}$ has non-increasing eigenvalues $\overline{\lambda}_1 \ge \cdots \ge \overline{\lambda}_r$. Under (ii) of Condition D.1, it holds $\mathbb{E}\|\boldsymbol{u}\|_2^2 = \mathrm{Tr}(\boldsymbol{\Sigma}_u) \le \|\boldsymbol{\Sigma}_\varepsilon\|\sum_{\ell=1}^{r}\overline{\lambda}_\ell^{-1} \lesssim rd^{-1}$ and thus $\|\boldsymbol{u}\|_2 = O_{\mathbb{P}}(d^{-1/2})$. Putting the pieces together, we get

$$\frac{1}{d}\sum_{k=1}^{d}\Delta_{k3} = O_{\mathbb{P}}(d^{-1/2}\sqrt{n}\,\|\boldsymbol{\mu}\|_2 + d^{-1/2}). \tag{D.52}$$

Combining (D.47), (D.49), (D.51) and (D.52), we reach

$$\frac{1}{d}\sum_{k=1}^{d}\Phi\Big(\frac{-z+\widehat{\boldsymbol{b}}_k^\top\widehat{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\Big) = \frac{1}{d}\sum_{k=1}^{d}\Phi\Big(\frac{-z+\sqrt{n}\,\boldsymbol{b}_k^\top\overline{\boldsymbol{f}}}{\sqrt{\sigma_{\varepsilon,kk}}}\Big) + O_{\mathbb{P}}(w_{n,d} + d^{-1/2}\sqrt{n}\,\|\boldsymbol{\mu}\|_2).$$

Using the same argument, it can similarly derived that

$$\frac{1}{d}\sum_{k=1}^{d}\Phi\Big(\frac{-z-\widehat{\boldsymbol{b}}_k^\top\widehat{\boldsymbol{u}}}{\sqrt{\widehat{\sigma}_{\varepsilon,kk}}}\Big) = \frac{1}{d}\sum_{k=1}^{d}\Phi\Big(\frac{-z-\sqrt{n}\,\boldsymbol{b}_k^\top\overline{\boldsymbol{f}}}{\sqrt{\sigma_{\varepsilon,kk}}}\Big) + O_{\mathbb{P}}(w_{n,d} + d^{-1/2}\sqrt{n}\,\|\boldsymbol{\mu}\|_2).$$

Together, the last two displays lead to the stated result (D.12). $\qquad\square$

# References

ANDERSON, T. W. and RUBIN, H. (1956). Statistical inference for factor analysis. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, 1954–1955, *Vol. V* 111–150. Univ. California Press, Berkeley.

BAI, J. and LI, K. (2012). Statistical analysis of factor models of high dimension. *Ann. Statist.* **40** 436–465.

BARRAS, L., SCAILLET, O. and WERMERS, R. (2010). False discoveries in mutual fund performance: Measuring luck in estimated alphas. *J. Finance* **65** 179–216.

BERK, J. B. and GREEN, R. C. (2004). Mutual fund flows and performance in rational markets. *J. Polit. Econ.* **112** 1269–1295.

BENTKUS, V. (2005). A Lyapunov-type bound in $R^d$. *Theory Probab. Appl.* **49** 311–323.

CAI, T. T. and LIU, W. (2011). Adaptive thresholding for sparse covariance matrix estimation. *J. Amer. Statist. Assoc.* **106** 672–684.

CATONI, O. (2012). Challenging the empirical mean and empirical variance: A deviation study. *Ann. Inst. Henri Poincaré Probab. Stat.* **48** 1148–1185.

CHAMBERLAIN, G. and ROTHSCHILD, M. (1983). Arbitrage, factor structure and mean-variance analysis in large asset markets. *Econometrica* **51** 1305–1324.

FAMA, E. F. and FRENCH, K. R. (1993). Common risk factors in the returns on stocks and bonds. *J. Financ. Econ.* **33** 3–56.

FAN, J. and HAN, X. (2017). Estimation of the false discovery proportion with unknown dependence. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **79** 1143–1164.

FAN, J., KE, Y., SUN, Q. and ZHOU, W.-X. (2019). FarmTest: Factor-adjusted robust multiple testing with approximate false discovery control. *J. Amer. Statist. Assoc.* To appear.

FAN, J., LI, Q. and WANG, Y. (2017). Estimation of high dimensional mean regression in the absence of symmetry and light tail assumptions. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **79** 247–265.

FAN, J., LIAO, Y. and MINCHEVA, M. (2013). Large covariance estimation by thresholding principal orthogonal complements. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **75** 603–680.

FAN, J., LIU, H. and WANG, W. (2018). Large covariance estimation through elliptical factor models. *Ann. Statist.* **46** 1383–1414.

FAN, J., WANG, W. and ZHONG, Y. (2018). An $\ell_\infty$ eigenvector perturbation bound and its application to robust covariance estimation. *J. Mach. Learn. Res.* **18**(207) 1–42.

HOEFFDING, W. (1963). Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* **58** 13–30.

HSU, D., KAKADE, S. M. and ZHANG, T. (2012). A tail inequality for quadratic forms of subgaussian random vectors. *Electron. Commun. Probab.* **17**(52) 1–6.

LAN, W. and DU, L. (2019). A factor-adjusted multiple testing procedure with application to mutual fund selection. *J. Bus. Econom. Statist.* **37** 147–157.

LI, W. V. and SHAO, Q.-M. (2002). A normal comparison inequality and its applications. *Probab. Theory Relat. Fields* **122** 494–508.

MINSKER, S. (2018). Sub-Gaussian estimators of the mean of a random matrix with heavy-tailed entries. *Ann. Statist.* **46** 2871–2903.

MINSKER, S. and STRAWN, N. (2017). Distributed statistical estimation and rates of convergence in normal approximation. *arXiv preprint arXiv:1704.02658.*

MINSKER, S. and WEI, X. (2018). Robust modifications of U-statistics and applications to covariance estimation problems. *arXiv preprint arXiv:1801.05565.*

PINELIS, I. and MOLZON, R. (2016). Optimal-order bounds on the rate of convergence to normality in the multivariate delta method. *Electron. J. Statist.* **10** 1001–1063.

ROTHMAN, A. J., LEVINA, E. and ZHU, J. (2009). Generalized thresholding of large covariance matrices. *J. Amer. Statist. Assoc.* **104** 177–186.

SHARPE, W. F. (1964). Capital asset prices: A theory of market equilibrium under conditions of risk. *J. Finance* **19** 425–442.

VERSHYNIN, R. (2012). Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing* 210–268. Cambridge Univ. Press, Cambridge.

YU, Y., WANG, T. and SAMWORTH, R. J. (2015). A useful variant of the Davis–Kahan theorem for statisticians. *Biometrika* **102** 315–323.

ZHOU, W.-X., BOSE, K., FAN, J. and LIU, H. (2018). A new perspective on robust $M$-estimation: Finite sample theory and applications to dependence-adjusted multiple testing. *Ann. Statist.* **46** 1904–1931.