

# Supplementary materials

## Appendix A: Proofs of Theorems

For any two real number-valued  $p \times q$  matrices  $A = (A_{ij})$  and  $B = (B_{ij})$ , define their Frobenius inner product as  $\langle A, B \rangle_F = \sum_{i,j} A_{ij} B_{ij}$ . For any matrix  $\mathbf{M}$ , we use  $\|\mathbf{M}\|_F$  and  $\|\mathbf{M}\|_{op}$  to denote its Frobenius norm and operator norm, respectively. For the sake of clarity, we define  $\Omega_{r,s} = \{\mathbf{C} : \text{rank}(\mathbf{C}) = r, \|\mathbf{C}\|_{2,0} \leq s\}$  as the set of the low-rank and row-sparse matrices and  $B_F(r_0, \mathbf{C}^*)$  as a ball centered at  $\mathbf{C}^*$  with radius  $r_0$ , which will be used later.

### A.1. Proof of Proposition 1

To begin with, we consider the second-order expansion of the loss function  $L(\mathbf{C}; \mathcal{D})$  by applying Taylor's theorem with Lagrange remainder. That is, fix the current estimate  $\mathbf{C}^m$ , and suppose  $\mathbf{C}$  is in a neighborhood of  $\mathbf{C}^m$ , then there exists  $0 < \theta < 1$  such that

$$L(\mathbf{C}; \mathcal{D}) = L(\mathbf{C}^m; \mathcal{D}) - \frac{1}{n} \left\langle \mathbf{X}(\mathbf{C} - \mathbf{C}^m), (\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^m}) \right\rangle_F + \frac{1}{2n} \sum_{i=1}^n \mathbf{x}_i^\top (\mathbf{C} - \mathbf{C}^m) \mathbf{H}_i^{\mathbf{C}'} (\mathbf{C} - \mathbf{C}^m)^\top \mathbf{x}_i,$$

where  $\mathbf{H}_i^{\mathbf{C}'}$  is the within observation Hessian matrix for the coefficient matrix  $\mathbf{C}' = \theta \mathbf{C} + (1 - \theta) \mathbf{C}^m$ ,  $i = 1, \dots, n$ . By the definition of  $S(\mathbf{C}; \mathcal{D} | \mathbf{C}^{(m)})$ , we have

$$S(\mathbf{C}; \mathcal{D} | \mathbf{C}^{(m)}) - L(\mathbf{C}; \mathcal{D}) = \frac{1}{2n} \sum_{i=1}^n \mathbf{x}_i^\top (\mathbf{C} - \mathbf{C}^m) \left( \frac{1}{2} I_{q-1} - \mathbf{H}_i^{\mathbf{C}'} \right) (\mathbf{C} - \mathbf{C}^m)^\top \mathbf{x}_i.$$

In [Simon et al. \(2013\)](#) Lemma 3.3 has proved, for any coefficient matrix  $\mathbf{C}'$ ,

$$\mathbf{H}_i^{\mathbf{C}'} \preceq 2 \max_{ij} \{ \mathbf{P}_{ij}^{\mathbf{C}'} (1 - \mathbf{P}_{ij}^{\mathbf{C}'}) \} I \preceq \frac{1}{2} I_{q-1}.$$

Therefore, we have

$$S(\mathbf{C}; \mathcal{D} | \mathbf{C}^{(m)}) \geq L(\mathbf{C}; \mathcal{D})$$

for all  $\mathbf{C}$  around the current estimate  $\mathbf{C}^m$ . Finally, the equality  $S(\mathbf{C}^m; \mathcal{D} | \mathbf{C}^{(m)}) = L(\mathbf{C}^m; \mathcal{D})$  holds by applying the definition of  $S(\mathbf{C}; \mathcal{D} | \mathbf{C}^{(m)})$ .

### A.2. Proof of Theorem 1

Recall that  $\hat{\mathbf{C}}$  is the minimizer of (4), then by the definition of surrogate function, we have

$$\hat{\mathbf{C}} = \arg \min_{\mathbf{C} \in \Omega_r} S(\mathbf{C}; \mathcal{D} | \hat{\mathbf{C}}).$$

Therefore, we can derive an explicit expression for  $\hat{\mathbf{C}}$  as follows.

$$\begin{aligned} \hat{\mathbf{C}}_{\hat{\mathcal{A}}} &= (\mathbf{X}_{\hat{\mathcal{A}}}^\top \mathbf{X}_{\hat{\mathcal{A}}})^{-1} \mathbf{X}_{\hat{\mathcal{A}}}^\top \mathbf{M}^{\hat{\mathbf{C}}} \mathbf{V}^{\hat{\mathbf{C}}} (\mathbf{V}^{\hat{\mathbf{C}}})^\top, \\ \hat{\mathbf{C}}_{(\hat{\mathcal{A}})^c} &= \mathbf{0}, \end{aligned}$$

where  $M^{\hat{C}} = \mathbf{X}\hat{C} + 2(\tilde{Y} - \tilde{P}^{\hat{C}})$  and  $\mathbf{V}^{\hat{C}} \in \mathbb{R}^{q \times r}$  consists of the eigenvectors that corresponds to the top  $r$ -th largest eigenvalues of  $(M^{\hat{C}})^\top \mathbf{X}_{\cdot \hat{A}} (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top M^{\hat{C}}$ . It can be seen that  $(\mathbf{V}^{\hat{C}})^\top \mathbf{V}^{\hat{C}} = I_r$ . Firstly, we take difference between  $\hat{C}_{\hat{A}}$  and  $C_{\hat{A}}^*$ , i.e.,

$$\begin{aligned} \|\hat{C}_{\hat{A}} - C_{\hat{A}}^*\|_F &= \|(\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top M^{\hat{C}} \mathbf{V}^{\hat{C}} (\mathbf{V}^{\hat{C}})^\top - C_{\hat{A}}^*\|_F \\ &= \|\hat{C}_{\hat{A}} \mathbf{V}^{\hat{C}} (\mathbf{V}^{\hat{C}})^\top - C_{\hat{A}}^* + (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\tilde{Y} - \tilde{P}^{\hat{C}}) \mathbf{V}^{\hat{C}} (\mathbf{V}^{\hat{C}})^\top\|_F \\ &= \|\hat{C}_{\hat{A}} \mathbf{V}^{\hat{C}} - C_{\hat{A}}^* \mathbf{V}^{\hat{C}} + (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\tilde{Y} - \tilde{P}^{\hat{C}}) \mathbf{V}^{\hat{C}}\|_F \\ &= \|\hat{C}_{\hat{A}} \mathbf{V}^{\hat{C}} - C_{\hat{A}}^* \mathbf{V}^{\hat{C}} + (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\tilde{Y} - \tilde{P}^{C^*} + \tilde{P}^{C^*} - \tilde{P}^{\hat{C}}) \mathbf{V}^{\hat{C}}\|_F \\ &\leq \|\hat{C}_{\hat{A}} - C_{\hat{A}}^*\|_F + (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2\|\tilde{Y} - \tilde{P}^{C^*} + \tilde{P}^{C^*} - \tilde{P}^{\hat{C}}\|_F. \end{aligned}$$

Using the triangle inequality gives

$$\|\hat{C}_{\hat{A}} - C_{\hat{A}}^*\| \leq \|\hat{C}_{\hat{A}} - C_{\hat{A}}^* - (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\tilde{P}^{\hat{C}} - \tilde{P}^{C^*})\|_F + \|(\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\tilde{Y} - \tilde{P}^{C^*})\|_F. \quad (1)$$

We know that the  $(i, k)$ -th element of  $\tilde{P}^{\hat{C}}$  is actually a function of  $c_k^\top \mathbf{x}_i$  with other elements fixed. Let  $f(\cdot)$  denote this function, i.e.,  $f(x) = \frac{\exp(x)}{\alpha + \exp(x)}$ , where  $\alpha$  is a constant. Then it is easy to show that  $f(x)$  has the Lipschitz continuity, that is,

$$\tilde{c} \|\mathbf{X}(\hat{C} - C^*)\|_F \leq \|\tilde{P}^{\hat{C}} - \tilde{P}^{C^*}\|_F \leq \frac{1}{4} \|\mathbf{X}(\hat{C} - C^*)\|_F.$$

Therefore, each element of  $\tilde{P}^{\hat{C}} - \tilde{P}^{C^*}$ , denoted by  $\Delta_{i,j}$ , satisfies  $\frac{1}{4} |(\mathbf{X}\hat{C} - \mathbf{X}C^*)_{i,j}| \geq |\Delta_{i,j}| \geq \tilde{c} |(\mathbf{X}\hat{C} - \mathbf{X}C^*)_{i,j}|$ . For convenience, we denote  $\Delta = (\Delta_{i,j})$  with  $\Delta_{i,j} = \alpha_{i,j} (\mathbf{X}\hat{C} - \mathbf{X}C^*)_{i,j}$  and  $\alpha_{i,j}$  lies in  $[\tilde{c}, \frac{1}{4}]$ .

Next we divide the columns of  $\mathbf{X}$  by the sets of  $\hat{A}$  and  $\mathcal{A}^*$ . In particular, define three instrumental matrices  $\mathbf{X}_1, \mathbf{X}_2$  and  $\mathbf{X}_3$  as following

$$\begin{aligned} (\mathbf{X}_1)_{\cdot \hat{A}} &= \mathbf{X}_{\cdot \hat{A}}, \\ (\mathbf{X}_1)_{\cdot \hat{A}^c} &= \mathbf{0}, \\ (\mathbf{X}_2)_{\cdot \mathcal{A}^* - \hat{A}} &= \mathbf{X}_{\cdot \mathcal{A}^* - \hat{A}}, \\ (\mathbf{X}_2)_{\cdot (\mathcal{A}^* - \hat{A})^c} &= \mathbf{0}, \\ (\mathbf{X}_3)_{\cdot (\mathcal{A}^* \cup \hat{A})^c} &= \mathbf{X}_{\cdot (\mathcal{A}^* - \hat{A})^c}, \\ (\mathbf{X}_3)_{\cdot \mathcal{A}^* \cup \hat{A}} &= \mathbf{0}. \end{aligned}$$

Similarly, we can separate the matrix  $\Delta$  into three instrumental matrices,  $\Delta_1, \Delta_2$ , and  $\Delta_3$ , in the following way.

$$\begin{aligned} (\Delta_1)_{i,j} &= \alpha_{i,j} (\mathbf{X}_1 \hat{C} - \mathbf{X}_1 C^*)_{i,j}, \\ (\Delta_2)_{i,j} &= \alpha_{i,j} (\mathbf{X}_2 \hat{C} - \mathbf{X}_2 C^*)_{i,j}, \\ (\Delta_3)_{i,j} &= \alpha_{i,j} (\mathbf{X}_3 \hat{C} - \mathbf{X}_3 C^*)_{i,j}. \end{aligned}$$

Notice that  $\Delta_3 \equiv \mathbf{0}$  and we define  $\Delta_3$  to provide a symmetric definition for  $\Delta_1$ . Since  $\Delta = \Delta_1 + \Delta_2 + \Delta_3$ , we can rewrite the inequality in (1) as

$$\begin{aligned} \|\hat{C}_{\hat{A}} - C_{\hat{A}}^*\| &\leq \|\hat{C}_{\hat{A}} - C_{\hat{A}}^* - (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\Delta_1 + \Delta_2 + \Delta_3)\|_F \\ &\quad + \|(\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\tilde{Y} - \tilde{P}^{C^*})\|_F \\ &\leq \|\hat{C}_{\hat{A}} - C_{\hat{A}}^* - (\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2\Delta_1\|_F + \|(\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2\Delta_2\|_F \\ &\quad + \|(\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2\Delta_3\|_F + \|(\mathbf{X}_{\cdot \hat{A}}^\top \mathbf{X}_{\cdot \hat{A}})^{-1} \mathbf{X}_{\cdot \hat{A}}^\top 2(\tilde{Y} - \tilde{P}^{C^*})\|_F, \end{aligned} \quad (2)$$

where a triangle inequality is used in the last inequality. For the first term, by the definition  $\Delta_1$ , we know that each element of the matrix  $(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top \Delta_1$  is actually  $\alpha_{i,j}$ . Thus, we can its upper bound by

$$\|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^* - (\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2\Delta_1\|_F \leq (1 - 2\tilde{c}) \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F.$$

For the second term, applying Conditions 3.1-3.2 as well as the definition of  $\Delta_2$ , we can obtain

$$\begin{aligned} & \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2\Delta_2\|_F \\ & \leq 2 \times \frac{1}{4} \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top \mathbf{X}_{A^*-\hat{A}} (\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*)\|_F \\ & \leq \frac{\theta_s}{2c_-(s)} \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F. \end{aligned}$$

Therefore, the inequality (2) becomes,

$$\begin{aligned} \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\| & \leq (1 - 2\tilde{c}) \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F + \frac{\theta_s}{2c_-(s)} \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F + \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F \\ & \leq (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)}) \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F + \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F, \end{aligned}$$

where the last inequality holds due to the fact that  $\hat{\mathbf{C}}$  lies in the cone set  $\Lambda(s)$ . Based on the above result, we can obtain the upper bound of  $\|\hat{\mathbf{C}} - \mathbf{C}^*\|_F^2$  as follows.

$$\begin{aligned} \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F^2 & = \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F^2 + \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F^2 \\ & \leq (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2 \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F^2 + \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F^2 \\ & \quad + (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)}) \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F \\ & \quad + \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F^2 \\ & \leq (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2 \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F^2 + \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F^2 \\ & \quad + (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)}) \frac{\bar{c}}{\sqrt{1 + \bar{c}^2}} \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F \\ & \quad + \|(\mathbf{X}_{\hat{A}}^\top \mathbf{X}_{\hat{A}})^{-1} \mathbf{X}_{\hat{A}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F^2 \end{aligned}$$

Since  $1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)} < 1$ , we can obtain

$$\begin{aligned} & (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2 \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F^2 + \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F^2 \\ & < \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F^2 + \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F^2 \\ & = \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F^2. \end{aligned}$$

Let  $\alpha_d = \frac{(1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2 \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F^2 + \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F^2}{\|\hat{\mathbf{C}} - \mathbf{C}^*\|_F^2} < 1$ . After some simplifications, we get

$$(\alpha_d - (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2) \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F = (1 - \alpha_d) \|\hat{\mathbf{C}}_{A^*-\hat{A}} - \mathbf{C}_{A^*-\hat{A}}^*\|_F \leq \bar{c}(1 - \alpha_d) \|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F.$$

Canceling the same terms  $\|\hat{\mathbf{C}}_{\hat{A}} - \mathbf{C}_{\hat{A}}^*\|_F$ , we get

$$\alpha_d \leq \frac{\bar{c} + (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2}{\bar{c} + 1} < 1.$$

Without loss of generality, we set

$$\alpha_d = \frac{\bar{c} + (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2}{\bar{c} + 1}.$$

Therefore,

$$\begin{aligned} \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F^2 &\leq \alpha_d \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F^2 + (1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)}) \frac{\bar{c}}{\sqrt{1 + \bar{c}^2}} \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F \|(\mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top \mathbf{X}_{\cdot \hat{\mathcal{A}}})^{-1} \mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F \\ &\quad + \|(\mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top \mathbf{X}_{\cdot \hat{\mathcal{A}}})^{-1} \mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_F^2. \end{aligned}$$

Solving the above quadratic inequality in  $\|\hat{\mathbf{C}} - \mathbf{C}^*\|_F$  yields

$$\begin{aligned} \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F &\leq \frac{-(1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)}) + \sqrt{(1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2 + 4(1 - \alpha_d)}}{2(1 - \alpha_d)} \left\| (\mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top \mathbf{X}_{\cdot \hat{\mathcal{A}}})^{-1} \mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*}) \right\|_F, \\ &\leq \frac{-(1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)}) + \sqrt{(1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2 + 4(1 - \alpha_d)}}{2(1 - \alpha_d)} \sqrt{r} \left\| (\mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top \mathbf{X}_{\cdot \hat{\mathcal{A}}})^{-1} \mathbf{X}_{\cdot \hat{\mathcal{A}}}^\top 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*}) \right\|_{op}, \end{aligned}$$

Then applying Condition 3.1 and Lemma 5, with probability at least  $1 - \delta$  for any  $\delta > 0$ , we have

$$\|\hat{\mathbf{C}} - \mathbf{C}^*\|_F \leq \frac{-(1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)}) + \sqrt{(1 - 2\tilde{c} + \bar{c} \frac{\theta_s}{2c_-(s)})^2 + 4(1 - \alpha_d)}}{(1 - \alpha_d)c_-(s)} \sqrt{\frac{3rqsP^2}{n} \log \frac{2pq}{\delta}}.$$

This completes the proof.

### A.3. Proof of Theorem 2

Recall that  $\mathbf{C}^*$  is the true coefficient matrix, and  $\mathbf{C}^{m,k}$  is the estimate at the  $k$ -th iteration of the inner loop of the  $m$ -th step of the outer loop in **Algorithm 2**. To prove Theorem 2, we first define several immediate variables. Let  $\hat{\mathbf{C}}^0 = \mathbf{C}^0$  denote the initialized value of  $\mathbf{C}$  in **Algorithm 2**. Then at the  $m$ -th step of **Algorithm 2** ( $m \geq 1$ ), define

$$\hat{\mathbf{C}}^m = \arg \min_{\mathbf{C} \in \Omega_r} S(\mathbf{C}; \mathcal{D} | \hat{\mathbf{C}}^{m-1}).$$

Denote the minimizer of (4) as  $\hat{\mathbf{C}}$ , then by the definition of surrogate function, we have

$$\hat{\mathbf{C}} = \arg \min_{\mathbf{C} \in \Omega_r} S(\mathbf{C}; \mathcal{D} | \hat{\mathbf{C}}).$$

Next we define the noise-free version of the surrogate function. For any coefficient matrix  $\mathbf{C}'$  in  $\Omega_r$ , define the noise-free surrogate function as

$$\tilde{S}(\mathbf{C}; \mathcal{D} | \mathbf{C}') = L(\mathbf{C}'; \mathcal{D}) + \frac{1}{4n} \|\mathbf{X}(\mathbf{C} - \mathbf{C}') - 2(\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\mathbf{C}'})\|_F^2 - \frac{1}{n} \|\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\mathbf{C}'}\|_F. \quad (3)$$

It is easy to show that the true coefficient matrix  $\mathbf{C}^*$  is actually the minimizer of the noise-free surrogate function around  $\mathbf{C}^*$ , that is,

$$\mathbf{C}^* = \arg \min_{\mathbf{C} \in \Omega_r} \tilde{S}(\mathbf{C}; \mathcal{D} | \mathbf{C}^*).$$

Similarly, we can define the “true” coefficient matrix  $\mathbf{C}^{m*}$  at the  $m$ -th step of **Algorithm 2**,

$$\mathbf{C}^{m*} = \arg \min_{\mathbf{C} \in \Omega_r} \tilde{S}(\mathbf{C}; \mathcal{D} | \hat{\mathbf{C}}^{m-1}).$$

The proof of Theorem 2 is based on three following lemmas.

**LEMMA 1.** *Assume that Conditions 3.1-3.2 hold. Then for the output estimate  $\mathbf{C}^{m,k}$  at the  $k$ -th inner loop iteration of the  $m$ -th step of **Algorithm 2** with  $s \geq s^*$ ,  $r \geq r^*$ , we have*

$$\begin{aligned} \|\mathbf{C}^{m,k} - \mathbf{C}^{m*}\|_{op} &\leq \alpha_1 \gamma^k \|\mathbf{C}^{m*}\|_{op} + (\alpha_2 + \alpha_3 \sqrt{r}) h(s), \\ \|\mathbf{C}^{m,k} - \mathbf{C}^{m*}\|_F &\leq \alpha_1 \sqrt{2r} \gamma^k \|\mathbf{C}^{m*}\|_{op} + \sqrt{2r} (\alpha_2 + \alpha_3 \sqrt{r}) h(s), \end{aligned} \quad (4)$$

where

$$h(s) = \frac{1}{n} \max_{|\mathcal{A}| \leq s} \|\mathbf{X}_{\cdot \mathcal{A}}^\top \text{Err}^m\|_{op},$$

$$\text{Err}^m = 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*}).$$

In addition, with probability at least  $1 - \delta$  for any constant  $\alpha > 0$  and sufficiently large  $n$ , we have

$$h(s) \leq 5 \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} (\|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F + (1 - 2\zeta(1 - \zeta)) \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F) + 5c_+(s) \|\mathbf{C}^* - \hat{\mathbf{C}}\|_F + 2\sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}}.$$

LEMMA 2. Assume conditions 3.1-3.2 hold, we have

$$\|\hat{\mathbf{C}}^m - \hat{\mathbf{C}}\|_F \leq (1 - 2\zeta(1 - \zeta))^m \frac{c_+(s)}{c_-(s)} \|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F,$$

where

$$\zeta = \min_{\mathbf{C} \in B_F(r_0, \mathbf{C}^*), \mathbf{C} \in \Omega_r, i, j} \tilde{\mathbf{P}}_{i, j}^{\mathbf{C}}.$$

is a constant lying between 0 and 1.

LEMMA 3. Assume conditions 3.1-3.2 hold, we have

$$\|\mathbf{C}^{m*} - \mathbf{C}^*\|_F \leq (1 - 2\zeta(1 - \zeta))^m \frac{c_+(s)}{c_-(s)} \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F + (1 - 2\zeta(1 - \zeta)) \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F,$$

where  $\zeta$  is defined as Lemma 2.

It follows from Lemma 1 and the triangle inequality that

$$\begin{aligned} \|\mathbf{C}^{m, k} - \mathbf{C}^*\|_F &\leq \|\mathbf{C}^{m, k} - \mathbf{C}^{m*}\|_F + \|\mathbf{C}^{m*} - \mathbf{C}^*\|_F \\ &\leq \alpha_1 \sqrt{2r} \gamma^k \|\mathbf{C}^{m*}\|_{op} + \sqrt{2r} (\alpha_2 + \alpha_3 \sqrt{r}) h(s) + \|\mathbf{C}^{m*} - \mathbf{C}^*\|_F \\ &\leq \alpha_1 \sqrt{2r} \gamma^k \|\mathbf{C}^*\|_{op} + \sqrt{2r} (\alpha_2 + \alpha_3 \sqrt{r}) h(s) + (\alpha_1 \sqrt{2r} \gamma^{k+1} + 1) \|\mathbf{C}^{m*} - \mathbf{C}^*\|_F. \end{aligned}$$

Now, we consider bounding the last two terms in the right hand side of the above inequality. In specific, we bound the second stochastic term by Lemma 1, and give an upper bound for the third term by Lemma 3. Thus, with probability at least  $1 - \delta$  for any constant  $\alpha > 0$  and sufficiently large  $n$ , we have

$$\begin{aligned} \|\mathbf{C}^{m, k} - \mathbf{C}^*\|_F &\leq \alpha_1 \sqrt{2r} \gamma^{k+1} \|\mathbf{C}^*\|_{op} + 5\sqrt{2r} (\alpha_2 + \alpha_3 \sqrt{r}) \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} \|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F \\ &\quad + \left( 5\sqrt{2r} (\alpha_2 + \alpha_3 \sqrt{r}) c_+(s) + \alpha_1 \sqrt{2r} \gamma^k + 1 \right) \frac{c_+(s)}{c_-(s)} (1 - 2\zeta(1 - \zeta))^m \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F \quad (5) \\ &\quad + 5\sqrt{2r} (\alpha_2 + \alpha_3 \sqrt{r}) c_+(s) \|\mathbf{C}^* - \hat{\mathbf{C}}\|_F + 2\sqrt{2r} (\alpha_2 + \alpha_3 \sqrt{r}) \sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}}. \end{aligned}$$

For the term  $\|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F$ , applying the triangle inequality and Theorem 1 gives the following result,

$$\begin{aligned} \|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F &\leq \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F + \|\hat{\mathbf{C}} - \mathbf{C}^*\|_F \\ &\leq \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F + \frac{-(1 - 2\tilde{c} + \tilde{c} \frac{\theta_s}{2c_-(s)}) + \sqrt{(1 - 2\tilde{c} + \tilde{c} \frac{\theta_s}{2c_-(s)})^2 + 4(1 - \alpha_d)}}{(1 - \alpha_d) c_-(s)} \sqrt{\frac{3qsP^2}{n} \log \left( \frac{2pq}{\delta} \right)} \end{aligned}$$

with probability at least  $1 - \delta$  for any constant  $\alpha > 0$  and sufficiently large  $n$ . Substituting the bound into (5), we have

$$\begin{aligned}
\|\mathbf{C}^{m,k} - \mathbf{C}^*\|_F &\leq \alpha_1 \sqrt{2r} \gamma^{k+1} \|\mathbf{C}^*\|_{op} + \left[ 5\sqrt{2} \left( \frac{\alpha_2}{\sqrt{r}} + \alpha_3 \right) \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} \right. \\
&\quad + \left. \left( 5\sqrt{2} \left( \frac{\alpha_2}{\sqrt{r}} + \alpha_3 \right) c_+(s) + \alpha_1 \frac{\sqrt{2}}{\sqrt{r}} \gamma^k + \frac{1}{r} \right) \frac{c_+(s)}{c_-(s)} (1 - 2\zeta(1 - \zeta))^m \right] r \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F \\
&\quad + \left( 2 + \left( \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} + c_+(s) \right) 5\alpha \right) \sqrt{2} \left( \frac{\alpha_2}{\sqrt{r}} + \alpha_3 \right) r \sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}} \\
&\leq \alpha_1 \sqrt{2r} \gamma^{k+1} \|\mathbf{C}^*\|_{op} + \left[ 5\sqrt{2} (\alpha_2 + \alpha_3) \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} \right. \\
&\quad + \left. \left( 5\sqrt{2} (\alpha_2 + \alpha_3) c_+(s) + \alpha_1 \sqrt{2} + 1 \right) \frac{c_+(s)}{c_-(s)} (1 - 2\zeta(1 - \zeta))^m \right] r \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F \\
&\quad + \left( 2 + \left( \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} + c_+(s) \right) 5\alpha \right) \sqrt{2} (\alpha_2 + \alpha_3) r \sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}}.
\end{aligned}$$

Taking  $c_2 = \sqrt{2}\alpha_1$ ,  $c_3 = 5\sqrt{2}(\alpha_2 + \alpha_3) \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} + (5\sqrt{2}(\alpha_2 + \alpha_3)c_+(s) + \alpha_1\sqrt{2} + 1) \frac{c_+(s)}{c_-(s)} (1 - 2\zeta(1 - \zeta))^m$ , and leaving  $c_4 = \left( 2 + \left( \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} + c_+(s) \right) 5\alpha \right) \sqrt{2}(\alpha_2 + \alpha_3) \sqrt{3 \log(2)P^2}$  gives Theorem 2.

#### A.4. Proof of Corollary 1

The first part of Corollary 1 follows immediately from Theorem 2. In particular, when  $k \geq \log_{\gamma} \frac{c_4 \sqrt{r}}{2c_2 \|\mathbf{C}^*\|_F} \sqrt{\frac{qs \log pq}{n}}$ , we have

$$\gamma^k \leq \frac{c_4 \sqrt{r}}{2c_2 \|\mathbf{C}^*\|_F} \sqrt{\frac{qs \log pq}{n}},$$

which is equal to

$$\gamma^k c_2 \|\mathbf{C}^*\|_F \leq \frac{c_4 \sqrt{r}}{2} \sqrt{\frac{qs \log pq}{n}}.$$

Similarly, when  $m \geq \log_{\zeta'} \frac{c_4 r}{2c_3 \|\mathbf{C}^0 - \mathbf{C}^*\|_F} \sqrt{\frac{qs \log pq}{n}}$ , we have

$$(\zeta')^m \leq \frac{c_4}{2c_3 \|\mathbf{C}^0 - \mathbf{C}^*\|_F} \sqrt{\frac{qs \log pq}{n}}.$$

which is equal to

$$(1 - 2\zeta(1 - \zeta))^m c_3 r \|\mathbf{C}^0 - \mathbf{C}^*\|_F \leq \frac{c_4 r}{2} \sqrt{\frac{qs \log pq}{n}}.$$

Next, combining the inequality  $\min_{i \in \text{supp}(\mathbf{C}^*)} \|\mathbf{C}_{i \cdot}^*\|_2 \geq 2c_4 r \sqrt{\frac{qs \log(pq)}{n}}$ , we have

$$\text{supp}(\mathbf{C}^*) \subseteq \text{supp}(\mathbf{C}^{m,k}).$$

If not, we have at least one true nonzero being missed in the implementation. Without loss of generality, we assume that the  $i$ -th row is screened out, i.e.,  $i \in \text{supp}(\mathbf{C}^*) \setminus \text{supp}(\mathbf{C}^{m,k})$ . Then

$$\|\mathbf{C}^* - \mathbf{C}^{m,k}\|_F \geq \|\mathbf{C}_{i \cdot}^*\|_2 \geq 2c_4 r \sqrt{\frac{qs \log(pq)}{n}},$$

which leads to a contradiction to the inequality in the first part.

Owing to the same size of  $\text{supp}(\mathbf{C}^*)$  and  $\text{supp}(\mathbf{C}^{m,k})$ , we conclude that

$$\text{supp}(\mathbf{C}^*) = \text{supp}(\mathbf{C}^{m,k}).$$

## Appendix B: Proofs of lemmas

### B.1. Proof of Lemma 1

The proof of Lemma 1 follows similar lines of the proof of Theorem 3.1 of Wen et al. (2022), with some minor modifications. Thus, we introduce the main difference and present an outline of the proof for the first part of Lemma 1. The main difference is the model setting: while a reduced rank regression model is considered in Wen et al. (2022), our method is built on the reduced rank multinomial logistic regression model. This difference leads to the difference in error term. In particular, the error term  $Err^m$  here is defined as

$$Err^m = 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*}).$$

In contrast, the error term in Wen et al. (2022) is  $\mathbf{Y} - \mathbf{X}\mathbf{C}^*$ . The involvement of  $\hat{\mathbf{C}}^{m-1}$  in the definition of error leads a much more complicated expression.

Before presenting the proof, we introduce some notations and lemmas needed. Let the row support set of  $\mathbf{C}^{m,k}$  be  $\mathcal{A}^{m,k}$ , and the set of zero rows be  $\mathcal{I}^{m,k} = (\mathcal{A}^{m,k})^c$ . We define the auxiliary index sets by

$$\begin{aligned} \mathcal{A}_{11}^{m,k} &= \mathcal{A}^* \cap \mathcal{A}^{m,k} \cap \mathcal{I}^{m,k+1}, & \mathcal{A}_{22}^{m,k} &= \mathcal{A}^* \cap \mathcal{I}^{m,k} \cap \mathcal{I}^{m,k+1}, \\ \mathcal{I}_{11}^{m,k} &= \mathcal{I}^* \cap \mathcal{A}^{m,k} \cap \mathcal{A}^{m,k+1}, & \mathcal{I}_{22}^{m,k} &= \mathcal{I}^* \cap \mathcal{I}^{m,k} \cap \mathcal{A}^{m,k+1}. \end{aligned}$$

We can see that  $\mathcal{A}_{11}^{m,k}$  is the subset of the true nonzero rows that are kicked out at the  $(m, k+1)$ -th iteration and selected at  $(m, k)$ -th iteration at the same time, and  $\mathcal{A}_{11}^{m,k}$  is the subset of the true nonzero rows that are kicked out at both the  $(m, k+1)$ -th iteration and the  $(m, k)$ -th iteration. Similarly,  $\mathcal{A}_{11}^{m,k}$  is the subset of the true nonzero rows that are kicked out at the  $(m, k+1)$ -th iteration and selected at  $(m, k)$ -th iteration at the same time, and  $\mathcal{A}_{11}^{m,k}$  is the subset of the true nonzero rows that are kicked out at both the  $(m, k+1)$ -th iteration and the  $(m, k)$ -th iteration. In fact, they are generalization of false positive and false negative for the switch from the  $(m, k)$ -th to the  $(m, k+1)$ -th iteration. We further define the difference in  $\mathbf{C}^*$  as  $D(\mathcal{A}^k) = \|\mathbf{C}_{\mathcal{A}^* \cap \mathcal{I}^{m,k}}^*\|_{op}$ . For convenience, we introduce some more definition as follows.

$$\Delta(\mathbf{C}^{m,k+1}) = \mathbf{C}_{\mathcal{A}^{m,k}}^{m,k+1} - \mathbf{C}_{\mathcal{A}^{m,k}}^{m*}, \quad \tilde{\Gamma}_{\mathcal{A}_{22}^{m,k}}^{k+1} = \frac{1}{n} \mathbf{X}_{\mathcal{A}_{22}^{m,k}}^\top (2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}\hat{\mathbf{C}}^{m-1} - \mathbf{X}\mathbf{C}^{m,k+1}).$$

An outline of the proof is given below.

1. Construct inequality to bound  $\|\mathbf{C}^{m,k} - \mathbf{C}^{m*}\|_{op}$  by  $D(\mathcal{A}^k)$  as follows.

$$\begin{aligned} \|\mathbf{C}_{\mathcal{A}^{m,k}}^{m,k+1} - \mathbf{C}_{\mathcal{A}^{m,k}}^{m*}\|_{op} &\leq \frac{\theta_s}{c_-(s)} D(\mathcal{A}^k) + \frac{2}{nc_-(s)} \|\mathbf{X}_{\mathcal{A}^{m,k}} Err^m\|_{op}, \\ \|\mathbf{C}^{m,k+1} - \mathbf{C}^{m*}\|_{op} &\leq (1 + \frac{\theta_s}{c_-(s)}) D(\mathcal{A}^k) + \frac{2}{nc_-(s)} \|\mathbf{X}_{\mathcal{A}^{m,k}} Err^m\|_{op}. \end{aligned}$$

This is the same with Lemma B.1 in Wen et al. (2022) except the definition of  $Err^m$ .

2. Develop the iterative inequality for  $D(\mathcal{A}^{k+1})$ ,

$$D(\mathcal{A}^{k+1}) \leq \gamma D(\mathcal{A}^k) + \alpha h(s).$$

To get the above result, we first need to prove three following results. The first is Lemma B.2 of Wen et al. (2022), which are listed below.

$$\begin{aligned} D(\mathcal{A}^{k+1}) &\leq \|\mathbf{C}_{\mathcal{A}_{11}^k}^{m*}\|_{op} + \|\mathbf{C}_{\mathcal{A}_{22}^k}^{m*}\|_{op}, \\ \|\mathbf{C}_{\mathcal{A}_{11}^k}^{m*}\|_{op} &\leq \|\Delta(\mathbf{C}^{m,k+1})_{\mathcal{A}_{11}^k}\|_{op} + \|\mathbf{C}_{\mathcal{A}_{11}^k}^{m,k+1}\|_{op}, \\ \|\mathbf{C}_{\mathcal{A}_{22}^k}^{m*}\|_{op} &\leq \frac{1}{c_-(s)} (\|\tilde{\Gamma}_{\mathcal{A}_{22}^k}^{k+1}\|_{op} + \theta_s \|\Delta(\mathbf{C}^{m,k+1})\|_{op} + \theta_s D(\mathcal{A}^k) + h(s)). \end{aligned}$$

The second is Lemma B.3 of [Wen et al. \(2022\)](#) stated as

$$\|\mathbf{B}_{\mathcal{A}_{11}^{m,k}}^{k+1}\|_{op} + \|\Gamma_{\mathcal{A}_{22}^{m,k}}^{k+1}\|_{op} \leq \sqrt{r}(\|\mathbf{C}_{\mathcal{I}_{11}^{m,k}}^{k+1}\|_{op} + \|\tilde{\Gamma}_{\mathcal{I}_{22}^{m,k}}^{k+1}\|_{op}).$$

The last one is Lemma 4, which is a modified version of Lemma B.4 of [Wen et al. \(2022\)](#). The proof of Lemma 4 is presented in Appendix B.4.

LEMMA 4. *Suppose Condition 3.1 and 3.2 hold, we have*

$$\|\tilde{\Gamma}_{\mathcal{A}_{22}^{m,k}}^{k+1}\|_{op} \leq \|\Gamma_{\mathcal{A}_{22}^{m,k}}^{k+1}\|_{op} + \frac{\alpha_0(1+c_-(s)c_+(s))+1}{n} \|\mathbf{X}_{\cdot\mathcal{A}^k}^\top \text{Err}^m\|_{op}, \quad (6)$$

where  $\alpha_0 > 0$  is some constant.

3. Recursively apply the above iterative inequality to derive

$$D(\mathcal{A}^{k+1}) \leq \gamma^{k+1} \|\mathbf{C}^{m*}\|_{op} + \alpha h(s).$$

4. Combining the results in Step 1 and Step 3, we have the results in Lemma 1.

Next, we prove the results regarding of  $h(s)$ . By definition, for any  $\mathcal{A}$  satisfying  $|\mathcal{A}| \leq s$ ,

$$\begin{aligned} \frac{1}{n} \|\mathbf{X}_{\cdot\mathcal{A}}^\top \text{Err}^m\|_{op} &= \frac{1}{n} \|\mathbf{X}_{\cdot\mathcal{A}}^\top \left( 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*}) \right)\|_{op} \\ &= \frac{1}{n} \|\mathbf{X}_{\cdot\mathcal{A}}^\top \left( 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*}) + 2(\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*}) \right)\|_{op} \\ &\leq \frac{1}{n} \|2\mathbf{X}_{\cdot\mathcal{A}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})\|_{op} + \frac{1}{n} \|2\mathbf{X}_{\cdot\mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}_{\cdot\mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_{op}. \end{aligned} \quad (7)$$

where the last inequality holds by applying the triangle inequality. From Lemma 5, the first term of the right hand side of (7) is bounded above as

$$\max_{\mathcal{A}: |\mathcal{A}| \leq s} \frac{1}{n} \|(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{C}^*})^\top \mathbf{X}_{\cdot\mathcal{A}}\|_{op} \leq \sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}}.$$

with probability at least  $1 - \delta$ .

As for the second term of (7), using the fact that  $\|A\|_{op} \leq \|A\|_F$  and the triangle inequality gives

$$\begin{aligned} &\frac{1}{n} \|2\mathbf{X}_{\cdot\mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}_{\cdot\mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_{op} \\ &\leq \frac{1}{n} \|2\mathbf{X}_{\cdot\mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}_{\cdot\mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_F \\ &\leq \frac{1}{n} \|2\mathbf{X}_{\cdot\mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}})\|_F + \frac{1}{n} \|\mathbf{X}_{\cdot\mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_F \end{aligned} \quad (8)$$

To bound the term  $\|2\mathbf{X}_{\cdot\mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}})\|_F/n$ , we bound the absolute value of each element in terms of  $\mathbf{X}\mathbf{C}$ . In particular, for any  $\mathbf{C}$  and other elements fixed, the  $(i, k)$ -th element of  $\tilde{\mathbf{P}}^{\mathbf{C}}$  is actually a function of  $\mathbf{c}_k^\top \mathbf{x}_i$ . Let  $f(\cdot)$  denote this function, i.e.,  $f(x) = \frac{\exp(x)}{\alpha + \exp(x)}$ , where  $\alpha$  is a constant. With simple calculation, we know that  $f(\cdot)$  is a contraction function with the Lipschitz constant less than 1. Hence,

$$\begin{aligned} &\frac{1}{n} \|2\mathbf{X}_{\cdot\mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}})\|_F + \frac{1}{n} \|\mathbf{X}_{\cdot\mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_F \\ &\leq \frac{2}{n} \|\mathbf{X}_{\cdot\mathcal{A}}^\top \mathbf{X}(\mathbf{C}^* - \hat{\mathbf{C}}^{m-1})\|_F + \frac{1}{n} \|\mathbf{X}_{\cdot\mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_F. \end{aligned} \quad (9)$$

The two terms of the right hand size of (9) are similar. Thus, we only present the proof for the upper bound of the first term, and that of the second term is omitted. Denote the row support set of  $\mathbf{C}^*$  and  $\hat{\mathbf{C}}^{m-1}$  as  $\mathcal{A}^*$  and  $\hat{\mathcal{A}}^{m-1}$ , respectively. Applying the triangle inequality again gives

$$\begin{aligned}
 & \frac{2}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top \mathbf{X}(\mathbf{C}^* - \hat{\mathbf{C}}^{m-1})\|_F \\
 &= \frac{2}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top \mathbf{X}_{\cdot \mathcal{A}^* \cup \hat{\mathcal{A}}^{m-1}}(\mathbf{C}_{\mathcal{A}^* \cup \hat{\mathcal{A}}^{m-1}}^* - \hat{\mathbf{C}}_{\mathcal{A}^* \cup \hat{\mathcal{A}}^{m-1}}^{m-1})\|_F \\
 &= \frac{2}{n} \left\| \mathbf{X}_{\cdot \mathcal{A}}^\top \left( \mathbf{X}_{\cdot \mathcal{A}^*}(\mathbf{C}_{\mathcal{A}^*}^* - \hat{\mathbf{C}}_{\mathcal{A}^*}^{m-1}) + \mathbf{X}_{\cdot \hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}(\mathbf{C}_{\hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}^* - \hat{\mathbf{C}}_{\hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}^{m-1}) \right) \right\|_F \\
 &\leq \frac{2}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top \mathbf{X}_{\cdot \mathcal{A}^*}(\mathbf{C}_{\mathcal{A}^*}^* - \hat{\mathbf{C}}_{\mathcal{A}^*}^{m-1})\|_F + \frac{2}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top \mathbf{X}_{\cdot \hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}(\mathbf{C}_{\hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}^* - \hat{\mathbf{C}}_{\hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}^{m-1})\|_F \\
 &\leq \frac{2}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top\|_F \|\mathbf{X}_{\cdot \mathcal{A}^*}\|_F \|\mathbf{C}_{\mathcal{A}^*}^* - \hat{\mathbf{C}}_{\mathcal{A}^*}^{m-1}\|_F + \frac{2}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top\|_F \|\mathbf{X}_{\cdot \hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}\|_F \|\mathbf{C}_{\hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}^* - \hat{\mathbf{C}}_{\hat{\mathcal{A}}^{m-1} \setminus \mathcal{A}^*}^{m-1}\|_F \\
 &\leq 4c_+(s) \|\mathbf{C}^* - \hat{\mathbf{C}}^{m-1}\|_F,
 \end{aligned}$$

where the last equality holds by using Condition 3.1. Similar procedure can be forwarded to show that

$$\frac{1}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_F \leq c_+(s) \|\mathbf{C}^{m*} - \hat{\mathbf{C}}^{m-1}\|_F.$$

Substituting the above inequalities into (8) and (9), we can obtain

$$\begin{aligned}
 & \frac{1}{n} \|2\mathbf{X}_{\cdot \mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}_{\cdot \mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_{op} \\
 &\leq 4c_+(s) \|\mathbf{C}^* - \hat{\mathbf{C}}^{m-1}\|_F + c_+(s) \|\mathbf{C}^{m*} - \hat{\mathbf{C}}^{m-1}\|_F \\
 &\leq 5c_+(s) \|\mathbf{C}^* - \hat{\mathbf{C}}\|_F + 5c_+(s) \|\hat{\mathbf{C}}^{m-1} - \hat{\mathbf{C}}\|_F + c_+(s) \|\mathbf{C}^{m*} - \mathbf{C}^*\|_F.
 \end{aligned}$$

Then we can bound above the right hand size by applying Lemmas 2-3,

$$\begin{aligned}
 & \frac{1}{n} \|2\mathbf{X}_{\cdot \mathcal{A}}^\top (\tilde{\mathbf{P}}^{\mathbf{C}^*} - \tilde{\mathbf{P}}^{\hat{\mathbf{C}}^{m-1}}) + \mathbf{X}_{\cdot \mathcal{A}}^\top \mathbf{X}(\hat{\mathbf{C}}^{m-1} - \mathbf{C}^{m*})\|_{op} \\
 &\leq 5 \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} (\|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F + (1 - 2\zeta(1 - \zeta)) \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F) + 5c_+(s) \|\mathbf{C}^* - \hat{\mathbf{C}}\|_F
 \end{aligned}$$

Therefore, with probability at least  $1 - \delta$ , we have

$$\begin{aligned}
 & \max_{\mathcal{A}: |\mathcal{A}| \leq s} \frac{1}{n} \|\mathbf{X}_{\cdot \mathcal{A}}^\top \text{Err}^m\|_{op} \leq 5 \frac{c_+(s)^2}{c_-(s)} (1 - 2\zeta(1 - \zeta))^{m-1} (\|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F + (1 - 2\zeta(1 - \zeta)) \|\hat{\mathbf{C}}^0 - \mathbf{C}^*\|_F) \\
 & + 5c_+(s) \|\mathbf{C}^* - \hat{\mathbf{C}}\|_F + 2\sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}},
 \end{aligned}$$

which finishes the proof.

## B.2. Proof of Lemma 2

The first part of the proof proceeds along similar lines of the proof of Theorem 1 of Balakrishnan et al. (2017). Denote  $\mathbf{G} = \mathbf{X}\mathbf{C}$  and  $\mathbf{G}' = \mathbf{X}\mathbf{C}'$ . Then, we can rewrite the surrogate function  $S(\mathbf{C}; \mathcal{D}|\mathbf{C}')$  with respect to  $\mathbf{G}$  as  $S_X(\mathbf{G}|\mathbf{G}')$ . For convenience, we introduce the mapping  $M$  given by

$$M(\mathbf{G}) := \arg \min_{\mathbf{G}'} S_X(\mathbf{G}'|\mathbf{G}). \tag{10}$$

Let  $\hat{\mathbf{G}} = \mathbf{X}\hat{\mathbf{C}}$ , then we know  $\hat{\mathbf{G}} = M(\hat{\mathbf{G}})$ , by the definition of  $\hat{\mathbf{C}}$  and equation (10). In the first part, we show that the mapping  $M$ , in  $B_F(r_0, \mathbf{C}^*)$  satisfies

$$\|M(\mathbf{G}) - \hat{\mathbf{G}}\|_F \leq (1 - 2\zeta(1 - \zeta)) \|\mathbf{G} - \hat{\mathbf{G}}\|_F \quad \text{for all } \mathbf{G}, \tag{11}$$

where  $0 < \zeta < 1$ ,  $\zeta = \min_{\mathbf{C} \in B_F(r_0, \mathbf{C}^*), \mathbf{C} \in \Omega_r, i, j} \tilde{\mathbf{P}}_{i,j}^{\mathbf{C}}$ . By the fact that  $1/4 \geq c_{\mathbf{C}^*} > 0$  (Condition 3.3), we know that  $1/4 \geq \zeta(1 - \zeta) > 0$ .

Denote the gradient of  $S_X(\mathbf{G}|\mathbf{G}')$  with respect to  $\mathbf{G}$  as  $\nabla S_X(\mathbf{G}|\mathbf{G}')$ . Then by the definition of surrogate function, we have

$$\nabla S_X(\mathbf{G}|\mathbf{G}') = \frac{1}{2n} \left( \mathbf{G} - \mathbf{G}' - 2(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{\mathbf{G}'}) \right),$$

where  $\tilde{\mathbf{P}}^{\mathbf{G}'}$  represents the term  $\tilde{\mathbf{P}}^{\mathbf{C}'}$  being rewritten as a function of  $\mathbf{G}' = \mathbf{X}\mathbf{C}'$ . Making use the fact that

$$\nabla S_X(\hat{\mathbf{G}}|\hat{\mathbf{G}}) - \nabla S_X(M(\mathbf{G})|\hat{\mathbf{G}}) = \frac{1}{2n} \left( \hat{\mathbf{G}} - M(\mathbf{G}) \right),$$

we can show that the following equation hold for all  $\mathbf{G}$ ,

$$\langle \nabla S_X(\hat{\mathbf{G}}|\hat{\mathbf{G}}) - \nabla S_X(M(\mathbf{G})|\hat{\mathbf{G}}), \hat{\mathbf{G}} - M(\mathbf{G}) \rangle_F = \beta \|\hat{\mathbf{G}} - M(\mathbf{G})\|_F^2 \quad (12)$$

with  $\beta = \frac{1}{2n}$ .

Next, we show that the following inequality hold for all  $\mathbf{G}$ ,

$$\|\nabla S_X(M(\mathbf{G})|\hat{\mathbf{G}}) - \nabla S_X(M(\mathbf{G})|\mathbf{G})\|_F \leq \tau \|\mathbf{G} - \hat{\mathbf{G}}\|_F, \quad (13)$$

where  $\tau = \frac{1}{2n}(1 - 2\zeta(1 - \zeta))$  and  $0 < \zeta < 1$  is a constant related to  $\hat{\mathbf{G}}$ . Noting that

$$\nabla S_X(M(\mathbf{G})|\hat{\mathbf{G}}) - \nabla S_X(M(\mathbf{G})|\mathbf{G}) = \frac{1}{2n} \left( \mathbf{G} - \hat{\mathbf{G}} - 2(\tilde{\mathbf{P}}^{\mathbf{G}} - \tilde{\mathbf{P}}^{\hat{\mathbf{G}}}) \right) \triangleq f(\hat{\mathbf{G}}).$$

Taking gradient of  $f(\cdot)$  with respect of  $\hat{\mathbf{G}}$  leads to

$$\nabla f(\hat{\mathbf{G}}) = \frac{1}{2n} \left( 1 - 2\tilde{\mathbf{P}}^{\hat{\mathbf{G}}}(1 - \tilde{\mathbf{P}}^{\hat{\mathbf{G}}}) \right).$$

Let  $P_{\min}^{\hat{\mathbf{G}}}$  denote the minimal element in the matrix  $\tilde{\mathbf{P}}^{\hat{\mathbf{G}}}$  and assume that  $P_{\min}^{\mathbf{G}} > 0$ . By the convexity of function  $f$ , we can bound the Frobenius norm of  $f(\hat{\mathbf{G}})$  as

$$\|\nabla S_X(M(\mathbf{G})|\hat{\mathbf{G}}) - \nabla S_X(M(\mathbf{G})|\mathbf{G})\|_F \leq \|\nabla f(\hat{\mathbf{G}})(\mathbf{G} - \hat{\mathbf{G}})\|_F \leq \tau \|\mathbf{G} - \hat{\mathbf{G}}\|_F$$

with  $\tau = \frac{1}{2n} \left( 1 - 2P_{\min}^{\hat{\mathbf{G}}}(1 - P_{\min}^{\hat{\mathbf{G}}}) \right)$  in  $B_F(r_0, \mathbf{C}^*)$ .

Finally, the Cauchy-Schwartz inequality implies that the left had side of (12) is upper bounded by

$$\langle \nabla S_X(\hat{\mathbf{G}}|\hat{\mathbf{G}}) - \nabla S_X(M(\mathbf{G})|\hat{\mathbf{G}}), \hat{\mathbf{G}} - M(\mathbf{G}) \rangle_F \leq \|\nabla S_X(M(\mathbf{G})|\hat{\mathbf{G}}) - \nabla S_X(M(\mathbf{G})|\mathbf{G})\|_F \|\mathbf{G} - \hat{\mathbf{G}}\|_F.$$

Substituting the above inequality into (12) and (13) yields

$$\beta \|\hat{\mathbf{C}} - M(\mathbf{C})\|_F^2 \leq \tau \|\hat{\mathbf{C}} - M(\mathbf{C})\|_F \|\mathbf{C} - \hat{\mathbf{C}}\|_F.$$

Canceling the term  $\beta \|\hat{\mathbf{C}} - M(\mathbf{C})\|_F$  on both sides and  $\zeta \leq \frac{\tau}{\beta}, \forall \mathbf{C} \in B_F(r_0, \mathbf{C}^*), \mathbf{C} \in \omega_r$ , we can obtain the inequality (11).

Let  $\hat{\mathbf{G}}^m = \mathbf{X}\hat{\mathbf{C}}^m$ , then we have  $\hat{\mathbf{G}}^m = M(\hat{\mathbf{G}}^{m-1})$ . By iteratively applying (11), we have for all  $m = 1, 2, \dots$ ,

$$\begin{aligned} \|\hat{\mathbf{G}}^m - \hat{\mathbf{G}}\|_F &\leq (1 - 2\zeta(1 - \zeta)) \|\hat{\mathbf{G}}^{m-1} - \hat{\mathbf{G}}\|_F \\ &\leq \dots \\ &\leq (1 - 2\zeta(1 - \zeta))^m \|\hat{\mathbf{G}}^0 - \hat{\mathbf{G}}\|_F. \end{aligned} \quad (14)$$

Applying Condition 3.1, we can obtain the following inequalities

$$\begin{aligned} \|\hat{\mathbf{G}}^{(m)} - \hat{\mathbf{G}}\|_F &= \|\mathbf{X}\hat{\mathbf{C}}^{(m)} - \mathbf{X}\hat{\mathbf{C}}\|_F \geq \sqrt{n}c_-(s) \|\hat{\mathbf{C}}^{(m)} - \hat{\mathbf{C}}\|_F, \\ \|\hat{\mathbf{G}}^{(0)} - \hat{\mathbf{G}}\|_F &= \|\mathbf{X}\mathbf{C}^0 - \mathbf{X}\hat{\mathbf{C}}\|_F \leq \sqrt{n}c_+(s) \|\hat{\mathbf{C}}^{(0)} - \hat{\mathbf{C}}\|_F. \end{aligned}$$

Substituting these two inequalities into (14), we have

$$\|\hat{\mathbf{C}}^{(m)} - \hat{\mathbf{C}}\|_F \leq (1 - 2\zeta(1 - \zeta))^m \frac{c_+(s)}{c_-(s)} \|\hat{\mathbf{C}}^0 - \hat{\mathbf{C}}\|_F.$$

### B.3. Proof of Lemma 3

Rewrite the noise-free surrogate function  $\tilde{S}(\mathbf{C}; \mathcal{D}|\mathbf{C}')$  with respect to  $G$  as  $\tilde{S}_X(\mathbf{G}|\mathbf{G}')$ . For convenience, we introduce the mapping  $\tilde{M}$  given by

$$\tilde{M}(G) := \arg \min_{\mathbf{G}'} \tilde{S}_X(\mathbf{G}'|\mathbf{G}). \quad (15)$$

The remaining proof follows similar lines of Lemma 2, and the details are omitted here.

### B.4. Proof of Lemma 4

Recall that  $\mathbf{C}^{m,k}$  is the output coefficient matrix at the  $k$ -th iteration of the inner loop and the  $m$ -th iteration of the outer loop in **Algorithm 2**. Let  $\mathcal{A}^{m,k}$  be the corresponding row support set, and  $\mathcal{I}^{m,k} = (\mathcal{A}^{m,k})^c$ . Similarly, denote  $\mathbf{B}^{m,k}$ ,  $\Gamma^{m,k}$ , and  $\mathbf{V}^{m,k}$  as the intermediate matrices at the  $k$ -th iteration of **Algorithm 1** within the  $m$ -th iteration of **Algorithm 2**. From step 5 in **Algorithm 1**, we have

$$\Gamma_{\mathcal{I}^{m,k}}^{m,k} = \mathbf{X}_{\mathcal{I}^{m,k}}^\top (\mathbf{M}\mathbf{V}^{m,k} - \mathbf{X}\mathbf{B}^{m,k})/n.$$

From the definition of  $\tilde{\Gamma}_{A_{22}^k}^{k+1}$  and  $\Gamma_{A_{22}^k}^{k+1}$ , we know

$$\tilde{\Gamma}_{A_{22}^k}^{k+1} = \Gamma_{A_{22}^k}^{k+1} (\mathbf{V}^{k+1})^\top + \mathbf{X}_{A^k}^\top \mathbf{M}^{(m)} (\mathbf{I} - \mathbf{P}_\mathbf{V})/n,$$

after some simplifications. Here we denote  $\mathbf{P}_\mathbf{V} = \mathbf{V}^{k+1}(\mathbf{V}^{k+1})^\top$  for convenience. Using Lemma B.5 in [Wen et al. \(2022\)](#), we can get

$$\|\tilde{\Gamma}_{A_{22}^k}^{k+1}\|_{op} \leq \|\Gamma_{A_{22}^k}^{k+1} (\mathbf{V}^{k+1})^\top\|_{op} + n^{-1} \|\mathbf{X}_{A_{22}^k}^\top \mathbf{M}^{(m)} (\mathbf{I} - \mathbf{P}_\mathbf{V})\|_{op}. \quad (16)$$

Assume that  $\text{rank}(\mathbf{X}_{A^k}^\top \mathbf{X} \mathbf{C}^{(m)*}) = r^*$  without loss of the generality. Moreover, we have  $\text{rank}(\mathbf{X}^\top \mathbf{X} \mathbf{C}^{(m)*}) = r^*$ . Thus there exists a linear transformer  $\mathbb{T} \in \mathbb{R}^{|A_{22}^k| \times |A^k|}$  satisfying  $\mathbf{X}_{A_{22}^k}^\top \mathbf{X} \mathbf{C}^{(m)*} = \mathbb{T}(\mathbf{X}_{A^k}^\top \mathbf{X} \mathbf{C}^{(m)*})$ . Because  $\mathbf{X}_{A_{22}^k}^\top \mathbf{X} \mathbf{C}^{(m)*}$  is a submatrix of the original matrix. It ensures that

$$\begin{aligned} n^{-1} \|\mathbf{X}_{A_{22}^k}^\top \mathbf{X} \mathbf{C}^{(m)*} (\mathbf{I} - \mathbf{P}_\mathbf{V})\|_{op} &\leq \frac{\|\mathbb{T}\|_{op}}{n} \|\mathbf{X}_{A^k}^\top \mathbf{X} \mathbf{C}^{(m)*} (\mathbf{I} - \mathbf{P}_\mathbf{V})\|_{op} \\ &\leq \frac{\alpha_0}{n} \|\mathbf{X}_{A^k}^\top \mathbf{X} \mathbf{C}^{(m)*} (\mathbf{I} - \mathbf{P}_\mathbf{V})\|_{op}. \end{aligned}$$

And we also have

$$\begin{aligned} \mathbf{X}_{A^k}^\top \mathbf{X} \mathbf{C}^* (\mathbf{I} - \mathbf{P}_\mathbf{V}) &= \mathbf{X}_{A^k}^\top \mathbf{M}^m (\mathbf{I} - \mathbf{P}_\mathbf{V}) - \mathbf{X}_{A^k}^\top \text{Err}^m (\mathbf{I} - \mathbf{P}_\mathbf{V}) \\ &= (\mathbf{X}_{A^k}^\top \mathbf{X}_{A^k})^{1/2} \sum_{l=1}^{q-1} \sigma_l^{k+1} \mathbf{u}_l^{k+1} (\mathbf{v}_l^{k+1})^\top (\mathbf{I} - \mathbf{P}_\mathbf{V}) - \mathbf{X}_{A^k}^\top \text{Err}^m (\mathbf{I} - \mathbf{P}_\mathbf{V}). \end{aligned} \quad (17)$$

The decomposition above comes from

$$(\mathbf{X}_{A^k}^\top \mathbf{X}_{A^k})^{-1/2} \mathbf{X}_{A^k}^\top \mathbf{M}^m = \sum_{l=1}^{q-1} \sigma_l^{k+1} \mathbf{u}_l^{k+1} (\mathbf{v}_l^{k+1})^\top,$$

where  $\sigma_l^{k+1}$ ,  $\mathbf{u}_l^{k+1}$  are the  $l$ -th largest singular value and corresponding left singular vector, respectively. Due to the orthogonality of  $\mathbf{v}_l^{k+1}$  and  $\mathbf{P}_\mathbf{V}$ , we have

$$\sum_{l=1}^{q-1} \sigma_l^{k+1} \mathbf{u}_l^{k+1} (\mathbf{v}_l^{k+1})^\top (\mathbf{I} - \mathbf{P}_\mathbf{V}) = \sum_{l=r+1}^{q-1} \sigma_l^{k+1} \mathbf{u}_l^{k+1} (\mathbf{v}_l^{k+1})^\top (\mathbf{I} - \mathbf{P}_\mathbf{V})$$

Meanwhile, by the definition of  $\sigma_{r+1}^{k+1}$ , we know that it is the  $(r+1)$ -th largest singular value of

$$(\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{-1/2} \mathbf{X}_{.A^k}^\top \mathbf{M}^m = (\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{-1/2} \mathbf{X}_{.A^k}^\top \mathbf{X} \mathbf{C}^{m*} + (\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{-1/2} \mathbf{X}_{.A^k}^\top \mathbf{E} r r^m.$$

Therefore, by the inequality (17), we have

$$\begin{aligned} \|\mathbf{X}_{.A^k}^\top \mathbf{X} \mathbf{C}^{m*} (\mathbf{I} - \mathbf{P}_V)\|_F &\leq \|(\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{1/2} \sum_{l=r+1}^{q-1} \sigma_l^{k+1} \mathbf{u}_l^{k+1} (\mathbf{v}_l^{k+1})^\top (\mathbf{I} - \mathbf{P}_V)\|_F + \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^m\|_F \\ &\leq \sigma_{r+1}^{k+1} \|(\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{1/2}\|_F + \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^m\|_F \\ &\leq \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^{(m)}\|_{op} \|(\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{-1/2}\|_{op} \|(\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{1/2}\|_{op} + \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^m\|_F. \end{aligned}$$

Combining (16) and applying the triangle inequality again, we can obtain

$$\begin{aligned} \|\tilde{\Gamma}_{A_{22}^k}^{k+1}\|_{op} &\leq \|\Gamma_{A_{22}^k}^{k+1}\|_{op} + n^{-1} \|\mathbf{X}_{.A_{22}^k}^\top \mathbf{X} \mathbf{C}^{(m)*} (\mathbf{I} - \mathbf{P}_V)\|_{op} + n^{-1} \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^m\|_{op} \\ &\leq \|\Gamma_{A_{22}^k}^{k+1}\|_{op} + \frac{\alpha_0}{n} \|\mathbf{X}_{.A^k}^\top \mathbf{X} \mathbf{C}^{(m)*} (\mathbf{I} - \mathbf{P}_V)\|_{op} + n^{-1} \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^m\|_{op} \\ &\leq \|\Gamma_{A_{22}^k}^{k+1}\|_{op} + \frac{\alpha_0}{n} \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^{(m)}\|_{op} \|(\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{-1/2}\|_{op} \|(\mathbf{X}_{.A^k}^\top \mathbf{X}_{.A^k})^{1/2}\|_{op} \\ &\quad + \frac{\alpha_0}{n} \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^{(m)}\|_{op} + n^{-1} \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^m\|_{op} \\ &\leq \|\Gamma_{A_{22}^k}^{k+1}\|_{op} + \frac{\alpha_0(1 + c_-(s)c_+(s)) + 1}{n} \|\mathbf{X}_{.A^k}^\top \mathbf{E} r r^m\|_{op}, \end{aligned}$$

where  $\alpha_0$  is the operator norm of  $\mathbb{T}$  and the last inequality comes from condition 3.1.

### B.5. Lemma 5 and its proof

LEMMA 5. Suppose each row of  $\mathbf{Y}$  comes from the multinomial logistic model with a deterministic  $n \times p$  design matrix  $\mathbf{X}$  and an  $p \times q$  coefficient matrix  $\mathbf{C}$ . Denote  $P^2 = \max_j \text{diag}(\frac{1}{n} \mathbf{X}^\top \mathbf{X}) \|\tilde{\mathbf{P}}_j^{C*}\|_F$  and the probability matrix as  $\mathbf{P}^{C*}$ . Let  $\tilde{\mathbf{Y}}$  and  $\tilde{\mathbf{P}}^{C*}$  denote the matrix by eliminating the first column of  $\mathbf{Y}$  and  $\mathbf{P}^{C*}$ . Then for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\max_{\mathcal{A}: |\mathcal{A}| \leq s} \frac{1}{n} \|(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C*})^\top \mathbf{X}_{.A}\|_{op} \leq \sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}}.$$

**Proof** For any  $\zeta > 0$ , by definition,

$$\begin{aligned} Pr \left( \max_{\mathcal{A}: |\mathcal{A}| \leq s} \frac{1}{n} \|(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C*})^\top \mathbf{X}_{.A}\|_{op} \geq \zeta \right) &= Pr \left( \max_{\mathcal{A}: |\mathcal{A}| \leq s} \frac{1}{n^2} \|(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C*})^\top \mathbf{X}_{.A}\|_{op}^2 \geq \zeta^2 \right) \\ &\leq Pr \left( \sup_{\|\mathbf{u}\|_2=1, \|\mathbf{u}\|_0 \leq s} n^{-2} \|(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C*})^\top \mathbf{X} \mathbf{u}\|_2^2 \geq \zeta^2 \right) \\ &\leq Pr \left( \sum_{i=1}^q \sup_{\|\mathbf{u}\|_2=1, \|\mathbf{u}\|_0 \leq s} n^{-2} ((\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C*})^\top \mathbf{X} \mathbf{u})_i^2 \geq \zeta^2 \right), \end{aligned}$$

where  $(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C^*})_{\cdot i}$  denotes the  $i$ -th column of matrix  $\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C^*}$ . Applying the union bound yields

$$\begin{aligned}
& Pr \left( \sum_{i=1}^q \sup_{\|\mathbf{u}\|_2=1, \|\mathbf{u}\|_0 \leq s} \frac{1}{n^2} ((\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C^*})_{\cdot i}^\top \mathbf{X} \mathbf{u})^2 \geq \zeta^2 \right) \\
& \leq \sum_{i=1}^q Pr \left( \sup_{\|\mathbf{u}\|_2=1, \|\mathbf{u}\|_0 \leq s} \frac{1}{n^2} ((\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C^*})_{\cdot i}^\top \mathbf{X} \mathbf{u})^2 \geq \frac{\zeta^2}{q} \right) \\
& = \sum_{i=1}^q Pr \left( \sup_{\|\mathbf{u}\|_2=1, \|\mathbf{u}\|_0 \leq s} \frac{1}{n} |(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C^*})_{\cdot i}^\top \mathbf{X} \mathbf{u}| \geq \frac{\zeta}{\sqrt{q}} \right) \\
& \leq \sum_{i=1}^q Pr \left( \max_{1 \leq j \leq p} \frac{1}{n} |(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C^*})_{\cdot i}^\top \mathbf{X}_{\cdot j}| \geq \frac{\zeta}{\sqrt{qs}} \right) \\
& \leq \sum_{i=1}^q \sum_{j=1}^p Pr \left( \frac{1}{\sqrt{n}} |(\tilde{\mathbf{Y}} - \tilde{\mathbf{P}}^{C^*})_{\cdot i}^\top \mathbf{X}_{\cdot j}| \geq \zeta \sqrt{\frac{n}{qs}} \right) \\
& \leq 2pq \exp \left( -\frac{n\zeta^2}{3qsP^2} \right),
\end{aligned}$$

where the Bernoulli Chernoff bound (Vershynin 2018) is applied in the last inequality. By taking  $\delta = 2pq \exp(-\frac{n\zeta^2}{2qsP^2})$ , we obtain  $\zeta = \sqrt{\frac{3qsP^2}{n} \log \frac{2pq}{\delta}}$ , which concludes the proof.

### Appendix C: HAM10000 data analysis

The HAM10000 dataset consists of 10,015 dermatoscopic images from different populations, acquired and stored by different modalities. It includes a representative collection of 7 important diagnostic categories in the realm of pigmented lesions: Actinic keratoses and intraepithelial carcinoma/Bowen’s disease (akiec), basal cell carcinoma (bcc), benign keratosis-like lesions (solar lentigines/seborrheic keratoses and lichen-planus like keratoses, bkl), dermatofibroma (df), melanoma (mel), melanocytic nevi (nv) and vascular lesions (angiomas, angiokeratomas, pyogenic granulomas and hemorrhage, vasc). Figure 1 shows five samples of the original images.

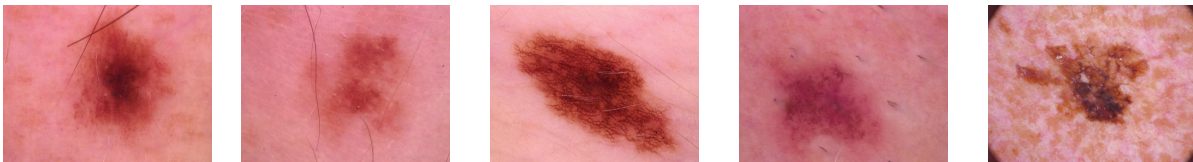


Figure 1 5 examples of HAM10000 dataset

Table 1 reports the results for each method. Notice that the minus log-likelihood value Pred is calculated on the whole data set since we do not have an independent test data set here. We can see that our method yields a sparse reduced rank model with the lowest sparsity. In contrast, the mLasso method identifies much more predictors and a higher Pred value. NNET has very poor performance for this dataset due to its high sparsity, spending most of the time wandering randomly. Due to the sensitivity and instability to data, RRVGLM does’t work for this specific data.

To compare the prediction accuracy and stability for these methods, we randomly split the whole data set into a training set with size 900 and a testing set with the rest samples. The optimal pair of rank and

**Table 1** The results on the whole data set.

Method	Pred	$\hat{r}$	$ \hat{\mathcal{A}} $
SRRMLR	16018.99	3	31
mLasso	18481.05	-	182
NNET	18611.97	-	192
RRVGLM	-	-	-

‘-’, these results can not be obtained.

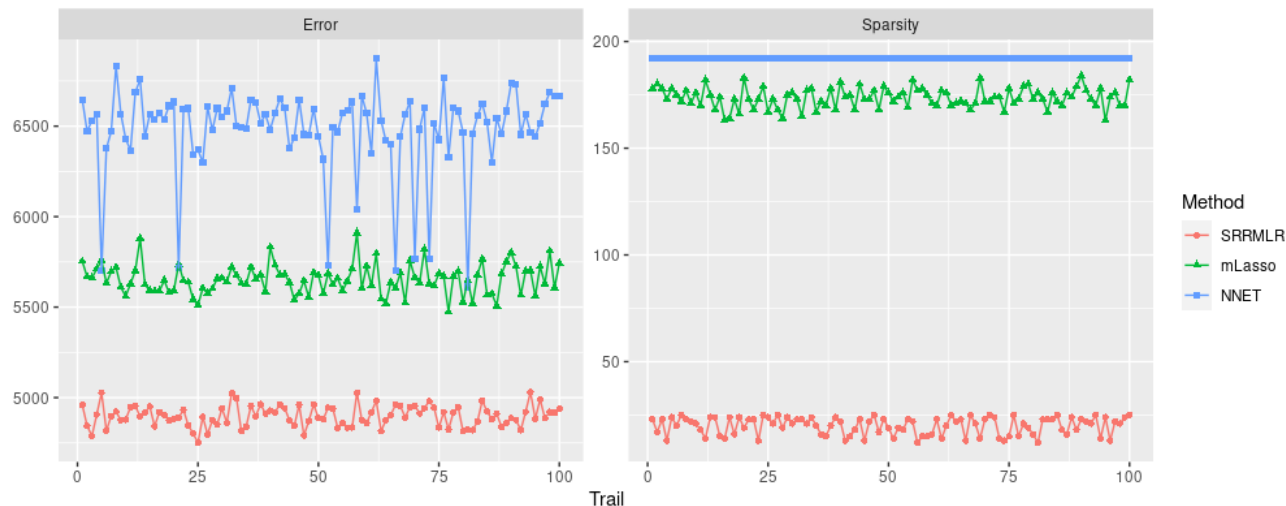
**Table 2** The average results on the test set of HAM10000 with the standard deviation in parentheses.

Method	Pred	$\hat{r}$	$ \hat{\mathcal{A}} $
SRRMLR	4919.94 (57.89)	3.84 (1.04)	29.09 (5.69)
mLasso	5646.84 (87.18)	-	173.59 (5.14)
NNET	6492.99 (232.37)	-	192 (0)
RRVGLM	-	-	-

‘-’, these results can not be obtained.

sparsity is determined by the strategy discussed above, but the 5-fold cross-validation procedure is carried out using the training set. According to the upper bounds developed in the simulation setup, we set the maximum range  $(r_{\max}, s_{\max})$  of these two tuning parameters to 6 and 38, respectively. In addition, we use the testing set to calibrate the predictive performance of each estimator  $\hat{C}$  by Pred in (11). The random splitting process is implemented 100 times to yields the average performance of the above methods.

Summary results are presented in Table 2. Figure 2 displays line charts of the 100 replications for the measurements of Pred and  $|\hat{\mathcal{A}}|$ . Overall, the results are accordance with those in Table 1. In particular, the proposed SRRMLR method yields significant better performance compared to other methods in all of the measurements. The mLasso achieves competitive prediction accuracy, but it produces a much larger model size with too many selected predictors. The RRVGLM cannot even be applied to this data set. The NNET leads to models with the largest model size and the highest value in Pred. Furthermore, the models from NNET are unstable in terms of predictive accuracy, which can be also seen from Figure 2. We conclude that the SRRMLR is feasible and effective when handling data in practice with the consideration of both prediction accuracy and model interpretability.



**Figure 2** The trails of prediction error (left) and row sparsity (right) among 100 replications.

## References

- Balakrishnan S, Wainwright MJ, Yu B (2017) Statistical guarantees for the EM algorithm: From population to sample-based analysis. *Annals of Statistics* 45(1):77 – 120.
- Simon N, Friedman J, Hastie T (2013) A blockwise descent algorithm for group-penalized multiresponse and multinomial regression. *arXiv preprint arXiv:1311.6529* .
- Vershynin R (2018) *High-dimensional probability: An introduction with applications in data science*, volume 47 (Cambridge university press).
- Wen C, Dong R, Wang X, Li W, Zhang H (2022) Simultaneous best subset selection and dimension reduction via primal-dual iterations. *arXiv:2211.15889* .