

Supplementary Material

6 Auxiliary Lemmas: Proof of Lemma 3

Proof. We can rewrite (6) as an optimization problem over the ℓ_1/ℓ_2 ball of radius C for some $C(\lambda_n) < \infty$. Since $\lambda_n > 0$, by KKT conditions, $\|\tilde{\Theta}_{\setminus r}\|_{1,2} = C$ for all optimal primal solution $\tilde{\Theta}_{\setminus r}$.

By definition of the ℓ_1/ℓ_2 subdifferential, we know that for any column $u \in V \setminus \{r\}$, we have $\|(\hat{Z}_{\setminus r})_u\|_2 \leq 1$. Considering the necessary optimality condition $\nabla \ell(\hat{\Theta}_{\setminus r}) + \lambda_n \hat{Z}_{\setminus r} = 0$, by complementary slackness condition, we have $\langle \tilde{\Theta}_{\setminus r}, \hat{Z}_{\setminus r} \rangle - C = \langle \tilde{\Theta}_{\setminus r}^T, \hat{Z}_{\setminus r} \rangle - \|\tilde{\Theta}_{\setminus r}\|_{1,2} = 0$. Now if for an arbitrary column $u \in V \setminus \{r\}$, we have $\|(\hat{Z}_{\setminus r})_u\|_2 < 1$ and $(\tilde{\Theta}_{\setminus r})_u \neq 0$ then this would contradict the condition that $\langle \tilde{\Theta}_{\setminus r}, \hat{Z}_{\setminus r} \rangle = \|\tilde{\Theta}_{\setminus r}\|_{1,2}$.

For this restricted problem, if the Hessian sub-matrix is positive definite, then the problem is strictly convex and it has a unique solution. \square

7 Derivatives of the Log-Likelihood Function

In this section, we point out the key properties of the gradient, Hessian and derivative of the Hessian for the log-likelihood function. These properties are used to prove the concentration lemmas.

7.1 Gradient

By simple derivation, we have

$$\begin{aligned} & \frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell^{(i)}(\Theta_{\setminus r}; D) \\ &= \mathcal{I}[x_t^{(i)} = k] \left(\mathcal{I}[x_r^{(i)} = \ell] - \mathbb{P}_{\Theta_{\setminus r}^*}[X_r = \ell | X_{\setminus r} = x_{\setminus r}^{(i)}] \right). \end{aligned}$$

It is easy to show that $\mathbb{E}_{\Theta_{\setminus r}^*} \left[\frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell^{(i)}(\Theta_{\setminus r}; D) \right] = 0$ and $\text{Var} \left(\frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell^{(i)}(\Theta_{\setminus r}; D) \right) \leq \frac{1}{4}$. With i.i.d assumption on drawn samples, we have $\text{Var} \left(\frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell(\Theta_{\setminus r}; D) \right) = \text{Var} \left(\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell^{(i)}(\Theta_{\setminus r}; D) \right) \leq \frac{1}{4n}$. Hence, for a

fixed $t \in V \setminus \{r\}$ by Jensen's inequality,

$$\begin{aligned} & \mathbb{E}_{\Theta_{\setminus r}^*} \left[\left\| \frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell(\Theta_{\setminus r}; D) \right\|_2 \right] \\ & \leq \sqrt{\mathbb{E}_{\Theta_{\setminus r}^*} \left[\left\| \frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell(\Theta_{\setminus r}; D) \right\|_2^2 \right]} \\ & \leq \frac{m-1}{2\sqrt{n}}. \end{aligned}$$

Considering the terms associated with $\theta_{rt;\ell k}^*$'s in the gradient vector of the log-likelihood function, for a fixed $t \in V \setminus \{r\}$, only $m-1$ (out of $(m-1)^2$) values are non-zero. By a simple calculation, we get

$$\max_{t \in V \setminus \{r\}} \left\| \frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell^{(i)}(\Theta_{\setminus r}; D) \right\|_2 \leq \sqrt{2} \quad \forall i.$$

By Azuma-Hoeffding inequality, we get

$$\mathbb{P} \left[\left\| \frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell(\Theta_{\setminus r}; D) \right\|_2 > \frac{m-1}{2\sqrt{n}} + \epsilon \right] \leq 2 \exp \left(-\frac{\epsilon^2}{4} n \right),$$

for all $t \in V \setminus \{r\}$. Using the union bound, we get

$$\begin{aligned} & \mathbb{P} \left[\max_{t \in V \setminus \{r\}} \left\| \frac{\partial}{\partial \theta_{rt;\ell k}^*} \ell(\Theta_{\setminus r}; D) \right\|_2 > \frac{m-1}{2\sqrt{n}} + \epsilon \right] \\ & \leq 2 \exp \left(-\frac{\epsilon^2}{4} n + \log(p-1) \right). \end{aligned} \quad (12)$$

7.2 Hessian

For the Hessian of the log-likelihood function, we have

$$\frac{\partial^2 \ell^{(i)}(\Theta_{\setminus r}; D)}{\partial \theta_{rt_2;\ell_2 k_2}^* \partial \theta_{rt_1;\ell_1 k_1}^*} = \mathcal{I}[x_{t_1}^{(i)} = k_1] \mathcal{I}[x_{t_2}^{(i)} = k_2] \eta_{\ell_1 \ell_2}(x^{(i)}),$$

where,

$$\begin{aligned} \eta_{\ell_1 \ell_2}(x^{(i)}) &:= \mathbb{P}_{\Theta_{\setminus r}^*} [X_r = \ell_1 | X_{\setminus r} = x_{\setminus r}^{(i)}] \\ & \left(\mathcal{I}[x_r^{(i)} = \ell_1] \mathcal{I}[x_r^{(i)} = \ell_2] - \mathbb{P}_{\Theta_{\setminus r}^*} [X_r = \ell_2 | X_{\setminus r} = x_{\setminus r}^{(i)}] \right). \end{aligned}$$

Consider the zero-mean random variable

$$\begin{aligned} Z_{t_1 \ell_1 k_1; t_2 \ell_2 k_2}^{(i)} &:= \\ & \frac{\partial^2 \ell^{(i)}(\Theta_{\setminus r}; D)}{\partial \theta_{rt_2;\ell_2 k_2}^* \partial \theta_{rt_1;\ell_1 k_1}^*} - \mathbb{E} \left[\frac{\partial^2 \ell(\Theta_{\setminus r}; D)}{\partial \theta_{rt_2;\ell_2 k_2}^* \partial \theta_{rt_1;\ell_1 k_1}^*} \right]. \end{aligned}$$

Notice that $\text{Var} \left(Z_{t_1 \ell_1 k_1; t_2 \ell_2 k_2}^{(i)} \right) \leq 1$ and consequently, by i.i.d assumption, $\text{Var} \left(\frac{1}{n} \sum_{i=1}^n Z_{t_1 \ell_1 k_1; t_2 \ell_2 k_2}^{(i)} \right) \leq \frac{1}{n}$.

Hence, for fixed values t_1, ℓ_1, k_1 and $t_2 \in S_2 \subseteq V \setminus \{r\}$, we have

$$\begin{aligned} & \mathbb{E}_{\Theta_{\setminus r}^*} \left[\left\| \frac{1}{n} \sum_{i=1}^n Z_{t_1 \ell_1 k_1; t_2 \ell_2 k_2}^{(i)} \right\|_2 \right] \\ & \leq \sqrt{\mathbb{E}_{\Theta_{\setminus r}^*} \left[\left\| \frac{1}{n} \sum_{i=1}^n Z_{t_1 \ell_1 k_1; t_2 \ell_2 k_2}^{(i)} \right\|_2^2 \right]} \\ & \leq \sqrt{\frac{|S_2|}{n}}. \end{aligned} \quad (13)$$

This random variable, for fixed values t_1, ℓ_1, k_1 and a fixed t_2 , is bounded and in particular, $\left\| \frac{1}{n} \sum_{i=1}^n Z_{t_1 \ell_1 k_1; t_2 \ell_2 k_2}^{(i)} \right\|_2 \leq 2$. By Azuma-Hoeffding inequality and the union bound,

$$\begin{aligned} & \mathbb{P} \left[\left\| Q_{S_r S_r}^n - Q_{S_r S_r}^* \right\|_{\infty, 2} > \frac{\sqrt{d_r}}{\sqrt{n}} + \epsilon \right] \\ & \leq 2 \exp \left(-\frac{\epsilon^2}{8} n + \log((m-1)^2 d_r) \right). \\ & \mathbb{P} \left[\left\| Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right\|_{\infty, 2} > \frac{\sqrt{d_r}}{\sqrt{n}} + \epsilon \right] \\ & \leq 2 \exp \left(-\frac{\epsilon^2}{8} n + \log((m-1)^2 (p-d_r-1)) \right). \end{aligned} \quad (14)$$

Similar analysis as (13) combined with the inequality $\Lambda_{\max}(\cdot) \leq \|\cdot\|_{\infty, 2}$, shows that

$$\begin{aligned} & \mathbb{P} \left[\Lambda_{\max} (Q_{S_r S_r}^n - Q_{S_r S_r}^*) > \frac{\sqrt{d_r}}{\sqrt{n}} + \epsilon \right] \\ & \leq 2 \exp \left(-\frac{\epsilon^2}{8} n + \log((m-1)^2 d_r) \right). \end{aligned} \quad (15)$$

We also need a control over the deviation of the inverse sample Fisher information matrix from the inverse of its mean. We have

$$\begin{aligned} & \Lambda_{\max} \left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right) \\ & = \Lambda_{\max} \left((Q_{S_r S_r}^*)^{-1} (Q_{S_r S_r}^* - Q_{S_r S_r}^n) (Q_{S_r S_r}^n)^{-1} \right) \\ & \leq \Lambda_{\max} \left((Q_{S_r S_r}^*)^{-1} \right) \Lambda_{\max} (Q_{S_r S_r}^* - Q_{S_r S_r}^n) \\ & \quad \Lambda_{\max} \left((Q_{S_r S_r}^n)^{-1} \right) \\ & \leq \frac{\sqrt{d_r}}{C_{\min} \sqrt{n}} \Lambda_{\max} \left((Q_{S_r S_r}^n)^{-1} \right). \end{aligned}$$

By part (B1) in Lemma 1, we have

$$\begin{aligned} & \mathbb{P} \left[\Lambda_{\max} \left((Q_{S_r S_r}^n)^{-1} \right) > \frac{1}{C_{\min}} + \epsilon \right] \\ & \leq 2 \exp \left(-\frac{\left(\frac{C_{\min} \epsilon \sqrt{n}}{1+C_{\min} \epsilon} - \sqrt{d_r} \right)^2}{8} + \log((m-1)^2 d_r) \right). \end{aligned} \quad (16)$$

Hence, we get,

$$\begin{aligned} & \mathbb{P} \left[\Lambda_{\max} \left((Q_{S_r S_r}^n)^{-1} (Q_{S_r S_r}^*)^{-1} \right) > \frac{\sqrt{d_r}}{C_{\min} \sqrt{n}} + \epsilon \right] \\ & \leq 4 \exp \left(-\frac{\left(\frac{C_{\min} \epsilon \sqrt{n}}{1+C_{\min} \epsilon} - \sqrt{d_r} \right)^2}{8} + \log((m-1)^2 d_r) \right). \end{aligned} \quad (17)$$

7.3 Derivative of Hessian

We want to bound the rate of the change for the elements of Hessian matrix. Let

$$\begin{aligned} & \nabla Q_{t_2 \ell_2 k_2; t_1 \ell_1 k_1}^{(i)} \\ & := \frac{\partial}{\partial \Theta_{\setminus r}} \frac{\partial^2 \ell^{(i)}(\Theta_{\setminus r}; D)}{\partial \theta_{rt_2; \ell_2 k_2}^* \partial \theta_{rt_1; \ell_1 k_1}^*} \\ & = \mathcal{I} \left[x_{t_1}^{(i)} = k_1 \right] \mathcal{I} \left[x_{t_2}^{(i)} = k_2 \right] \frac{\partial}{\partial \Theta_{\setminus r}} \eta_{\ell_1 \ell_2} \left(x^{(i)} \right). \end{aligned}$$

Recall the definition of $\eta(\cdot)$ from section 7.2. We have

$$\begin{aligned} & \frac{\partial \eta_{\ell_1 \ell_2} \left(x^{(i)} \right)}{\partial \theta_{rt_3; \ell_3 k_3}} = \mathcal{I} \left[x_{t_3}^{(i)} = k_3 \right] \mathbb{P}_{\Theta_{\setminus r}^*} \left[X_r = \ell_1 \mid X_{\setminus r} = x_{\setminus r}^{(i)} \right] \\ & \quad \left(\eta_{\ell_2 \ell_3} \left(x^{(i)} \right) - \frac{\eta_{\ell_1 \ell_2} \left(x^{(i)} \right) \eta_{\ell_1 \ell_3} \left(x^{(i)} \right)}{\mathbb{P}_{\Theta_{\setminus r}^*} \left[X_r = \ell_1 \mid X_{\setminus r} = x_{\setminus r}^{(i)} \right]^2} \right). \end{aligned}$$

For any $t_3 \in V \setminus \{r\}$, each entry is bounded by $\frac{1}{2}$ and there are only $m-1$ non-zero entries for each k_3 . Hence, for any t_3 , one can calculate that $\left\| \frac{\partial}{\partial \theta_{rt_3; \ell_3 k_3}} \eta_{\ell_1 \ell_2} \left(x^{(i)} \right) \right\|_2 \leq \frac{m-1}{\sqrt{2}}$ for all i . Finally, for all ℓ_1 and ℓ_2 we have

$$\max_{t_3 \in V \setminus \{r\}} \left\| \frac{\partial}{\partial \theta_{rt_3; \ell_3 k_3}} \eta_{\ell_1 \ell_2} \left(x^{(i)} \right) \right\|_2 \leq \frac{m-1}{\sqrt{2}}. \quad (18)$$

8 Proof of Lemma 1

(B1) By variational representation of the smallest eigenvalue, we have

$$\begin{aligned} \Lambda_{\min} (Q_{S_r S_r}^*) & = \min_{\|x\|_2=1} x^T Q_{S_r S_r}^* x \\ & \leq y^T Q_{S_r S_r}^n y + y^T (Q_{S_r S_r}^* - Q_{S_r S_r}^n) y, \end{aligned}$$

for all $y \in \mathbb{R}^{(m-1)^2 d_r}$ with $\|y\|_2 = 1$ and in particular for the unit-norm minimal eigenvalue of $Q_{S_r S_r}^n$. Hence,

$$\Lambda_{\min}(Q_{S_r S_r}^n) \geq \Lambda_{\min}(Q_{S_r S_r}^*) - \Lambda_{\max}(Q_{S_r S_r}^* - Q_{S_r S_r}^n).$$

By (15), we get

$$\begin{aligned} \mathbb{P}[\Lambda_{\min}(Q_{S_r S_r}^n) < C_{\min} - \epsilon] &\leq \mathbb{P}[\Lambda_{\max}(Q_{S_r S_r}^* - Q_{S_r S_r}^n) > \epsilon] \\ &\leq 2 \exp\left(-\frac{(\epsilon\sqrt{n} - \sqrt{d_r})^2}{8} + \log((m-1)^2 d_r)\right). \end{aligned}$$

(B2) We can write

$$\begin{aligned} Q_{S_r^c S_r}^n (Q_{S_r S_r}^n)^{-1} &= \underbrace{Q_{S_r^c S_r}^* (Q_{S_r S_r}^*)^{-1}}_{T_0} \\ &\quad + \underbrace{Q_{S_r^c S_r}^* \left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right)}_{T_1} \\ &\quad + \underbrace{\left(Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right) (Q_{S_r S_r}^*)^{-1}}_{T_2} \\ &\quad + \underbrace{\left(Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right) \left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right)}_{T_3}. \end{aligned}$$

Considering assumption (A3), $\|T_0\|_{\infty,2} < \frac{1-2\alpha}{\sqrt{d_r}}$ and hence, it suffices to show that $\|T_i\|_{\infty,2} < \frac{\alpha}{3\sqrt{d_r}}$ for $i = 1, 2, 3$. For the first term, we have

$$\begin{aligned} &\left\| Q_{S_r^c S_r}^* \left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right) \right\|_{\infty,2} \\ &= \left\| Q_{S_r^c S_r}^* (Q_{S_r S_r}^*)^{-1} (Q_{S_r S_r}^* - Q_{S_r S_r}^n) (Q_{S_r S_r}^n)^{-1} \right\|_{\infty,2} \\ &\leq \left\| Q_{S_r^c S_r}^* (Q_{S_r S_r}^*)^{-1} \right\|_{\infty,2} \Lambda_{\max}(Q_{S_r S_r}^* - Q_{S_r S_r}^n) \\ &\quad \Lambda_{\max}\left((Q_{S_r S_r}^n)^{-1} \right) \\ &\leq \frac{1-2\alpha}{\sqrt{d_r}} \frac{\sqrt{d_r}}{\sqrt{n}} \frac{1}{C_{\min}}. \end{aligned}$$

The last inequality follows from (14) and (16) with high probability. Setting $\bar{C}_{\min} = \min(C_{\min}, 1)$, by applying the union bound,

$$\begin{aligned} \mathbb{P}\left[\left\| Q_{S_r^c S_r}^* \left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right) \right\|_{\infty,2} > \epsilon\right] \\ \leq 4 \exp\left(-\frac{\left(\bar{C}_{\min} \epsilon \sqrt{n} - \sqrt{d_r} - \frac{1-2\alpha}{C_{\min}}\right)^2}{8} + \log((m-1)^2 d_r)\right). \end{aligned}$$

For the second term, we have

$$\begin{aligned} &\left\| \left(Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right) (Q_{S_r S_r}^*)^{-1} \right\|_{\infty,2} \\ &\leq \left\| Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right\|_{\infty,2} \Lambda_{\max}\left((Q_{S_r S_r}^*)^{-1} \right) \\ &\leq \frac{\sqrt{d_r}}{\sqrt{n}} \frac{1}{C_{\min}}. \end{aligned}$$

The last inequality follows from (14) with high probability. Hence, we have

$$\begin{aligned} \mathbb{P}\left[\left\| \left(Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right) (Q_{S_r S_r}^*)^{-1} \right\|_{\infty,2} > \epsilon\right] \\ \leq 2 \exp\left(-\frac{\left(\epsilon\sqrt{n} - \frac{(1+C_{\min})\sqrt{d_r}}{C_{\min}}\right)^2}{8} + \log((m-1)^2(p-1-d_r))\right). \end{aligned}$$

For the third term, we have

$$\begin{aligned} &\left\| \left(Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right) \left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right) \right\|_{\infty,2} \\ &\leq \left\| Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right\|_{\infty,2} \Lambda_{\max}\left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right) \\ &\leq \frac{\sqrt{d_r}}{\sqrt{n}} \frac{\sqrt{d_r}}{C_{\min}^2 \sqrt{n}} = \frac{d_r}{C_{\min}^2 n} \end{aligned}$$

The last inequality follows from (14) and (17). Hence, we have

$$\begin{aligned} \mathbb{P}\left[\left\| \left(Q_{S_r^c S_r}^n - Q_{S_r^c S_r}^* \right) \left((Q_{S_r S_r}^n)^{-1} - (Q_{S_r S_r}^*)^{-1} \right) \right\|_{\infty,2} > \epsilon\right] \\ \leq 6 \exp\left(-\frac{\left(\bar{C}_{\min} \epsilon \sqrt{n} - \left(1 + \frac{\sqrt{d_r}}{C_{\min}^2 \sqrt{n}}\right) \sqrt{d_r}\right)^2}{8} \right. \\ \left. + \log((m-1)^2(p-1-d_r))\right). \end{aligned}$$

The result follows by substituting ϵ with $\frac{\alpha}{3\sqrt{d_r}}$.

(B3) We can write

$$\begin{aligned} \mathbb{P}[\Lambda_{\max}(\mathcal{J}^n) > D_{\max} + \epsilon] \\ \leq \mathbb{P}\left[\left\| \frac{1}{n} \sum_{i=1}^n (\mathcal{J}^{(i)} - \mathcal{J}^*) \right\|_F > \epsilon\right]. \end{aligned}$$

Consequently, same analysis as part (B1) gives the result.

This concludes the proof of the Lemma.

9 Sufficiency Lemmas for Pairwise Dependencies

Lemma 5. *The constructed candidate primal-dual pair $(\hat{\Theta}_{\setminus r}, \hat{Z}_{\setminus r})$ satisfy the conditions of the Lemma 3*

with probability $1 - c_1 \exp(-c_2 n)$ for some positive constants $c_1, c_2 \in \mathbb{R}$.

Proof. Using the mean-value theorem, for some $\bar{\Theta}_{\setminus r}$ in the convex combination of $\hat{\Theta}_{\setminus r}$ and $\Theta_{\setminus r}^*$, we have

$$\begin{aligned} & \nabla^2 \ell(\Theta_{\setminus r}^*; D) \left[\hat{\Theta}_{\setminus r} - \Theta_{\setminus r}^* \right] \\ &= \nabla \ell(\hat{\Theta}_{\setminus r}; D) - \nabla \ell(\Theta_{\setminus r}^*; D) \\ & \quad + \left(\nabla^2 \ell(\Theta_{\setminus r}^*; D) - \nabla^2 \ell(\bar{\Theta}_{\setminus r}; D) \right) \left[\hat{\Theta}_{\setminus r} - \Theta_{\setminus r}^* \right] \\ &= -\lambda_n \hat{Z}_{\setminus r} - \underbrace{\nabla \ell(\Theta_{\setminus r}^*; D)}_{W_{\setminus r}^n} \\ & \quad + \underbrace{\left(\nabla^2 \ell(\Theta_{\setminus r}^*; D) - \nabla^2 \ell(\bar{\Theta}_{\setminus r}; D) \right)}_{R_{\setminus r}^n} \left[\hat{\Theta}_{\setminus r} - \Theta_{\setminus r}^* \right]. \end{aligned}$$

We can rewrite these set of equations as two sets of equations over S_r and S_r^c . By Lemma 1, the Hessian sub-matrix on S_r is invertible with high probability and thus we get

$$\begin{aligned} & Q_{S_r^c S_r}^n (Q_{S_r S_r}^n)^{-1} \left(-\lambda_n (\hat{Z}_{\setminus r})_{S_r} - (W_{\setminus r}^n)_{S_r} + (R_{\setminus r}^n)_{S_r} \right) \\ &= -\lambda_n (\hat{Z}_{\setminus r})_{S_r^c} - (W_{\setminus r}^n)_{S_r^c} + (R_{\setminus r}^n)_{S_r^c}. \end{aligned}$$

Equivalently, we get

$$\begin{aligned} (\hat{Z}_{\setminus r})_{S_r^c} &= \frac{1}{\lambda_n} \left[(W_{\setminus r}^n)_{S_r^c} - (R_{\setminus r}^n)_{S_r^c} \right] \\ & \quad - \frac{1}{\lambda_n} Q_{S_r^c S_r}^n (Q_{S_r S_r}^n)^{-1} \left((W_{\setminus r}^n)_{S_r} - (R_{\setminus r}^n)_{S_r} \right) \\ & \quad + Q_{S_r^c S_r}^n (Q_{S_r S_r}^n)^{-1} (\hat{Z}_{\setminus r})_{S_r}. \end{aligned}$$

Notice that $\left\| (\hat{Z}_{\setminus r})_{S_r} \right\|_{\infty, 2} = 1$. Thus, we can establish the following bound

$$\begin{aligned} & \left\| (\hat{Z}_{\setminus r})_{S_r^c} \right\|_{\infty, 2} \\ & \leq \left(1 + \left\| Q_{S_r^c S_r}^n (Q_{S_r S_r}^n)^{-1} \right\|_{\infty, 2} \sqrt{d_r} \right) \\ & \quad \left[\frac{\left\| (W_{\setminus r}^n)_{S_r^c} \right\|_{\infty, 2}}{\lambda_n} + \frac{\left\| (R_{\setminus r}^n)_{S_r^c} \right\|_{\infty, 2}}{\lambda_n} + 1 \right] - 1 \\ & \leq (2 - \alpha) \left(\frac{\alpha}{4(2 - \alpha)} + \frac{\alpha}{4(2 - \alpha)} + 1 \right) - 1 \\ & = 1 - \frac{\alpha}{2} < 1. \end{aligned}$$

The second inequality holds with high probability according to Lemma 1 and Lemma 6. \square

Lemma 6. For quantities defined in the proof of Lemma 5, the following inequalities hold:

$$\begin{aligned} & \mathbb{P} \left[\frac{\left\| W_{\setminus r}^n \right\|_{\infty, 2}}{\lambda_n} \geq \frac{\alpha}{4(2 - \alpha)} \right] \\ & \leq 2 \exp \left(- \frac{\left(\frac{\alpha}{4(2 - \alpha)} \lambda_n \sqrt{n} - \frac{m-1}{2} \right)^2}{4} + \log(p-1) \right) \\ & \mathbb{P} \left[\frac{\left\| R_{\setminus r}^n \right\|_{\infty, 2}}{\lambda_n} > \frac{\alpha}{4(2 - \alpha)} \right] \\ & \leq 2 \exp \left(- \frac{\left(\frac{\alpha}{4(2 - \alpha)} \lambda_n \sqrt{n} - \frac{m-1}{2} \right)^2}{4} + \log(p-1) \right). \end{aligned}$$

Proof. The first inequality follows directly from (12), for $\epsilon = \frac{\alpha}{4(2 - \alpha)} \lambda_n - \frac{m-1}{2\sqrt{n}}$, provided that $\lambda_n \geq \frac{2(2 - \alpha)}{\alpha} \frac{m-1}{\sqrt{n}}$. This probability goes to zero, if $\lambda_n \geq \frac{8(2 - \alpha)}{\alpha} \left(\sqrt{\frac{\log(p-1)}{n}} + \frac{m-1}{4\sqrt{n}} \right)$.

Before we proceed, we want to point out a technical fact that we will use it through the rest of the proof. For λ_n achieves the lower bound mentioned above, any positive value K and $n \geq \frac{1}{K^2} \frac{64(2 - \alpha)^2}{\alpha^2} \left(\sqrt{\log(p-1)} + \frac{m-1}{4} \right)^2 d_r^2$, we have $\lambda_n d_r \leq K$. Hence, we can assume $\lambda_n d_r$ is less than any fixed constant K for sufficiently large n .

In order to bound $R_{\setminus r}^n$, we need to bound $\left\| (\hat{\Theta}_{\setminus r})_{S_r} - (\Theta_{\setminus r}^*)_{S_r} \right\|_{\infty, 2}$, using the technique used in Rothman et al. [28]. Let $G : \mathbb{R}^{(m-1)^2 d_r} \rightarrow \mathbb{R}$ be a function defined as

$$\begin{aligned} G((U)_{S_r}) &:= \ell((\Theta_{\setminus r}^*)_{S_r} + (U)_{S_r}; D) - \ell((\Theta_{\setminus r}^*)_{S_r}; D) \\ & \quad + \lambda_n \left(\left\| (\Theta_{\setminus r}^*)_{S_r} + (U)_{S_r} \right\|_{1, 2} - \left\| (\Theta_{\setminus r}^*)_{S_r} \right\|_{1, 2} \right). \end{aligned}$$

By optimality of $\hat{\Theta}_{\setminus r}$, it is clear that $(\hat{U})_{S_r} = (\hat{\Theta}_{\setminus r})_{S_r} - (\Theta_{\setminus r}^*)_{S_r}$ minimizes G . Since $G(\mathbf{0}) = 0$ by construction, we have $G((\hat{U})_{S_r}) \leq 0$. Suppose there exist an ℓ_∞/ℓ_2 ball with radius B_r such that for any $\left\| (U)_{S_r} \right\|_{\infty, 2} = B_r$, we have that $G((U)_{S_r}) > 0$. Then, we can claim that $\left\| (\hat{U})_{S_r} \right\|_{\infty, 2} \leq B_r$; because if, in contrary, we assume that $(\hat{U})_{S_r}$ is outside the ball,

then for an appropriate choice of $t \in (0, 1)$, the point $t \left(\hat{U} \right)_{S_r} + (1-t)\mathbf{0}$ lies on the boundary of the ball. By convexity of G , we have

$$G \left(t \left(\hat{U} \right)_{S_r} + (1-t)\mathbf{0} \right) \leq tG \left(\left(\hat{U} \right)_{S_r} \right) + (1-t)G(\mathbf{0}) \leq 0.$$

This is a contradiction to the assumption of the positivity of G on the boundary of the ball.

Let $(U)_{S_r} \in \mathbb{R}^{(m-1)d_r}$ be an arbitrary vector with $\left\| (U)_{S_r} \right\|_{\infty,2} = \frac{5}{C_{\min}} \lambda_n$. Applying mean value theorem to the log likelihood function, for some $\beta \in [0, 1]$, we get

$$\begin{aligned} G \left((U)_{S_r} \right) &= \left\langle (W_{\setminus r})_{S_r}, (U)_{S_r} \right\rangle \\ &+ \left\langle (U)_{S_r}, \nabla^2 \ell \left(\left(\Theta_{\setminus r}^* \right)_{S_r} + \beta (U)_{S_r}; D \right) (U)_{S_r} \right\rangle \\ &+ \lambda_n \left(\left\| \left(\Theta_{\setminus r}^* \right)_{S_r} + (U)_{S_r} \right\|_{1,2} - \left\| \left(\Theta_{\setminus r}^* \right)_{S_r} \right\|_{1,2} \right). \end{aligned} \quad (19)$$

We bound each of these three terms individually. By Cauchy-Schwartz inequality, we have

$$\begin{aligned} \left| \left\langle (W_{\setminus r})_{S_r}, (U)_{S_r} \right\rangle \right| &\leq \left\| (W_{\setminus r})_{S_r} \right\|_{\infty,2} \left\| (U)_{S_r} \right\|_{1,2} \\ &\leq \frac{\alpha}{4(2-\alpha)} \lambda_n d_r \frac{5}{C_{\min}} \lambda_n \\ &\leq \frac{5}{4C_{\min}} d_r \lambda_n^2. \end{aligned}$$

Moreover, by triangle inequality,

$$\begin{aligned} \lambda_n \left(\left\| \left(\Theta_{\setminus r}^* \right)_{S_r} + (U)_{S_r} \right\|_{1,2} - \left\| \left(\Theta_{\setminus r}^* \right)_{S_r} \right\|_{1,2} \right) \\ \geq -\lambda_n \left\| (U)_{S_r} \right\|_{1,2} \\ \geq -\frac{5}{C_{\min}} d_r \lambda_n^2. \end{aligned}$$

To bound the other term, notice that by Taylor expansion,

we get

$$\begin{aligned} &\Lambda_{\min} \left(\nabla^2 \ell \left(\left(\Theta_{\setminus r}^* \right)_{S_r} + \beta (U)_{S_r}; D \right) \right) \\ &\geq \min_{\beta \in [0,1]} \Lambda_{\min} \left(\nabla^2 \ell \left(\left(\Theta_{\setminus r}^* \right)_{S_r} + \beta (U)_{S_r}; D \right) \right) \\ &\geq \Lambda_{\min} (Q_{S_r, S_r}^*) \\ &\quad - \max_{\beta \in [0,1]} \Lambda_{\max} \left(\left\langle \frac{\partial \nabla^2 \ell (\Theta_{S_r}; D)}{\partial \Theta_{S_r}} \Big|_{\left(\Theta_{\setminus r}^* \right)_{S_r} + \beta (U)_{S_r}}, (U)_{S_r} \right\rangle \right) \\ &\geq C_{\min} - \left(\max_{t_3 \in V \setminus \{r\}} \left\| \frac{\partial}{\partial \theta_{rt_3; \ell_3 k_3}} \eta_{\ell_1 \ell_2} \left(x^{(i)} \right) \right\|_2 \sqrt{d_r} \right) \\ &\quad \Lambda_{\max}(\mathfrak{S}^*) \sqrt{d_r} \left\| (U)_{S_r} \right\|_{\infty,2}, \end{aligned} \quad (20)$$

where, $\eta(\cdot)$ is defined in Section 7.2. We know that $\Lambda_{\max}(\mathfrak{S}^*) = \Lambda_{\max}(\mathcal{J}^*)$ as a property of Kronecher product. By (18) and assumption on the maximum eigenvalue of \mathcal{J}^* , we have

$$\begin{aligned} &\Lambda_{\min} \left(\nabla^2 \ell \left(\left(\Theta_{\setminus r}^* \right)_{S_r} + \beta (U)_{S_r}; D \right) \right) \\ &\geq C_{\min} - \frac{m-1}{\sqrt{2}} d_r D_{\max} \left\| (U)_{S_r} \right\|_{\infty,2} \\ &\geq C_{\min} - \frac{m-1}{\sqrt{2}} d_r D_{\max} \frac{5}{C_{\min}} \lambda_n \\ &\geq \frac{C_{\min}}{2} \left(\lambda_n d_r \leq \frac{C_{\min}^2}{\sqrt{50}(m-1)D_{\max}} \right). \end{aligned}$$

Hence, from (19), we get

$$G \left((U)_{S_r} \right) \geq d_r \frac{5}{C_{\min}} \lambda_n^2 \left(-\frac{1}{4} + \frac{5}{2} - 1 \right) > 0.$$

We can conclude that

$$\left\| \left(\hat{\Theta}_{\setminus r} \right)_{S_r} - \left(\Theta_{\setminus r}^* \right)_{S_r} \right\|_{\infty,2} \leq \frac{5}{C_{\min}} \lambda_n. \quad (21)$$

with high probability. With similar analysis on the maximum eigenvalue of the derivative of Hessian as in (20), it is easy to show that

$$\begin{aligned} &\frac{\left\| R_{\setminus r}^n \right\|_{\infty,2}}{\lambda_n} \\ &\leq \frac{1}{\lambda_n} \frac{m-1}{\sqrt{2}} d_r D_{\max} \left\| \left(\hat{\Theta}_{\setminus r} \right)_{S_r} - \left(\Theta_{\setminus r}^* \right)_{S_r} \right\|_{\infty,2}^2 \\ &\leq \frac{m-1}{\sqrt{2}} d_r D_{\max} \frac{25}{C_{\min}^2} \lambda_n \\ &\leq \frac{\alpha}{4(2-\alpha)}, \end{aligned}$$

provided that $\lambda_n d_r \leq \frac{C_{\min}^2}{50\sqrt{2}(m-1)D_{\max}} \frac{\alpha}{2-\alpha}$. \square

10 Proof of Lemma 4

(D1) By variational representation of the smallest eigenvalue, we have

$$\begin{aligned}
 & \Lambda_{\min} \left(\left[\nabla^2 \ell(\bar{\Theta}_P^*; D) \right]_{S_r S_r} \right) \\
 & \geq \Lambda_{\min} \left(\left[\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D) \right]_{S_r S_r} \right) \\
 & \quad - \Lambda_{\max} \left(\left[\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D) \right]_{S_r S_r} - \left[\nabla^2 \ell(\bar{\Theta}_P^*; D) \right]_{S_r S_r} \right) \\
 & \geq C_{\min}(1 + \gamma) \\
 & \quad - \Lambda_{\max} \left(\left[\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D) \right]_{S_r S_r} - \left[\nabla^2 \ell(\bar{\Theta}_P^*; D) \right]_{S_r S_r} \right).
 \end{aligned}$$

In the second inequality, we used the result of Lemma 1, i.e., the inequality holds with probability stated in Lemma 4. By Taylor expansion, for some $\beta \in [0, 1]$, and by (23), we get

$$\begin{aligned}
 & \Lambda_{\max} \left(\left[\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D) \right]_{S_r S_r} - \left[\nabla^2 \ell(\bar{\Theta}_P^*; D) \right]_{S_r S_r} \right) \\
 & \leq \Lambda_{\max} \left(\left\langle \frac{\partial \left[\nabla^2 \ell(\bar{\Theta}; D) \right]_{S_r S_r}}{\partial \bar{\Theta}} \bigg|_{\bar{\Theta}_{\setminus r}^* - \beta \bar{\Theta}_{P^c}^*}, \bar{\Theta}_{P^c}^* \right\rangle \right) \\
 & \leq \left\| \nabla \eta_{\ell_1 \ell_2}(x^{(i)}) \right\|_{\infty} D_{\max} \|\bar{\Theta}_{P^c}^*\|_1 \\
 & = \gamma C_{\min}.
 \end{aligned}$$

Note that $\|\nabla \eta_{\ell_1 \ell_2}(x^{(i)})\|_{\infty} \leq 1$ for $\eta(\cdot)$ defined in section 7.3. The last inequality holds as a result of Lemma 1 with the probability stated in Lemma 4. Hence, the result follows.

(D2) We can write

$$\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r^c S_r} \left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r S_r} \right)^{-1} = \sum_{i=0}^3 T_i,$$

where,

$$\begin{aligned}
 T_0 &= \nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r^c S_r} \left(\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r S_r} \right)^{-1} \\
 T_1 &= \nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r^c S_r} \\
 & \quad \left(\left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r S_r} \right)^{-1} - \left(\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r S_r} \right)^{-1} \right) \\
 T_2 &= \left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r^c S_r} - \nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r^c S_r} \right) \\
 & \quad \left(\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r S_r} \right)^{-1} \\
 T_3 &= \left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r^c S_r} - \nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r^c S_r} \right) \\
 & \quad \left(\left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r S_r} \right)^{-1} - \left(\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r S_r} \right)^{-1} \right).
 \end{aligned}$$

By Lemma 1, we have that $\|T_0\|_{\infty, 1} \leq \frac{1-\tau}{\sqrt{d_r}}$ with the probability stated in Lemma 4. For the second term, we have

$$\begin{aligned}
 & \|T_1\|_{\infty, 2} \\
 & \leq \|T_0\|_{\infty, 2} \Lambda_{\max} \left(\underbrace{\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r S_r} - \nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r S_r}}_{T_{12}} \right) \\
 & \quad \Lambda_{\max} \left(\underbrace{\left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r S_r} \right)^{-1}}_{T_{13}} \right) \\
 & \leq \frac{1-\tau}{\sqrt{d_r}} \gamma C_{\min} \frac{1}{C_{\min}} = \frac{1-\tau}{\sqrt{d_r}} \gamma.
 \end{aligned}$$

We used the result of (D1) for $\Lambda_{\max}(T_{13}) \leq \frac{1}{C_{\min}}$.

For the third term, we have

$$\begin{aligned}
 & \|T_2\|_{\infty, 2} \leq \left\| \underbrace{\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r^c S_r} - \nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r^c S_r}}_{T_{21}} \right\|_{\infty, 2} \\
 & \quad \Lambda_{\max} \left(\underbrace{\left(\nabla^2 \ell(\bar{\Theta}_{\setminus r}^*; D)_{S_r S_r} \right)^{-1}}_{T_{22}} \right) \\
 & \leq \gamma C_{\min} \frac{1}{C_{\min}(1 + \gamma)} \\
 & = \frac{\gamma}{1 + \gamma}.
 \end{aligned}$$

For the fourth term, we have

$$\begin{aligned}
 & \|T_3\|_{\infty, 2} \leq \|T_{21}\|_{\infty, 2} \Lambda_{\max}(T_{22}) \Lambda_{\max}(T_{12}) \Lambda_{\max}(T_{13}) \\
 & \leq \gamma C_{\min} \frac{1}{C_{\min}(1 + \gamma)} \gamma C_{\min} \frac{1}{C_{\min}} \\
 & \leq \frac{\gamma^2}{1 + \gamma}.
 \end{aligned}$$

Putting all pieces together, we get the result.

(D3) The result follows directly from Lemma 1.

This concludes the proof of Lemma.

11 Sufficiency Lemmas for Higher Order Dependencies

Lemma 7. *The constructed candidate primal-dual pair $(\hat{\Theta}_{\setminus r}, \hat{Z}_{\setminus r})$ satisfy the conditions of the Lemma 3*

with probability $1 - c_1 \exp(-c_2 n)$ for some positive constants $c_1, c_2 \in \mathbb{R}$.

Proof. Using the mean-value theorem, for some $\bar{\Theta}_{\setminus r}$ in the convex combination of $\hat{\Theta}_{\setminus r}$ and $\bar{\Theta}_P^*$, we have

$$\begin{aligned} & \nabla^2 \ell(\bar{\Theta}_P^*; D) \left[\hat{\Theta}_{\setminus r} - \bar{\Theta}_P^* \right] \\ &= \nabla \ell(\hat{\Theta}_{\setminus r}; D) - \nabla \ell(\bar{\Theta}_P^*; D) \\ & \quad + (\nabla^2 \ell(\bar{\Theta}_P^*; D) - \nabla^2 \ell(\bar{\Theta}_{\setminus r}; D)) \left[\hat{\Theta}_{\setminus r} - \bar{\Theta}_P^* \right] \\ &= -\lambda_n \hat{Z}_{\setminus r} - \underbrace{\nabla \ell(\bar{\Theta}_P^*; D)}_{\bar{W}_{\setminus r}^n} \\ & \quad + \underbrace{(\nabla^2 \ell(\bar{\Theta}_P^*; D) - \nabla^2 \ell(\bar{\Theta}_{\setminus r}; D))}_{\bar{R}_{\setminus r}^n} \left[\hat{\Theta}_{\setminus r} - \bar{\Theta}_P^* \right]. \end{aligned}$$

We can rewrite these set of equations as two sets of equations over S_r and S_r^c . By Lemma 4, the Hessian sub-matrix on S_r is invertible with high probability and thus we get

$$\begin{aligned} & \nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r^c S_r} \left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r S_r} \right)^{-1} \\ & \left(-\lambda_n \left(\hat{Z}_{\setminus r} \right)_{S_r} - \left(\bar{W}_{\setminus r}^n \right)_{S_r} + \left(\bar{R}_{\setminus r}^n \right)_{S_r} \right) \\ & \quad = -\lambda_n \left(\hat{Z}_{\setminus r} \right)_{S_r^c} - \left(\bar{W}_{\setminus r}^n \right)_{S_r^c} + \left(\bar{R}_{\setminus r}^n \right)_{S_r^c}. \end{aligned}$$

Notice that $\left\| \left(\hat{Z}_{\setminus r} \right)_{S_r} \right\|_{\infty, 2} = 1$ and hence, we get

$$\begin{aligned} & \left\| \left(\hat{Z}_{\setminus r} \right)_{S_r^c} \right\|_{\infty, 2} \\ & \leq \left(1 + \left\| \nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r^c S_r} \left(\nabla^2 \ell(\bar{\Theta}_P^*; D)_{S_r S_r} \right)^{-1} \right\|_{\infty, 2} \sqrt{d_r} \right) \\ & \quad \left[\frac{\left\| \bar{W}_{\setminus r}^n \right\|_{\infty, 2}}{\lambda_n} + \frac{\left\| \bar{R}_{\setminus r}^n \right\|_{\infty, 2}}{\lambda_n} + 1 \right] - 1 \\ & \leq (2 - \alpha) \left(\frac{\alpha}{4(2 - \alpha)} + \frac{\alpha}{4(2 - \alpha)} + 1 \right) - 1 \\ & = 1 - \frac{\alpha}{2} < 1. \end{aligned}$$

The second inequality holds with high probability according to Lemma 4 and Lemma 8. \square

Lemma 8. For quantities defined in the proof of

Lemma 7, the following inequalities hold:

$$\begin{aligned} & \mathbb{P} \left[\frac{\left\| \bar{W}_{\setminus r}^n \right\|_{\infty, 2}}{\lambda_n} > \frac{\alpha}{4(2 - \alpha)} \right] \\ & \leq 2 \exp \left(- \frac{\left(\left(\frac{\alpha}{4(2 - \alpha)} \lambda_n^{-\frac{1}{2}} \left\| \bar{\Theta}_{P^c}^* \right\|_1 \right) \sqrt{n - \frac{m-1}{2}} \right)^2}{4} \right. \\ & \quad \left. + \log(p - 1) \right) \\ & \mathbb{P} \left[\frac{\left\| \bar{R}_{\setminus r}^n \right\|_{\infty, 2}}{\lambda_n} > \frac{\alpha}{4(2 - \alpha)} \right] \\ & \leq 2 \exp \left(- \frac{\left(\left(\frac{\alpha}{4(2 - \alpha)} \lambda_n^{-\frac{1}{2}} \left\| \bar{\Theta}_{P^c}^* \right\|_1 \right) \sqrt{n - \frac{m-1}{2}} \right)^2}{4} \right. \\ & \quad \left. + \log(p - 1) \right). \end{aligned}$$

Proof. By simple derivation, we have

$$\begin{aligned} & \frac{\partial}{\partial \theta_{rt; \ell k}^*} \ell^{(i)}(\bar{\Theta}_P; D) = \mathcal{I} \left[x_t^{(i)} = k \right] \\ & \quad \left(\mathcal{I} \left[x_r^{(i)} = \ell \right] - \mathbb{P}_{\bar{\Theta}_P^*} \left[X_r = \ell \mid X_{\setminus r} = x_{\setminus r}^{(i)} \right] \right). \end{aligned}$$

It is easy to show that

$$\begin{aligned} & \mathbb{E}_{\bar{\Theta}_{\setminus r}^*} \left[\frac{\partial}{\partial \theta_{rt; \ell k}^*} \ell^{(i)}(\bar{\Theta}_P; D) \right] \\ &= \mathbb{P}_{\bar{\Theta}_{\setminus r}^*} \left[X_r = \ell \mid X_t = k, X_{\setminus r, t} = x_{\setminus r, t} \right] \\ & \quad - \mathbb{P}_{\bar{\Theta}_P^*} \left[X_r = \ell \mid X_t = k, X_{\setminus r, t} = x_{\setminus r, t} \right] \\ & \leq \left\| \bar{\Theta}_{P^c}^* \right\|_1 \\ & \quad \max_{\beta \in [0, 1]} \left\| \nabla \mathbb{P}_{\bar{\Theta}_{\setminus r}^* - \beta \bar{\Theta}_{P^c}^*} \left[X_r = \ell \mid X_t = k, X_{\setminus r, t} = x_{\setminus r, t} \right] \right\|_{\infty} \\ & \leq \frac{1}{4} \left\| \bar{\Theta}_{P^c}^* \right\|_1, \end{aligned}$$

where, with abuse of notation $\bar{\Theta}_{\setminus r}^* - \beta \bar{\Theta}_{P^c}^*$ represents the matrix $\bar{\Theta}_{\setminus r}^*$ perturbed only on the entries corresponding to $\bar{\Theta}_{P^c}^*$. Also, one can show that $\text{Var} \left(\frac{\partial}{\partial \theta_{rt; \ell k}^*} \ell^{(i)}(\bar{\Theta}_{\setminus r}; D) \right) \leq \frac{1}{4}$. Consequently, with i.i.d assumption on drawn samples, we have $\text{Var} \left(\frac{\partial}{\partial \theta_{rt; \ell k}^*} \ell(\bar{\Theta}_{\setminus r}; D) \right) \leq \frac{1}{4n}$. For a fixed $t \in V \setminus \{r\}$

by Jensen's inequality,

$$\begin{aligned} \mathbb{E}_{\Theta_{\setminus r}^*} \left[\left\| \frac{\partial}{\partial \bar{\theta}_{rt;lk}^*} \ell(\Theta_{\setminus r}; D) \right\|_2 \right] &\leq \sqrt{\mathbb{E}_{\Theta_{\setminus r}^*} \left[\left\| \frac{\partial}{\partial \bar{\theta}_{rt;lk}^*} \ell(\Theta_{\setminus r}; D) \right\|_2^2 \right]} \\ &\leq \frac{1}{2} \sqrt{\frac{(m-1)^2}{n} + \|\bar{\Theta}_{P^c}^*\|_1^2} \\ &\leq \frac{m-1}{2\sqrt{n}} + \frac{1}{2} \|\bar{\Theta}_{P^c}^*\|_1. \end{aligned}$$

We have $\max_{t \in V \setminus \{r\}} \left\| \frac{\partial}{\partial \bar{\theta}_{rt;lk}^*} \ell^{(i)}(\Theta_{\setminus r}; D) \right\|_2 \leq \sqrt{2}$ for all i and hence, by Azuma-Hoeffding inequality and the union bound, we get

$$\mathbb{P} \left[\left\| \frac{\partial}{\partial \bar{\theta}_{rt;lk}^*} \ell(\Theta_{\setminus r}; D) \right\|_{\infty, 2} > \frac{m-1}{2\sqrt{n}} + \frac{1}{2} \|\bar{\Theta}_{P^c}^*\|_1 + \epsilon \right] \leq 2 \exp \left(-\frac{\epsilon^2}{4} n + \log(p-1) \right).$$

For $\lambda_n \geq \frac{8(2-\alpha)}{\alpha} \left(\frac{m-1}{4\sqrt{n}} + \frac{1}{4} \|\bar{\Theta}_{P^c}^*\|_1 \right)$, the result follows.

In order to bound $\bar{R}_{\setminus r}^n$, we need to control the estimation error $\left(\hat{\Theta}_{\setminus r} \right)_{S_r} - \left(\bar{\Theta}_P^* \right)_{S_r}$. Let $H : \mathbb{R}^{(m-1)^2 d_r} \rightarrow \mathbb{R}$ be a function defined as

$$\begin{aligned} H(U_{S_r}) &:= \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r} + U_{S_r}; D \right) - \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r}; D \right) \\ &\quad + \lambda_n \left(\left\| \left(\bar{\Theta}_P^* \right)_{S_r} + U_{S_r} \right\|_{1,2} - \left\| \left(\bar{\Theta}_P^* \right)_{S_r} \right\|_{1,2} \right). \end{aligned}$$

By optimality of $\hat{\Theta}_{\setminus r}$, it is clear that $U^* = \left(\hat{\Theta}_{\setminus r} \right)_{S_r} - \left(\bar{\Theta}_P^* \right)_{S_r}$ minimizes H . Since $H(\mathbf{0}) = 0$ by construction, we have $H(U^*) \leq 0$. Suppose there exist an ℓ_∞/ℓ_2 ball with radius B_r such that for any $\|U\|_{\infty, 2} = B_r$, we have that $H(U) > 0$. Then, we can claim that $\|U^*\|_{\infty, 2} \leq B_r$. See proof of Lemma 6 for more discussion on this proof technique. Let $U_0 \in \mathbb{R}^{(m-1)^2 d_r}$ be an arbitrary vector with $\|U_0\|_{\infty, 2} = \frac{5}{C_{\min}} \lambda_n$. We have

$$\begin{aligned} H(U_0) &:= \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r} + U_0; D \right) - \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r}; D \right) \\ &\quad + \lambda_n \left(\left\| \left(\bar{\Theta}_P^* \right)_{S_r} + U_0 \right\|_{1,2} - \left\| \left(\bar{\Theta}_P^* \right)_{S_r} \right\|_{1,2} \right). \end{aligned} \tag{22}$$

We bound each of these three terms individually. Applying mean value theorem to the log likelihood func-

tion, for some $\beta \in [0, 1]$, we get

$$\begin{aligned} &\ell \left(\left(\bar{\Theta}_P^* \right)_{S_r} + U_0; D \right) - \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r}; D \right) \\ &= \left\langle \left(\bar{W}_{\setminus r}^n \right)_{S_r}, U_0 \right\rangle + \left\langle U_0, \nabla^2 \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r} + \beta U_0; D \right) U_0 \right\rangle. \end{aligned}$$

Note that $\frac{\alpha}{4(2-\alpha)} \lambda_n \leq \frac{1}{4} \lambda_n$ and hence, by our bound on $\bar{W}_{\setminus r}^n$ and Cauchy-Schwartz inequality, we have

$$\begin{aligned} \left| \left\langle \left(\bar{W}_{\setminus r}^n \right)_{S_r}, U_0 \right\rangle \right| &\leq \left\| \left(\bar{W}_{\setminus r}^n \right)_{S_r} \right\|_{\infty, 2} \|U_0\|_{1,2} \\ &\leq \frac{\lambda_n}{4} d_r \|U_0\|_{\infty, 2} \\ &\leq \frac{5}{4C_{\min}} \lambda_n^2 d_r. \end{aligned}$$

To bound the other term, by Taylor expansion, we get

$$\begin{aligned} &\Lambda_{\min} \left(\nabla^2 \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r} + \beta U_0; D \right) \right) \\ &\geq \min_{\beta \in [0, 1]} \Lambda_{\min} \left(\nabla^2 \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r} + \beta U_0; D \right) \right) \\ &\geq \Lambda_{\min} \left(\nabla^2 \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r}; D \right) \right) \\ &\quad - \max_{\beta \in [0, 1]} \Lambda_{\max} \left(\left\langle \frac{\partial \nabla^2 \ell \left(\left(\bar{\Theta}_P^* \right)_{S_r}; D \right)}{\partial \left(\bar{\Theta}_P^* \right)_{S_r}} \Big|_{\left(\bar{\Theta}_P^* \right)_{S_r} + \beta U_0}, U_0 \right\rangle \right) \\ &\geq C_{\min} \\ &\quad - \max_{t_3 \in V \setminus \{r\}} \left\| \frac{\partial \eta_{\ell_1 \ell_2} (x^{(i)})}{\partial \bar{\theta}_{rt_3; \ell_3 k_3}} \right\|_2 d_r \Lambda_{\max}(\mathfrak{S}^*) \|U_0\|_{\infty, 2} \\ &\geq C_{\min} - \frac{m-1}{\sqrt{2}} d_r D_{\max} \|U_0\|_{\infty, 2} \\ &\geq \frac{C_{\min}}{2} \left(\lambda_n d_r \leq \frac{C_{\min}^2}{\sqrt{50}(m-1)D_{\max}} \right). \end{aligned} \tag{23}$$

Here, we used the fact that $\Lambda_{\max}(\mathfrak{S}^*) = \Lambda_{\max}(\mathcal{J}^*)$ as a property of Kronecher product and also our assumption on the maximum eigenvalue of \mathcal{J}^* . By triangle inequality,

$$\begin{aligned} \lambda_n \left(\left\| \left(\bar{\Theta}_P^* \right)_{S_r} + U_0 \right\|_{1,2} - \left\| \left(\bar{\Theta}_P^* \right)_{S_r} \right\|_{1,2} \right) &\geq -\lambda_n \|U_0\|_{1,2} \\ &\geq -\lambda_n d_r \|U_0\|_{\infty, 2} \\ &\geq -\frac{5\lambda_n^2 d_r}{C_{\min}}. \end{aligned}$$

Hence, from (22), we get $H(U_0) \geq \frac{5\lambda_n^2 d_r}{4C_{\min}} > 0$ and hence,

$$\left\| \left(\hat{\Theta}_{\setminus r} \right)_{S_r} - \left(\bar{\Theta}_P^* \right)_{S_r} \right\|_{\infty, 2} \leq \frac{5}{C_{\min}} \lambda_n, \tag{24}$$

with high probability. With similar analysis as in 23,

we have

$$\begin{aligned}
 & \frac{\|\bar{R}_{\setminus r}^n\|_{\infty,2}}{\lambda_n} \\
 & \leq \frac{1}{\lambda_n} \frac{m-1}{\sqrt{2}} d_r D_{\max} \left\| \left(\hat{\Theta}_{\setminus r} \right)_{S_r} - \left(\Theta_{\setminus r}^* \right)_{S_r} \right\|_{\infty,2}^2 \\
 & \leq \frac{m-1}{\sqrt{2}} d_r D_{\max} \frac{25}{C_{\min}^2} \lambda_n \\
 & \leq \frac{\alpha}{4(2-\alpha)},
 \end{aligned}$$

provided that $\lambda_n d_r \leq \frac{C_{\min}^2}{50\sqrt{2}(m-1)D_{\max}} \frac{\alpha}{2-\alpha}$.

□