

A Locally Adaptive Shrinkage Approach to False Selection Rate Control in High-Dimensional Classification

Bowen Gang¹, Yuantao Shi² and Wenguang Sun³

¹*Fudan University*, ²*University of Chicago* and ³*Zhejiang University*

Supplementary Material

This supplement contains the proofs of main theorems, propositions and corollaries (Section S1), proofs of technical lemmas (Section S2), an argument establishing the asymptotic equivalence of FSR and mFSR (Section S3), numerical illustrations of the effects of shrinkage factor (Section S4), additional simulation result on a “boarderline sparse” setting (Section S5), additional numerical results on a lung cancer data (Section S6), an example showing the advantage LASS has over LPD (Section S7) and a proof showing class-specific FSR control implies global FSR control (Section S8).

S1 Proofs of main theorems, propositions and corollaries

We swap the proofs of Theorems 4 and 3 as the latter is simpler. Other proofs are arranged according to the orders in the text.

S1.1 Proof of Theorem 1

We only need to show $\delta_{OR}^j = 2\mathbb{I}(T^j < t_{OR}^2)$ controls mFSR^2 at level α_2 and maximizes $\mathbb{E}\{\sum_{i=1}^m \mathbb{I}(\theta_j = 2, \delta_j = 2)\}$. The same argument can be used to show that $\delta_{OR}^j = \mathbb{I}(1 - T^j < t_{OR}^1)$ controls mFSR^1 at level α_1 and maximizes $\mathbb{E}\{\sum_{i=1}^m \mathbb{I}(\theta_j = 1, \delta_j = 1)\}$. The theorem is proved when we combine these two statements.

The proof is divided into two parts. In part (a), we establish two properties of the classification rule $\delta^2(t) = \{2\mathbb{I}(T^j < t) : 1 \leq j \leq m\}$. We show that it produces $\text{mFSR}^2 < t$ for all $t \in (0, 1)$ and that its mFSR^2 is monotonic in t . In part (b) we show that when the threshold is t_{OR}^2 , the classification rule has $\text{mFSR}^2 = \alpha_2$ and maximizes $\mathbb{E}\{\sum_{i=1}^m \mathbb{I}(\theta_j = 2, \delta_j = 2)\}$ amongst all valid classification rule with $\text{mFSR}^2 \leq \alpha_2$.

Part(a). Consider classification rule $\{2 \cdot \mathbb{I}(T^j < t) : 1 \leq j \leq m\}$. Let $Q^2(t) = \alpha_t$ be its the mFSR^2 level. We first show that $\alpha_t < t$. Since $T^j = P(\theta_j = 1 | \mathbf{W}_j)$, then

$$\mathbb{E}\left\{\sum_{j=1}^m \mathbb{I}(\theta_j = 1, \delta_j = 2)\right\} = \mathbb{E}_{\mathbf{W}}\left[\left\{\sum_{j=1}^m \mathbb{E}_{\theta|\mathbf{W}}\mathbb{I}(\theta_j = 1, \delta_j = 2)\right\}\right] = \mathbb{E}_{\mathbf{W}}\left\{\sum_{i=1}^m T^j \mathbb{I}(\delta_j = 2)\right\},$$

where \mathbb{E} is the expectation over (\mathbf{W}, θ) , $E_{\mathbf{W}}$ is the expectation over \mathbf{W} , and $E_{\theta|\mathbf{W}}$ is the expectation over θ holding \mathbf{W} fixed. Recall $Q_{OR}^2(t) = \alpha_t$ is the

FSR² level of the classification rule $\{2 \cdot \mathbb{I}(T^j < t) : 1 \leq i \leq m\}$, we have

$$\mathbb{E}_{\mathbf{W}} \left\{ \sum_{j=1}^m (T^j - \alpha_t) \mathbb{I}(T^j < t) \right\} = 0. \quad (\text{S1.1})$$

This implies that $\alpha_t < t$. To see this, if $\alpha_t \geq t$, then $(T^j - \alpha_t) \mathbb{I}(T^j < t) < 0$, which contradicts the right hand side.

Next, we show that $Q^2(t)$ is nondecreasing in t . That is, letting $Q^2(t_j) = \alpha_{t_j}$, if $t_1 < t_2$, then $\alpha_{t_1} \leq \alpha_{t_2}$. We argue by contradiction. Suppose that $t_1 < t_2$ but $\alpha_{t_1} > \alpha_{t_2}$. Then

$$\begin{aligned} & (T^j - \alpha_{t_2}) \mathbb{I}(T^j < t_2) & (\text{S1.2}) \\ = & (T^j - \alpha_{t_1}) \mathbb{I}(T^j < t_1) + (\alpha_{t_1} - \alpha_{t_2}) \mathbb{I}(T^j < t_1) + (T^j - \alpha_{t_2}) \mathbb{I}(t_1 \leq T^j < t_2) \\ \geq & (T^j - \alpha_{t_1}) \mathbb{I}(T^j < t_1) + (\alpha_{t_1} - \alpha_{t_2}) \mathbb{I}(T^j < t_1) + (T^j - \alpha_{t_1}) \mathbb{I}(t_1 \leq T^j < t_2). \end{aligned}$$

By (S1.1) we have $\mathbb{E} \left\{ \sum_{j=1}^m (T^j - \alpha_{t_1}) \mathbb{I}(T^j < t_1) \right\} = 0$, together with the fact that $\alpha_{t_1} < t_1$ we have $\mathbb{E} \left\{ \sum_{j=1}^m (T^j - \alpha_{t_2}) \mathbb{I}(T^j < t_2) \right\} > 0$, contradicting (S1.1).

Part(b). Define $t_{OR}^2 = \sup_t \{t \in (0, 1) : Q^2(t) \leq \alpha\}$. By part (a), $Q^2(t)$ is non-decreasing in t . By continuity, $Q^2(T) = \alpha_2$. Next, consider the oracle rule $\delta_{OR}^2 = (\delta_{OR}^{2,1}, \dots, \delta_{OR}^{2,m}) = \{2\mathbb{I}(T^j < t_{OR}^2) : 1 \leq j \leq m\}$ and an arbitrary rule $\delta_* = (\delta_*^1, \dots, \delta_*^m)$ such that $\text{mFSR}_{\delta_*}^2 \leq \alpha_2$. Using the result in part

(a), we have

$$\mathbb{E} \left\{ \sum_{j=1}^m (T^j - \alpha) \mathbb{I}(\delta_{OR}^{2,j} = 2) \right\} = 0 \quad \text{and} \quad \mathbb{E} \left\{ \sum_{j=1}^m (T^j - \alpha) \mathbb{I}(\delta_*^j = 2) \right\} \leq 0. \quad (\text{S1.3})$$

Taking the difference of the two equations in (S1.3), we have

$$\mathbb{E} \left\{ \sum_{j=1}^m (T^j - \alpha) \mathbb{I}(\delta_{OR}^{2,j} = 2) - (T^j - \alpha) \mathbb{I}(\delta_*^j = 2) \right\} \geq 0. \quad (\text{S1.4})$$

Next consider the transformation $f(x) = (x - \alpha)/(1 - x)$. Note that $f'(x) = (1 - \alpha)/(1 - x)^2 > 0$, $f(x)$ is monotonically increasing, the order is preserved by this transformation: if $T^i < t_{OR}^2$ then $f(T^i) < f(t_{OR}^2)$. This means we can rewrite the oracle rule as

$$\delta_{OR}^{2,i} = 2 \mathbb{I} \left[\left\{ (T^i - \alpha)/(1 - T^i) \right\} < \lambda_{OR}^2 \right],$$

where $\lambda_{OR}^2 = (t_{OR}^2 - \alpha)/(1 - t_{OR}^2)$. It will be useful to note that, from part (a), we have $\alpha_{t_2} < t_{OR}^2 < 1$, which implies that $\lambda_{OR}^2 > 0$. Note that

$$\mathbb{E} \left[\sum_{j=1}^m \left\{ \mathbb{I}(\delta_{OR}^{2,j} = 2) - \mathbb{I}(\delta_*^j = 2) \right\} \left\{ (T^j - \alpha) - \lambda_{OR}^2 (1 - T^j) \right\} \right] \leq 0 \quad (\text{S1.5})$$

To see this, consider that if $\mathbb{I}(\delta_{OR}^{2,j} = 2) - \mathbb{I}(\delta_*^j = 2) \neq 0$, then either (i) $\mathbb{I}(\delta_{OR}^{2,j} = 2) - \mathbb{I}(\delta_*^j = 2) > 0$ or (ii) $\mathbb{I}(\delta_{OR}^{2,j} = 2) - \mathbb{I}(\delta_*^j = 2) < 0$ holds. If

(i) holds, then $\delta_{OR}^{2,j} = 2$ and it follows that $\{(T^j - \alpha)/(1 - T^j)\} < \lambda_{OR}^2$. If (ii) holds, then $\delta_{OR}^j \neq 2$ and $\{(T^j - \alpha)/(1 - T^j)\} \geq \lambda_{OR}^2$. In both cases, we have

$$\{\mathbb{I}(\delta_{OR}^{2,j} = 2) - \mathbb{I}(\delta_*^j = 2)\} \{(T^j - \alpha) - \lambda_{OR}^2(1 - T^j)\} \leq 0.$$

Summing over all m terms and taking the expectation yields (S1.5).

Combining (S1.4) and (S1.5), we obtain

$$0 \leq \lambda_{OR}^2 \mathbb{E} \left[\sum_{j=1}^m \{\mathbb{I}(\delta_{OR}^{2,j} = 2) - \mathbb{I}(\delta_*^j = 2)\} (1 - T^j) \right]$$

Finally, since $\lambda_{OR}^2 > 0$, it follows that

$$\mathbb{E} \left[\sum_{i=1}^m \{\mathbb{I}(\delta_{OR}^j = 2) - \mathbb{I}(\delta_*^j = 2)\} (1 - T^j) \right] \geq 0.$$

□

S1.2 Proof of Theorem 2

As in the proof of Theorem 1, we will show $\boldsymbol{\delta}_{OR}^{*2} = \{\delta_{OR}^{*2,j} = 2\mathbb{I}(T^j < \hat{t}_{OR}^2), j = 1, \dots, m\}$ controls mFSR² and FSR² at level α_2 . Let $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)$. Note

that

$$\text{FSR}^2 = \mathbb{E} \left[\frac{\sum_{j=1}^m \delta_{OR}^{*2,j} \mathbb{I}(\theta_j \neq 2)}{\{\sum_{j=1}^m \delta_{OR}^{*2,j}\} \vee 1} \right] = \mathbb{E}_{\delta_{OR}^{*2}} \mathbb{E}_{\boldsymbol{\theta} | \delta_{OR}^{*2}} \left[\frac{\sum_{j=1}^m \delta_{OR}^{*2,j} \mathbb{I}(\theta_j \neq 2)}{\{\sum_{j=1}^m \delta_{OR}^{*2,j}\} \vee 1} \middle| \delta_{OR}^{*2} \right].$$

By definition of the rule δ_{OR}^{*2} , if $\sum_{j=1}^m \delta_{OR}^{*2,j} = k \geq 1$ then

$$\mathbb{E}_{\boldsymbol{\theta} | \delta_{OR}^{*2}} \sum_{j=1}^m \delta_{OR}^{*2,j} \mathbb{I}(\theta_j \neq 2) = \sum_{i=1}^m \mathbb{P}(\theta_j \neq 2) \delta_{OR}^{*2,j} \leq k \alpha_2.$$

It follows that $\mathbb{E}_{\boldsymbol{\theta} | \delta_{OR}^{*2}} \left[\frac{\sum_{j=1}^m \delta_{OR}^{*2,j} \mathbb{I}(\theta_j \neq 2)}{\{\sum_{j=1}^m \delta_{OR}^{*2,j}\} \vee 1} \middle| \delta_{OR}^{*2} \right] \leq \alpha$, and $\text{FSR}^2 \leq \alpha_2$.

Let $\mathcal{W} = (\mathbf{W}_1, \dots, \mathbf{W}_m)$. For mFSR^2 we have

$$\begin{aligned} \frac{\mathbb{E} \left(\sum_{j=1}^m \mathbb{I}(\theta_j \neq 2) \delta_{OR}^{*2,j} \right)}{\mathbb{E} \left(\sum_{j=1}^m \delta_{OR}^{*2,j} \right)} &= \frac{\mathbb{E}_{\mathcal{W}} \mathbb{E}_{\boldsymbol{\theta} | \mathcal{W}} \left\{ \sum_{j=1}^m \mathbb{I}(\theta_j \neq 2) \delta_{OR}^{*2,j} \right\}}{\mathbb{E} \left(\sum_{j=1}^m \delta_{OR}^{*2,j} \right)} \\ &= \frac{\mathbb{E}_{\mathcal{W}} \mathbb{E}_{\boldsymbol{\theta} | \mathcal{W}} \left\{ \left(\sum_{j=1}^m \delta_{OR}^{*2,j} \right) \left(\sum_{j=1}^m \mathbb{I}(\theta_j \neq 2) \delta_{OR}^{*2,j} / \sum_{j=1}^m \delta_{OR}^{*2,j} \right) \right\}}{\mathbb{E} \left(\sum_{j=1}^m \delta_{OR}^{*2,j} \right)} \\ &= \frac{\mathbb{E}_{\mathcal{W}} \left(\sum_{j=1}^m \delta_{OR}^{*2,j} \right) \mathbb{E}_{\mathcal{W}} \left\{ \mathbb{E}_{\boldsymbol{\theta} | \mathcal{W}} \left(\sum_{j=1}^m \mathbb{I}(\theta_j \neq 2) \delta_{OR}^{*2,j} / \sum_{j=1}^m \delta_{OR}^{*2,j} \right) \right\}}{\mathbb{E}_{\mathcal{W}} \left(\sum_{j=1}^m \delta_{OR}^{*2,j} \right)} \\ &= \mathbb{E}_{\mathcal{W}} \left\{ \mathbb{E}_{\boldsymbol{\theta} | \mathcal{W}} \left(\sum_{j=1}^m \mathbb{I}(\theta_j \neq 2) \delta_{OR}^{*2,j} / \sum_{j=1}^m \delta_{OR}^{*2,j} \right) \right\}. \end{aligned}$$

The fourth equality holds because δ_{OR}^{*2} depends on $\boldsymbol{\theta}$ only through \mathcal{W} . Now

given \mathcal{W} and that $\sum_i \delta_{OR}^{*2,j} = k$, we have

$$\begin{aligned} \mathbb{E}_{\mathcal{W}} \left\{ E_{\theta|\mathcal{W}} \left(\frac{\sum_{j=1}^m \mathbb{I}(\theta_j \neq 2) \delta_{OR}^{*2,j}}{\sum_{j=1}^m \delta_{OR}^{*2,j}} \right) \right\} &= \mathbb{E}_{\mathcal{W}} \left\{ \frac{\mathbb{E}_{\theta|\mathcal{W}} \left(\sum_{j=1}^m \mathbb{I}(\theta_j \neq 2) \right) \delta_{OR}^{*2,j}}{k} \right\} \\ &= \mathbb{E}_{\mathcal{W}} \left(\frac{\sum_{j=1}^m T^j \delta_{OR}^{*2,j}}{k} \right). \end{aligned}$$

By the definition of δ_{OR}^{*2} , we have

$$\mathbb{E}_{\mathcal{W}} \left(\frac{\sum_{j=1}^m T^j \delta_{OR}^{*2,j}}{k} \right) \leq \alpha_2.$$

Hence δ_{OR}^{*2} satisfies $\text{mFSR}^2 \leq \alpha_2$. □

S1.3 Proof of Proposition 1

Let $q_k(x) := \tilde{g}_{1k}(|x|) / \{g_0(|x|) + \tilde{g}_{1k}(|x|)\}$, where g_0 and \tilde{g}_{1k} are the density function of $\mathcal{N}\left(0, \frac{n_1+n_2}{n_1 n_2}\right)$ and $\mathcal{N}\left(\left\{(2+b)\sqrt{\sigma_{kk}} + \sqrt{(2+b)^2 \sigma_{kk} + 4}\right\} \sqrt{\frac{(n_1+n_2)}{2n_1 n_2} \log p}, \frac{n_1+n_2}{n_1 n_2}\right)$ respectively. We will assume without loss of generality that $d_k > 0$.

We shall first prove the result when the true variance σ_{kk} is known and then argue that with probability greater than $1 - pe^{-O(n)}$ the result still holds when σ_{kk} is replaced by $\hat{\sigma}_{kk}$. Let g_{2k} be the density function of $\mathcal{N}\left(d_k, \sigma_{kk} \frac{n_1+n_2}{n_1 n_2}\right)$. Assume d_k is strong, we investigate the asymptotic

behavior of

$$\begin{aligned}
1 - \mathbb{E}\{g_k(X)|X \sim g_{2k}\} &= \int_{-\infty}^{\infty} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx & (S1.6) \\
&= \int_{-\infty}^{d_k - \frac{\epsilon}{4} \sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx \\
&\quad + \int_{d_k - \frac{\epsilon}{4} \sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}}^{\infty} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx.
\end{aligned}$$

Consider the first term, it is easy to see that

$$\int_{-\infty}^{d_k - \frac{\epsilon}{4} \sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx \leq \int_{-\infty}^{d_k - \frac{\epsilon}{4} \sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}} g_{2k}(x) dx = O(p^{-\epsilon^2/(64\sigma_{kk})}).$$

Since $d_k \in \mathcal{G}_1$, we have $x > (a_k/2 + 3\epsilon/4) \sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}$ if $\left(d_k - (\epsilon/4) \sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}, \infty\right)$;

on this interval we have

$$\begin{aligned}
\frac{\tilde{g}_{1k}(|x|)}{g_0(|x|)} &= \exp\left(\frac{2a_k|x|\sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p} - a_k^2 \frac{(n_1+n_2)}{2n_1n_2} \log p}{2 \frac{n_1+n_2}{n_1n_2}}\right) \\
&\geq \exp\left(\frac{1}{2}a_k(a_k/2 + 3\epsilon/4) \log p - \frac{1}{4}a_k^2 \log p\right) \\
&\geq \exp\left(\frac{3a_k\epsilon}{8} \log p\right) \\
&\geq p^{3a_k\epsilon/8}.
\end{aligned}$$

It follows that

$$\begin{aligned} \int_{(d_k - \epsilon/4)\sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}}^{\infty} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx &\leq \sup_{x > (d_k - \epsilon/4)\sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}} \frac{1}{1 + \tilde{g}_{1k}(x)/g_0(x)} \\ &\leq p^{-3a_k\epsilon/8}. \end{aligned}$$

Using (S1.6), we have $\mathbb{E}(q_k | \theta_k = 1) = O(p^{-\epsilon_1})$, where $\epsilon_1 > 0$ is some constant.

Next, when $d_k \in \mathcal{G}_3$, we can similarly show that

$$\mathbb{E}\{q_k(X) | X \sim g_{2k}\} = \left| 1 - \int_{-\infty}^{\infty} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx \right| = O(p^{-\epsilon_2}).$$

Finally, consider the the case when $d_k \in \mathcal{G}_2$. We have

$$\begin{aligned} \int_{-\infty}^{\infty} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx &= \int_{-\infty}^{\infty} \frac{g_{2k}(x)}{1 + \exp\left(\frac{2a_k\sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}|x| - a_k^2\frac{(n_1+n_2)}{2n_1n_2} \log p}{2\frac{n_1+n_2}{n_1n_2}}\right)} dx \\ &= \int_{-\infty}^{\infty} \frac{g_{2k}(x)}{1 + \exp\left(\frac{1}{\sqrt{2}}a_k|x|\sqrt{\frac{n_1n_2}{n_1+n_2} \log p} - \frac{a_k^2}{4} \log p\right)} dx. \end{aligned}$$

$$\begin{aligned}
\text{Let } \mathcal{A}_t &= \left\{ x : -t/\sqrt{\frac{2n_1n_2}{n_1+n_2}} + d_k \leq x \leq t/\sqrt{\frac{2n_1n_2}{n_1+n_2}} + d_k \right\}, \\
& \int_{-\infty}^{\infty} \frac{g_{2k}(x)}{1 + \exp\left(\frac{1}{\sqrt{2}}a_k|x|\sqrt{\frac{n_1n_2}{n_1+n_2}}\log p - \frac{a_k^2}{4}\log p\right)} - g_{2k}(x)dx \\
&= \int_{x \in \mathcal{A}_t^c} \frac{g_{2k}(x)}{1 + \exp\left(\frac{1}{\sqrt{2}}a_k|x|\sqrt{\frac{n_1n_2}{n_1+n_2}}\log p - \frac{a_k^2}{4}\log p\right)} - g_{2k}(x)dx \\
&+ \int_{x \in \mathcal{A}_t} \frac{g_{2k}(x)}{1 + \exp\left(\frac{1}{\sqrt{2}}a_k|x|\sqrt{\frac{n_1n_2}{n_1+n_2}}\log p - \frac{a_k^2}{4}\log p\right)} - g_{2k}(x)dx.
\end{aligned}$$

Some algebra shows that

$$\begin{aligned}
& \left| \int_{x \in \mathcal{A}_t^c} \frac{g_{2k}(x)}{1 + \exp\left(\frac{1}{\sqrt{2}}a_k|x|\sqrt{\frac{n_1n_2}{n_1+n_2}}\log p - \frac{a_k^2}{4}\log p\right)} - g_{2k}(x)dx \right| \\
& \leq \int_{x \in \mathcal{A}_t^c} g_{2k}(x)dx = O\left(e^{-\frac{t^2}{4\sigma_{kk}}}\right), \text{ and} \tag{S1.7}
\end{aligned}$$

$$\begin{aligned}
& \left| \int_{x \in \mathcal{A}_t} \frac{g_{2k}(x)}{1 + \exp\left(\frac{1}{\sqrt{2}}a_k|x|\sqrt{\frac{n_1n_2}{n_1+n_2}}\log p - \frac{a_k^2}{4}\log p\right)} dx \right| \\
& \leq \left| \int_{x \in \mathcal{A}_t} \frac{\phi_{2k}(x)}{1 + \exp\left(\frac{1}{2}a_k t \sqrt{\log p} + o(\log p) - \frac{a_k^2}{4}\log p\right)} - g_{2k}(x)dx \right|. \tag{S1.8}
\end{aligned}$$

Take $t = (2 + b/2)\sqrt{\sigma_{kk}\log p}$, we can see that (S1.7) is bounded by $o(1/p^{1+b})$.

For (S1.8) to be bounded by $o(1/p^{1+a})$ for some constant $a > 0$, we need

$a_k^2/4 - (1 + b/4)a_k\sqrt{\sigma_{kk}} > 1 + c$ for some constant $c > 0$. Some computa-

tion shows that $a_k > (2 + b/2)\sqrt{\sigma_{kk}} + \sqrt{(2 + b/2)^2\sigma_{kk} + 4}$ can satisfy both

requirements. Since we take $a_k = (2 + b)\sqrt{\sigma_{kk}} + \sqrt{(2 + b)^2\sigma_{kk} + 4}$ we have

(S1.8) is bounded by:

$$\left| \int_{x \in \mathcal{A}_t} \frac{g_{2k}(x)}{1 + o(1/p^{1+\epsilon_3})} - g_{2k}(x) dx \right| = O(1/p^{1+\epsilon_3}) \int g_{2k}(x) dx = O(1/p^{1+\epsilon_3}),$$

where $\epsilon_3 > 0$ is some constant. Use the fact that $\int \phi_{2k} dx = 1$ we conclude that

$$\mathbb{E}\{q_k(X) | X \sim g_{2k}\} = \left| 1 - \int_{-\infty}^{\infty} \frac{g_0(|x|)g_{2k}(x)}{g_0(|x|) + \tilde{g}_{1k}(|x|)} dx \right| = O(1/p^{1+\epsilon_3}).$$

Take $\gamma = \min(\epsilon_1, \epsilon_2, \epsilon_3)$, the results follow.

Define $\hat{a}_k = (2 + b)\sqrt{\hat{\sigma}_{kk}} + \sqrt{(2 + b)^2\hat{\sigma}_{kk} + 4}$. From the above proof, we can see that if we would like the results to hold when we replace a_k by \hat{a}_k , we need $x > (\hat{a}_k/2 + 3\epsilon'/4) \sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}$ on the interval $(d_k - (\epsilon/4)\sqrt{\frac{(n_1+n_2)}{2n_1n_2} \log p}, \infty)$ when d_k is strong and $\hat{a}_k > (2+b/2)\sqrt{\sigma_{kk}} + \sqrt{(2 + b/2)^2\sigma_{kk} + 4}$ when d_i is weak. Here ϵ' is some positive constant less than ϵ and depends solely on ϵ . As a_k satisfies these two conditions, we can choose a constant d which depends solely on ϵ' and b such that \hat{a}_k satisfies these two conditions when $|\hat{a}_k - a_k| < d$.

When the true variance is unknown, we use the pooled sample variance

$$\hat{\sigma}_{kk} = \frac{n_1 - 1}{n_1 + n_2 - 2} \sum_{i=1}^{n_1} (X_{ik} - \bar{X}_k)^2 + \frac{n_2 - 1}{n_1 + n_2 - 2} \sum_{i=1}^{n_2} (Y_{ik} - \bar{Y}_k)^2$$

to estimate σ_{kk} , and use \hat{a}_k to estimate a_k . Since \hat{a}_k is a continuous function of $\hat{\sigma}_{kk}$ and $\hat{a}_k = a_k$ when $\hat{\sigma}_{kk} = \sigma_{kk}$, we can choose a constant λ small enough such that $|\hat{a}_k - a_k| < d$ when $|\hat{\sigma}_{kk} - \sigma_{kk}| < \epsilon_0^{-1}\lambda$. (Recall that ϵ_0^{-1} is the upper bound of σ_{kk} .) The constant λ here depends solely on d and ϵ_0 , which are fixed and independent of n_1 , n_2 and p . Recall that since X_{ik} are *i.i.d* for $1 \leq i \leq n_1$ and Y_{ik} are *i.i.d* for $1 \leq i \leq n_2$, we have $(n_1 + n_2 - 2)\hat{\sigma}_{kk}/\sigma_{kk} \sim \chi_{n_1+n_2-2}^2$. Lemma 1 in Foygel and Drton (2010) and Lemma 4 in Cai (2002) proved the following concentration inequality for χ^2 random variable. For $n \geq 4\lambda^{-2} + 1$, we have

$$\begin{aligned} P\{\chi_n^2 > n(1+\lambda)\} &\leq \frac{1}{\lambda\sqrt{\pi n}} e^{-\frac{n}{2}(\lambda - \log(1+\lambda))}; \\ P\{\chi_n^2 < n(1-\lambda)\} &\leq \frac{1}{\lambda\sqrt{\pi(n-1)}} e^{-\frac{n-1}{2}(\lambda + \log(1-\lambda))}. \end{aligned}$$

Noting that λ is a constant independent of n_1 , n_2 and p , we apply the above results to conclude that

$$\begin{aligned} &P\left(\left|\frac{\hat{\sigma}_{kk}}{\sigma_{kk}} - 1\right| < \lambda\right) \\ &\geq 1 - \frac{1}{\lambda\sqrt{\pi(n_1+n_2-2)}} e^{-\frac{n_1+n_2-2}{2}(\lambda - \log(1+\lambda))} - \frac{1}{\lambda\sqrt{\pi(n_1+n_2-3)}} e^{-\frac{n_1+n_2-3}{2}(\lambda + \log(1-\lambda))} \\ &= 1 - e^{-\frac{n_1+n_2-2}{2}(\lambda - \log(1+\lambda)) + O(\log(n_1+n_2))} - e^{-\frac{n_1+n_2-3}{2}(\lambda + \log(1-\lambda)) + O(\log(n_1+n_2))}. \\ &= 1 - e^{-O(n_1+n_2)} \end{aligned}$$

Since $P\left(\bigcup_{i=1}^p |\hat{\sigma}_{kk} - \sigma_{kk}| > \epsilon_0^{-1}\lambda\right) \leq \sum_{i=1}^p P(|\hat{\sigma}_{kk} - \sigma_{kk}| > \epsilon_0^{-1}\lambda) \leq \sum_{i=1}^p P(|\hat{\sigma}_{kk} - \sigma_{kk}| > \sigma_{kk}\lambda)$, we have with probability greater than $1 - pe^{-O(n_1+n_2)}$, $|\hat{\sigma}_{kk} - \sigma_{kk}| < \epsilon_0^{-1}\lambda$ for all $1 \leq i \leq p$.

Combining with the fact that $\mathbb{E}(q_k)$ is bounded, the theorem follows. \square

S1.4 Proof of Theorem 4

Define $\text{ECC}_{\delta}^2 = \mathbb{E}\left\{\sum_{j=1}^m \mathbb{I}(\delta_j = 2, \theta_j = 2)\right\}$. We will prove $\text{mFSR}_{\delta}^2 = \text{mFSR}_{\delta_{OR}}^2 + o(1)$ and $\text{ECC}_{\delta}^2/\text{ECC}_{\delta_{OR}}^2 = 1 + o(1)$. Then by the same argument one can show $\text{mFSR}_{\delta}^1 = \text{mFSR}_{\delta_{OR}}^1 + o(1)$ and $\text{ECC}_{\delta}^1/\text{ECC}_{\delta_{OR}}^1 = 1 + o(1)$, then the theorem follows.

We begin with a summary of notation used throughout the proof:

- $Q(t) = m^{-1} \left\{ \sum_{j=1}^m (T^j - \alpha_2) \mathbb{I}(T^j < t) \right\}$.
- $\widehat{Q}(t) = m^{-1} \left\{ \sum_{j=1}^m (\widehat{T}^j - \alpha_2) \mathbb{I}(\widehat{T}^j < t) \right\}$.
- $Q_{\infty}(t) = \mathbb{E}\{(T - \alpha_2) \mathbb{I}(T < t)\}$.
- $t_{\infty} = \min\{\sup\{t : Q_{\infty}(t) \leq 0\}, 0.5\}$ is the “ideal” threshold.

Note that T^j is a function of \mathbf{W}_j only. Since each \mathbf{W}_j are iid, it follows that T^j are also iid. Without loss of generality, assume the first s_1 signals are strong, the next s_2 signals are moderate and the rest are weak. We break the proof into 2 cases:

Case 1: $\sum_{k=1}^{s_1} d_k^2 < \infty$ as $n, p \rightarrow \infty$.

Case 2: $\sum_{k=1}^{s_1} d_k^2 \rightarrow \infty$ as $n, p \rightarrow \infty$.

Proof of case 1

We first show that $\text{mFSR}_{\hat{\delta}}^2 = \text{mFSR}_{\delta_{OR}}^2 + o(1)$. We define a continuous version of $\hat{Q}(t)$ using the following procedure: If t_1 and t_2 are two adjacent points of discontinuity, on the interval $[t_1, t_2]$,

$$\hat{Q}_C(t) = \frac{t - t_2}{t_1 - t_2} \hat{Q}(t_1) + \frac{t - t_1}{t_2 - t_1} \hat{Q}(t_2).$$

As t_∞ is proved to be larger than α_1 in section S1.1, it is easy to verify that $\hat{Q}_C(t)$ is continuous and monotone on the interval $[\alpha_1, 1)$. Hence, its inverse \hat{Q}_C^{-1} is well-defined, continuous, and monotone.

Next, we shall show the following two results in turn: (i) $\hat{Q}(t) \xrightarrow{p} Q_\infty(t)$ and (ii) $\hat{Q}_C^{-1}(0) \xrightarrow{p} t_\infty$. We shall need the following two lemmas, which are proved later:

Lemma 1. $\mathbb{E}(\hat{T} - T)^2 \rightarrow 0$.

Lemma 2. Let $V_j = (T^j - \alpha_2)\mathbb{I}(T^j < t)$ and $\hat{V}_j = (\hat{T}^j - \alpha_2)\mathbb{I}(\hat{T}^j < t)$. Then $\mathbb{E}(\hat{V}_j - V_j)^2 = o(1)$.

To show (i), note that $Q(t) \xrightarrow{p} Q_\infty(t)$ by the WLLN, so that we only

need to establish that $\widehat{Q}(t) \xrightarrow{P} Q(t)$. Let $S_m = \sum_{j=1}^m (\widehat{V}_j - V_j)$. By Lemma 2 and the Cauchy-Schwartz inequality, $\mathbb{E} \left\{ (\widehat{V}_i - V_i) (\widehat{V}_j - V_j) \right\} = o(1)$. It follows that

$$\begin{aligned} \text{Var} (m^{-1}S_m) &= m^{-2}\text{Var}(S_m) \leq m^{-2} \sum_{j=1}^m \mathbb{E} \left\{ (\widehat{V}_j - V_j)^2 \right\} \\ &\quad + O \left(\frac{1}{m^2} \sum_{i,j:i \neq j} \mathbb{E} \left\{ (\widehat{V}_i - V_i) (\widehat{V}_j - V_j) \right\} \right) \\ &= o(1). \end{aligned}$$

By Lemma 2, $\mathbb{E}(m^{-1}S_m) \rightarrow 0$, applying Chebyshev's inequality, we obtain $m^{-1}S_m = \widehat{Q}(t) - Q(t) \xrightarrow{P} 0$. Hence (i) is proved.

Next, we show (ii). Since $\widehat{Q}_C(t)$ is continuous, for any $\varepsilon > 0$, we can find $\eta > 0$ such that $\left| \widehat{Q}_C^{-1}(Q_\infty(t_\infty)) - \widehat{Q}_C^{-1} \left\{ \widehat{Q}_C(t_\infty) \right\} \right| < \varepsilon$ if $\left| \widehat{Q}_C(t_\infty) - Q_\infty(t_\infty) \right| < \eta$. It follows that

$$P \left\{ \left| \widehat{Q}_C(t_\infty) - Q_\infty(t_\infty) \right| > \eta \right\} \geq P \left\{ \left| \widehat{Q}_C^{-1}(Q_\infty(t_\infty)) - \widehat{Q}_C^{-1} \left\{ \widehat{Q}_C(t_\infty) \right\} \right| > \varepsilon \right\}.$$

Lemma 1 and the WLLN imply that $\widehat{Q}_C(t) \xrightarrow{P} Q_\infty(t)$. Then,

$$P \left(\left| \widehat{Q}_C(t_\infty) - Q_\infty(t_\infty) \right| > \eta \right) \rightarrow 0.$$

Hence, we have

$$\widehat{Q}_C^{-1}(Q_\infty(t_\infty)) \xrightarrow{p} \widehat{Q}_C^{-1} \left\{ \widehat{Q}_C(t_\infty) \right\} = t_\infty, \quad (\text{S1.1})$$

completing the proof of (ii). Notice that $Q_\infty(t)$ is continuous by construction, by (i) we also have $\widehat{Q}(t) \xrightarrow{p} \widehat{Q}_C(t)$.

We can similarly define the continuous version of $Q(t)$ as $Q_C(t)$ and the corresponding threshold as $Q_C^{-1}(0)$. Write $\widehat{\boldsymbol{\delta}} = \widehat{\boldsymbol{\delta}}^1 + \widehat{\boldsymbol{\delta}}^2$, where $\widehat{\boldsymbol{\delta}}^1$ is of the form $\mathbb{I} \left\{ 1 - \widehat{T}^j \leq \beta_1 \right\}$ and $\widehat{\boldsymbol{\delta}}^2$ is of the form $2\mathbb{I} \left\{ \widehat{T}^j \leq \beta_2 \right\}$ for some $\beta_1, \beta_2 > 0$. Similarly, write $\boldsymbol{\delta} = \boldsymbol{\delta}^1 + \boldsymbol{\delta}^2$, where $\boldsymbol{\delta}^1$ is of the form $\mathbb{I} \left\{ 1 - T^j \leq t_{OR}^1 \right\}$ and $\boldsymbol{\delta}_{OR}^2$ is of the form $2\mathbb{I} \left\{ T^j \leq t_{OR}^2 \right\}$ for some $\beta_1, \beta_2 > 0$. Then by construction, we have

$$\widehat{\boldsymbol{\delta}}^2 = \left[2\mathbb{I} \left\{ \widehat{T}^j \leq \widehat{Q}_C^{-1}(Q_\infty(t_\infty)) \right\} : 1 \leq j \leq m \right]$$

$$\text{and } \boldsymbol{\delta}_{OR}^2 = \left[2\mathbb{I} \left\{ T^j \leq Q_C^{-1}(Q_\infty(t_\infty)) \right\} : 1 \leq i \leq m \right].$$

Also, following the previous arguments, we can show that

$$Q_C^{-1}(Q_\infty(t_\infty)) \xrightarrow{p} t_\infty. \quad (\text{S1.2})$$

According to (S1.1) and (S1.2), we have

$$\widehat{Q}_C^{-1}(Q_\infty(t_\infty)) = Q_C^{-1}(Q_\infty(t_\infty)) + o_p(1). \quad (\text{S1.3})$$

Note that the mFSR² level of δ_{OR} and $\hat{\delta}$ are

$$\text{mFSR}_{\delta_{OR}}^2 = \frac{P_{\theta_j=1} \{T^j \leq Q_C^{-1}(Q_\infty(t_\infty))\}}{P \{T^j \leq Q_C^{-1}(Q_\infty(t_\infty))\}} \quad \text{and} \quad \text{mFSR}_{\hat{\delta}}^2 = \frac{P_{\theta_j=1} \{\hat{T}^j \leq Q_C^{-1}(Q_\infty(t_\infty))\}}{P \{\hat{T}^j \leq Q_C^{-1}(Q_\infty(t_\infty))\}}.$$

From Lemma 1, $\hat{T}^j \xrightarrow{P} T^j$, (S1.2), and by the continuous mapping theorem, $\text{mFSR}_{\hat{\delta}}^2 = \text{mFSR}_{\delta_{OR}}^2 + o(1) \leq \alpha_2 + o(1)$. By the asymptotic equivalence between mFSR and FSR (Supplement S3), the desired result follows.

Next we show that $\text{ECC}_{\hat{\delta}}^2 / \text{ECC}_{\delta_{OR}}^2 = 1 + o(1)$. By definition,

$$\text{ECC}_{\hat{\delta}}^2 = \mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(\hat{T}^j \leq \hat{Q}_C^{-1}(Q_\infty(t_\infty)))(1 - T^j) \right\}$$

$$\text{and} \quad \text{ECC}_{\delta_{OR}}^2 = \mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(T^j \leq Q_C^{-1}(Q_\infty(t_\infty)))(1 - T^j) \right\}.$$

Using the fact that $\hat{T}^j \xrightarrow{P} T^j$ and $\hat{Q}_C^{-1}(Q_\infty(t_\infty)) \xrightarrow{P} Q_C^{-1}(Q_\infty(t_\infty))$, the result follows.

Proof of case 2

From the proof of Lemma 1, we have:

$$\mathbb{E} \left(S_j^\pi - \hat{S}_j \right)^2 = O\left(s_1 \frac{n_1 + n_2}{n_1 n_2}\right) + o\left(\sum_{k=1}^{s_1} d_k^2\right) + O\left(\sum_{k=1}^{s_1} d_k \frac{n_1 + n_2}{n_1 n_2}\right) + o(1). \quad (\text{S1.4})$$

Since $\sum_{k=1}^{s_1} d_k^2 \geq cs_1 \frac{n_1+n_2}{n_1 n_2} \log p$ for some $c > 0$ and $\mathbb{E}(S_j^\pi) \propto \sum_{k=1}^{s_1} d_k^2$, by (S1.4)

$$\mathbb{E} \left(S_j^\pi - \hat{S}_j \right)^2 / \sum_{k=1}^{s_1} d_k^2 = O \left(\frac{1}{\log p} \right).$$

It follows that $\mathbb{E} \sum_{k=1}^{s_1} \{q_k(\bar{X}_k - \bar{Y}_k)\}^2 \propto \mathbb{E}(|S_j^\pi|) \propto \sum_{k=1}^{s_1} d_k^2$. If $\sum_{k=1}^{s_k} d_i^2 \rightarrow \infty$, then $\hat{T}^j \rightarrow 0$ or $\hat{T}^j \rightarrow \infty$. In either case, we have perfect separation between the two classes asymptotically. The theorem follows. \square

S1.5 Proof of Theorem 3

Again, without loss of generality, assume the first s_1 signals are strong, the next s_2 signals are moderate and the rest are weak. Consider the same model as described in section 2, but replace $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ by $\tilde{\boldsymbol{\mu}}_1 = (\mu_{11}, \dots, \mu_{1s_1}, 0, \dots, 0)^t$ and $\tilde{\boldsymbol{\mu}}_2 = (\mu_{21}, \dots, \mu_{2s_1}, 0, \dots, 0)^t$ respectively. Let $\tilde{\mathbf{W}}_j = \mathbf{W}_j - \boldsymbol{\mu}_1 + \tilde{\boldsymbol{\mu}}_1$ when $\theta_j = 1$ and $\tilde{\mathbf{W}}_j = \mathbf{W}_j - \boldsymbol{\mu}_2 + \tilde{\boldsymbol{\mu}}_2$ when $\theta_j = 2$. It is then clear that

$$\tilde{\mathbf{W}}_j \sim \mathbb{I}(\theta_j = 1)\mathcal{N}(\tilde{\boldsymbol{\mu}}_1, \Sigma) + \mathbb{I}(\theta_j = 2)\mathcal{N}(\tilde{\boldsymbol{\mu}}_2, \Sigma).$$

Denote $\tilde{\mathbf{d}} = \tilde{\boldsymbol{\mu}}_1 - \tilde{\boldsymbol{\mu}}_2$, define

$$Z_j := Z(\tilde{\mathbf{W}}_j) := \left(\tilde{\mathbf{W}}_j - \frac{\tilde{\boldsymbol{\mu}}_1 + \tilde{\boldsymbol{\mu}}_2}{2} \right)^\top \Sigma^{-1} \tilde{\mathbf{d}},$$

$Z_{OR}^j := Z_{OR}(\tilde{\mathbf{W}}_j) = \exp(Z_j)/(\exp(Z_j) + 1)$. Consider the decision rule:

$$\delta_Z^j = 2\mathbb{I}\{Z_{OR}^j \leq \min(Z_{OR}^{(k_4)}, 0.5)\} + \mathbb{I}\{1 - Z_{OR}^j \leq \min(1 - Z_{OR}^{(m-k_3)}, 0.5)\},$$

where

$$k_3 = \inf \left\{ j : \frac{1}{j+1} \sum_{i=0}^j (1 - Z_{OR}^{(m-i)}) \leq \alpha_1 \right\} \quad \text{and} \quad k_4 = \sup \left\{ j : \frac{1}{j} \sum_{i=1}^j Z_{OR}^{(i)} \leq \alpha_2 \right\}.$$

By the proof of theorem 1 and theorem 3, $\delta_Z = (\delta_Z^1, \dots, \delta_Z^m)$ controls mFSR¹

and mFSR² at level α_1 and α_2 respectively. Now notice that

$$\begin{aligned} & \mathbb{E} \left(Z_j - \hat{S}_j \right)^2 \\ &= \mathbb{E} \left\{ \left(\tilde{\mathbf{W}}_j - \frac{\tilde{\boldsymbol{\mu}}_1 + \tilde{\boldsymbol{\mu}}_2}{2} \right)^\top \Sigma^{-1} \tilde{\mathbf{d}} - \left(\mathbf{W}_j - \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} \right)^\top \hat{\Sigma}^{-1} \hat{\mathbf{d}} \right\}^2 \\ &= \mathbb{E} \left[\mathbb{E} \left\{ \left(\mathbf{W}_i - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} (\tilde{\mathbf{d}} - \hat{\mathbf{d}}) \right\}^2 + \mathbb{E} \left\{ \left(\tilde{\mathbf{W}}_j - \mathbf{W}_j + \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} \hat{\mathbf{d}} \right\}^2 \right. \\ & \quad \left. + \mathbb{E} \left\{ \left(\mathbf{W}_i - \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} \right)^\top (\Sigma - \hat{\Sigma}^{-1}) \hat{\mathbf{d}} \right\}^2 \right] \\ &= I + II + III \end{aligned}$$

For II, note that

$$\begin{aligned}
 & \mathbb{E} \left\{ \left(\tilde{\mathbf{W}}_j - \mathbf{W}_j + \frac{\tilde{\boldsymbol{\mu}}_1 + \tilde{\boldsymbol{\mu}}_2}{2} - \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} \right)^\top \Sigma^{-1} \hat{\mathbf{d}} \right\}^2 \\
 &= O \left[\mathbb{E} \left\{ \hat{\mathbf{d}}^\top (\Sigma^{-1})^t \frac{1}{4} \frac{n_1 + n_2}{n_1 n_2} \Sigma \Sigma^{-1} \hat{\mathbf{d}} \right\} + \{ \lambda_{\max}(\Sigma^{-1}) p^{-\gamma} \log p \}^2 \right] \\
 &= O \left[\mathbb{E} \left\{ \hat{\mathbf{d}}^\top (\Sigma^{-1})^\top \frac{1}{4} \frac{n_1 + n_2}{n_1 n_2} \hat{\mathbf{d}} \right\} + \{ \lambda_{\max}(\Sigma^{-1}) p^{-\gamma} \log p \}^2 \right] \\
 &= O \left[\mathbb{E} \left(\frac{n_1 + n_2}{n_1 n_2} \lambda_{\max}(\Sigma^{-1}) \|\mathbf{d} - \hat{\mathbf{d}}\|^2 \right) + \mathbb{E} \left(\frac{n_1 + n_2}{n_1 n_2} \lambda_{\max}(\Sigma^{-1}) \|\mathbf{d}\|^2 \right) + \{ \lambda_{\max}(\Sigma^{-1}) p^{-\gamma} \log p \}^2 \right] \\
 &= O \left[\mathbb{E} \frac{n_1 + n_2}{n_1 n_2} \lambda_{\max}(\Sigma^{-1}) (\|\mathbf{d} - \hat{\mathbf{d}}\|^2 + \|\mathbf{d}\|^2) + \{ \lambda_{\max}(\Sigma^{-1}) p^{-\gamma} \log p \}^2 \right].
 \end{aligned}$$

By (A1) $\{ \lambda_{\max}(\Sigma^{-1}) p^{-\gamma} \log p \}^2 \rightarrow 0$, use the same computation in the proof

of lemma 1, we have

$$\mathbb{E} \left(Z_j - \hat{S}_j \right)^2 = O \left(\mathbb{E} \sum_{k=1}^p \{ q_k (\bar{X}_k - \bar{Y}_k) - \tilde{d}_k \}^2 \right) + o(1),$$

$$\mathbb{E} \sum_{k=1}^{s_1} \{ q_k (\bar{X}_k - \bar{Y}_k) - \tilde{d}_k \}^2 = O(s_1 \frac{n_1 + n_2}{n_1 n_2}) + o(\sum_{k=1}^{s_1} d_k^2) + O(\sum_{k=1}^{s_1} d_k \frac{n_1 + n_2}{n_1 n_2}) + o(1).$$

We know that

$$\mathbb{E} \sum_{i=s_1+1}^{s_1+s_2} \{ q_k (\bar{X}_k - \bar{Y}_k) \}^2 = O(s_2) O(p^{-\gamma}) O \left(\frac{n_1 + n_2}{n_1 n_2} \log p \right) = O(p^{-\gamma} \log p) \rightarrow 0,$$

$$\mathbb{E} \sum_{k=s_2+1}^p \{ q_k (\bar{X}_k - \bar{Y}_k) \}^2 = O(p) O(p^{-1-\gamma}) o \left(\frac{n_1 + n_2}{n_1 n_2} \log p \right) \rightarrow 0.$$

Thus,

$$\mathbb{E} \left(Z_j - \hat{S}_j \right)^2 = O\left(s_1 \frac{n_1 + n_2}{n_1 n_2}\right) + o\left(\sum_{k=1}^{s_1} d_k^2\right) + O\left(\sum_{k=1}^{s_1} d_k \frac{n_1 + n_2}{n_1 n_2}\right) + o(1). \quad (\text{S1.5})$$

Since $\sum_{k=1}^{s_1} d_k^2 \geq c s_1 \log p \frac{n_1 + n_2}{n_1 n_2}$ for some $c > 0$ and $\mathbb{E}(Z_j) \propto \sum_{k=1}^{s_1} d_k^2$, by

(S1.5)

$$\mathbb{E} \left(Z_j - \hat{S}_j \right)^2 \bigg/ \sum_{k=1}^{s_1} d_k^2 = O\left(\frac{1}{\log p}\right)$$

It follows that $\mathbb{E} \sum_{k=1}^{s_1} \{q_k(\bar{X}_k - \bar{Y}_k)\}^2 \propto \mathbb{E}(|Z_j|) \propto \sum_{k=1}^{s_1} d_k^2$. If $\sum_{k=1}^{s_1} d_k^2 \rightarrow$

∞ , then $\hat{T}^j \rightarrow 0$ or $\hat{T}^j \rightarrow 1$. In either case, we have perfect separation

between the two classes asymptotically. The theorem follows trivially. If

$\sum_{k=1}^{s_1} d_k^2 < \infty$, then

$$\left(\mathbb{E} \left| \hat{T}^j - Z_{OR}^j \right| \right)^2 = e^{O(\sum_{k=1}^{s_1} d_k^2)} \left\{ O\left(s_1 \frac{n_1 + n_2}{n_1 n_2}\right) + O\left(\sum_{k=1}^{s_1} d_k^2 p^{-2\gamma}\right) + O(p^{-\gamma} \log p) \right\}.$$

Since $\sum_{k=1}^{s_1} d_k^2 \geq c s_1 \log p \frac{n_1 + n_2}{n_1 n_2}$ and $\sum_{k=1}^{s_1} d_k^2 < \infty$, it follows that $e^{O(\sum_{k=1}^{s_1} d_k^2)} <$

∞

and $O\left(s_1 \frac{n_1 + n_2}{n_1 n_2}\right) + O(\sum_{k=1}^{s_1} d_k^2 p^{-2\gamma}) + O(p^{-\gamma} \log p) \rightarrow 0$. Thus, $\mathbb{E} \left(\hat{T}^j - Z_{OR}^j \right)^2 \rightarrow$

0. The theorem follows by using the same argument presented in the proof

of theorem 4. □

S1.6 Proof of Corollary 1

Note that $\min\{T^j, 1 - T^j\} \leq 0.5$. Hence, if we choose $\alpha_1 = \alpha_2 = 0.5$, $\boldsymbol{\delta}_{OR}$ makes no indecision. Since $\boldsymbol{\delta}_{OR}$ has the highest ECC among all rules that controls mFSR^1 and mFSR^2 at level 0.5, and $\boldsymbol{\delta}^F$ also controls mFSR^1 and mFSR^2 at level 0.5, it follows that $\text{ECC}_{\boldsymbol{\delta}_{OR}} \geq \text{ECC}_{\boldsymbol{\delta}^F}$. On the other hand, $\boldsymbol{\delta}^F$ has the lowest risk among all rules that do not make indecision and risk equals to $1 - \text{ECC}/m$. It follows that $\text{ECC}_{\boldsymbol{\delta}_{OR}} \leq \text{ECC}_{\boldsymbol{\delta}^F}$. Hence, $\text{ECC}_{\boldsymbol{\delta}_{OR}} = \text{ECC}_{\boldsymbol{\delta}^F}$. By Theorem 4, we have $\text{ECC}_{\boldsymbol{\delta}_{OR}}/\text{ECC}_{\hat{\boldsymbol{\delta}}} \rightarrow 1$, it follows that $R(\hat{\boldsymbol{\delta}}) \rightarrow R(\boldsymbol{\delta}^F)$. \square

S2 Proof of Technical Lemmas

Supplement S2 contains proofs of technical lemmas. We will assume without loss of generality, the first s_1 signals are strong, the next s_2 signals are moderate and the rest are weak.

S2.1 Proof of Lemma 1

Let $\hat{\mathbf{d}} = (q_1(\bar{X}_1 - \bar{Y}_1), \dots, q_p(\bar{X}_p - \bar{Y}_p))$. We have

$$\begin{aligned} \mathbb{E}\|\hat{\mathbf{d}} - \mathbf{d}\|^2 &= \sum_{k=1}^p \left\{ \mathbb{E}(q_k^2(\bar{X}_k - \bar{Y}_k)^2) - 2d_k \mathbb{E}\{q_k(\bar{X}_k - \bar{Y}_k)\} + d_k^2 \right\} \\ &= \sum_{k=1}^p \left\{ \text{Cov}\{q_k^2, (\bar{X}_k - \bar{Y}_k)^2\} + \mathbb{E}(q_k^2)\mathbb{E}(\bar{X}_k - \bar{Y}_k)^2 \right\} \\ &\quad - \sum_{k=1}^p \left\{ 2d_k [\text{Cov}\{q_k, (\bar{X}_k - \bar{Y}_k)\} + d_k \mathbb{E}q_k] + d_k^2 \right\}. \end{aligned}$$

Since q_k is bounded between 0 and 1, $\text{Var}(q_k^2)$ and $\text{Var}(q_i)$ are both bounded.

For $1 \leq k \leq s_1$, note that $(\bar{X}_k - \bar{Y}_k) \sim N(d_k, \sigma_{kk} \frac{n_1+n_2}{n_1 n_2})$, and $d_k^2 < \infty$ we

have $\text{Var}\{(\bar{X}_k - \bar{Y}_k)^2\} = O(\frac{n_1+n_2}{n_1 n_2})$. Use Cauchy-Schwarz inequality, for

$1 \leq k \leq s_1$, each summand has absolute value bounded by

$$\begin{aligned} &O\left(\frac{n_1 + n_2}{n_1 n_2}\right) + \mathbb{E} \left\{ q_k^2 \left(\frac{n_1 + n_2}{n_1 n_2} \sigma_{kk} + d_k^2 \right) \right\} - 2d_k^2 \mathbb{E}(q_k) + d_k^2 - 2d_k \text{Cov}(q_k, \bar{X}_k - \bar{Y}_k) \\ &= O\left(\frac{n_1 + n_2}{n_1 n_2}\right) + d_k^2 \mathbb{E}(q_k^2) + d_k^2 - 2d_k^2 \mathbb{E}(q_k) + O\left(d_k \frac{n_1 + n_2}{n_1 n_2}\right) + O\left(\frac{n_1 + n_2}{n_1 n_2}\right) \\ &= O\left(\frac{n_1 + n_2}{n_1 n_2}\right) + O\left(d_k \frac{n_1 + n_2}{n_1 n_2}\right) + O\{d_k^2(1 - \mathbb{E}(q_k))^2\} + O\{d_k^2(E(q_k^2) - (E q_k)^2)\} \\ &= O\left(\frac{n_1 + n_2}{n_1 n_2}\right) + O\left(d_k \frac{n_1 + n_2}{n_1 n_2}\right) + o(d_k^2). \end{aligned}$$

For $s_1 + 1 \leq i \leq p$,

$$\mathbb{E}\{q_k(\bar{X}_k - \bar{Y}_k) - d_k\}^2 \leq d_k^2.$$

Hence, by (A4)

$$\mathbb{E}\|\hat{\mathbf{d}} - \mathbf{d}\|^2 = O\left(s_1 \frac{n_1 + n_2}{n_1 n_2}\right) + o\left(\sum_{k=1}^{s_1} d_k^2\right) + O\left(\sum_{k=1}^{s_1} d_k \frac{n_1 + n_2}{n_1 n_2}\right) + o(1). \quad (\text{S2.6})$$

Since $s_1 \frac{n_1 + n_2}{n_1 n_2} = O(\sum_{k=1}^{s_1} d_k^2)$ and by assumption of case 1, $\sum_{k=1}^{s_1} d_k^2 = O(1)$ and $\sum_{k=1}^{s_1} d_k \frac{n_1 + n_2}{n_1 n_2} = o(1)$ we have $\mathbb{E}\|\hat{\mathbf{d}} - \mathbf{d}\|^2 = o(1)$. Now, since $\bar{\mathbf{X}} + \bar{\mathbf{Y}}$ and $\bar{\mathbf{X}} - \bar{\mathbf{Y}}$ are independent, and using Basu's Theorem we know that $\bar{\mathbf{X}} + \bar{\mathbf{Y}}$ and $\hat{\sigma}_{kk}$ are independent, thus $\bar{\mathbf{X}} + \bar{\mathbf{Y}}$ and $\hat{\mathbf{d}}$ are independent and we have

$$\begin{aligned} & \mathbb{E}(S_j^\pi - \hat{S}_j)^2 \\ &= \mathbb{E} \left\{ \left(\mathbf{W}_j - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} \mathbf{d} - \left(\mathbf{W}_j - \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} \right)^\top \hat{\Sigma}^{-1} \hat{\mathbf{d}} \right\}^2 \\ &= \mathbb{E} \left\{ \left(\mathbf{W}_j - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} (\mathbf{d} - \hat{\mathbf{d}}) + \left(\frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} \hat{\mathbf{d}} \right. \\ & \quad \left. + \left(\mathbf{W}_j - \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} \right)^\top (\Sigma^{-1} - \hat{\Sigma}^{-1}) \hat{\mathbf{d}} \right\}^2 \\ &= O \left[\mathbb{E} \left\{ \left(\mathbf{W}_j - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} (\mathbf{d} - \hat{\mathbf{d}}) \right\}^2 + \mathbb{E} \left\{ \left(\frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} \hat{\mathbf{d}} \right\}^2 \right. \\ & \quad \left. + \mathbb{E} \left\{ \left(\mathbf{W}_j - \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} \right)^\top (\Sigma - \hat{\Sigma}^{-1}) \hat{\mathbf{d}} \right\}^2 \right] \\ &= I + II + III. \end{aligned}$$

For term I , as in case 1 we have $\|\mathbf{d}\|^2$ bounded:

$$\mathbb{E} \left\{ \left(\mathbf{W}_j - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} (\mathbf{d} - \hat{\mathbf{d}}) \right\}^2 = O \left\{ \mathbb{E} \left(\lambda_{max}(\Sigma^{-1})^2 \|\mathbf{d} - \hat{\mathbf{d}}\|^2 \right) \right\},$$

where $\lambda_{max}(\Sigma^{-1})$ is the largest eigenvalue of Σ^{-1} , which in our case is bounded by some constant.

For term II , we have

$$\begin{aligned} & \mathbb{E} \left\{ \left(\frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)^\top \Sigma^{-1} \hat{\mathbf{d}} \right\}^2 \\ &= O \left\{ \mathbb{E} \left(\frac{n_1 + n_2}{n_1 n_2} \lambda_{max}(\Sigma^{-1}) \|\mathbf{d} - \hat{\mathbf{d}}\|^2 \right) + \mathbb{E} \left(\frac{n_1 + n_2}{n_1 n_2} \lambda_{max}(\Sigma^{-1}) \|\mathbf{d}\|^2 \right) \right\} \\ &= O \left\{ \frac{n_1 + n_2}{n_1 n_2} \mathbb{E} \left(\lambda_{max}(\Sigma^{-1}) (\|\mathbf{d} - \hat{\mathbf{d}}\|^2 + \|\mathbf{d}\|^2) \right) \right\}. \end{aligned}$$

For term *III*, we have

$$\begin{aligned}
 & \mathbb{E} \left\{ \left(\mathbf{W}_j - \frac{\bar{\mathbf{X}} + \bar{\mathbf{Y}}}{2} \right)^\top (\Sigma^{-1} - \hat{\Sigma}^{-1}) \hat{\mathbf{d}} \right\}^2 \\
 &= \mathbb{E} \left[\hat{\mathbf{d}}^\top (\Sigma - \hat{\Sigma}^{-1})^\top \left\{ \left(1 + \frac{1}{4} \frac{n_1 + n_2}{n_1 n_2} \right) \Sigma^{-1} + \frac{1}{4} \mathbf{d} \mathbf{d}^\top \right\} (\Sigma^{-1} - \hat{\Sigma}^{-1}) \hat{\mathbf{d}} \right] \\
 &\leq \mathbb{E} \left[\left\{ \left(1 + \frac{1}{4} \frac{n_1 + n_2}{n_1 n_2} \right) \frac{1}{\epsilon_0} + \frac{1}{4} \max\{d_i d_j, 1 \leq i, j \leq p\} \right\} \hat{\mathbf{d}}^\top (\Sigma^{-1} - \hat{\Sigma}^{-1})^\top (\Sigma^{-1} - \hat{\Sigma}^{-1}) \hat{\mathbf{d}} \right] \\
 &= O \left[\mathbb{E} \left\{ \|(\Sigma^{-1} - \hat{\Sigma}^{-1})\|_2^2 \|\hat{\mathbf{d}}\|^2 \right\} \right] \\
 &= O \left[\mathbb{E} \left\{ \|(\Sigma^{-1} - \hat{\Sigma}^{-1})\|_2^2 \|\mathbf{d}\|^2 \right\} \right] \\
 &= o(1).
 \end{aligned}$$

The last equality follows from assumption (A2) and the fact that $\|\mathbf{d}\|^2 < \infty$

in case 1. Hence,

$$\mathbb{E} \left(S_j^\pi - \hat{S}_j \right)^2 = O \left\{ \mathbb{E} \left(\|\hat{\mathbf{d}} - \mathbf{d}\|^2 \right) \right\} + O \left(\frac{n_1 + n_2}{n_1 n_2} \right) + o(1) = o(1).$$

Next we ask under what conditions do we have $\mathbb{E} \left| \hat{T}^j - T^j \right|^2 \rightarrow 0$. Let $\delta > 0$

be some constant, applying Chebyshev's inequality, we have $P \left(|\hat{S}_j - S_j^\pi| > \delta \right) =$

$O \left\{ \mathbb{E} \left(|S_j^\pi - \hat{S}_j| \right)^2 \right\} \rightarrow 0$. When $|\hat{S}_j - S_j^\pi| \leq \delta$, apply Cauchy-Schwartz in-

equality, we have:

$$\begin{aligned}
& \mathbb{E}_{|S_j^\pi - \hat{S}_j| < \delta} \left| \exp(\hat{S}_j) - \exp(S_j^\pi) \right|^2 \\
& \leq \mathbb{E} e^{2S_j^\pi + 2\delta} \mathbb{E} \left(S_j^\pi - \hat{S}_j \right)^2 \\
& = e^{O(\sum_{k=1}^{s_1} d_k^2)} \cdot o(1).
\end{aligned}$$

Under the assumption of case 1 the above goes to 0. Therefore, as \hat{T}^j, T^j are bounded above by 1, we have:

$$\begin{aligned}
\mathbb{E}(\hat{T}^j - T^j)^2 & \leq \mathbb{E} \left| \hat{T}^j - T^j \right| \\
& \leq 2P \left(\left| \hat{S}_j - S_j^\pi \right| > \delta \right) + E_{|S_j - S_j^\pi| < \delta} \left| \exp(\hat{S}_j) - \exp(S_j^\pi) \right|.
\end{aligned}$$

The lemma follows. □

S2.2 Proof of Lemma 2

Using the definitions of \hat{V}_j and V_j , we can show that

$$\begin{aligned}
\frac{1}{2} \left(\hat{V}_j - V_j \right)^2 & \leq \left(\hat{T}^j - T^j \right)^2 \mathbb{I} \left(\hat{T}^j \leq t, T^j \leq t \right) + \left(\hat{T}^j - \alpha_2 \right)^2 \mathbb{I} \left(\hat{T}^j \leq t, T^j > t \right) \\
& \quad + \left(T^j - \alpha_2 \right)^2 \mathbb{I} \left(\hat{T}^j > t, T^j \leq t \right).
\end{aligned}$$

Let us refer to the three summands on the right hand as (1), (2) and (3) respectively. By Lemma 1, (1) = $o(1)$. Then let $\varepsilon > 0$, and consider that

$$\begin{aligned} & P\left(\hat{T}^j \leq t, T^j > t\right) \\ & \leq P\left(\hat{T}^j \leq t, T^j \in (t, t + \varepsilon)\right) + P\left(\hat{T}^j \leq t, T^j \geq t + \varepsilon\right) \\ & \leq P\left\{T^j \in (t, t + \varepsilon)\right\} + P\left(\left|T^j - \hat{T}^j\right| > \varepsilon\right). \end{aligned}$$

The first term on the right hand is vanishingly small as $\varepsilon \rightarrow 0$ because T^j is a continuous random variable. The second term converges to 0 by Lemma 1. Noting that $0 \leq T^j \leq 1$, we conclude (2) = $o(1)$. In a similar fashion, we can show that (3) = $o(1)$, thus proving the lemma. \square

S3 Asymptotic Equivalence of FSR and mFSR

We show FSR_{δ}^c and mFSR_{δ}^c are asymptotically equivalent. Let $\mathcal{X}_{\delta}^c = \frac{1}{m} \sum_{j=1}^m \mathbb{I}(\delta_j = c, \theta_j \neq c)$ and $\mathcal{Y}_{\delta}^c = \frac{1}{m} \sum_{j=1}^m \mathbb{I}(\delta_j = c)$. The goal is to show $|\text{FSR}_{\delta}^c - \text{mFSR}_{\delta}^c| = o(1)$.

$$|\text{FSR}_{\delta}^c - \text{mFSR}_{\delta}^c| \leq \mathbb{E} \left\{ \left| \frac{\mathcal{X}_{\delta}^c}{\mathcal{Y}_{\delta}^c} - \frac{\mathcal{X}_{\delta}^c}{\mathbb{E}\mathcal{Y}_{\delta}^c} \right| \mathbb{I}(\mathcal{Y}_{\delta}^c > 0) \right\} = \mathbb{E} \left\{ \frac{\mathcal{X}_{\delta}^c}{\mathcal{Y}_{\delta}^c} \mathbb{I}(\mathcal{Y}_{\delta}^c > 0) \frac{|\mathcal{Y}_{\delta}^c - \mathbb{E}\mathcal{Y}_{\delta}^c|}{\mathbb{E}\mathcal{Y}_{\delta}^c} \right\}$$

Since $\mathcal{X}_\delta^c \leq \mathcal{Y}_\delta^c$, we have

$$\mathbb{E} \left\{ \frac{\mathcal{X}_\delta^c}{\mathcal{Y}_\delta^c} \mathbb{I}(\mathcal{Y}_\delta^c > 0) \frac{|\mathcal{Y}_\delta^c - \mathbb{E}\mathcal{Y}_\delta^c|}{\mathbb{E}\mathcal{Y}_\delta^c} \right\} \leq \mathbb{E} \left\{ \frac{|\mathcal{Y}_\delta^c - \mathbb{E}\mathcal{Y}_\delta^c|}{\mathbb{E}\mathcal{Y}_\delta^c} \right\} \leq \frac{(\mathbb{E}|\mathcal{Y}_\delta^c - \mathbb{E}\mathcal{Y}_\delta^c|^2)^{1/2}}{\mathbb{E}\mathcal{Y}_\delta^c} = \frac{(\text{Var}\mathcal{Y}_\delta^c)^{1/2}}{\mathbb{E}\mathcal{Y}_\delta^c}$$

By assumption (A3), we have $m\mathcal{Y}_{\delta_{OR}}^c \sim \text{Binom}(m, \eta)$ with $\eta > 0$. Therefore, $\mathbb{E}\mathcal{Y}_{\delta_{OR}}^c = \eta$ and $\text{Var}(\mathcal{Y}_{\delta_{OR}}^c)^{1/2} = \sqrt{\eta(1-\eta)/m}$, $|\text{FSR}_{\delta_{OR}}^c - \text{mFSR}_{\delta_{OR}}^c| = O(m^{-1/2})$. In the setting of Theorem 4, the data driven procedure is mimicking δ_{OR} . It can be seen from the proof of Theorem 4 that mFSR and FSR of the data driven procedure are also asymptotically equivalent. Similarly, for the setting of Theorem 3, consider the rule δ_Z defined in the proof of Theorem 3, by assumption (A3) we also have $|\text{FSR}_{\delta_Z}^c - \text{mFSR}_{\delta_Z}^c| = O(m^{-1/2})$. In the setting of Theorem 3, the data driven procedure is mimicking δ_Z . It can be seen from the proof of Theorem 3 that mFSR and FSR of the data driven procedure are also asymptotically equivalent. \square

S4 Adaptation to unknown sparsity: illustrations

We present numerical examples to provide insights on why the shrinkage rule (3.1) works well across the sparse and dense regimes.

We start with the sparse case. Suppose $n_1 = n_2 = n$. Let $n = 5$, $p = 625$, $\mathbf{d} = (d_1, \dots, d_p)$, where $d_k = 2.5$ for $k \in \{1, \dots, 50\}$ and $d_k = 0$ for

$k \in \{51, \dots, 625\}$. The observations are generated as $\mathbf{X}_i \sim \mathcal{N}(\mathbf{0}, \frac{1}{2}I_p)$ and $\mathbf{Y}_i \sim \mathcal{N}(\mathbf{d}, \frac{1}{2}I_p)$, $i = 1, \dots, n$, where $\mathbf{0}$ is a p -dimensional vector of zeros and I_p a p -dimensional identity matrix. It follows that $\bar{X}_k - \bar{Y}_k \sim \mathcal{N}(d_k, 1/n)$. We contrast the proposed shrinkage rule with hard/soft thresholding rules in Figure 1. Panels (a) and (b) plot d_k and observed $\bar{X}_k - \bar{Y}_k$, respectively, with coordinates corresponding to zero/nonzero d_k being marked in blue/red. Panel (c) plots the shrinkage estimates $(\bar{X}_k - \bar{Y}_k)q_k$, where q_k is defined in (3.2) with $b = 0$. Panels (d)–(f) show the hard/soft thresholding estimates with different thresholds λ , where $\rho_h(t, \lambda, \sigma) := t\mathbb{I}(|t| > \lambda\sqrt{\sigma})$ and $\rho_s(t, \lambda, \sigma) := \text{sgn}(t) \cdot \max(|t| - \lambda\sqrt{\sigma}, 0)$ are the hard and soft thresholding functions, respectively.

We can see from Panel (c) that almost all blue points are pulled towards and centered around the line of zero by our proposed shrinkage rule. The multiplicative factor q_k is as effective as existing thresholding methods for noise reduction in the sense that the patterns in Panel (c) is qualitatively similar to those in Panels (d) to (f), except that in (d) and (e) too many blue points have survived whereas in (f) too many nonzero signals have been killed. Panel (c) shrinks most noisy entries to zero while being capable of preserving a significant portion of nonzero signals. The bottom row compares the shrinkage functions ρ_h and ρ_s with q_k by setting $\hat{\sigma}_{kk} = 0.5$

S4. ADAPTATION TO UNKNOWN SPARSITY: ILLUSTRATIONS

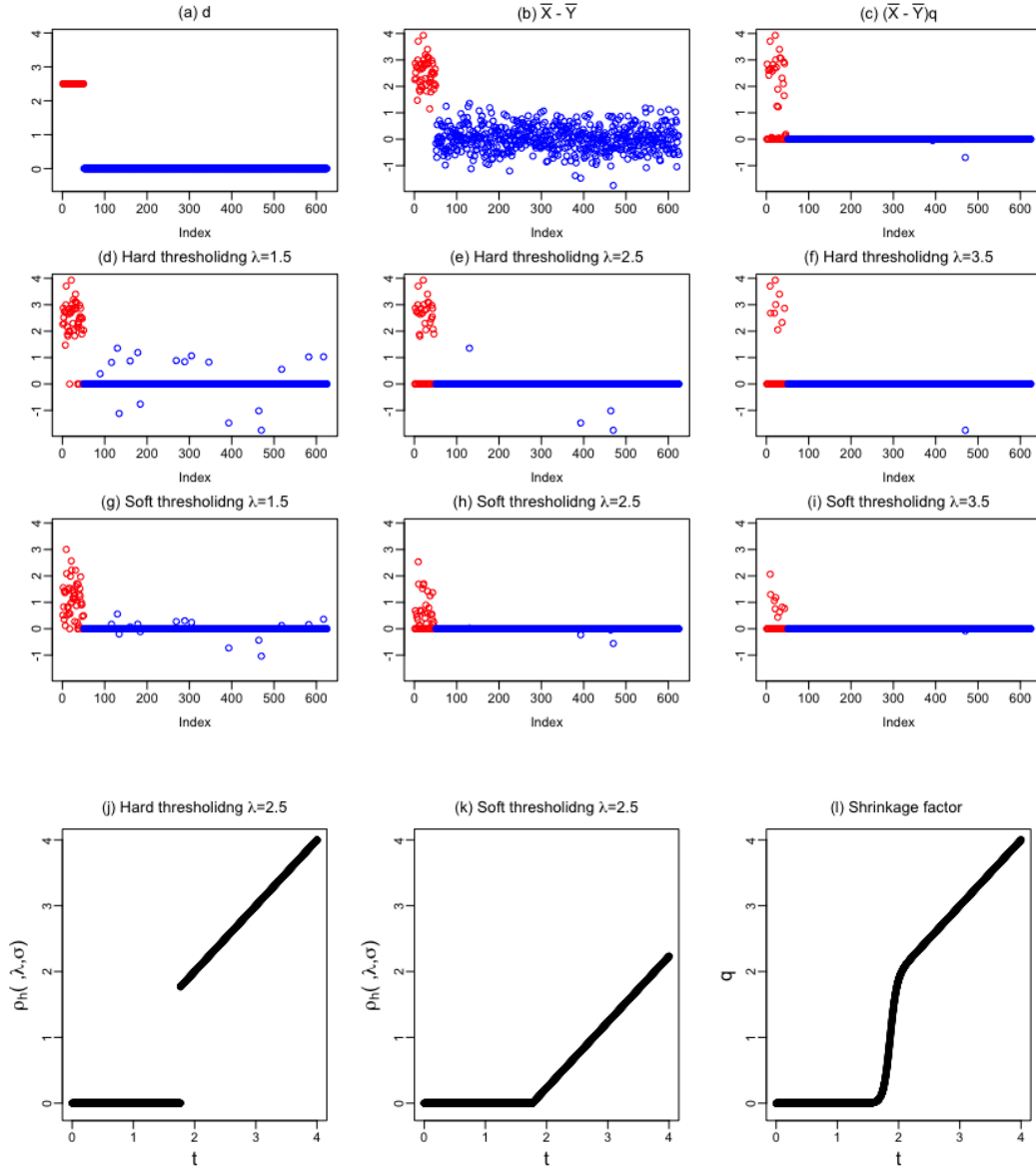


Figure 1: A comparison of our shrinkage rule (3.2) with hard/soft thresholding rules in the sparse case. Panels (a), (b) and (c) respectively plot the true d_k , observed $\bar{X}_k - \bar{Y}_k$ and $(\bar{X}_k - \bar{Y}_k)q_k$. Panels (d)-(f) [(g)-(i)] present results of hard-thresholding (soft-thresholding) rules. The effects of different shrinkage functions are provided at the bottom row.

and $\lambda = 2.5$. We can see that our shrinkage function is continuous, nearly unbiased for large signals and yields similar effects as that of the hard-thresholding function. The proposed shrinkage rule is desirably tuning-free in the sense that one can simply set $b = 0$ in practice; this merit is justified in our theoretical analysis and corroborated by our numerical results. By contrast, the performance of existing thresholding rules depends critically on the value of λ , which is nontrivial to choose.

Next we turn to the dense case. The data are generated in a similar way as before except that we let $d_k = 2.5(k - 1)/624$, $k = 1, \dots, 625$, i.e. d_k increases linearly from 0 to 2.5. We mark the coordinates of \mathbf{d} in three colors: red if $d_k > (\sqrt{0.5} + \sqrt{0.5 + 1})\sqrt{\log p/n} \approx 2.19$, black if $1.5 < d_k < (\sqrt{0.5} + \sqrt{0.5 + 1})\sqrt{\log p/n}$ and blue otherwise¹. We plot our shrinkage estimate and thresholding estimates in Figure 2. In this high-dimensional “dense” regime, the boundaries between weak, moderate and strong signals are blurred. Therefore the working assumptions (sparsity and dichotomy of \mathbf{d}) underpinning the use of thresholding rules can be problematic. In contrast with the (roughly) linear patterns of the surviving points in Panels (d) to (i), the curved pattern in Panel (c) provides differential amount of shrinkage for weak and strong signals, achieving a more desirable balance

¹There are no clearcut boundaries; the colors in this example are only chosen for illustration purpose.

S4. ADAPTATION TO UNKNOWN SPARSITY: ILLUSTRATIONS

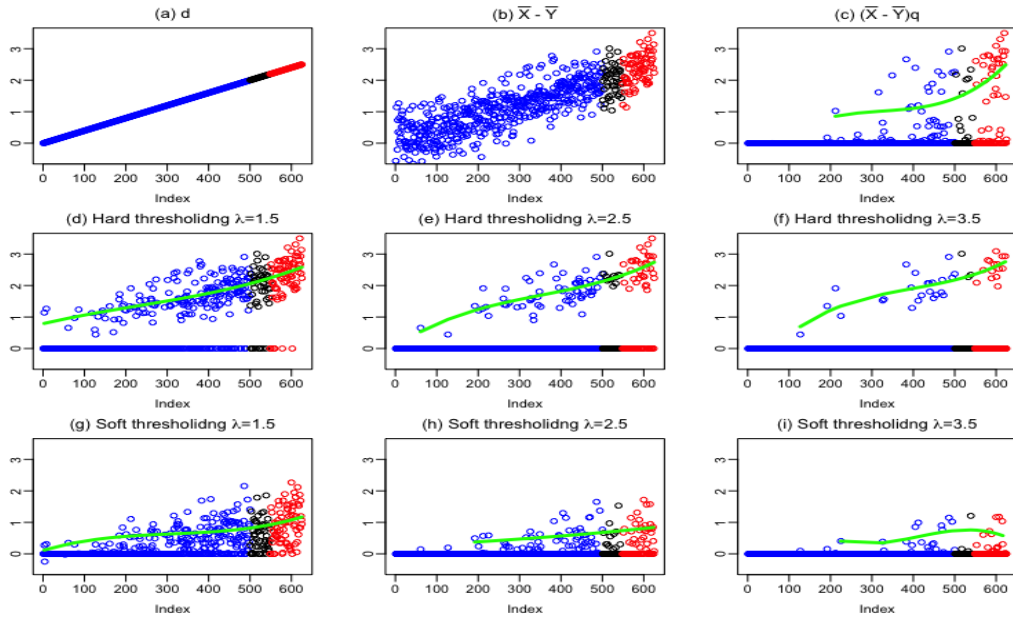


Figure 2: Comparison of the proposed shrinkage rule with thresholding rules. The proposed shrinkage rule exhibits a curved pattern [Panel (c)], which is more effective in eliminating weak signals and keeping strong signals. Green lines are cubic polynomial fits of the points above the horizontal line of 0.2.

between reducing the uncertainties and maintaining useful signals.

S5 Additional Simulations

Let $\boldsymbol{\mu}_1 = (0, \dots, 0)^\top \in \mathbb{R}^p$, and $\boldsymbol{\mu}_2$ be a vector with the first $p^{1/2}$ entries being 0.5 and the rest being 0. Consider the following three correlation structures (same as in the main text).

Model 1: Band graph. Let $\Sigma^{-1} = \Omega = (\omega_{ij})_{p \times p}$, where $\omega_{ii} = 1$, $\omega_{i,i+1} = \omega_{i+1,i} = 0.35$, $\omega_{i,i+2} = \omega_{i+2,i} = 0.175$, and $\omega_{ij} = 0$ if $|i - j| > 2$.

Model 2: AR(1) structure. Let $\Sigma^{-1} = \Omega = (\omega_{ij})_{p \times p}$, where $\omega_{ij} = 0.3^{|i-j|}$.

Model 3: Block structure. Let $\Sigma^{-1} = \Omega = (\mathbf{B} + \delta I_p)/(1 + \delta)$, where $b_{ij} = b_{ji} = 0.05 \cdot \text{Bernoulli}(0.1)$ for $1 \leq i \leq p/2, i < j \leq p$, $b_{ij} = b_{ji} = 0.05$ for $p/2 + 1 \leq i < j \leq p$, $b_{ii} = 1$ for $1 \leq i \leq p$, and $\delta = \max\{-\lambda_{\min}(\mathbf{B}), 0\} + 0.1$.

The size of the training set is $n = 400$, with p varying from 500 to 1000. The mis-classification rate is computed based on $m = 2000$ test points generated from $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma)$ or $\mathcal{N}(\boldsymbol{\mu}_2, \Sigma)$ with equal probability. We repeat the experiment for 100 times, and report the misclassification rates (in percentage) in Table 1. As p varies from 500 to 1000, the proportion of non-zero entries of \mathbf{d} varies from 0.045 to 0.031, so this can still be viewed as a sparse setting. It can be seen that LASS remains competitive.

S5. ADDITIONAL SIMULATIONS

p	Oracle	Naive	LASS	LPD	AdaLDA	Lasso	Ebay
Model 1							
500	4.97	17.97	5.32	8.01	5.73	6.91	5.70
600	4.17	22.11	4.58	7.48	5.02	6.21	4.86
700	3.58	28.44	3.91	6.78	4.29	5.43	4.16
800	3.01	45.61	3.41	8.14	4.83	4.96	3.63
900	2.67	32.09	3.08	8.84	3.45	4.49	3.24
1000	2.48	30.05	2.85	13.12	3.27	4.23	2.97
Model 2							
500	5.82	19.32	6.21	8.62	6.58	7.84	6.27
600	4.98	23.44	5.41	8.11	5.85	7.03	5.50
700	4.29	29.81	4.65	7.37	5.00	6.17	4.70
800	3.72	46.62	4.11	8.30	4.56	5.76	4.18
900	3.20	32.67	3.68	8.24	4.02	5.26	3.73
1000	3.12	30.47	3.44	10.93	3.87	4.93	3.47
Model 3							
500	11.07	26.61	13.13	15.27	14.32	14.76	13.21
600	9.79	29.58	12.29	14.94	13.63	14.05	12.33
700	9.01	35.20	11.76	14.59	13.21	13.48	11.78
800	8.03	47.23	11.19	14.42	12.62	13.08	11.25
900	7.27	36.62	10.80	15.67	12.40	12.56	10.85
1000	6.90	33.69	10.88	14.72	12.63	12.82	11.04

Table 1: Comparison of average misclassification rate in percentage. The smallest error rate (next to that by the oracle) in each setting has been boldfaced.

S6 Lung cancer data

We analyze the lung cancer data (Gordon et al., 2002, available at <https://leo.ugr.es/elvira/DBCRpository/LungCancer/LungCancer-Harvard2.html>), a benchmark data set in high-dimensional classification problems. The data set collects expression levels for $p = 12,533$ genes on 181 tissue samples. Among the 181 samples, 31 are from malignant pleural mesothelioma (MPM) group and 150 are from adenocarcinoma (ADCA) group. The training set contains 16 samples from MPM and 16 samples from ADCA. The testing set contains 134 samples from MPM and 15 samples from ADCA.

We follow the same data pre-processing steps in Cai and Liu (2011). First, the sample variances of individual genes are obtained based on the training data. Next, we drop 195 genes whose sample variances are greater than 10^2 or less than 10^{-2} after rescaling by a factor of 10^4 . Finally, to reduce the computational cost (which mainly comes from estimating the precision matrix), only top 200 genes (with the largest absolute values of the two sample t -statistics) are used for constructing different classification rules. Figure 3 illustrates the scores $\hat{T}^j = \frac{\exp(\hat{S}_j)}{1+\exp(\hat{S}_j)}$ estimated by LASS, Naive, Ebay and Lasso. For better illustration, the first 134 testing points are from ADCA group and the next 15 are from MPM group.

To achieve better separation, the values of the first 134 points (the

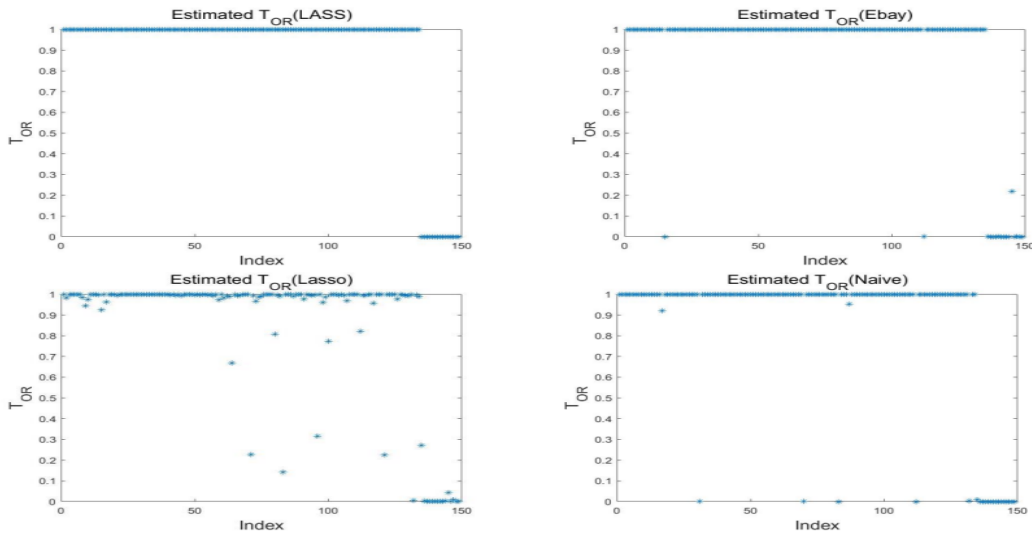


Figure 3: Comparison of estimated T : LASS separates the two classes almost perfectly. Ebay does better than Naive and Lasso but worse than LASS.

Table 2: Number of misclassifications by each method.

	LASS	Naive	LPD	AdaLDA	Lasso	Ebay
Misclassification	0 (0%)	5 (3.36%)	1 (0.67%)	1 (0.67%)	5 (3.36%)	3 (2.01%)

next 15 points) should be high (low). We can see that LASS shows perfect separation of the two classes. While Ebay provides a clearer separation than Naive and Lasso, it is less effective than LASS. For this particular data set, if we choose to control the FSR at level $\alpha = 0.1$, then LASS makes no indecision. Thus, it makes sense to compare the misclassification rates under the classical setup. The results are presented in Table 2, from which we can see that LASS has the best performance.

S7 An example comparing LASS and LPD

In this section we give an example to contrast our theory with existing theories in high-dimensional LDA, where the assumption on the sparsity of \mathbf{d} is commonly believed to be “indispensable”. For example, Cai and Liu (2011) requires $|\Sigma^{-1}\mathbf{d}|_0 = o(\sqrt{(n_1 + n_2)/\log p})$ to achieve risk consistency. This sparsity condition seems to be necessary if the scope of analysis is limited to the class of thresholding rules. It is not an artifact of the theoretical analysis, as we can see from the numerical results in Section 5.1.2 where the disadvantages of existing works are reflected. However, LASS still achieves risk consistency even when the condition in Cai and Liu (2011) is violated. We illustrate this through the next example.

Example 1. Consider an asymptotic setup where $n_1 = n_2 = n$, $p = n^2$, $\Sigma = I_p$ and $\boldsymbol{\mu}_1 = (0, \dots, 0)^\top$. Let $\boldsymbol{\mu}_2 = (1, \dots, 1, 0, \dots, 0)^\top$ be a vector with the first k entries being 1 and the rest 0. Let $k = n$ and $\Delta_p = \mathbf{d}^\top \Sigma^{-1} \mathbf{d}$, then in this setting $\Delta_p = n$. It is clear that $|\Sigma^{-1}\mathbf{d}|_0 = n$ and $|\Sigma^{-1}\mathbf{d}|_1 = n$. Denote $\hat{\boldsymbol{\delta}}_{LPD}$ the LPD rule in Cai and Liu (2011). To guarantee that $R(\hat{\boldsymbol{\delta}}_{LPD}) - R(\boldsymbol{\delta}^F) \rightarrow 0$, the theory in Cai and Liu (2011) requires that $|\Sigma^{-1}\mathbf{d}|_0 = o(\sqrt{n/\log p})$ or $\frac{|\Sigma^{-1}\mathbf{d}|_1}{\Delta_p^{1/2}} + \frac{|\Sigma^{-1}\mathbf{d}|_1^2}{\Delta_p^2} = o(\sqrt{n/\log p})$, both of which are violated. By contrast, note that all nonzero elements in \mathbf{d} are in \mathcal{G}_1 (strong signals), the conditions of Corollary 1 are satisfied for any $k \geq 1$,

which guarantees that LASS achieves risk consistency across sparse and dense settings.

As opposed to existing works that produce approximately the same amount of shrinkage for elements in both \mathcal{G}_1 (strong) and \mathcal{G}_3 (moderate), LASS adopts an adaptive strategy that makes it possible to provide differential amounts of shrinkage for the elements in \mathcal{G}_1 and \mathcal{G}_3 (Proposition 1 in Section 3.1 and Figure 2c in Section S4). The resulting shrinkage rule is more effective in keeping strong signals and eliminating noisy elements; this explains the superiority of LASS in both theory and numerical performance.

S8 Class-specific FSR control implies asymptotic global FSR control

We will show if mFSR^c is controlled for $c = 1, 2$ then mFSR is also controlled. The result then follows from Section S3. Recall

$$\text{mFSR}^c = \frac{\mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(\delta_j = c, \theta_j \neq c) \right\}}{\mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(\delta_j = c) \right\}}, \quad \text{mFSR} = \frac{\mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(\delta_j \neq \theta_j, \delta_j \neq 0) \right\}}{\mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(\delta_j \neq 0) \right\}}.$$

Denote

$$a^c = \mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(\delta_j = c, \theta_j \neq c) \right\}, \quad b^c = \mathbb{E} \left\{ \sum_{j=1}^m \mathbb{I}(\delta_j = c) \right\}.$$

Then $\text{mFSR}^c = a^c/b^c$, and

$$\text{mFSR} = \frac{a^1 + a^2}{b^1 + b^2} = \frac{a^1}{b^1} \frac{b^1}{b^1 + b^2} + \frac{a^2}{b^2} \frac{b^2}{b^2 + b^1} \leq \alpha \frac{b^1}{b^1 + b^2} + \alpha \frac{b^2}{b^1 + b^2} \leq \alpha.$$

S9 Discussion on condition (A2)

For all numerical experiments, we used the ACLIME estimator to estimate precision matrix Σ^{-1} . Suppose we have samples $X_1, \dots, X_n \sim N(\mu, \Sigma)$. In Cai et al. (2016) the authors proved that if the precision matrix $\Sigma^{-1} \in \mathbb{R}^{p \times p}$ belongs to $\mathcal{G}_q(c_{n,p}, M_{n,p})$:

$$\mathcal{G}_q(c_{n,p}, M_{n,p}) = \left\{ \begin{array}{l} \Omega = (\omega_{ij})_{1 \leq i, j \leq p} : \max_j \sum_{i=1}^p |\omega_{ij}|^q \leq c_{n,p}, \\ \|\Omega\|_1 \leq M_{n,p}, \lambda_{\max}(\Omega)/\lambda_{\min}(\Omega) \leq M_1, \Omega \succ 0 \end{array} \right\}.$$

Under some regularity conditions on p , n , $c_{n,p}$ the ACLIME estimator satisfies

$$\sup_{\mathcal{G}_q(c_{n,p}, M_{n,p})} \mathbb{E} \|\hat{\Sigma}^{-1} - \Sigma^{-1}\|_w^2 = M_{n,p}^{2-2q} c_{n,p} (\log p/n)^{1-q}$$

for all $1 \leq w \leq \infty$ and $0 \leq q < 1$. Hence if $M_{n,p}$ and $c_{n,p}$ are chosen such that

$$M_{n,p}^{2-2q} c_{n,p} (\log p/n)^{1-q} = o(1),$$

then the ACLIME estimator satisfies condition (A2).

We emphasize that our Condition (A2) is strictly weaker than the condition $\Sigma^{-1} \in \mathcal{G}_q(c_{n,p}, M_{n,p})$ considered in Cai et al. (2016). Moreover, LASS does not depend on any particular class of precision matrix estimators nor do we regard the parameter space $\mathcal{G}_q(c_{n,p}, M_{n,p})$ to be of any specific type. Instead, a range of estimators, as proposed in Yuan (2010), Liu and Luo (2015), Cai et al. (2016) and Avella-Medina et al. (2018), all have good empirical performances and are consistent under various conditions. These estimators can be employed to construct LASS classifiers too.

In practice we often have some prior knowledge on the covariance structure, and we can choose which estimator to use according to this prior knowledge. For example, if we know that the covariance matrix has an AR(1) structure (i.e. $\Sigma = (\sigma_{ij})$, with $\sigma_{ij} = \rho^{|i-j|}$) then the inverse of the MLE of the covariance matrix $\hat{\Sigma}_{mle}^{-1}$ satisfies (A2). Hence, we do not need to assume Σ^{-1} belongs to any particular parameter space.

Estimating precision matrix is still an active area of research. It is conceivable (and hopeful) that future works on this topic can further relax the sufficient conditions under which an existing estimator is consistent. It would also be of interest to propose new estimators that are consistent in a wider range of parameter spaces.

Bibliography

- Avella-Medina, M., H. S. Battey, J. Fan, and Q. Li (2018). Robust estimation of high-dimensional covariance and precision matrices. *Biometrika* 105(2), 271–284.
- Cai, T. and W. Liu (2011). A direct estimation approach to sparse linear discriminant analysis. *Journal of the American statistical association* 106(496), 1566–1577.
- Cai, T. T. (2002). On block thresholding in wavelet regression: Adaptivity, block size, and threshold level. *Statistica Sinica*, 1241–1273.
- Cai, T. T., W. Liu, H. H. Zhou, et al. (2016). Estimating sparse precision matrix: Optimal rates of convergence and adaptive estimation. *The Annals of Statistics* 44(2), 455–488.
- Foygel, R. and M. Drton (2010). Extended bayesian information criteria for gaussian graphical models. *Advances in neural information processing systems* 23.
- Gordon, G. J., R. V. Jensen, L.-L. Hsiao, S. R. Gullans, J. E. Blumenstock, S. Ramaswamy, W. G. Richards, D. J. Sugarbaker, and R. Bueno (2002). Translation of microarray data into clinically relevant cancer diagnostic

tests using gene expression ratios in lung cancer and mesothelioma. *Cancer research* 62(17), 4963–4967.

Liu, W. and X. Luo (2015). Fast and adaptive sparse precision matrix estimation in high dimensions. *Journal of multivariate analysis* 135, 153–162.

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