

Appendix

A. Proof of Proposition 1

As mentioned in the statement, $\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^* = \tilde{\boldsymbol{\alpha}} + \tilde{\boldsymbol{\beta}} - \bar{\boldsymbol{\theta}} = \tilde{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}} + \tilde{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}} = \boldsymbol{\Delta} + \boldsymbol{\Gamma}$. If (i) either $[\boldsymbol{\Delta}]_j$ or $[\boldsymbol{\Gamma}]_j$ is zero, or (ii) $\text{sign}([\boldsymbol{\Delta}]_j) = \text{sign}([\boldsymbol{\Gamma}]_j)$, then $([\boldsymbol{\Delta}]_j + [\boldsymbol{\Gamma}]_j)^2 \geq ([\boldsymbol{\Delta}]_j)^2 + ([\boldsymbol{\Gamma}]_j)^2$. Therefore, if this happens for every j , the inequality (10) holds. When either $[\boldsymbol{\alpha}^*]_j \neq 0$ or $[\boldsymbol{\beta}^*]_j \neq 0$ holds, $[\boldsymbol{\Delta}]_j = 0$ or $[\boldsymbol{\Gamma}]_j = 0$ is guaranteed by construction (rule 2 or 3 above). If both $[\boldsymbol{\alpha}^*]_j$ and $[\boldsymbol{\beta}^*]_j$ are zero (in case of rule 1), the following lemma ensures that $[\tilde{\boldsymbol{\alpha}}]_j$ and $[\tilde{\boldsymbol{\beta}}]_j$ (and therefore $[\boldsymbol{\Delta}]_j$ and $[\boldsymbol{\Gamma}]_j$ because $[\tilde{\boldsymbol{\alpha}}]_j = [\tilde{\boldsymbol{\beta}}]_j = 0$ in this case) have same signs.

Lemma 1. *The signs of $\tilde{\boldsymbol{\alpha}}$ and $\tilde{\boldsymbol{\beta}}$ are always consistent whenever both are not zeros: $\text{sign}([\tilde{\boldsymbol{\alpha}}]_j) = \text{sign}([\tilde{\boldsymbol{\beta}}]_j)$ for all j such that $[\tilde{\boldsymbol{\alpha}}]_j \neq 0$ and $[\tilde{\boldsymbol{\beta}}]_j \neq 0$.*

The proof of Lemma 1 is trivial. Suppose that $[\tilde{\boldsymbol{\alpha}}]_j$ and $[\tilde{\boldsymbol{\beta}}]_j$ have opposite signs; say $[\tilde{\boldsymbol{\alpha}}]_j > 0$ and $[\tilde{\boldsymbol{\beta}}]_j < 0$. Then $[\tilde{\boldsymbol{\alpha}}]_j - \epsilon$ and $[\tilde{\boldsymbol{\beta}}]_j + \epsilon$ with arbitrary small positive ϵ and all others fixed, will have the same loss by $\mathcal{L}(\cdot)$, but smaller values in the regularizers, which violates the stationary condition of local minimum.

Since the sign consistency is always guaranteed (or at least one of them is zero), the decomposability in the statement trivially holds.

Showing (11) also comes from the definitions of support sets. For any j such that $j \notin s^*$ but $j \in U$, we set $\bar{\boldsymbol{\alpha}}_j := \tilde{\boldsymbol{\alpha}}_j$ and hence $\boldsymbol{\Delta}_j = 0$ by definition. Therefore, the projection does not make any difference. The equality holds for $\boldsymbol{\Delta}_{\bar{s}}$ since $s^* \subseteq \bar{s} \subseteq U$. The same reasoning holds for $\boldsymbol{\Gamma}$ as well.

B. Proofs for ℓ_2 Error Bounds

B.1. Proof of Theorem 1

Recall that we get $(\bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}})$ by the transformation $\mathcal{T}(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*; \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}})$. For notational simplicity, we use $\boldsymbol{\Lambda}$ to denote the error vector on our estimation: $\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$. We also define the individual error vectors as $\boldsymbol{\Delta} := \tilde{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}$ and $\boldsymbol{\Gamma} := \tilde{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}$. Note that $\boldsymbol{\Lambda} = \boldsymbol{\Delta} + \boldsymbol{\Gamma}$ since $\tilde{\boldsymbol{\theta}} = \tilde{\boldsymbol{\alpha}} + \tilde{\boldsymbol{\beta}}$ and $\boldsymbol{\theta}^* = \bar{\boldsymbol{\alpha}} + \bar{\boldsymbol{\beta}}$ by definitions.

We first show that under (RSC), $\|\boldsymbol{\Lambda}\|_2 \leq 1$ is guaranteed, and hence the first inequality (12) for the case of $\|\boldsymbol{\Lambda}\|_2 \leq 1$ only matters. Toward this, given all the assumptions in the statement, suppose that $\|\boldsymbol{\Lambda}\|_2 \geq 1$.

Then setting $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$ as $(\boldsymbol{\theta}^*, \boldsymbol{\Lambda})$ in (13), we obtain

$$\kappa_2 \|\boldsymbol{\Lambda}\|_2 - \tau_2 \|\boldsymbol{\Lambda}\|_{\eta} \leq \langle \nabla \mathcal{L}(\boldsymbol{\theta}^* + \boldsymbol{\Lambda}) - \nabla \mathcal{L}(\boldsymbol{\theta}^*), \boldsymbol{\Lambda} \rangle. \quad (19)$$

At the same time, the stationary condition of (14) ensures that any local optimum satisfies

$$\langle \nabla \mathcal{L}(\tilde{\boldsymbol{\theta}}) + \nabla \|\tilde{\boldsymbol{\theta}}\|_{\lambda}, \boldsymbol{\theta} - \tilde{\boldsymbol{\theta}} \rangle \geq 0 \quad (20)$$

for any feasible $\boldsymbol{\theta}$. This is more general than the first-order stationary condition; for some local optima, the first-order stationary condition does not hold due to the inequality constraint of program (14). (see (Loh & Wainwright, 2014; 2015) for more details.) Since $\boldsymbol{\theta}^*$ is feasible by assumption, setting $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ yields

$$\langle \nabla \mathcal{L}(\tilde{\boldsymbol{\theta}}) + \nabla \|\tilde{\boldsymbol{\theta}}\|_{\lambda}, \boldsymbol{\Lambda} \rangle \leq 0. \quad (21)$$

Combining (19) and (21) yields

$$\kappa_2 \|\boldsymbol{\Lambda}\|_2 - \tau_2 \|\boldsymbol{\Lambda}\|_{\eta} \leq \langle -\nabla \mathcal{L}(\boldsymbol{\theta}^*) - \nabla \|\tilde{\boldsymbol{\theta}}\|_{\lambda}, \boldsymbol{\Lambda} \rangle.$$

By Hölder's inequality,

$$\begin{aligned} \langle -\nabla \mathcal{L}(\boldsymbol{\theta}^*) - \nabla \|\tilde{\boldsymbol{\theta}}\|_{\lambda}, \boldsymbol{\Lambda} \rangle &\leq \left\| \nabla \mathcal{L}(\boldsymbol{\theta}^*) + \nabla \|\tilde{\boldsymbol{\theta}}\|_{\lambda} \right\|_{\eta}^* \|\boldsymbol{\Lambda}\|_{\eta} \\ &\leq \left(\|\nabla \mathcal{L}(\boldsymbol{\theta}^*)\|_{\eta}^* + \bar{\eta} \right) \|\boldsymbol{\Lambda}\|_{\eta} \end{aligned} \quad (22)$$

where we utilize Lemma 3 with defining $\bar{\eta} := \max\{\frac{\lambda_1}{\eta_1}, \frac{\lambda_2}{\eta_2}\}$ in the second inequality. Moreover, since $\|\nabla\mathcal{L}(\boldsymbol{\theta}^*)\|_{\boldsymbol{\eta}}^* = \max\left(\frac{\|\nabla\mathcal{L}(\boldsymbol{\theta}^*)\|_{\infty}}{\eta_1}, \frac{\|\nabla\mathcal{L}(\boldsymbol{\theta}^*)\|_{\infty, a^*}}{\eta_2}\right) \leq \max\left(\frac{\lambda_1}{4\eta_1}, \frac{\lambda_2}{4\eta_2}\right)$ by assumption, we obtain

$$\kappa_2\|\boldsymbol{\Lambda}\|_2 - \tau_2\|\boldsymbol{\Lambda}\|_{\boldsymbol{\eta}} \leq \frac{5\bar{\eta}}{4}\|\boldsymbol{\Lambda}\|_{\boldsymbol{\eta}}. \quad (23)$$

Rearranging (23), we have

$$\|\boldsymbol{\Lambda}\|_2 \leq \frac{\|\boldsymbol{\Lambda}\|_{\boldsymbol{\eta}}}{\kappa_2} \left(\frac{5\bar{\eta}}{4} + \tau_2\right) \stackrel{(i)}{\leq} \frac{2r}{\kappa_2} \left(\frac{5\bar{\eta}}{4} + \tau_2\right) \stackrel{(ii)}{\leq} 1$$

where (i) follows the fact $\|\boldsymbol{\Lambda}\|_{\boldsymbol{\eta}} \leq \|\tilde{\boldsymbol{\theta}}\|_{\boldsymbol{\eta}} + \|\boldsymbol{\theta}^*\|_{\boldsymbol{\eta}} \leq 2r$ (since $\boldsymbol{\theta}^*$ is feasible) and (ii) holds under the conditions that $r \leq \frac{\kappa_2}{4\tau_2}$ and $r\bar{\eta} \leq \frac{\kappa_2}{5}$.

From now we revisit the stationary condition of the problem (14) and the RSC condition (12) for $\|\boldsymbol{\Lambda}\|_2 \leq 1$. Recall that \bar{s} is the union of $\text{supp}(\boldsymbol{\alpha}^*)$ and $\text{supp}(\tilde{\boldsymbol{\alpha}})$, and $\boldsymbol{\Delta}_{\bar{s}}$ is the projection of $\boldsymbol{\Delta}$ onto the space w.r.t. \bar{s} (simply meaning $[\boldsymbol{\Delta}_{\bar{s}}]_j = 0$ if $j \notin \bar{s}$). $\bar{\mathbf{b}}$ was similarly defined as $\text{supp}(\boldsymbol{\beta}^*) \cup \text{supp}(\tilde{\boldsymbol{\beta}})$. Rearranging the stationary condition (21), we obtain

$$\begin{aligned} \langle \nabla\mathcal{L}(\tilde{\boldsymbol{\theta}}), \boldsymbol{\Lambda} \rangle &\leq -\langle \nabla\|\tilde{\boldsymbol{\theta}}\|_{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^* \rangle \stackrel{(i)}{\leq} \|\tilde{\boldsymbol{\theta}}\|_{\boldsymbol{\lambda}} - \|\tilde{\boldsymbol{\theta}}\|_{\boldsymbol{\lambda}} \\ &\stackrel{(ii)}{\leq} \lambda_1\|\tilde{\boldsymbol{\alpha}}\|_1 + \lambda_2\|\tilde{\boldsymbol{\beta}}\|_{1,a} - \|\tilde{\boldsymbol{\theta}}\|_{\boldsymbol{\lambda}} \stackrel{(iii)}{\leq} \lambda_1\|\tilde{\boldsymbol{\alpha}}\|_1 + \lambda_2\|\tilde{\boldsymbol{\beta}}\|_{1,a} - \lambda_1\|\tilde{\boldsymbol{\alpha}}\|_1 - \lambda_2\|\tilde{\boldsymbol{\beta}}\|_{1,a} \\ &= \lambda_1\left(\|\tilde{\boldsymbol{\alpha}}\|_1 + \|[\boldsymbol{\Delta}]_{\bar{s}^c}\|_1 - \|[\boldsymbol{\Delta}]_{\bar{s}^c}\|_1 - \|\tilde{\boldsymbol{\alpha}}\|_1\right) \\ &\quad + \lambda_2\left(\|\tilde{\boldsymbol{\beta}}\|_{1,a} + \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a} - \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a} - \|\tilde{\boldsymbol{\beta}}\|_{1,a}\right) \end{aligned}$$

where we use (i) the convexity and (iii) the triangular inequality of $\|\cdot\|_{\boldsymbol{\lambda}}$ norm, and we obtain the inequality (ii) since $\|\tilde{\boldsymbol{\theta}}\|_{\boldsymbol{\lambda}}$ has the infimal sum of two regularizers by definition. Hence,

$$\begin{aligned} &\langle \nabla\mathcal{L}(\tilde{\boldsymbol{\theta}}), \boldsymbol{\Lambda} \rangle \\ &\leq \lambda_1\left(\|\tilde{\boldsymbol{\alpha}}\|_1 + \|[\boldsymbol{\Delta}]_{\bar{s}^c}\|_1 - \|[\boldsymbol{\Delta}]_{\bar{s}^c}\|_1 - \|\tilde{\boldsymbol{\alpha}}\|_1\right) \\ &\quad + \lambda_2\left(\|\tilde{\boldsymbol{\beta}}\|_{1,a} + \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a} - \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a} - \|\tilde{\boldsymbol{\beta}}\|_{1,a}\right) \\ &\stackrel{(i)}{=} \lambda_1\left(\|\tilde{\boldsymbol{\alpha}} + [\boldsymbol{\Delta}]_{\bar{s}^c}\|_1 - \|[\boldsymbol{\Delta}]_{\bar{s}^c}\|_1 - \|\tilde{\boldsymbol{\alpha}}\|_1\right) \\ &\quad + \lambda_2\left(\|\tilde{\boldsymbol{\beta}} + [\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a} - \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a} - \|\tilde{\boldsymbol{\beta}}\|_{1,a}\right) \\ &\stackrel{(ii)}{\leq} \lambda_1\left(\|\tilde{\boldsymbol{\alpha}} + [\boldsymbol{\Delta}]_{\bar{s}^c} + [\boldsymbol{\Delta}]_{\bar{s}}\|_1 + \|[\boldsymbol{\Delta}]_{\bar{s}}\|_1 - \|[\boldsymbol{\Delta}]_{\bar{s}^c}\|_1 - \|\tilde{\boldsymbol{\alpha}}\|_1\right) \\ &\quad + \lambda_2\left(\|\tilde{\boldsymbol{\beta}} + [\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c} + [\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}}\|_{1,a} + \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}}\|_{1,a} - \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a} - \|\tilde{\boldsymbol{\beta}}\|_{1,a}\right) \\ &= \lambda_1\left(\|[\boldsymbol{\Delta}]_{\bar{s}}\|_1 - \|[\boldsymbol{\Delta}]_{\bar{s}^c}\|_1\right) + \lambda_2\left(\|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}}\|_{1,a} - \|[\boldsymbol{\Gamma}]_{\bar{\mathbf{b}}^c}\|_{1,a}\right) \end{aligned} \quad (24)$$

where (i) and (ii) hold by the decomposability and triangular inequality of norms, respectively.

We also compute the upper bound for term $|\langle \nabla\mathcal{L}(\tilde{\boldsymbol{\theta}}), \boldsymbol{\Lambda} \rangle|$:

$$\begin{aligned} -\langle \nabla\mathcal{L}(\tilde{\boldsymbol{\theta}}), \boldsymbol{\Lambda} \rangle &= -\langle \nabla\mathcal{L}(\boldsymbol{\theta}^*), \boldsymbol{\Lambda} \rangle \stackrel{(i)}{=} -\langle \nabla\mathcal{L}(\boldsymbol{\theta}^*), \boldsymbol{\Delta} \rangle - \langle \nabla\mathcal{L}(\boldsymbol{\theta}^*), \boldsymbol{\Gamma} \rangle \\ &\stackrel{(ii)}{\leq} \|\nabla\mathcal{L}(\boldsymbol{\theta}^*)\|_{\infty} \|\boldsymbol{\Delta}\|_1 + \|\nabla\mathcal{L}(\boldsymbol{\theta}^*)\|_{\infty, a^*} \|\boldsymbol{\Gamma}\|_{1,a} \end{aligned} \quad (25)$$

since (i) holds from the fact $\boldsymbol{\Lambda} = \boldsymbol{\Delta} + \boldsymbol{\Gamma}$ by definition, and (ii) does by the two standard Hölder's inequalities.

We now combine (12) with (24) and (25), and obtain

$$\begin{aligned} \kappa_1 \|\mathbf{A}\|_2^2 - \tau_1 \|\mathbf{A}\|_{\eta}^2 &\leq \langle \nabla \mathcal{L}(\bar{\boldsymbol{\theta}} + \mathbf{A}) - \nabla \mathcal{L}(\bar{\boldsymbol{\theta}}), \mathbf{A} \rangle \\ &\leq \|\nabla \mathcal{L}(\bar{\boldsymbol{\theta}})\|_{\infty} \|\mathbf{A}\|_1 + \|\nabla \mathcal{L}(\bar{\boldsymbol{\theta}})\|_{\infty, a^*} \|\mathbf{\Gamma}\|_{1, a} \\ &\quad + \lambda_1 \left(\|\mathbf{[\Delta]_{\bar{s}}}\|_1 - \|\mathbf{[\Delta]_{\bar{s}^c}}\|_1 \right) + \lambda_2 \left(\|\mathbf{[\Gamma]_{\bar{b}}}\|_{1, a} - \|\mathbf{[\Gamma]_{\bar{b}^c}}\|_{1, a} \right). \end{aligned}$$

Moreover, since $\|\mathbf{A}\|_{\eta} \leq \eta_1 \|\mathbf{A}\|_1 + \eta_2 \|\mathbf{\Gamma}\|_{1, a}$ by definition of dirty regularizer $\|\cdot\|_{\eta}$, the above inequality can be rearranged as

$$\begin{aligned} \kappa_1 \|\mathbf{A}\|_2^2 &\leq \left(\|\nabla \mathcal{L}(\bar{\boldsymbol{\theta}})\|_{\infty} + 2\tau_1 \eta_1 r \right) \|\mathbf{A}\|_1 + \left(\|\nabla \mathcal{L}(\bar{\boldsymbol{\theta}})\|_{\infty, a^*} + 2\tau_1 \eta_2 r \right) \|\mathbf{\Gamma}\|_{1, a} \\ &\quad + \lambda_1 \left(\|\mathbf{[\Delta]_{\bar{s}}}\|_1 - \|\mathbf{[\Delta]_{\bar{s}^c}}\|_1 \right) + \lambda_2 \left(\|\mathbf{[\Gamma]_{\bar{b}}}\|_{1, a} - \|\mathbf{[\Gamma]_{\bar{b}^c}}\|_{1, a} \right) \\ &\leq \frac{3\lambda_1}{2} \|\mathbf{[\Delta]_{\bar{s}}}\|_1 + \frac{3\lambda_2}{2} \|\mathbf{[\Gamma]_{\bar{b}}}\|_{1, a}. \end{aligned}$$

Note that in the last inequality, we utilize the assumption on the setting λ_1 and λ_2 in the statement (that is, $4\|\nabla \mathcal{L}(\boldsymbol{\theta}^*)\|_{\infty} \leq \lambda_1$ and $8\tau_1 \eta_1 r \leq \lambda_1$ for \mathbf{A} case) and the fact that $\mathbf{A} = \mathbf{[\Delta]_{\bar{s}}} + \mathbf{[\Delta]_{\bar{s}^c}}$. Also note that we dropped minus terms at the end.

Now, in order to relate the ℓ_1 and ℓ_2 norms of projected error vector, $\mathbf{[\Delta]_{\bar{s}}}$, we need to consider the sparsity level of $\bar{\boldsymbol{\alpha}}$. During the construction of $\bar{\boldsymbol{\alpha}}$, the sparsity level of $\bar{\boldsymbol{\alpha}}$ might be possibly greater than that of $\boldsymbol{\alpha}^*$. However, by (11) of Proposition 1, we have $\boldsymbol{\Delta}_{\mathbf{s}^*} = \boldsymbol{\Delta}_{\bar{\mathbf{s}}}$ as well as $\boldsymbol{\Gamma}_{\mathbf{b}^*} = \boldsymbol{\Gamma}_{\bar{\mathbf{b}}}$. Hence,

$$\begin{aligned} \kappa_1 \|\mathbf{A}\|_2^2 &\leq \frac{3\lambda_1}{2} \sqrt{s} \|\mathbf{[\Delta]_{\bar{s}}}\|_2 + \frac{3\lambda_2}{2} \sqrt{s\mathcal{G}} \|\mathbf{[\Gamma]_{\bar{b}}}\|_2 \\ &\leq \frac{3\lambda_1}{2} \sqrt{s} \|\mathbf{A}\|_2 + \frac{3\lambda_2}{2} \sqrt{s\mathcal{G}} \|\mathbf{\Gamma}\|_2 \leq \frac{3\bar{\lambda}}{2} \left(\|\mathbf{A}\|_2 + \|\mathbf{\Gamma}\|_2 \right) \end{aligned} \quad (26)$$

where $\bar{\lambda} := \max\{\lambda_1 \sqrt{s}, \lambda_2 \sqrt{s\mathcal{G}}\}$.

Combining (26) with the result in (10) of Proposition 1, we have

$$\frac{\kappa_1}{2} \left(\|\mathbf{A}\|_2 + \|\mathbf{\Gamma}\|_2 \right)^2 \leq \kappa_1 \left(\|\mathbf{A}\|_2^2 + \|\mathbf{\Gamma}\|_2^2 \right) \leq \kappa_1 \|\mathbf{A}\|_2^2 \leq \frac{3\bar{\lambda}}{2} \left(\|\mathbf{A}\|_2 + \|\mathbf{\Gamma}\|_2 \right),$$

and finally the upper bound of ℓ_2 error can be computed as stated:

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_2 + \|\mathbf{\Gamma}\|_2 \leq \frac{3}{\kappa_1} \bar{\lambda},$$

which completes the proof.

B.2. Proof of Corollary 1

We first show the loss function of (3) satisfy RSC condition in the following proposition.

Proposition 2. Consider a design matrix $X \in \mathbb{R}^{n \times p}$ whose rows are independently sampled from $N(0, \Sigma)$. Then, with probability at least $1 - c_1 \exp(-c_2 n)$ for some constants c_1 and c_2 ,

$$\frac{1}{n} \|X\boldsymbol{\theta}\|_2^2 \geq \kappa'_1 \|\boldsymbol{\theta}\|_2^2 - \tau'_1 \|\boldsymbol{\theta}\|_{\eta'}^2 \quad \text{for any } \boldsymbol{\theta} \in \mathbb{R}^p. \quad (27)$$

where $\eta'_1 = \sqrt{\frac{\log p}{n}}$ and $\eta'_2 = \frac{\mathbb{E}(\|\boldsymbol{\varepsilon}\|_{\infty, a^*})}{\sqrt{n}}$ letting a standard normal vector $\boldsymbol{\varepsilon}$ for $\boldsymbol{\eta}'$, and κ'_1 and τ'_1 are some constants depending only on Σ .

As discussed in (Negahban et al., 2012), for the case of $a = 2$, $\frac{\mathbb{E}(\|\boldsymbol{\varepsilon}\|_{\infty, 2})}{\sqrt{n}} \leq \left\{ \sqrt{\frac{m}{n}} + \sqrt{\frac{3 \log q}{n}} \right\}$ by the standard tail bound of χ^2 variables.

The next step in order to appeal to Theorem 1, is to set the proper regularization parameters. Under the conditions in the statement, (Negahban et al., 2012) show that

$$4\|\nabla\mathcal{L}(\boldsymbol{\theta}^*)\|_\infty = 4\left\|\frac{1}{n}X^\top w\right\|_\infty \leq 8\sigma\sqrt{\frac{\log p}{n}}$$

with probability at least $1 - c_1 \exp(-c_2 n \lambda_1^2)$. Moreover,

$$4\|\nabla\mathcal{L}(\boldsymbol{\theta}^*)\|_{\infty, a^*} = 4 \max_{t=1,2,\dots,q} \left\|\frac{1}{n}X_{G_t}^\top w\right\|_{a^*} \leq 8\sigma \left\{ \frac{m^{1-1/a}}{\sqrt{n}} + \sqrt{\frac{\log q}{n}} \right\}$$

with probability at least $1 - 2 \exp(-2 \log q)$.

Since $\|\boldsymbol{\theta}\|_{\boldsymbol{\eta}'}^2 \leq \frac{3}{64} \|\boldsymbol{\theta}\|_\lambda^2$ for the values of $\boldsymbol{\eta}'$ and λ specified, (12) holds with $\kappa_1 = \kappa_1'$ and $\tau_1 = \frac{3}{64} \tau_1'$ by (27). In addition, (13) holds with $\kappa_2 = \kappa_1$ and $\tau_2 = \frac{3}{32} \tau_1'$, by Lemma 8.

Therefore, combining all pieces, for a constant $r := \min \left\{ \sqrt{\frac{8\kappa_1}{3\tau_1}}, \frac{\kappa_1}{5}, \frac{8}{3\tau_1} \right\}$, we can guarantee that

$$\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 \leq \frac{24\sigma}{\kappa_1} \max \left\{ \sqrt{\frac{s \log p}{n}}, \sqrt{\frac{s_G m}{n}} + \sqrt{\frac{s_G \log q}{n}} \right\}.$$

B.3. Proof of Proposition 2

The proof of Proposition 2 is the simple extension of proofs given by (Raskutti et al., 2010; Negahban et al., 2012): with probability greater than $1 - c_1 \exp(-\lambda_2 n)$ for some constants c_1 and c_2 ,

$$\frac{\|X\boldsymbol{\theta}\|_2}{\sqrt{n}} \geq \frac{\lambda_{\min}(\Sigma^{1/2})}{4} \|\boldsymbol{\theta}\|_2 - 3 \frac{\mathbb{E}(\|w\|_\eta^*)}{\sqrt{n}} \|\boldsymbol{\theta}\|_\eta \quad \text{for all } \boldsymbol{\theta} \in \mathbb{R}^p \quad (28)$$

where $w \sim N(0, \Sigma)$, $\lambda_{\min}(\Sigma^{1/2})$ is the minimum eigenvalue of $\Sigma^{1/2}$. Here we use $\|\cdot\|_\eta$ and its dual norm to arrive at the statement rather than ℓ_1 (and ℓ_∞) as in (Raskutti et al., 2010) or $\|\cdot\|_{1,a}$ (and its dual) in (Negahban et al., 2012).

In particular, (Raskutti et al., 2010) obtains the result in terms of ℓ_1 from the fact $\frac{\mathbb{E}(\|w\|_\infty)}{\sqrt{n}} \leq 3(\max_i \Sigma_{ii}) \sqrt{\frac{\log p}{n}}$ while (Negahban et al., 2012) does for group lasso case from $\frac{\mathbb{E}(\|w\|_{\infty, a^*})}{\sqrt{n}} \leq \max_{t=1,\dots,N_G} \|\Sigma^{1/2}\|_{G_t} \|a^*\|_{a^*} \frac{\mathbb{E}(\|\varepsilon\|_{\infty, a^*})}{\sqrt{n}}$ where ε is a standard normal vector, and $\|A\|_{a^*} := \max_{\|\boldsymbol{\theta}\|_{a^*}=1} \|A\boldsymbol{\theta}\|_{a^*}$.

Recalling $\mathbb{E}(\|w\|_\eta^*) = \mathbb{E}(\max\{\frac{\|w\|_\infty}{\eta_1}, \frac{\|w\|_{\infty, a^*}}{\eta_2}\})$,

$$\mathbb{E}(\|w\|_\eta^*) \leq \mathbb{E}\left(\frac{\|w\|_\infty}{\eta_1} + \frac{\|w\|_{\infty, a^*}}{\eta_2}\right) = \mathbb{E}\left(\frac{\|w\|_\infty}{\eta_1}\right) + \mathbb{E}\left(\frac{\|w\|_{\infty, a^*}}{\eta_2}\right) \quad (29)$$

since both $\frac{\|w\|_\infty}{\eta_1}$ and $\frac{\|w\|_{\infty, a^*}}{\eta_2}$ are always greater than or equal to zero. Setting $\eta_1 = \sqrt{\frac{\log p}{n}}$ and $\eta_2 = \frac{\mathbb{E}(\|\varepsilon\|_{\infty, a^*})}{\sqrt{n}}$ for $\boldsymbol{\eta}$, we have

$$\mathbb{E}(\|w\|_\eta^*) \leq \left(3(\max_i \Sigma_{ii}) + \max_{t=1,\dots,N_G} \|\Sigma^{1/2}\|_{G_t} \|a^*\|\right) \sqrt{n} \quad (30)$$

and we can establish the bound

$$\frac{\|X\boldsymbol{\theta}\|_2}{\sqrt{n}} \geq \frac{\lambda_{\min}(\Sigma^{1/2})}{4} \|\boldsymbol{\theta}\|_2 - 3 \left(3(\max_i \Sigma_{ii}) + \max_{t=1,\dots,N_G} \|\Sigma^{1/2}\|_{G_t} \|a^*\|\right) \|\boldsymbol{\theta}\|_\eta. \quad (31)$$

Finally, (27) can be obtained with constants $\kappa_1 = \frac{\lambda_{\min}(\Sigma)}{32}$ and $\tau_1 = 9 \left(3(\max_i \Sigma_{ii}) + \max_{t=1,\dots,N_G} \|\Sigma^{1/2}\|_{G_t} \|a^*\|\right)^2$ from the fact that $a \geq c - b$ implies $2(a^2 + b^2) \geq (a + b)^2 \geq c^2$ for positive real a, b and c .

B.4. Proof of Corollary 2

As shown in (Loh & Wainwright, 2015), it can be easily shown that the RSC condition (12) holds with $\kappa_1 = (\|\Theta^*\|_2 + 1)^{-2}$ and $\tau_1 = 0$. Hence by Lemma 7, (13) also holds with $\kappa_2 = \kappa_1$ and $\tau_2 = 0$. Hence, the condition on r in Theorem 1 is reduced to $r \leq \frac{1}{5(\|\Theta^*\|_2 + 1)^2}$. Moreover, since $\nabla \mathcal{L}(\theta^*) = \hat{\Sigma} - (\Theta^*)^{-1}$, the choices of regularization parameters satisfy the condition of Theorem 1, and its error bound holds.

C. Proofs for Support Set Recovery Guarantees

Before providing the actual proof of Theorem 2, we briefly review the primal-dual witness (PDW) proof technique (Wainwright, 2009; Jalali et al., 2010; Yang et al., 2015; Loh & Wainwright, 2014) for our setting:

(i) Solve the *restricted* problem

$$\underset{\theta \in \mathbb{R}^U, \|\theta\|_\eta \leq r}{\text{minimize}} \quad \mathcal{L}(\theta) + \mathcal{R}(\theta; \lambda) \quad (32)$$

where $U := \text{supp}(\alpha^*) \cup \text{supp}(\beta^*)$ as defined earlier. We set $\hat{\theta}_U$ as the local minimum of this problem.

(ii) Define $\hat{z}_1 \in [\nabla \|\theta\|_\lambda]_U$ and $q_\lambda(\theta) := \|\theta\|_\lambda - \mathcal{R}(\theta; \lambda)$. Choose \hat{z}_2 such that

$$\nabla \mathcal{L}(\hat{\theta}) - \nabla q_\lambda(\hat{\theta}) + \hat{z} = \mathbf{0} \quad (33)$$

where $\hat{z} := (\hat{z}_1, \hat{z}_2)$, $\hat{\theta} := (\hat{\theta}_U, \mathbf{0})$ and $[\nabla \|\theta\|_\lambda]_j = [\nabla \mathcal{R}(\theta; \lambda)]_j = [\hat{z}_2]_j$ for $j \notin U$ so that $[\nabla q_\lambda(\hat{\theta})]_j = 0$ at $\hat{\theta}_j = 0$. Establish strict dual feasibility of $\|\hat{z}_2\|_\lambda^* \leq 1 - \delta$ for some $\delta \in (0, 1]$.

(iii) Show that all stationary points of (17) are supported on U .

C.1. Proof of Theorem 2

Since $\hat{\theta}$ has the same sparsity pattern with θ^* , we can begin with very similar analysis as the proof of Theorem 1 by handling $\hat{\theta}$ as θ^* . Hence, in this proof we re-define $(\tilde{\alpha}, \tilde{\beta}) := \mathcal{T}(\hat{\alpha}, \hat{\beta}; \tilde{\alpha}, \tilde{\beta})$ where $(\tilde{\alpha}, \tilde{\beta})$ is the (local) minimizer of non-convex dirty regularizer such that $\tilde{\alpha} + \tilde{\beta} = \tilde{\theta}$. Therefore, by construction $\tilde{\theta} = \hat{\theta}$. We also re-use the notation Λ , Δ and Γ to represent again $\Lambda := \tilde{\theta} - \hat{\theta}$, $\tilde{\alpha} - \hat{\alpha}$ and $\tilde{\beta} - \hat{\beta}$ respectively but for newly defined $\tilde{\alpha}$ and $\tilde{\beta}$.

We begin with the first-order stationary condition of (17):

$$\langle \nabla \mathcal{L}(\tilde{\theta}) + \nabla \mathcal{R}(\tilde{\theta}; \lambda), \Lambda \rangle \leq 0, \quad (34)$$

which is in the form discussed in (21).

As before, we first show that $\|\Lambda\|_2 \leq 1$ under (RSC). In order to show by contradiction, suppose that $\|\Lambda\|_2 \geq 1$.

Then setting $(\theta_1, \theta_2) = (\hat{\theta}, \Lambda)$ in (13), we obtain

$$\kappa_2 \|\Lambda\|_2 - \tau_2 \|\Lambda\|_\eta \leq \langle \nabla \mathcal{L}(\hat{\theta} + \Lambda) - \nabla \mathcal{L}(\hat{\theta}), \Lambda \rangle. \quad (35)$$

Combining (34) and (35) yields

$$\kappa_2 \|\Lambda\|_2 - \tau_2 \|\Lambda\|_\eta \leq \langle -\nabla \mathcal{L}(\hat{\theta}) - \nabla \mathcal{R}(\tilde{\theta}; \lambda), \Lambda \rangle. \quad (36)$$

Since we set \hat{z} to satisfy (33),

$$\nabla \mathcal{L}(\hat{\theta}) = \nabla q_\lambda(\hat{\theta}) - \hat{z} = \nabla \|\hat{\theta}\|_\lambda - \nabla \mathcal{R}(\hat{\theta}; \lambda) - \hat{z} = -\nabla \mathcal{R}(\hat{\theta}; \lambda),$$

and hence by Hölder's inequality and Lemma 4,

$$\begin{aligned} \kappa_2 \|\Lambda\|_2 - \tau_2 \|\Lambda\|_\eta &\leq \langle \nabla \mathcal{R}(\hat{\theta}; \lambda) - \nabla \mathcal{R}(\tilde{\theta}; \lambda), \Lambda \rangle \\ &\leq \left(\|\nabla \mathcal{R}(\hat{\theta}; \lambda)\|_\eta^* + \|\nabla \mathcal{R}(\tilde{\theta}; \lambda)\|_\eta^* \right) \|\Lambda\|_\eta \leq 2\bar{\eta} \|\Lambda\|_\eta \end{aligned}$$

where $\bar{\eta} := \max\{\frac{\lambda_1}{\eta_1}, \frac{\lambda_2}{\eta_2}\}$. Since $\|\mathbf{A}\|_\lambda = \|\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}\|_\lambda \leq 2r$ by the constraint of (17),

$$\kappa_2 \|\mathbf{A}\|_2 \leq (\tau_2 + 2\bar{\eta}) \|\mathbf{A}\|_\eta \leq 2r(\tau_2 + 2\bar{\eta}).$$

Therefore, if $2r(\tau_2 + 2\bar{\eta}) \leq 1$ as assumed, we should have $\|\mathbf{A}\|_2 \leq 1$.

Now we focus on the RSC condition in (12):

$$\kappa_1 \|\mathbf{A}\|_2^2 - \tau_1 \|\mathbf{A}\|_\eta^2 \leq \langle \nabla \mathcal{L}(\tilde{\boldsymbol{\theta}}) - \nabla \mathcal{L}(\hat{\boldsymbol{\theta}}), \mathbf{A} \rangle. \quad (37)$$

Let $\bar{\mathcal{L}}(\boldsymbol{\theta}) := \mathcal{L}(\boldsymbol{\theta}) - q_\lambda(\boldsymbol{\theta})$. (Recalling $q_\lambda(\boldsymbol{\theta}) := \|\boldsymbol{\theta}\|_\lambda - \mathcal{R}(\boldsymbol{\theta}; \lambda)$, $\bar{\mathcal{L}}(\boldsymbol{\theta})$ actually can be rewritten as $\mathcal{L}(\boldsymbol{\theta}) + \mathcal{R}(\boldsymbol{\theta}; \lambda) - \|\boldsymbol{\theta}\|_\lambda$). Then, from (37),

$$\kappa_1 \|\mathbf{A}\|_2^2 - \tau_1 \|\mathbf{A}\|_\eta^2 \leq \langle \nabla \bar{\mathcal{L}}(\tilde{\boldsymbol{\theta}}) - \nabla \bar{\mathcal{L}}(\hat{\boldsymbol{\theta}}), \mathbf{A} \rangle + \langle \nabla q_\lambda(\tilde{\boldsymbol{\theta}}) - \nabla q_\lambda(\hat{\boldsymbol{\theta}}), \mathbf{A} \rangle. \quad (38)$$

In order to upper-bound the second term in the RHS of (38), we apply the mean-value theorem:

$$\nabla q_\lambda(\tilde{\boldsymbol{\theta}}) - \nabla q_\lambda(\hat{\boldsymbol{\theta}}) = \nabla^2 q_\lambda(\underline{\boldsymbol{\theta}})(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})$$

where $\underline{\boldsymbol{\theta}}$ is a parameter vector on the line between $\tilde{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$. By the Cauchy-Schwarz inequality

$$|\langle \nabla q_\lambda(\tilde{\boldsymbol{\theta}}) - \nabla q_\lambda(\hat{\boldsymbol{\theta}}), \mathbf{A} \rangle| \leq \|\nabla^2 q_\lambda(\underline{\boldsymbol{\theta}})\|_2 \|\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}\|_2^2 \quad (39)$$

where $\|\cdot\|_2$ denotes the spectral norm of the matrix.

Combining (38) and (39) with the assumption $\|\nabla^2 q_\lambda(\underline{\boldsymbol{\theta}})\|_2 \leq \mu$ (which holds for non-convex regularizers considered in this paper, as shown in Lemma 5) yields

$$(\kappa_1 - \mu) \|\mathbf{A}\|_2^2 - \tau_1 \|\mathbf{A}\|_\eta^2 \leq \langle \nabla \bar{\mathcal{L}}(\tilde{\boldsymbol{\theta}}) - \nabla \bar{\mathcal{L}}(\hat{\boldsymbol{\theta}}), \mathbf{A} \rangle. \quad (40)$$

At the same time, we represent the stationary condition (34) using the notation $q_\lambda(\tilde{\boldsymbol{\theta}})$:

$$0 \leq \langle \nabla \bar{\mathcal{L}}(\tilde{\boldsymbol{\theta}}), \hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}} \rangle + \langle \tilde{\mathbf{z}}, \hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}} \rangle \quad (41)$$

where $\tilde{\mathbf{z}} \in \partial \|\tilde{\boldsymbol{\theta}}\|_\lambda$. From our construction of $\tilde{\mathbf{z}}$ in (33), we have $\langle \nabla \bar{\mathcal{L}}(\hat{\boldsymbol{\theta}}) + \tilde{\mathbf{z}}, \hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}} \rangle = 0$, and when combined with (41) this implies

$$\begin{aligned} 0 &\leq \langle \nabla \bar{\mathcal{L}}(\hat{\boldsymbol{\theta}}) - \nabla \bar{\mathcal{L}}(\tilde{\boldsymbol{\theta}}), \tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}} \rangle - \langle \tilde{\mathbf{z}}, \hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}} \rangle - \langle \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}} \rangle \\ &= \langle \nabla \bar{\mathcal{L}}(\hat{\boldsymbol{\theta}}) - \nabla \bar{\mathcal{L}}(\tilde{\boldsymbol{\theta}}), \tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}} \rangle + \langle \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}} \rangle - \|\hat{\boldsymbol{\theta}}\|_\lambda + \langle \tilde{\mathbf{z}}, \hat{\boldsymbol{\theta}} \rangle - \|\tilde{\boldsymbol{\theta}}\|_\lambda. \end{aligned} \quad (42)$$

With (40), this inequality implies

$$\begin{aligned} (\kappa_1 - \mu) \|\mathbf{A}\|_2^2 - \tau_1 \|\mathbf{A}\|_\eta^2 &\leq \langle \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}} \rangle - \|\hat{\boldsymbol{\theta}}\|_\lambda + \langle \tilde{\mathbf{z}}, \hat{\boldsymbol{\theta}} \rangle - \|\tilde{\boldsymbol{\theta}}\|_\lambda \\ &\leq \langle \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}} \rangle - \|\tilde{\boldsymbol{\theta}}\|_\lambda \end{aligned} \quad (43)$$

where we use the fact $\langle \tilde{\mathbf{z}}, \hat{\boldsymbol{\theta}} \rangle \leq \|\tilde{\mathbf{z}}\|_\lambda^* \|\hat{\boldsymbol{\theta}}\|_\lambda \leq \|\tilde{\boldsymbol{\theta}}\|_\lambda$ by Lemma 3.

Assume for now that

$$\|\mathbf{A}\|_\eta \leq \max\{\lambda_1 \sqrt{s}, \lambda_2 \sqrt{s\mathcal{G}}\} \max\left\{\frac{\eta_1}{\lambda_1}, \frac{\eta_2}{\lambda_2}\right\} \sqrt{2} \left(\frac{4}{\delta} + 2\right) \|\mathbf{A}\|_2 \quad (44)$$

for some $\delta \in (0, 1]$. Then, defining $C := \max\{\lambda_1 \sqrt{s}, \lambda_2 \sqrt{s\mathcal{G}}\} \max\left\{\frac{\eta_1}{\lambda_1}, \frac{\eta_2}{\lambda_2}\right\} \sqrt{2} \left(\frac{4}{\delta} + 2\right)$ so that $\|\mathbf{A}\|_\eta \leq C \|\mathbf{A}\|_2$, we have

$$(\kappa_1 - \mu) \|\mathbf{A}\|_2^2 - \tau_1 \|\mathbf{A}\|_\eta^2 \geq \kappa_1 \|\mathbf{A}\|_2^2 - \tau_1 C^2 \|\mathbf{A}\|_2^2. \quad (45)$$

Therefore, as long as $\kappa_1 > \tau_1 C^2$,

$$0 \leq (\kappa_1 - \mu) \|\mathbf{\Lambda}\|_2^2 - \tau_1 \|\mathbf{\Lambda}\|_\eta^2 \leq \langle \widehat{\mathbf{z}}, \widetilde{\boldsymbol{\theta}} \rangle - \|\widetilde{\boldsymbol{\theta}}\|_\lambda \quad (46)$$

implying $\|\widetilde{\boldsymbol{\theta}}\|_\lambda \leq \langle \widehat{\mathbf{z}}, \widetilde{\boldsymbol{\theta}} \rangle$. Actually, since $\langle \widehat{\mathbf{z}}, \widetilde{\boldsymbol{\theta}} \rangle \leq \|\widehat{\mathbf{z}}\|_\lambda^* \|\widetilde{\boldsymbol{\theta}}\|_\lambda \leq \|\widetilde{\boldsymbol{\theta}}\|_\lambda$ by Hölder's inequality and Lemma 3, it should hold that $\|\widetilde{\boldsymbol{\theta}}\|_\lambda = \langle \widehat{\mathbf{z}}, \widetilde{\boldsymbol{\theta}} \rangle$. As discussed in previous works (Wainwright, 2009; Jalali et al., 2010; Yang et al., 2015; Loh & Wainwright, 2014) using PDW approach, this equality guarantees that for all $j \notin U$, $\widetilde{\boldsymbol{\theta}}_j = 0$ under the strict dual feasibility.

The remaining procedure is to show (44), which is proved in the following lemma.

Lemma 2. Suppose $\|\widehat{\mathbf{z}}_{U^c}\|_\lambda^* \leq 1 - \delta$ for some $\delta \in (0, 1]$, and $2r\tau_1\eta_1 \leq \frac{\delta}{2}\lambda_1$ and $2r\tau_1\eta_2 \leq \frac{\delta}{2}\lambda_2$. Then

$$\|\mathbf{\Lambda}\|_\eta \leq \max\{\lambda_1\sqrt{s}, \lambda_2\sqrt{s_G}\} \max\left\{\frac{\eta_1}{\lambda_1}, \frac{\eta_2}{\lambda_2}\right\} \sqrt{2}\left(\frac{4}{\delta} + 2\right) \|\mathbf{\Lambda}\|_2. \quad (47)$$

From the first inequality of (42), the following inequality can be easily derived in a similar way of constructing (43):

$$(\kappa_1 - \mu) \|\mathbf{\Lambda}\|_2^2 - \tau_1 \|\mathbf{\Lambda}\|_\eta^2 \leq \underbrace{\langle \widehat{\mathbf{z}}, \widetilde{\boldsymbol{\theta}} \rangle - \|\widetilde{\boldsymbol{\theta}}\|_\lambda}_{(I)} + \underbrace{\langle \widehat{\mathbf{z}}, \mathbf{\Lambda} \rangle}_{(II)}. \quad (48)$$

By the same reasoning in (24),

$$(I) \leq \|\widehat{\boldsymbol{\theta}}\|_\lambda - \|\widetilde{\boldsymbol{\theta}}\|_\lambda \leq \lambda_1 \left(\|\mathbf{[\Delta]_{\bar{s}}}\|_1 - \|\mathbf{[\Delta]_{\bar{s}^c}}\|_1 \right) + \lambda_2 \left(\|\mathbf{[\Gamma]_{\bar{b}}}\|_{1,a} - \|\mathbf{[\Gamma]_{\bar{b}^c}}\|_{1,a} \right). \quad (49)$$

The second term in (48) can be upper bounded as follows: if $\|\widehat{\mathbf{z}}\|_\lambda^* \leq 1 - \delta$ as assumed,

$$\begin{aligned} (II) &= \langle \widehat{\mathbf{z}}_1, \mathbf{\Lambda}_U \rangle + \langle \widehat{\mathbf{z}}_2, \mathbf{\Lambda}_{U^c} \rangle \leq \|\widehat{\mathbf{z}}_1\|_\lambda^* \|\mathbf{\Lambda}_U\|_\lambda + \|\widehat{\mathbf{z}}_2\|_\lambda^* \|\mathbf{\Lambda}_{U^c}\|_\lambda \\ &\leq \|\mathbf{\Lambda}_U\|_\lambda + (1 - \delta) \|\mathbf{\Lambda}_{U^c}\|_\lambda \\ &\leq \lambda_1 \left(\|\mathbf{[\Delta]_U}\|_1 + (1 - \delta) \|\mathbf{[\Delta]_{U^c}}\|_1 \right) + \lambda_2 \left(\|\mathbf{[\Gamma]_U}\|_{1,a} + (1 - \delta) \|\mathbf{[\Gamma]_{U^c}}\|_{1,a} \right) \\ &= \lambda_1 \left(\|\mathbf{[\Delta]_{\bar{s}}}\|_1 + (1 - \delta) \|\mathbf{[\Delta]_{\bar{s}^c}}\|_1 \right) + \lambda_2 \left(\|\mathbf{[\Gamma]_{\bar{b}}}\|_{1,a} + (1 - \delta) \|\mathbf{[\Gamma]_{\bar{b}^c}}\|_{1,a} \right) \end{aligned} \quad (50)$$

where the last equality follows the result in Proposition 1.

By (49) and (50),

$$\begin{aligned} -\tau_1 \|\mathbf{\Lambda}\|_\eta^2 &\leq (\kappa_1 - \mu) \|\mathbf{\Lambda}\|_2^2 - \tau_1 \|\mathbf{\Lambda}\|_\eta^2 \\ &\leq \lambda_1 \left(2\|\mathbf{[\Delta]_{\bar{s}}}\|_1 - \delta\|\mathbf{[\Delta]_{\bar{s}^c}}\|_1 \right) + \lambda_2 \left(2\|\mathbf{[\Gamma]_{\bar{b}}}\|_{1,a} - \delta\|\mathbf{[\Gamma]_{\bar{b}^c}}\|_{1,a} \right). \end{aligned} \quad (51)$$

Furthermore,

$$\tau_1 \|\mathbf{\Lambda}\|_\eta^2 \leq 2r\tau_1 \|\mathbf{\Lambda}\|_\eta \leq 2r\tau_1 \left(\eta_1 \|\mathbf{\Delta}\|_1 + \eta_2 \|\mathbf{\Gamma}\|_{1,a} \right). \quad (52)$$

Hence, if $2r\tau_1\eta_1 \leq \frac{\delta}{2}\lambda_1$ and $2r\tau_1\eta_2 \leq \frac{\delta}{2}\lambda_2$ as stated, combining 51 and 52 establishes

$$\begin{aligned} -\frac{\delta}{2} \left(\lambda_1 \|\mathbf{\Delta}\|_1 + \lambda_2 \|\mathbf{\Gamma}\|_{1,a} \right) &\leq -2r\tau_1 \left(\eta_1 \|\mathbf{\Delta}\|_1 + \eta_2 \|\mathbf{\Gamma}\|_{1,a} \right) \\ &\leq \lambda_1 \left(2\|\mathbf{[\Delta]_{\bar{s}}}\|_1 - \delta\|\mathbf{[\Delta]_{\bar{s}^c}}\|_1 \right) + \lambda_2 \left(2\|\mathbf{[\Gamma]_{\bar{b}}}\|_{1,a} - \delta\|\mathbf{[\Gamma]_{\bar{b}^c}}\|_{1,a} \right) \end{aligned}$$

or equivalently

$$\frac{\delta}{2} \left(\lambda_1 \|\mathbf{[\Delta]_{\bar{s}^c}}\|_1 + \lambda_2 \|\mathbf{[\Gamma]_{\bar{b}^c}}\|_{1,a} \right) \leq \left(2 + \frac{\delta}{2} \right) \left(\lambda_1 \|\mathbf{[\Delta]_{\bar{s}}}\|_1 + \lambda_2 \|\mathbf{[\Gamma]_{\bar{b}}}\|_{1,a} \right).$$

Finally, we can have

$$\begin{aligned}
 \|\mathbf{\Delta}\|_{\eta} &\leq \eta_1 \|\mathbf{\Delta}\|_1 + \eta_2 \|\mathbf{\Gamma}\|_{1,a} \leq \bar{\eta}' \left(\lambda_1 \|\mathbf{\Delta}\|_1 + \lambda_2 \|\mathbf{\Gamma}\|_{1,a} \right) \\
 &= \bar{\eta}' \left(\lambda_1 \left\| [\mathbf{\Delta}]_{\bar{s}} \right\|_1 + \lambda_1 \left\| [\mathbf{\Delta}]_{\bar{s}^c} \right\|_1 + \lambda_2 \left\| [\mathbf{\Gamma}]_{\bar{b}} \right\|_{1,a} + \lambda_2 \left\| [\mathbf{\Gamma}]_{\bar{b}^c} \right\|_{1,a} \right) \\
 &\leq \bar{\eta}' \left(\frac{4}{\delta} + 2 \right) \left(\lambda_1 \left\| [\mathbf{\Delta}]_{\bar{s}} \right\|_1 + \lambda_2 \left\| [\mathbf{\Gamma}]_{\bar{b}} \right\|_{1,a} \right) \leq \bar{\eta}' \bar{\lambda} \left(\frac{4}{\delta} + 2 \right) \left(\|\mathbf{\Delta}\|_2 + \|\mathbf{\Gamma}\|_2 \right) \\
 &\leq \sqrt{2} \bar{\eta}' \bar{\lambda} \left(\frac{4}{\delta} + 2 \right) \|\mathbf{\Lambda}\|_2
 \end{aligned} \tag{53}$$

where $\bar{\eta}' := \max\{\frac{\eta_1}{\lambda_1}, \frac{\eta_2}{\lambda_2}\}$, $\bar{\lambda} := \max\{\lambda_1 \sqrt{s}, \lambda_2 \sqrt{s_G}\}$, and we have $\|\mathbf{\Delta}\|_2 + \|\mathbf{\Gamma}\|_2 \leq \sqrt{2} \|\mathbf{\Lambda}\|_2$ (since $\frac{1}{2}(\|\mathbf{\Delta}\|_2 + \|\mathbf{\Gamma}\|_2)^2 \leq \|\mathbf{\Delta}\|_2^2 + \|\mathbf{\Gamma}\|_2^2 \leq \|\mathbf{\Lambda}\|_2^2$) by Proposition 1 in the last inequality.

C.2. Proof of Corollary 3

The statement can be shown to hold by combining the result of Theorem 2 and lemmas in (Loh & Wainwright, 2014)

(a) Uniqueness. The proof of Lemma 2 in (Loh & Wainwright, 2014) shows that under (12)

$$v^\top \nabla^2 \mathcal{L}(\boldsymbol{\theta}) v \geq \kappa_1 \|v\|_2^2 - \tau_1 \|v\|_{\eta}^2, \quad \forall v \in \{v \in \mathbb{R}^p \mid \text{supp}(v) \subseteq U, \|v\|_2 = 1\}.$$

By definition of $\|\cdot\|_{\lambda}$, for any $v \in \mathbb{R}^U$,

$$\|v\|_{\eta} \leq \eta_1 \|v_S\|_1 + \eta_2 \|v_B\|_{1,a} \leq \eta_1 \sqrt{s} \|v_S\|_2 + \eta_2 \sqrt{s_G} \|v_B\|_2 \leq \max\{\eta_1 \sqrt{s}, \eta_2 \sqrt{s_G}\} \|v\|_2.$$

Hence,

$$v^\top \nabla^2 \mathcal{L}(\boldsymbol{\theta}) v \geq \kappa_1 \|v\|_2^2 - \tau_1 \left(\max\{\eta_1 \sqrt{s}, \eta_2 \sqrt{s_G}\} \right)^2 \|v\|_2^2, \quad \forall v \in \{v \in \mathbb{R}^p \mid \text{supp}(v) \subseteq U, \|v\|_2 = 1\}.$$

which implies

$$\left(\nabla^2 \mathcal{L}(\boldsymbol{\theta}) \right)_{UU} \succeq \left(\kappa - \tau_1 \left(\max\{\eta_1 \sqrt{s}, \eta_2 \sqrt{s_G}\} \right)^2 \right) I$$

for all $\boldsymbol{\theta} \in \mathbb{R}^U$. Therefore if $\frac{\kappa_1 - \mu}{2} \geq \tau_1 \left(\max\{\eta_1 \sqrt{s}, \eta_2 \sqrt{s_G}\} \right)^2$, it is guaranteed that

$$\left(\nabla^2 \mathcal{L}(\boldsymbol{\theta}) - \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2 \right)_{UU} \succeq \left(\kappa - \mu - \tau_1 \left(\max\{\eta_1 \sqrt{s}, \eta_2 \sqrt{s_G}\} \right)^2 \right) I \succeq \tau_1 \left(\max\{\eta_1 \sqrt{s}, \eta_2 \sqrt{s_G}\} \right)^2 I$$

and hence $(\mathcal{L}(\boldsymbol{\theta}) - \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2) + (\frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2 + \mathcal{R}(\boldsymbol{\theta}; \boldsymbol{\lambda}))$ is strictly convex over \mathbb{R}^U as the sum of strictly convex and convex function. By Theorem 2, any stationary point $\tilde{\boldsymbol{\theta}}$ of the program (17) should be in the form of $(\tilde{\boldsymbol{\theta}}_U, \mathbf{0}_{U^c})$. Moreover, the restricted program is strictly convex as just shown, $\tilde{\boldsymbol{\theta}}_U$ is unique.

(b) ℓ_{∞} error bound. By the definition of \hat{Q} , $\hat{Q}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) = \nabla \mathcal{L}(\hat{\boldsymbol{\theta}}) - \nabla \mathcal{L}(\boldsymbol{\theta}^*)$. As shown in the proof of Theorem 2 of (Loh & Wainwright, 2014),

$$\begin{aligned}
 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_{\infty} &\leq \left\| \left(\hat{Q}_{UU} \right)^{-1} \left(\nabla \mathcal{L}(\boldsymbol{\theta}^*)_U - \nabla q_{\lambda}(\hat{\boldsymbol{\theta}})_U + \hat{\boldsymbol{z}}_1 \right) \right\|_{\infty} \\
 &\leq \left\| \left(\hat{Q}_{UU} \right)^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}^*)_U \right\|_{\infty} + \min\{\lambda_1, \lambda_2\} \left\| \left(\hat{Q}_{UU} \right)^{-1} \right\|_{\infty}
 \end{aligned} \tag{54}$$

where we use Lemma 6 so that $\|\nabla q_{\lambda}(\hat{\boldsymbol{\theta}})_U - \hat{\boldsymbol{z}}_1\|_{\lambda}^* = \max\{\|\nabla q_{\lambda}(\hat{\boldsymbol{\theta}})_U - \hat{\boldsymbol{z}}_1\|_{\infty} / \lambda_1, \|\nabla q_{\lambda}(\hat{\boldsymbol{\theta}})_U - \hat{\boldsymbol{z}}_1\|_{\infty, a^*} / \lambda_2\} \leq 1$ implying $\|\nabla q_{\lambda}(\hat{\boldsymbol{\theta}})_U - \hat{\boldsymbol{z}}_1\|_{\infty} \leq \min\{\lambda_1, \lambda_2\}$.

(c) ℓ_∞ error bound for (μ, γ) -amenable regularizers. By the assumption on θ_{\min}^* , it is guaranteed $|\widehat{\theta}_j| \geq |\theta_{\min}^*| - \|\widehat{\theta} - \theta^*\|_\infty \geq 2 \max\{\lambda_1, \lambda_2\}\gamma$ for all $j \in U$. Then, at least either $\widehat{\alpha}_j$ or $\widehat{\beta}_j$ is larger than $\gamma \max\{\lambda_1, \lambda_2\}$, and, as a result,

$$\nabla q_\lambda(\widehat{\theta})_U = \widehat{z}_1 \quad (55)$$

and (54) reduces to $\|\widehat{\theta} - \theta^*\|_\infty \leq \left\| (\widehat{Q}_{UU})^{-1} \nabla \mathcal{L}(\theta^*)_U \right\|_\infty$.

In order to show the dual feasibility, we combine the zero-subgradient condition (33) and the definition of \widehat{Q} as described in (Loh & Wainwright, 2014):

$$\begin{bmatrix} \widehat{Q}_{UU} & \widehat{Q}_{UU^c} \\ \widehat{Q}_{U^cU} & \widehat{Q}_{U^cU^c} \end{bmatrix} \begin{bmatrix} \widehat{\theta}_U - \theta_U^* \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} \nabla \mathcal{L}(\theta^*)_U - \nabla q_\lambda(\widehat{\theta})_U \\ \nabla \mathcal{L}(\theta^*)_{U^c} - \nabla q_\lambda(\widehat{\theta})_{U^c} \end{bmatrix} + \begin{bmatrix} \widehat{z}_1 \\ \widehat{z}_2 \end{bmatrix} = \mathbf{0},$$

and rearranging it for \widehat{z}_2 with (55) and the fact $\nabla q_\lambda(\widehat{\theta})_{U^c} = 0$, yields

$$\widehat{z}_2 = -\nabla \mathcal{L}(\theta^*)_{U^c} + \widehat{Q}_{U^cU} (\widehat{Q}_{UU})^{-1} \nabla \mathcal{L}(\theta^*)_U \quad (56)$$

Therefore, if $\|\nabla \mathcal{L}(\theta^*)_{U^c}\|_\lambda^* \leq \frac{1-\delta}{2}$ and $\|\widehat{Q}_{U^cU} (\widehat{Q}_{UU})^{-1} \nabla \mathcal{L}(\theta^*)_U\|_\lambda^* \leq \frac{1-\delta}{2}$ additionally hold, the strict dual feasibility holds:

$$\|\widehat{z}_2\|_\lambda^* \leq \|\nabla \mathcal{L}(\theta^*)_{U^c}\|_\lambda^* + \|\widehat{Q}_{U^cU} (\widehat{Q}_{UU})^{-1} \nabla \mathcal{L}(\theta^*)_U\|_\lambda^* \leq \frac{1-\delta}{2} + \frac{1-\delta}{2} \leq 1 - \delta.$$

C.3. Proof of Corollary 4

For the problem, if $(j, k) \in U^c$, then $[\widehat{z}_2]_j^{(k)}$ in (56), can be written as

$$\frac{1}{n} \left\langle X_j^{(k)}, I - \frac{1}{n} X_{U_k}^{(k)} \left(\frac{1}{n} \langle X_{U_k}^{(k)}, X_{U_k}^{(k)} \rangle \right)^{-1} \left(X_{U_k}^{(k)} \right)^\top \right\rangle w^{(k)}. \quad (57)$$

In order to show the strict dual feasibility $\|\widehat{z}_2\|_\lambda^* \leq 1 - \delta$, we need to show both $\max_{(j,k) \in U^c} |[\widehat{z}_2]_j^{(k)}| \leq (1 - \delta)\lambda_1$ and $\max_j \sum_{(j,k) \in U^c} |[\widehat{z}_2]_j^{(k)}| \leq (1 - \delta)\lambda_2$.

Setting $t = 4\sigma \sqrt{\frac{\log(pm)}{n}}$ for (61) of Lemma 9 yields $\max_{(j,k) \in U^c} |[\widehat{z}_2]_j^{(k)}| \leq 4\sigma \sqrt{\frac{\log(pm)}{n}}$ with probability at least $1 - 2 \exp(-3 \log(pm))$. Similarly, setting $t = 4\sigma \sqrt{\frac{\log p + m \log 2}{n}}$ for (62) yields $\max_j \sum_{(j,k) \in U^c} |[\widehat{z}_2]_j^{(k)}| \leq 4\sigma \sqrt{\frac{\log p + m \log 2}{n}}$ with probability at least $1 - 2 \exp(-3(\log p + m \log 2))$. In addition, by the similar reasoning in the proof of Proposition 2, we can easily show that the RSC condition holds w.h.p. for η such that $\max\{\lambda_1/\eta_1, \lambda_2/\eta_2\}$ is some constant depending on σ (the only difference is $a = \infty$ in this example, but $\mathbb{E}(\|\varepsilon\|_\infty)$ for a standard normal vector $\varepsilon \sim N(0, I_{p \times p})$ scales as the same rates of λ_2 by (62) of Lemma 9).

Therefore, the strict dual feasibility holds for the selection of parameters, and the support set recovery is guaranteed w.h.p. ℓ_∞ bound is also trivially derived from the combination of (63) of Lemma 9 and Corollary 3 (we can compute the upper bound of $\left\| (\widehat{Q}_{UU})^{-1} \nabla \mathcal{L}(\theta^*)_U \right\|_\infty$ by setting $t = \sigma \sqrt{\frac{100 \log(pm)}{nC_{\min}}}$ in (63)).

D. Useful Lemmas for Proofs

Lemma 3. At any θ , $\|\nabla \|\theta\|_\lambda\|_\eta^* \leq \max\left\{\frac{\lambda_1}{\eta_1}, \frac{\lambda_2}{\eta_2}\right\}$.

Proof. For any fixed θ , let s be the sub-gradient of $\|\theta\|_\lambda$ at θ , $\nabla \|\theta\|_\lambda$, in \mathbb{R}^p . By definition of the sub-gradient, s satisfies

$$\|\theta + v\|_\lambda \geq \|\theta\|_\lambda + \langle s, v \rangle \quad \text{for all } v \in \mathbb{R}^p.$$

Taking this inequality the supremum over all v such that $\|v\|_\lambda = 1$, we have

$$\sup_{\|v\|_\lambda=1} \|\theta + v\|_\lambda \geq \|\theta\|_\lambda + \sup_{\|v\|_\lambda=1} \langle s, v \rangle. \quad (58)$$

By rearranging (58),

$$\|s\|_\lambda^* \stackrel{(i)}{=} \sup_{\|v\|_\lambda=1} \langle s, v \rangle \leq \sup_{\|v\|_\lambda=1} \|\theta + v\|_\lambda - \|\theta\|_\lambda \leq \sup_{\|v\|_\lambda=1} \|v\|_\lambda \stackrel{(ii)}{}$$

where (i) is the definition of dual norm of $\|\cdot\|_\lambda$, and (ii) follows the triangular inequality of the norm $\|\cdot\|_\lambda$. Since $\sup_{\|v\|_\lambda=1} \|v\|_\lambda \leq 1$ by definitions, we obtain $\|s\|_\lambda^* = \max \left\{ \frac{\|\nabla\|\theta\|_\lambda\|_\infty}{\lambda_1}, \frac{\|\nabla\|\theta\|_\lambda\|_{\infty,a}}{\lambda_2} \right\} \leq 1$, finally implying

$$\|\nabla\|\theta\|_\lambda\|_\eta^* \leq \max \left\{ \frac{\|\nabla\|\theta\|_\lambda\|_\infty}{\eta_1}, \frac{\|\nabla\|\theta\|_\lambda\|_{\infty,a}}{\eta_2} \right\} \leq \max \left\{ \frac{\lambda_1}{\eta_1}, \frac{\lambda_2}{\eta_2} \right\}.$$

□

Lemma 4. At any θ , $\|\nabla\mathcal{R}(\theta; \lambda)\|_\eta^* \leq \max \left\{ \frac{\lambda_1}{\eta_1}, \frac{\lambda_2}{\eta_2} \right\}$.

Proof. At any θ , if we compute the derivative, $[\nabla\mathcal{R}(\theta; \lambda)]_j$ is upper bounded by $\max\{\partial\rho_{\lambda_1}([\theta]_j), \partial\phi_{\lambda_2,a}[\theta]_{g_j}\}$ where g_j is the group of indices that j belongs to. As shown in Lemma 8 of (Loh & Wainwright, 2014), $\partial\rho_{\lambda_1}([\theta]_j) \leq \partial\|[\theta]_j\|_1$. Similarly, by definition of ϕ , we have $\partial\phi_{\lambda_2,a}([\theta]_{g_j}) \leq \partial\|[\theta]_{g_j}\|_{1,a}$. Therefore, for every index j , $|\nabla\mathcal{R}(\theta; \lambda)_j| \leq |\nabla\|\theta\|_\lambda|_j$, and $\|\nabla\mathcal{R}(\theta; \lambda)\|_\eta^* \leq \|\nabla\|\theta\|_\lambda\|_\eta^*$. The final result comes by Lemma 3. □

Lemma 5. Consider the non-convex penalty functions in (C1). Then, for any θ ,

$$\|\nabla^2 q_\lambda(\theta)\|_2 \leq \mu. \quad (59)$$

Proof. For notational simplicity, we define the function $F: \mathbb{R}^p \rightarrow \mathbb{R}^p$ as $F(\theta; \lambda) := \nabla R(\theta; \lambda)$, hence the i -th coordinate of F is

$$F_i := \frac{\partial R}{\partial \alpha} \cdot \frac{\partial \alpha}{\partial \theta_i} + \frac{\partial R}{\partial \beta} \cdot \frac{\partial \beta}{\partial \theta_i}$$

where we suppress the dependency on θ and λ of $F(\theta; \lambda)$ and $R(\theta; \lambda)$ for compact presentation. By applying another chain rule to compute ∇F , we immediately obtain

$$\begin{aligned} \frac{\partial F_i}{\partial \theta_j} &= \sum_{k,l} \frac{\partial^2 R}{\partial \alpha_k \partial \alpha_l} \cdot \frac{\partial \alpha_k}{\partial \theta_i} \frac{\partial \alpha_l}{\partial \theta_j} + \sum_k \frac{\partial R}{\partial \alpha_k} \cdot \frac{\partial^2 \alpha_k}{\partial \theta_i \partial \theta_j} \\ &\quad + \sum_{k,l} \frac{\partial^2 R}{\partial \beta_k \partial \beta_l} \cdot \frac{\partial \beta_k}{\partial \theta_i} \frac{\partial \beta_l}{\partial \theta_j} + \sum_k \frac{\partial R}{\partial \beta_k} \cdot \frac{\partial^2 \beta_k}{\partial \theta_i \partial \theta_j}. \end{aligned}$$

Furthermore, given penalty functions defined in (C1), it can be shown that

1. $\frac{\partial^2 R}{\partial \alpha_k \partial \alpha_l} = \frac{\partial^2 R}{\partial \beta_k \partial \beta_l} = 0$ for all $k \neq l$, and $|\frac{\partial^2 R}{\partial \alpha_k^2}| \leq \mu$ for all k .
2. If $\frac{\partial R}{\partial \alpha_k} = 0$ for some k , then $\frac{\partial \alpha_k}{\partial \theta_i}$ (and hence $\frac{\partial \beta_k}{\partial \theta_i}$) can have any value, but if $\frac{\partial R}{\partial \alpha_k} \neq 0$ then $\frac{\partial \alpha_k}{\partial \theta_i} = 0$ for $i \neq k$, and moreover $\frac{\partial \alpha_k}{\partial \theta_k}$ can be either 1 ($\frac{\partial \beta_k}{\partial \theta_k}$ should be 0 in this case) or 0 ($\frac{\partial \alpha_k}{\partial \theta_k}$ should be 1).
3. If $\frac{\partial R}{\partial \alpha_k} \neq 0$, then $\frac{\partial^2 \alpha_k}{\partial \theta_i \partial \theta_j} = 0$.

Based on these facts, ∇F can be reduce to the diagonal matrix whose the maximum absolute element is upper bounded by μ , implying (59) (since $\nabla^2\|\theta\|_\lambda$ can be shown to be zero matrix by the same reason above). □

Lemma 6. $\|\nabla\|\theta\|_\lambda - \nabla\mathcal{R}(\theta; \lambda)\|_\eta^* \leq \max\left\{\frac{\lambda_1}{\eta_1}, \frac{\lambda_2}{\eta_2}\right\}$.

Proof. By definitions of $\|\theta\|_\lambda$ and $R(\theta; \lambda)$, every pair of coordinates $(\nabla\|\theta\|_\lambda, \nabla\mathcal{R}(\theta; \lambda))$ has the same signs, and moreover based on Lemma 3 and 4, the statement holds. \square

Lemma 7 (Lemma 8 of (Loh & Wainwright, 2015)). *If \mathcal{L} is convex and (12) holds, then (13) holds as well, with $\kappa_2 = \kappa_1$ and $\tau_2 = 2r\tau_1$.*

Proof. As shown in the proof of Lemma 8 of (Loh & Wainwright, 2015), for any $\theta_2 \in \mathbb{R}^p$ such that $\|\theta_2\|_2 \geq 1$,

$$\begin{aligned} \langle \nabla\mathcal{L}(\theta_1 + \theta_2) - \nabla\mathcal{L}(\theta_1), \theta_2 \rangle &\geq \|\theta_2\|_2 \left(\kappa_1 - \tau_1 \frac{\|\theta_2\|_\eta^2}{\|\theta_2\|_2^2} \right) \\ &\geq \|\theta_2\|_2 \left(\kappa_1 - 2r\tau_1 \frac{\|\theta_2\|_\eta}{\|\theta_2\|_2} \right) \\ &\geq \|\theta_2\|_2 \left(\kappa_1 - 2r\tau_1 \frac{\|\theta_2\|_\eta}{\|\theta_2\|_2} \right) \\ &\geq \kappa_1 \|\theta_2\|_2 - 2r\tau_1 \|\theta_2\|_\eta. \end{aligned}$$

\square

Lemma 8 (Lemma 9 of (Loh & Wainwright, 2015)). *If (12) holds globally for all $\theta_2 \in \mathbb{R}^p$, then (13) holds as well, with $\kappa_2 = \kappa_1$ and $\tau_2 = 2r\tau_1$.*

Proof. If $\|\theta_2\|_2 \geq 1$,

$$\begin{aligned} \langle \nabla\mathcal{L}(\theta_1 + \theta_2) - \nabla\mathcal{L}(\theta_1), \theta_2 \rangle &\geq \kappa_1 \|\theta_2\|_2^2 - \tau_1 \|\theta_2\|_\eta^2 \\ &\geq \kappa_1 \|\theta_2\|_2 - 2r\tau_1 \|\theta_2\|_\eta. \end{aligned}$$

\square

Lemma 9 (From the proof of Lemma 9 of (Jalali et al., 2010)). *Let $W_j^{(k)}$ be*

$$\frac{1}{n} \left\langle X_j^{(k)}, I - \frac{1}{n} X_{U_k}^{(k)} \left(\frac{1}{n} \langle X_{U_k}^{(k)}, X_{U_k}^{(k)} \rangle \right)^{-1} \left(X_{U_k}^{(k)} \right)^\top \right\rangle w^{(k)} \quad (60)$$

where U_k^c denotes the support of k -th column of Θ^* . For all $t \geq 0$,

$$\mathbb{P} \left[\max_{(j,k) \in U_k^c} |W_j^{(k)}| \geq t \right] \leq 2 \exp \left(- \frac{t^2 n}{4\sigma^2} + \log(pm) \right), \quad (61)$$

and

$$\mathbb{P} \left[\max_j \sum_{(j,k) \in U_k^c} |W_j^{(k)}| \geq t \right] \leq 2 \exp \left(- \frac{t^2 n}{4\sigma^2} + m \log 2 + \log p \right). \quad (62)$$

In addition, for all $t \geq 0$,

$$\mathbb{P} \left[\max_{k=1,2,\dots,m} \left\| \left(\frac{1}{n} \langle X_{U_k}^{(k)}, X_{U_k}^{(k)} \rangle \right)^{-1} \frac{1}{n} \left(X_{U_k}^{(k)} \right)^\top w^{(k)} \right\|_\infty \geq t \right] \leq 2 \exp \left(- \frac{t^2 n C_{\min}}{50\sigma^2} + \log(pm) \right) \quad (63)$$

where $C_{\min} := \min_{k=1,2,\dots,m} \lambda_{\min} \left(\Sigma_{U_k}^{(k)} \right) > 0$.

Note that (Jalali et al., 2010) originally consider only $j \in \cap_{k=1}^m U_k^c$, but (61) and (62) hold by the same reasoning.

E. Additional details on the simulation experiments

Figure 4 provides an example of how the computing time of convex and non-convex dirty models vary with the regularization parameters λ_1, λ_2 . The iteration number needed for each method under the same convergence tolerance exhibit a similar profile. In many cases non-convex dirty models converge faster and take less computing time.

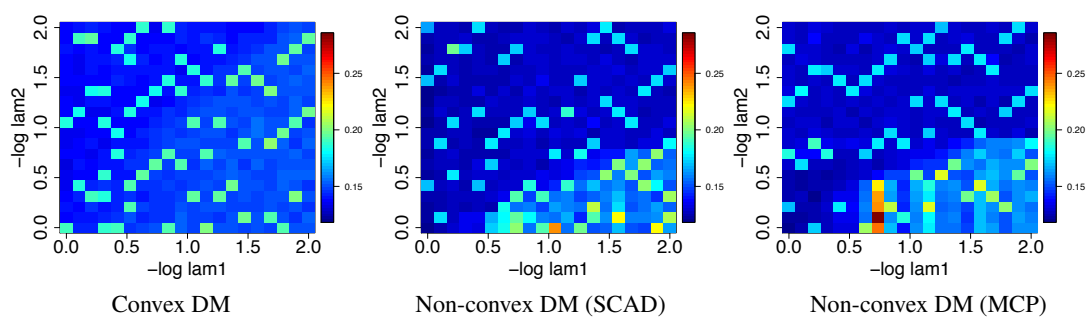


Figure 4. Timing for comparison methods for varying penalty parameters.