

## 8. Appendix

We prove Theorem 1 in this section. Since the proof is technical and lengthy, for the clarity of presentation, we organize the proof as follows. To begin with, in Section 8.1, we review two standard concentration inequalities, the Chernoff inequality and the Hoeffding inequality, which will be used to prove some technical lemmas. We then present and prove these technical lemmas in Section 8.2. These technical lemmas are subsequently used to validate some auxiliary results, which are presented in Section 8.3. Finally, we prove Theorem 1 based on these auxiliary results.

### 8.1. Concentration Inequalities

**Lemma 1** (Hoeffding Inequality). *Let  $X_1, X_2, \dots, X_n$  be  $n$  i.i.d. random variables drawn from the distribution  $\mathcal{D}$ , with  $0 \leq X_i \leq a$ ,  $\forall i \in \{1, 2, \dots, n\}$ . Let  $\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i$ . Then, for any  $t > 0$ ,*

$$P(|\bar{X} - \mathbb{E}[\bar{X}]| \geq t) \leq 2 \exp\left(-\frac{2nt^2}{a^2}\right).$$

**Lemma 2** (Chernoff Inequality). *Let  $X_1, X_2, \dots, X_n$  be  $n$  random variables and let  $X := \sum_{i=1}^n X_i$ . Then, for any  $t > 0$ ,*

$$P(X \geq \epsilon) \leq \exp(-t\epsilon) \mathbb{E} \left[ \exp \left( \sum_{i=1}^n tX_i \right) \right]. \quad (14)$$

Furthermore, if  $X_i$ 's are independent, then

$$P(X \geq \epsilon) \leq \min_{t>0} \exp(-t\epsilon) \prod_{i=1}^n \mathbb{E} [\exp(tX_i)]. \quad (15)$$

### 8.2. Technical Lemmas

We use  $\|\cdot\|_{\max}$  to represent the max norm of a matrix, which is equal to the maximum of the absolute value of all the elements in the matrix.

**Lemma 3.** *Let  $\mathbb{X}$  be given. Suppose that  $0 < \max_{i,i' \in \{1,2,\dots,n\}} \|\mathbf{x}_i \mathbf{x}_{i'}^\top\|_{\max} < \epsilon^2$ . Then,*

$$P \left( \max_{j \neq j', j \neq j', j, j' \in \{1,2,\dots,p\}} |\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]| \geq \epsilon^2 \sqrt{\frac{\log p}{n}} \right) \leq 2 \exp(-2 \log p).$$

*Proof.* Since  $0 < \max_{i,i' \in \{1,2,\dots,n\}} \|\mathbf{x}_i \mathbf{x}_{i'}^\top\|_{\max} < \epsilon^2$ , we let  $a = \epsilon^2$  and  $t = \epsilon^2 \sqrt{\frac{\log p}{n}}$  in Lemma 1 to yield the result.  $\square$

**Lemma 4.** *Let  $\mathbb{X}$  be given. Suppose that  $0 < \max_{i \in \{1,2,\dots,n\}} \|\mathbf{x}_i\|_{\infty} < \epsilon$ . Then,*

$$P \left( \max_{j \in \{1,2,\dots,p\}} |\mathbb{E}_{\mathbb{X}}[X_j] - \mathbb{E}[X_j]| \geq \epsilon \sqrt{\frac{\log p}{n}} \right) \leq 2 \exp(-2 \log p).$$

*Proof.* Since  $0 < \max_{i \in \{1,2,\dots,n\}} \|\mathbf{x}_i\|_{\infty} < \epsilon$ , we let  $a = \epsilon$  and  $t = \epsilon \sqrt{\frac{\log p}{n}}$  in Lemma 1 to yield the result.  $\square$

**Lemma 5.** *Let  $\mathbb{X}$  be given. Suppose that  $0 < \max_{i \in \{1,2,\dots,n\}} \|\mathbf{x}_i\|_{\infty} < \epsilon$ . Then,*

$$P \left( \max_{j,j' \in \{1,2,\dots,p\}} |\mathbb{E}_{\mathbb{X}}[\mathbb{E}[X_j X_{j'} | \mathbf{X}_{-j}] - \mathbb{E}[\mathbb{E}[X_j X_{j'} | \mathbf{X}_{-j}]]| \geq C_3 \epsilon \sqrt{\frac{\log p}{n}} \right) \leq 2 \exp(-2 \log p).$$

*Proof.* Since  $0 < \max_{i \in \{1,2,\dots,n\}} \|\mathbf{x}_i\|_{\infty} < \epsilon$  and  $\mathbb{E}[X_j | \mathbf{x}_{i,-j}] \leq C_3$  by Assumption 2, we have that  $0 < \mathbb{E}[X_j X_{j'} | \mathbf{x}_{i,-j}] \leq C_3 \epsilon$ . Therefore, we let  $a = C_3 \epsilon$  and  $t = C_3 \epsilon \sqrt{\frac{\log p}{n}}$  in Lemma 1 to yield the result.  $\square$

**Remark**

The subtlety of the definitions of  $C_3$  and  $C_4$  in Assumption 2, as well as the notion of  $\epsilon$  in Lemma 3, Lemma 4, and Lemma 5 should be noted. Formally, the  $n$  data points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  in  $\mathbb{X}$  can be viewed as assignments to the corresponding random variables  $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(n)}$  following the PSQR parameterized by  $\Theta^*$ . In Assumption 2, we are interested in a set  $\mathcal{X} \subseteq \mathbb{N}^p$ , such that  $\forall i \in \{1, 2, \dots, n\}$  and  $\forall j \in \{1, 2, \dots, p\}$ ,

$$\max_{\mathbf{X}^{(i)} \in \mathcal{X}} \mathbb{E} [X_j | \mathbf{X}_{-j}^{(i)}] \leq C_3 \quad \text{and} \quad \max_{\mathbf{X}^{(i)} \in \mathcal{X}} |\lambda_{ij}^* - \mathbb{E} [X_j | \mathbf{X}_{-j}^{(i)}]| \leq C_4.$$

In Lemma 3, Lemma 4, and Lemma 5, we are interested in a set  $\mathcal{X} \subseteq \mathbb{N}^p$ , such that  $\forall i, i' \in \{1, 2, \dots, n\}$ , where  $i \neq i'$ ,

$$0 < \max_{\mathbf{X}^{(i)}, \mathbf{X}^{(i')} \in \mathcal{X}} \|\mathbf{X}^{(i)} \mathbf{X}^{(i')\top}\|_{\max} < \epsilon^2 \quad \text{and} \quad 0 < \max_{\mathbf{X}^{(i)} \in \mathcal{X}} \|\mathbf{X}^{(i)}\|_{\infty} < \epsilon.$$

Also, implicitly, we have that  $\mathbf{x}_i \in \mathcal{X}$ ,  $\forall i \in \{1, 2, \dots, n\}$ .

**Lemma 6.** Let  $\mathbb{X}$  be given. Then,

$$\mathbb{P} \left( \max_{j \in \{1, 2, \dots, p\}} |\mathbb{E}_{\mathbb{X}}[\mathbb{E}[X_j | \mathbf{X}_{-j}]] - \mathbb{E}[\mathbb{E}[X_j | \mathbf{X}_{-j}]]| \geq C_3 \sqrt{\frac{\log p}{n}} \right) \leq 2 \exp(-2 \log p).$$

*Proof.* Since  $\mathbb{E}[X_j | \mathbf{x}_{i,-j}] \leq C_3$  by Assumption 2, we have that  $0 < \mathbb{E}[X_j | \mathbf{x}_{i,-j}] \leq C_3$ . Therefore, we let  $a = C_3$  and  $t = C_3 \sqrt{\frac{\log p}{n}}$  in Lemma 1 to yield the result.  $\square$

**Lemma 7.** Let  $\mathbf{X}$  be a random vector drawn from a PSQR distribution parameterized by  $\Theta^*$ . Suppose that  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}^\top$  is the set of  $n$  i.i.d. samples of  $\mathbf{X}$ . Given  $j \in \{1, 2, \dots, p\}$ ,  $\epsilon_1 := 3 \log p + \log n$ , and  $\epsilon_2 := C_1 + \sqrt{\frac{2 \log p}{n}}$ ,

$$\mathbb{P}(X_j \geq \epsilon_1) \leq \exp(C_1 + C_2/2 - \epsilon_1), \quad \text{and} \quad \mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n x_{ij} \geq \epsilon_2 \right) \leq \exp \left[ -\frac{n(\epsilon_2 - C_1)^2}{2C_2} \right].$$

*Proof.* We start with proving the first inequality. To this end, consider the following equation due to Taylor expansion:

$$\begin{aligned} \log \mathbb{E} [\exp(X_j)] &= B(\Theta^*, \mathbf{0} + \mathbf{e}_j) - B(\Theta^*, \mathbf{0}) = \nabla^\top B(\Theta^*, \mathbf{0}) \mathbf{e}_j + \frac{1}{2} \mathbf{e}_j^\top \nabla^2 B(\Theta^*, k \mathbf{e}_j) \mathbf{e}_j \\ &= \mathbb{E}[X_j] + \frac{1}{2} \frac{\partial^2}{\partial b_j^2} B(\Theta^*, \mathbf{0} + k \mathbf{e}_j) \leq C_1 + C_2/2, \end{aligned} \tag{16}$$

where  $k \in [0, 1]$ ,  $\mathbf{e}_j$  is a vector whose  $j^{\text{th}}$  component is one and zeros elsewhere, and the last inequality is due to Assumption 1. Then, let  $t = 1$  and  $X = X_j$  in Lemma 2,

$$\mathbb{P}(X_j \geq \epsilon_1) = \exp(-\epsilon_1) \mathbb{E} [\exp(X_j)] \leq \exp(C_1 + C_2/2 - \epsilon_1).$$

Now, we prove the second bound. For any  $0 < a < 1$  and some  $k \in [0, 1]$ , with Taylor expansion,

$$\begin{aligned} \log \mathbb{E} [\exp(aX_i)] &= B(\Theta^*, \mathbf{0} + a \mathbf{e}_j) - B(\Theta^*, \mathbf{0}) = a \nabla^\top B(\Theta^*, \mathbf{0}) \mathbf{e}_j + \frac{a^2}{2} \mathbf{e}_j^\top \nabla^2 B(\Theta^*, \mathbf{0} + a k \mathbf{e}_j) \mathbf{e}_j \\ &= a \mathbb{E}(X_j) + \frac{a^2}{2} \frac{\partial^2}{\partial b_j^2} B(\Theta^*, \mathbf{0} + a k \mathbf{e}_j) \leq a C_1 + \frac{a^2}{2} C_2, \end{aligned} \tag{17}$$

where the last inequality is due to Assumption 1. Then, following the proof technique above, we have

$$\mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n X_i \geq \epsilon_2 \right) = \mathbb{P} \left( \sum_{i=1}^n X_i \geq n \epsilon_2 \right) \leq \min_{t>0} \exp(-t n \epsilon_2) \prod_{i=1}^n \mathbb{E} [\exp(t X_i)]$$

$$\begin{aligned}
 &\leq \min_{t>0} \exp(-tn\epsilon_2) \prod_{i=1}^n \exp\left(C_1 t + \frac{C_2}{2} t^2\right) \\
 &= \min_{t>0} \exp\left[(C_1 - \epsilon_2)nt + \frac{nC_2}{2} t^2\right] \leq \exp\left[-\frac{n(\epsilon_2 - C_1)^2}{2C_2}\right],
 \end{aligned}$$

where the minimum is obtained when  $t = \frac{\epsilon_2 - C_1}{C_2}$ , and we have used the fact that  $\epsilon_2 > C_1$ .  $\square$

### 8.3. Auxiliary Results

**Lemma 8.** *Let  $r := 4C_5\lambda$ . Then with probability of at least  $1 - ((\exp(C_1 + C_2/2) + 8)p^{-2} + p^{-1/C_2})$ , the following two inequalities simultaneously hold:*

$$\|\nabla F(\Theta^*)\|_\infty \leq 2[C_3(3\log p + \log n) + (3\log p + \log n)^2] \sqrt{\frac{\log p}{n}} + 2C_4 \left(C_1 + \sqrt{\frac{2\log p}{n}}\right), \quad (18)$$

$$\|\tilde{\Theta}_S - \Theta_S^*\|_\infty \leq r. \quad (19)$$

*Proof.* We prove (18) and (19) in turn.

#### PROOF OF (18)

To begin with, we prove (18). By the definition of  $F$  in (13), for  $j < j'$ , the derivative of  $F(\Theta^*)$  is:

$$\frac{\partial F(\Theta^*)}{\partial \theta_{jj'}} = \frac{1}{n} \sum_{i=1}^n [-x_{ij'}x_{ij} + \lambda_{ij}^*x_{ij'} - x_{ij}x_{ij'} + \lambda_{ij'}^*x_{ij}] = -2\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] + \frac{1}{n} \sum_{i=1}^n \lambda_{ij}^*x_{ij'} + \frac{1}{n} \sum_{i=1}^n \lambda_{ij'}^*x_{ij}. \quad (20)$$

and

$$\frac{\partial}{\partial \theta_{jj}} F(\Theta^*) = \frac{1}{n} \sum_{i=1}^n [-x_{ij} + \lambda_{ij}^*] = -\mathbb{E}_{\mathbb{X}}[X_j] + \frac{1}{n} \sum_{i=1}^n \lambda_{ij}^*, \quad (21)$$

where  $\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] := \frac{1}{n} \sum_{i=1}^n x_{ij}x_{ij'}$  and  $\mathbb{E}_{\mathbb{X}}[X_j] := \frac{1}{n} \sum_{i=1}^n x_{ij}$  are the expectations of  $X_j X_{j'}$  and  $X_j$  over the empirical distribution given by the dataset  $\mathbb{X}$ .

Then, by defining  $\mathbb{E}[X_j X_{j'}]$  as the expectation of the multiplication of two components of an multivariate square root Poisson random vector whose distribution is parameterized by  $\Theta^*$ , and by Assumption 2, (20) can be controlled via

$$\begin{aligned}
 \left| \frac{\partial}{\partial \theta_{jj'}} F(\Theta^*) \right| &= \left| \frac{1}{n} \sum_{i=1}^n \lambda_{ij}^*x_{ij'} - \mathbb{E}[X_j X_{j'}] + \frac{1}{n} \sum_{i=1}^n \lambda_{ij'}^*x_{ij} - \mathbb{E}[X_j X_{j'}] + 2\mathbb{E}[X_j X_{j'}] - 2\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] \right| \\
 &\leq \left| \frac{1}{n} \sum_{i=1}^n \lambda_{ij}^*x_{ij'} - \mathbb{E}[X_j X_{j'}] \right| + \left| \frac{1}{n} \sum_{i=1}^n \lambda_{ij'}^*x_{ij} - \mathbb{E}[X_j X_{j'}] \right| + 2|\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]| \\
 &= \left| \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}] + \lambda_{ij}^* - \mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}]) x_{ij'} - \mathbb{E}[X_j X_{j'}] \right| \\
 &\quad + \left| \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[X_{j'} | \mathbf{X}_{-j'} = \mathbf{x}_{i,-j'}] + \lambda_{ij'}^* - \mathbb{E}[X_{j'} | \mathbf{X}_{-j'} = \mathbf{x}_{i,-j'}]) x_{ij} - \mathbb{E}[X_j X_{j'}] \right| \\
 &\quad + 2|\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]| \\
 &\leq \left| \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}]) x_{ij'} - \mathbb{E}[X_j X_{j'}] \right| + \frac{1}{n} \sum_{i=1}^n |\lambda_{ij}^* - \mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}]| x_{ij'} \\
 &\quad + \left| \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[X_{j'} | \mathbf{X}_{-j'} = \mathbf{x}_{i,-j'}]) x_{ij} - \mathbb{E}[X_j X_{j'}] \right| + \frac{1}{n} \sum_{i=1}^n |\lambda_{ij'}^* - \mathbb{E}[X_{j'} | \mathbf{X}_{-j'} = \mathbf{x}_{i,-j'}]| x_{ij} \\
 &\quad + 2|\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]|
 \end{aligned}$$

$$\begin{aligned}
 &\leq \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}] x_{ij'} - \mathbb{E}[X_j X_{j'}] \right| + \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_{j'} | \mathbf{X}_{-j'} = \mathbf{x}_{i,-j'}] x_{ij} - \mathbb{E}[X_j X_{j'}] \right| \\
 &+ 2|\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]| + C_4(\mathbb{E}_{\mathbb{X}}[X_j] + \mathbb{E}_{\mathbb{X}}[X_{j'}]) \\
 &= \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_j X_{j'} | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}] - \mathbb{E}[X_j X_{j'}] \right| + \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_j X_{j'} | \mathbf{X}_{-j'} = \mathbf{x}_{i,-j'}] - \mathbb{E}[X_j X_{j'}] \right| \\
 &+ 2|\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]| + C_4(\mathbb{E}_{\mathbb{X}}[X_j] + \mathbb{E}_{\mathbb{X}}[X_{j'}]) \\
 &= 2|\mathbb{E}_{\mathbb{X}}[\mathbb{E}[X_j X_{j'} | \mathbf{X}_{-j}] - \mathbb{E}[X_j X_{j'}]| + 2|\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]| + C_4(\mathbb{E}_{\mathbb{X}}[X_j] + \mathbb{E}_{\mathbb{X}}[X_{j'}]) \\
 &= 2|\mathbb{E}_{\mathbb{X}}[\mathbb{E}[X_j X_{j'} | \mathbf{X}_{-j}] - \mathbb{E}[\mathbb{E}[X_j X_{j'} | \mathbf{X}_{-j}]]| + 2|\mathbb{E}_{\mathbb{X}}[X_j X_{j'}] - \mathbb{E}[X_j X_{j'}]| + C_4(\mathbb{E}_{\mathbb{X}}[X_j] + \mathbb{E}_{\mathbb{X}}[X_{j'}]),
 \end{aligned}$$

where we have used the law of total expectation in the last equality.

Similarly, (21) can be controlled via

$$\begin{aligned}
 \left| \frac{\partial}{\partial \theta_{jj}} F(\Theta^*) \right| &= \left| -\mathbb{E}_{\mathbb{X}}[X_j] + \frac{1}{n} \sum_{i=1}^n \lambda_{ij}^* \right| = \left| -\mathbb{E}_{\mathbb{X}}[X_j] + \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}] + \lambda_{ij}^* - \mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}]) \right| \\
 &= \left| -\mathbb{E}_{\mathbb{X}}[X_j] + \mathbb{E}[X_j] - \mathbb{E}[X_j] + \mathbb{E}_{\mathbb{X}}[\mathbb{E}[X_j | \mathbf{X}_{-j}]] + \frac{1}{n} \sum_{i=1}^n (\lambda_{ij}^* - \mathbb{E}[X_j | \mathbf{X}_{-j} = \mathbf{x}_{i,-j}]) \right| \\
 &\leq |\mathbb{E}_{\mathbb{X}}[\mathbb{E}[X_j | \mathbf{X}_{-j}]] - \mathbb{E}[X_j]| + |\mathbb{E}_{\mathbb{X}}[X_j] - \mathbb{E}[X_j]| + C_4 \\
 &= |\mathbb{E}_{\mathbb{X}}[\mathbb{E}[X_j | \mathbf{X}_{-j}]] - \mathbb{E}[\mathbb{E}[X_j | \mathbf{X}_{-j}]]| + |\mathbb{E}_{\mathbb{X}}[X_j] - \mathbb{E}[X_j]| + C_4.
 \end{aligned}$$

We define four events:

$$\begin{aligned}
 E_1 &:= \left\{ \max_{j \neq j', j, j' \in \{1, 2, \dots, p\}} \left| \frac{\partial}{\partial \theta_{jj'}} F(\Theta^*) \right| \geq 2(C_3 \epsilon_1 + \epsilon_1^2) \sqrt{\frac{\log p}{n}} + 2C_4 \epsilon_2 \right\}, \\
 E_2 &:= \left\{ \max_{j \in \{1, 2, \dots, p\}} \left| \frac{\partial}{\partial \theta_{jj}} F(\Theta^*) \right| \geq (C_3 + \epsilon_1) \sqrt{\frac{\log p}{n}} + C_4/n \right\}, \\
 E_3 &:= \left\{ 0 < \max_{i \in \{1, 2, \dots, n\}} \|\mathbf{x}_i\|_{\infty} < \epsilon_1 \right\}, \quad \text{and} \quad E_4 := \left\{ 0 < \max_{j \in \{1, 2, \dots, p\}} \mathbb{E}_{\mathbb{X}}[X_j] < \epsilon_2 \right\},
 \end{aligned}$$

where  $\epsilon_1 := 3 \log p + \log n$  and  $\epsilon_2 := C_1 + \sqrt{\frac{2 \log p}{n}}$  are defined in Lemma 7. By Lemma 3, Lemma 4, Lemma 5 and Lemma 6, it follows that

$$\mathbb{P}(E_1 | E_3, E_4) \leq 4 \exp(-2 \log p) \quad \text{and} \quad \mathbb{P}(E_2 | E_3, E_4) \leq 4 \exp(-2 \log p). \quad (22)$$

Therefore,

$$\begin{aligned}
 \mathbb{P}(E_1 \cup E_2) &= \mathbb{P}(E_1 \cup E_2 | E_3, E_4) \mathbb{P}(E_3, E_4) + \mathbb{P}(E_1 \cup E_2 | E_3^c, E_4) \mathbb{P}(E_3^c, E_4) \\
 &+ \mathbb{P}(E_1 \cup E_2 | E_3, E_4^c) \mathbb{P}(E_3, E_4^c) + \mathbb{P}(E_1 \cup E_2 | E_3^c, E_4^c) \mathbb{P}(E_3^c, E_4^c) \\
 &\leq \mathbb{P}(E_1 | E_3, E_4) + \mathbb{P}(E_2 | E_3, E_4) + \mathbb{P}(E_3^c, E_4) + \mathbb{P}(E_3, E_4^c) + \mathbb{P}(E_3^c, E_4^c) \\
 &\leq \mathbb{P}(E_1 | E_3, E_4) + \mathbb{P}(E_2 | E_3, E_4) + \mathbb{P}(E_3^c) + \mathbb{P}(E_4^c) \\
 &\leq 8 \exp(-2 \log p) + \exp(C_1 + C_2/2 - \epsilon_1) np + \exp \left[ -\frac{n(\epsilon_2 - C_1)^2}{2C_2} \right] \\
 &= 8 \exp(-2 \log p) + \frac{\exp(C_1 + C_2/2)}{p^2} + p^{-\frac{1}{c_2}},
 \end{aligned} \quad (23)$$

where the superscript  $c$  over an event represents the complement of that event, and the last inequality is due to (22) and Lemma 7. Also notice that by the definitions of  $E_1$  and  $E_2$ ,

$$2(C_3 \epsilon_1 + \epsilon_1^2) \sqrt{\frac{\log p}{n}} + 2C_4 \epsilon_2 > (C_3 + \epsilon_1) \sqrt{\frac{\log p}{n}} + C_4/n.$$

Therefore, with probability of  $1 - \mathbb{P}(E_1 \cup E_2) \geq 1 - ((\exp(C_1 + C_2/2) + 8)p^{-2} + p^{-1/C_2})$ , neither  $E_1$  nor  $E_2$  occurs, and hence

$$\begin{aligned} \|\nabla F(\Theta^*)\|_\infty &\leq 2(C_3\epsilon_1 + \epsilon_1^2)\sqrt{\frac{\log p}{n}} + 2C_4\epsilon_2 \\ &= 2[C_3(3\log p + \log n) + (3\log p + \log n)^2]\sqrt{\frac{\log p}{n}} + 2C_4\left(C_1 + \sqrt{\frac{2\log p}{n}}\right). \end{aligned}$$

PROOF OF (19)

Then, we study (19). We consider a map defined as  $G(\Delta_S) := -\mathbf{H}_{SS}^{-1}[\nabla_S F(\Theta^* + \Delta_S) + \lambda \hat{\mathbf{Z}}_S] + \Delta_S$ . If  $\|\Delta\|_\infty \leq r$ , by Taylor expansion of  $\nabla_S F(\Theta^* + \Delta)$  centered at  $\nabla_S F(\Theta^*)$ ,

$$\begin{aligned} \|G(\Delta_S)\|_\infty &= \left\| -\mathbf{H}_{SS}^{-1}[\nabla_S F(\Theta^*) + \mathbf{H}_{SS}\Delta_S + \mathbf{R}_S(\Delta) + \lambda \hat{\mathbf{Z}}_S] + \Delta_S \right\|_\infty = \left\| -\mathbf{H}_{SS}^{-1}(\nabla_S F(\Theta^*) + \mathbf{R}_S(\Delta) + \lambda \hat{\mathbf{Z}}_S) \right\|_\infty \\ &\leq \left\| \mathbf{H}_{SS}^{-1} \right\|_\infty (\|\nabla_S F(\Theta^*)\|_\infty + \|\mathbf{R}_S(\Delta)\|_\infty + \lambda \|\hat{\mathbf{Z}}_S\|_\infty) \leq (C_5(\lambda + C_6 r^2 + \lambda)) = C_5 C_6 r^2 + 2C_5 \lambda, \end{aligned}$$

where the inequality is due to  $\|\nabla_S F(\Theta^*)\|_\infty \leq \lambda$  conditioning on  $E_1^c \cap E_2^c$  and according to (18). Then, based on the definition of  $r$ , we can derive the upper bound of  $\|G(\Delta_S)\|_\infty$  as  $\|G(\Delta_S)\|_\infty \leq r/2 + r/2 = r$ .

Therefore, according to the fixed point theorem (Ortega and Rheinboldt, 2000; Yang and Ravikumar, 2011), there exists  $\Delta_S$  satisfying  $G(\Delta_S) = \Delta_S$ , which indicates  $\nabla_S F(\Theta^* + \Delta) + \lambda \hat{\mathbf{Z}}_S = \mathbf{0}$ . Considering that the optimal solution to (25) is unique,  $\hat{\Delta}_S = \Delta_S$ , whose infinite norm is bounded by  $\|\Delta_S\|_\infty \leq r$ , with probability larger than  $1 - ((\exp(C_1 + C_2/2) + 8)p^{-2} + p^{-1/C_2})$ .  $\square$

**Lemma 9.** *Let  $\hat{\Theta}$  be an optimal solution to (12), and  $\hat{\mathbf{Z}}$  be the corresponding dual solution. If  $\hat{\mathbf{Z}}$  satisfies  $\|\hat{\mathbf{Z}}_I\|_\infty < 1$ , then any given optimal solution to (12)  $\tilde{\Theta}$  satisfies  $\tilde{\Theta}_I = \mathbf{0}$ . Moreover, if  $\mathbf{H}_{SS}$  is positive definite, then the solution to (12) is unique.*

*Proof.* Specifically, following the same rationale as Lemma 1 in Wainwright 2009, Lemma 1 in Ravikumar et al. 2010, and Lemma 2 in Yang and Ravikumar 2011, we can derive Lemma 9 characterizing the optimal solution of (12).  $\square$

#### 8.4. Proof of Theorem 1

The proof follows the primal-dual witness (PDW) technique, which is widely used in this line of research (Wainwright, 2009; Ravikumar et al., 2010; Yang and Ravikumar, 2011; Yang et al., 2015a). Specifically, by Lemma 9, we can prove the sparsistency by building an optimal solution to (12) satisfying  $\|\hat{\mathbf{Z}}_I\|_\infty < 1$ , which is summarized as *strict dual feasibility* (SDF). To this end, we apply PDW to build a qualified optimal solution with the assumption that  $\mathbf{H}_{SS}$  is positively definite.

##### SOLVE A RESTRICTED PROBLEM

First of all, we derive the KKT condition of (12):

$$\nabla F(\hat{\Theta}) + \lambda \hat{\mathbf{Z}} = \mathbf{0}. \quad (24)$$

To construct an optimal primal-dual pair solution, we define  $\tilde{\Theta}$  as an optimal solution to the restricted problem:

$$\tilde{\Theta} := \min_{\Theta} F(\Theta) + \lambda \|\Theta\|_1, \quad (25)$$

with  $\Theta_I = \mathbf{0}$ , where  $\tilde{\Theta}$  is unique according to Lemma 9 with the assumption that  $\mathbf{H}_{SS} \succ \mathbf{0}$ . Denote the subgradient corresponding to  $\tilde{\Theta}$  as  $\tilde{\mathbf{Z}}$ . Then  $(\tilde{\Theta}, \tilde{\mathbf{Z}})$  is optimal for the restricted problem (25). Therefore,  $\tilde{\mathbf{Z}}_S$  can be determined according to the values of  $\tilde{\Theta}_S$  via the KKT conditions of (25). As a result,

$$\nabla_S F(\tilde{\Theta}) + \lambda \tilde{\mathbf{Z}}_S = \mathbf{0}, \quad (26)$$

where  $\nabla_S$  represents the gradient components with respect to  $S$ . Furthermore, by letting  $\hat{\Theta} = \tilde{\Theta}$ , we determine  $\tilde{\mathbf{Z}}_I$  according to (24). It remains to show that  $\tilde{\mathbf{Z}}_I$  satisfies SDF.

CHECK SDF

Now, we demonstrate that  $\tilde{\Theta}$  and  $\tilde{\mathbf{Z}}$  satisfy SDF. By (26), and by the Taylor expansion of  $\nabla_S F(\tilde{\Theta})$ , we have that

$$\mathbf{H}_{SS}\tilde{\Delta}_S + \nabla_S F(\Theta^*) + \mathbf{R}_S(\tilde{\Delta}) + \lambda\tilde{\mathbf{Z}}_S = \mathbf{0} \Rightarrow \tilde{\Delta}_S = \mathbf{H}_{SS}^{-1} \left[ -\nabla_S F(\Theta^*) - \mathbf{R}_S(\tilde{\Delta}) - \lambda\tilde{\mathbf{Z}}_S \right], \quad (27)$$

where  $\tilde{\Delta} := \tilde{\Theta} - \Theta^*$ ,  $\mathbf{R}_S(\tilde{\Delta})$  represents the components of  $\mathbf{R}(\tilde{\Delta})$  corresponding to  $S$ , and we have used the fact that  $\mathbf{H}_{SS}$  is positive definite and hence invertible. By the definition of  $\tilde{\Theta}$  and  $\tilde{\mathbf{Z}}$ ,

$$\nabla F(\tilde{\Theta}) + \lambda\tilde{\mathbf{Z}} = \mathbf{0} \Rightarrow \nabla F(\Theta^*) + \mathbf{H}\tilde{\Delta} + \mathbf{R}(\tilde{\Delta}) + \lambda\tilde{\mathbf{Z}} = \mathbf{0} \Rightarrow \nabla_I F(\tilde{\Theta}) + \mathbf{H}_{IS}\tilde{\Delta}_S + \mathbf{R}_I(\tilde{\Delta}) + \lambda\tilde{\mathbf{Z}}_I = \mathbf{0}, \quad (28)$$

where  $\mathbf{R}_I(\tilde{\Delta})$  represents the components of  $\mathbf{R}(\tilde{\Delta})$  corresponding to  $I$ , and we have used the fact that  $\tilde{\Delta}_I = \mathbf{0}$  because  $\tilde{\Theta}_I = \Theta^* = \mathbf{0}$ . As a result,

$$\begin{aligned} \lambda\|\tilde{\mathbf{Z}}_I\|_\infty &= \|\mathbf{H}_{IS}\tilde{\Delta}_S + \nabla_I F(\Theta^*) + \mathbf{R}_I(\tilde{\Delta})\|_\infty \\ &\leq \left\| \mathbf{H}_{IS}\mathbf{H}_{SS}^{-1} \left[ -\nabla_S F(\Theta^*) - \mathbf{R}_S(\tilde{\Delta}) - \lambda\tilde{\mathbf{Z}}_S \right] \right\|_\infty + \|\nabla_I F(\Theta^*) + \mathbf{R}_I(\tilde{\Delta})\|_\infty \\ &\leq \|\mathbf{H}_{IS}\mathbf{H}_{SS}^{-1}\|_\infty \left\| \nabla_S F(\Theta^*) + \mathbf{R}_S(\tilde{\Delta}) \right\|_\infty + \|\mathbf{H}_{IS}\mathbf{H}_{SS}^{-1}\|_\infty \|\lambda\tilde{\mathbf{Z}}_S\|_\infty + \|\nabla_I F(\Theta^*) + \mathbf{R}_I(\tilde{\Delta})\|_\infty \\ &\leq (1-\alpha) \left( \|\nabla_S F(\Theta^*)\|_\infty + \|\mathbf{R}_S(\tilde{\Delta})\|_\infty \right) + (1-\alpha)\lambda + \left( \|\nabla_I F(\Theta^*)\|_\infty + \|\mathbf{R}_I(\tilde{\Delta})\|_\infty \right) \\ &\leq (2-\alpha) \left( \|\nabla F(\Theta^*)\|_\infty + \|\mathbf{R}(\tilde{\Delta})\|_\infty \right) + (1-\alpha)\lambda, \end{aligned} \quad (29)$$

where we have used (27) in the first inequality, and the third inequality is due to Assumption 3.

With (29), it remains to control  $\|\nabla F(\Theta^*)\|_\infty$  and  $\|\mathbf{R}(\tilde{\Delta})\|_\infty$ . On one hand, according to Lemma 8 and the assumption on  $\lambda$  in Theorem 1,  $\|\nabla F(\Theta^*)\|_\infty \leq 2 \left[ 3C_3 \log p + C_3 \log n + (3 \log p + \log n)^2 \right] \sqrt{\frac{\log p}{n}} + 2C_4 \left( C_1 + \sqrt{\frac{2 \log p}{n}} \right) \leq \frac{\alpha\lambda}{4}$ , with probability larger than  $1 - ((\exp(C_1 + C_2/2) + 8)p^{-2} + p^{-1/C_2})$ .

On the other hand, according to Assumption 4 and Lemma 8,

$$\|\mathbf{R}(\tilde{\Delta})\|_\infty \leq C_6 \|\tilde{\Delta}\|_\infty^2 \leq C_6 r^2 \leq C_6 (4C_5\lambda)^2 = \lambda \frac{64C_5^2 C_6}{\alpha} \frac{\alpha\lambda}{4} \leq \left( C_7 \sqrt{\frac{\log^5 p}{n}} \right) \frac{64C_5^2 C_6}{\alpha} \frac{\alpha\lambda}{4}, \quad (30)$$

where in the last inequality we have used the assumption  $\lambda \propto \sqrt{\frac{\log^5 p}{n}}$  in Theorem 1, and hence there exists  $C_7$  satisfying  $\lambda \leq C_7 \sqrt{\frac{\log^5 p}{n}}$ . Therefore, when we choose  $n \geq (64C_7 C_5^2 C_6 / \alpha)^2 \log^5 p$  as assumed in Theorem 1, then from (30), we can conclude that  $\|\mathbf{R}(\tilde{\Delta})\|_\infty \leq \frac{\alpha\lambda}{4}$ . As a result,  $\lambda\|\tilde{\mathbf{Z}}_I\|_\infty$  can be bounded by  $\lambda\|\tilde{\mathbf{Z}}_I\|_\infty < \alpha\lambda/2 + \alpha\lambda/2 + (1-\alpha)\lambda = \lambda$ . Combined with Lemma 9, we demonstrate that any optimal solution of (12) satisfies  $\tilde{\Theta}_I = \mathbf{0}$ . Furthermore, (19) controls the difference between the optimal solution of (12) and the real parameter by  $\|\tilde{\Delta}_S\|_\infty \leq r$ , by the fact that  $r \leq \|\Theta_S^*\|_\infty$  in Theorem 1,  $\tilde{\Theta}_S$  shares the same sign with  $\Theta_S^*$ .