

## Multivariate Varying-coefficient Models via Tensor Decomposition

Fengyu Zhang, Ya Zhou, Kejun He, and Raymond K. W. Wong

*Renmin University of China, Chinese Academy of Medical Sciences and Peking Union Medical College,*

*Renmin University of China, and Texas A&M University*

### Supplementary Material

## S.1 Computation Details

### S.1.1 Details on Updating $\mathbf{S}$

In this subsection, we provide more details on solving (4.1) of the main paper through the MM algorithm. To construct the majorized function for  $P_{\mathbf{A}^{(t)}}(\mathbf{S})$ , we first let  $\mathbf{S}_k$  and  $\mathbf{S}_{(1),k}$  be the  $k$ -th iterations of  $\mathbf{S}$  and  $\mathbf{S}_{(1)}$  in the MM algorithm respectively. Using the concave property of the square root function, we have

$$\|\mathbf{x}\|_2 = \sqrt{\|\mathbf{x}\|_2^2} \leq \sqrt{\|\mathbf{x}_0\|_2^2} + \frac{\|\mathbf{x}\|_2^2 - \|\mathbf{x}_0\|_2^2}{2\sqrt{\|\mathbf{x}_0\|_2^2}} := g(\mathbf{x}|\mathbf{x}_0).$$

It can be verified that  $g(\mathbf{x}|\mathbf{x}_0)$  is a majorization function for  $\|\mathbf{x}\|_2$ . Since  $P_{\mathbf{A}^{(t)}}(\mathbf{S})$  is a summation over  $\ell_2$  norms, a majorized function for  $P_{\mathbf{A}^{(t)}}(\mathbf{S})$  at  $\mathbf{S}_k$  is

$$\Phi_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k) = \lambda \sum_{j=1}^p \left\{ \sqrt{\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2^2} + \frac{\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1)}\|_2^2 - \|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2^2}{2\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2} \right\}. \quad (\text{S.1})$$

Denoting  $\lambda_{j,k}^{(t)} = 1/\{2\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2\}$ ,  $j = 1, \dots, p$ , and  $\boldsymbol{\Sigma}_k^{(t)} = \text{diag}(\lambda_{1,k}^{(t)}, \dots, \lambda_{p,k}^{(t)})$ , the majorized function  $\Phi_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k)$  in (S.1) can be rewritten as

$$\begin{aligned} \Phi_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k) = P_{\mathbf{A}^{(t)}}(\mathbf{S}_k) + \lambda \left[ \text{vec}(\mathbf{S}_{(1)})^\top \{ \mathbf{I}_{R_2 R_3} \otimes (\mathbf{A}^{(t)})^\top \boldsymbol{\Sigma}_k^{(t)} \mathbf{A}^{(t)} \} \text{vec}(\mathbf{S}_{(1)}) \right. \\ \left. - \text{vec}(\mathbf{S}_{(1),k})^\top \{ \mathbf{I}_{R_2 R_3} \otimes (\mathbf{A}^{(t)})^\top \boldsymbol{\Sigma}_k^{(t)} \mathbf{A}^{(t)} \} \text{vec}(\mathbf{S}_{(1),k}) \right]. \end{aligned}$$

Note that  $\Phi_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k)$  is well-defined except for  $\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2 = 0$ . When  $\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2 = 0$ ,  $j = 1, 2, \dots, p$ , we make a further discussion in two cases: (i)  $\|\mathbf{a}_j^{(t)}\|_2 = 0$ ; (ii)  $\|\mathbf{a}_j^{(t)}\|_2 \neq 0$ , but  $\mathbf{a}_j^{(t)}$  belongs to the null space of  $\mathbf{S}_{(1),k}^\top$ .

For case (i), denote  $\tilde{\mathcal{J}} = \{j : \|\mathbf{a}_j^{(t)}\|_2 = 0, j = 1, \dots, p\}$ . If  $j \in \tilde{\mathcal{J}}$ , the  $j$ -th component in  $P_{\mathbf{A}^{(t)}}(\mathbf{S})$  equals zero, which allows us to set the  $j$ -th diagonal element in  $\tilde{\Sigma}_k^{(t)}$  to be zero. Define  $\tilde{\Sigma}_k^{(t)} = \text{diag}(\tilde{\lambda}_{1,k}^{(t)}, \dots, \tilde{\lambda}_{p,k}^{(t)})$  and  $\tilde{\lambda}_{j,k}^{(t)} = \lambda_{j,k}^{(t)} \cdot \mathbf{I}\{j \in \tilde{\mathcal{J}}\}$ ,  $j = 1, \dots, p$ , where  $\mathbf{I}\{\cdot\}$  is the indicator function. Denote the refined quadratic function

$$\tilde{\Phi}_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k) = P_{\mathbf{A}^{(t)}}(\mathbf{S}_k) + \text{vec}(\mathbf{\Delta}_{(1),k})^\top \{ \mathbf{I}_{R_2 R_3} \otimes (\mathbf{A}^{(t)})^\top \tilde{\Sigma}_k^{(t)} \mathbf{A}^{(t)} \} \text{vec}(\mathbf{\Delta}_{(1),k}).$$

Similar to  $\Phi_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k)$ ,  $\tilde{\Phi}_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k)$  is also a majorized function for  $P_{\mathbf{A}^{(t)}}(\mathbf{S})$ . Now we proceed to the minimizing step of the MM algorithm, and the corresponding minimization problem is

$$\min_{\mathbf{S}} \{ H_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{S}) + \tilde{\Phi}_{\mathbf{A}^{(t)}}(\mathbf{S}|\mathbf{S}_k) \}.$$

Taking gradient of the objective function and setting it to zero gives

$$\text{vec}(\mathbf{S}_{(1),k+1}) = \left[ \sum_{i=1}^n (\mathbf{U}_{\mathbf{S},i}^{(t)})^\top \mathbf{U}_{\mathbf{S},i}^{(t)} + \mathbf{I}_{R_2 R_3} \otimes \{ (\mathbf{A}^{(t)})^\top \tilde{\Sigma}_k^{(t)} \mathbf{A}^{(t)} \} \right]^+ \left( \sum_{i=1}^n (\mathbf{U}_{\mathbf{S},i}^{(t)})^\top \mathbf{y}_i \right), \quad (\text{S.2})$$

where  $\mathbf{U}_{\mathbf{S},i}^{(t)} = \mathbf{C}^{(t)} \otimes \{ \mathbf{b}^\top(t_i) \mathbf{B}^{(t)} \} \otimes (\mathbf{x}_i^\top \mathbf{A}^{(t)})$  and  $(\cdot)^+$  denotes the Moore-Penrose inverse.

For case (ii), an improved version of local quadratic approximation has been suggested in Hunter and Li (2005). To handle the ill-posedness of (S.1) when  $2\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2 = 0$ , Hunter and Li (2005) proposed to replace  $2\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2$  by  $2(\varepsilon + \|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2)$  for some  $\varepsilon > 0$ . After applying such replacement to all  $(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}$  in (S.1),  $j = 1, \dots, p$ , the adjusted version is no longer a majorizer of  $P_{\mathbf{A}^{(t)}}(\mathbf{S})$  as required by the MM algorithm. Nonetheless, Hunter and Li (2005) showed that it majorizes a perturbed version of  $P_{\mathbf{A}^{(t)}}(\mathbf{S})$ . This would lead to an objective function that is similar to the one in case (i). In addition, the minimizer of this function should be close to the minimizer of the original function as long as  $\varepsilon$  is small enough and the original objective function is not extremely flat in the neighborhood of the minimizer. We thus define  $\tilde{\Sigma}_{\varepsilon,k}^{(t)} = \text{diag}(\tilde{\lambda}_{1,k}^{(t)}(\varepsilon), \dots, \tilde{\lambda}_{p,k}^{(t)}(\varepsilon))$  with  $\tilde{\lambda}_{j,k}^{(t)}(\varepsilon) = 1/\{2\varepsilon + 2\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2\}$ . We get the perturbed version of the approximation function

$$\tilde{\Phi}_{\mathbf{A}^{(t)}}^\varepsilon(\mathbf{S}|\mathbf{S}_k) = P_{\mathbf{A}^{(t)}}(\mathbf{S}_k) + \text{vec}(\mathbf{\Delta}_{(1),k})^\top \{ \mathbf{I}_{R_2 R_3} \otimes (\mathbf{A}^{(t)})^\top \tilde{\Sigma}_{\varepsilon,k}^{(t)} \mathbf{A}^{(t)} \} \text{vec}(\mathbf{\Delta}_{(1),k}).$$

Similar to case (i), we replace  $\tilde{\Sigma}_k^{(t)}$  by  $\tilde{\Sigma}_{\varepsilon,k}^{(t)}$  in (S.2) and obtain an updating rule for case (ii):

$$\text{vec}(\mathbf{S}_{(1),k+1}) = \left[ \sum_{i=1}^n (\mathbf{U}_{\mathbf{S},i}^{(t)})^\top \mathbf{U}_{\mathbf{S},i}^{(t)} + \mathbf{I}_{R_2 R_3} \otimes \{ (\mathbf{A}^{(t)})^\top \tilde{\Sigma}_{\varepsilon,k}^{(t)} \mathbf{A}^{(t)} \} \right]^+ \left( \sum_{i=1}^n (\mathbf{U}_{\mathbf{S},i}^{(t)})^\top \mathbf{y}_i \right). \quad (\text{S.3})$$

The choice of  $\varepsilon$  can be

$$\varepsilon = \frac{\tau}{2n} \min \left\{ \|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),0}\|_2 : (\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),0} \neq 0, j = 1, \dots, p \right\}$$

with  $\tau = 10^{-4}$  as suggested in Hunter and Li (2005), where  $\mathbf{S}_{(1),0}$  is the initial value of  $\mathbf{S}_{(1)}$  for the MM algorithm to solve (4.1) of the main paper. We summarize the MM algorithm in Algorithm S.1.

---

**Algorithm S.1:** MM Algorithm for Updating  $\mathbf{S}$ .

---

**Input** : Data set  $\{\mathbf{y}_i, \mathbf{X}_i, t_i\}_{i=1}^n$ ; Fixed  $\mathbf{A}^{(t)}$ ,  $\mathbf{B}^{(t)}$ , and  $\mathbf{C}^{(t)}$ ; Random initial  $\mathbf{S}_{(1),0} \in \mathbb{R}^{R_1 \times R_2 R_3}$ ;

$\delta > 0$ .

**Output:**  $\tilde{\mathbf{S}}_{(1)}^{(t+1)}$  (so is  $\tilde{\mathbf{S}}^{(t+1)}$ ).

**for**  $k$  from 0, 1, ... **do**

**if** For any  $j = 1, \dots, p$ ,  $\|(\mathbf{a}_j^{(t)})^\top \mathbf{S}_{(1),k}\|_2 = 0$  and  $\|\mathbf{a}_j^{(t)}\|_2 \neq 0$  **then**

Update  $\text{vec}(\mathbf{S}_{(1),k+1})$  by (S.3) and calculate

$$d_{\mathbf{S}}^\varepsilon = H_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{S}_k) + \Phi_{\mathbf{A}^{(t)}}^\varepsilon(\mathbf{S}_k | \mathbf{S}_k) H_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{S}_{k+1}) - \Phi_{\mathbf{A}^{(t)}}^\varepsilon(\mathbf{S}_{k+1} | \mathbf{S}_k).$$

**else**

Update  $\text{vec}(\mathbf{S}_{(1),k+1})$  by (S.2) and calculate

$$d_{\mathbf{S}} = H_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{S}_k) + \Phi_{\mathbf{A}^{(t)}}(\mathbf{S}_k | \mathbf{S}_k) - H_{\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{S}_{k+1}) - \Phi_{\mathbf{A}^{(t)}}(\mathbf{S}_{k+1} | \mathbf{S}_k).$$

**end**

If  $d_{\mathbf{S}} \leq \delta$  or  $d_{\mathbf{S}}^\varepsilon \leq \delta$ , then output  $\mathbf{S}_{(1),k+1}$  as  $\tilde{\mathbf{S}}_{(1)}^{(t+1)}$ . If not, set  $k \leftarrow k + 1$  and continue the loop.

**end**

---

### S.1.2 Details on Updating $\mathbf{A}$

In this subsection, we provide more details on solving (4.4) of the main paper through the Alternating Direction Method of Multipliers (ADMM, Gabay and Mercier, 1976) algorithm. Let  $\mathbf{A}_{(k)}$ ,  $\mathbf{\Gamma}_{(k)}$ , and  $\boldsymbol{\nu}_{(k)}$  be the  $k$ -th iteration of the ADMM algorithm. With a penalty parameter  $\rho$ , the ADMM solves the subproblem (4.4) by generating the

following iterates at the  $(k+1)$ -th step

$$\begin{aligned}\mathbf{A}_{(k+1)} &= \arg \min_{\mathbf{A}} \left\{ H_{\tilde{\mathbf{S}}^{(t+1)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{A}) + \rho/2 \left\| \mathbf{A} \tilde{\mathbf{S}}_{(1)}^{(t+1)} - \mathbf{\Gamma}_{(k)} + \boldsymbol{\nu}_{(k)}/\rho \right\|_2^2 \right\}, \\ \mathbf{\Gamma}_{(k+1)} &= \arg \min_{\mathbf{\Gamma}} \left\{ \lambda \sum_{j=1}^p g(\gamma_j) + \rho/2 \left\| \mathbf{A}_{(k+1)} \tilde{\mathbf{S}}_{(1)}^{(t+1)} - \mathbf{\Gamma} + \boldsymbol{\nu}_{(k)}/\rho \right\|_2^2 \right\}, \\ \boldsymbol{\nu}_{(k+1)} &= \boldsymbol{\nu}_{(k)} + \rho (\mathbf{A}_{(k+1)} \tilde{\mathbf{S}}_{(1)}^{(t+1)} - \mathbf{\Gamma}_{(k+1)}).\end{aligned}\tag{S.4}$$

For simplicity, denote  $\bar{\mathbf{y}} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_n^\top)^\top \in \mathbb{R}^{nq}$  and  $\mathbf{U}_{\mathbf{A}} = (\mathbf{U}_{\mathbf{A},1}^\top, \dots, \mathbf{U}_{\mathbf{A},n}^\top)^\top \in \mathbb{R}^{nq \times pR_1}$ , where  $\mathbf{U}_{\mathbf{A},i} = \{[\mathbf{C}^{(t)} \otimes \{\mathbf{b}^\top(t_i) \mathbf{B}^{(t)}\}] (\tilde{\mathbf{S}}_{(1)}^{(t+1)})^\top\} \otimes \mathbf{x}_i^\top \in \mathbb{R}^{q \times pR_1}$ ,  $i = 1, 2, \dots, n$ .  $H_{\tilde{\mathbf{S}}^{(t+1)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{A})$  can be reformulated as

$$H_{\tilde{\mathbf{S}}^{(t+1)}, \mathbf{B}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{A}) = \frac{1}{2} \left\| \bar{\mathbf{y}} - \mathbf{U}_{\mathbf{A}} \text{vec}(\mathbf{A}) \right\|_2^2.$$

We first note that the first line of (S.4) has a closed-form solution

$$\begin{aligned}\text{vec}(\mathbf{A}) &= \left[ \mathbf{U}_{\mathbf{A}}^\top \mathbf{U}_{\mathbf{A}} + \rho \left\{ \tilde{\mathbf{S}}_{(1)}^{(t+1)} \otimes \mathbf{I}_p \right\} \left\{ (\tilde{\mathbf{S}}_{(1)}^{(t+1)})^\top \otimes \mathbf{I}_p \right\} \right]^{-1} \\ &\quad \times \left[ \mathbf{U}_{\mathbf{A}}^\top \bar{\mathbf{y}} + \rho \left\{ \tilde{\mathbf{S}}_{(1)}^{(t+1)} \otimes \mathbf{I}_p \right\} \left\{ \text{vec}(\mathbf{\Gamma}_{(k)}) - \frac{1}{\rho} \text{vec}(\boldsymbol{\nu}_{(k)}) \right\} \right].\end{aligned}\tag{S.5}$$

As for the second line of (S.4), we can split it into  $p$  subproblems:

$$\gamma_{j,(k+1)} = \arg \min_{\gamma \in \mathbb{R}^{R_2 R_3}} \left\{ \lambda g(\gamma) + \frac{\rho}{2} \left\| (\tilde{\mathbf{S}}_{(1)}^{(t+1)})^\top \mathbf{a}_{j,(k+1)} - \gamma + \frac{1}{\rho} \boldsymbol{\nu}_{j,(k)} \right\|_2^2 \right\}, \quad j = 1, 2, \dots, p,$$

where  $\mathbf{a}_{j,(k+1)}^\top$  and  $\boldsymbol{\nu}_{j,(k)}^\top$  are the  $j$ -th rows of  $\mathbf{A}_{(k+1)}$  and  $\boldsymbol{\nu}_{(k)}$  respectively. We introduce a proximal operator Boyd et al. (2011) to simplify our problem. For a given function  $h(\cdot)$ , its proximal map is defined by

$$\text{prox}_{h(\cdot)}(u) = \arg \min_v \left\{ h(v) + \frac{1}{2} \|u - v\|_2^2 \right\}.$$

The proximal map  $\text{prox}_{h(\cdot)}(u)$  exists and is unique for all  $u$  if  $h(\cdot)$  is a closed and convex function. Using this notation, we observe that  $\gamma_{j,(k+1)} = \text{prox}_{\lambda g(\cdot)/\rho}(\mathbf{u}_j)$  with  $\mathbf{u}_j = (\tilde{\mathbf{S}}_{(1)}^{(t+1)})^\top \mathbf{a}_{j,(k+1)} + \boldsymbol{\nu}_{j,(k)}/\rho$ . Furthermore, we can obtain the analytic form of the proximal operator with respect to  $g(\cdot)$ ,

$$\gamma_{j,(k+1)} = \left( 1 - \frac{\lambda}{\rho \|\mathbf{u}_j\|_2} \right)_+ \mathbf{u}_j,\tag{S.6}$$

where  $(x)_+ = \max(0, x)$ .

We use the standard stopping criterion for the ADMM algorithm based on primal and dual residuals as in Section 3.3 of Boyd et al. (2011). To be specific, we let  $r_{(k+1)}$  and  $s_{(k+1)}$  denote the primal and dual residuals, which

are defined as

$$\begin{cases} r_{(k+1)} = \|\mathbf{A}_{(k+1)} \tilde{\mathbf{S}}_{(1)}^{(t+1)} - \mathbf{\Gamma}_{(k+1)}\|_2 \\ s_{(k+1)} = \|\rho(\mathbf{\Gamma}_{(k+1)} - \mathbf{\Gamma}_{(k)}) (\tilde{\mathbf{S}}_{(1)}^{(t+1)})^\top\|_2. \end{cases} \quad (\text{S.7})$$

The ADMM algorithm is terminated when  $r_{(k+1)} \leq \epsilon^{\text{pri}}$  and  $s_{(k+1)} \leq \epsilon^{\text{dual}}$  are satisfied, where the primal and dual feasibility  $\epsilon^{\text{pri}}$  and  $\epsilon^{\text{dual}}$  are specified as suggested by Boyd et al. (2011). Theoretically, the primal and dual residuals will converge to zero for any fixed  $\rho > 0$ . However, different choices of  $\rho$  may lead to different convergence speeds. An efficient and stable way of varying penalty strategy is suggested by Zhu (2017). The updating rule for  $\rho$  is set to be

$$\rho = \begin{cases} \eta\rho, & \text{if } r_{(k)}/\epsilon^{\text{pri}} \geq \mu s_{(k)}/\epsilon^{\text{dual}} \\ \eta^{-1}\rho, & \text{if } s_{(k)}/\epsilon^{\text{dual}} \geq \mu r_{(k)}/\epsilon^{\text{pri}}, \end{cases} \quad (\text{S.8})$$

where  $\mu$  and  $\eta$  are set to be 10 and 2, respectively, as suggested by Boyd et al. (2011). The above discussions lead to Algorithm S.2 for updating  $\mathbf{A}$ .

### S.1.3 Details on Updating $\mathbf{B}$

In this section, we provide detailed steps to update  $\mathbf{B}$  using the manifold gradient algorithm. Recall that

$$\mathbf{B}^{(t+1)} = \arg \min_{\mathbf{B} \in \text{St}(R_2, K)} H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}),$$

where  $\text{St}(R_2, K)$  denotes the Stiefel manifold

$$\text{St}(R_2, K) = \{\mathbf{B} \in \mathbb{R}^{K \times R_2} : \mathbf{B}^\top \mathbf{B} = \mathbf{I}_{R_2}\}.$$

Note that  $H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})$  has the form as

$$H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}) = \sum_{i=1}^n \left\| \mathbf{y}_i - ([\mathbf{C}^{(t)} \otimes (\mathbf{x}_i^\top \mathbf{A}^{(t+1)})] \{\mathbf{S}_{(2)}^{(t+1)}\}^\top] \otimes \mathbf{b}^\top(t_i)) \text{vec}(\mathbf{B}) \right\|_2^2.$$

Let  $\mathbf{B}_{(k)}$  denote the  $k$ -th iteration of the manifold gradient algorithm. First observe that the Euclidean gradient of

$H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})$  can be directly obtained as

$$\nabla H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}) = -2 \text{Mat}_{K, R_2} \left[ \sum_{i=1}^n \{(\mathbf{U}_{\mathbf{B}, i}^{(t)})^\top \mathbf{y}_i - (\mathbf{U}_{\mathbf{B}, i}^{(t)})^\top \mathbf{U}_{\mathbf{B}, i}^{(t)} \text{vec}(\mathbf{B})\} \right],$$

---

**Algorithm S.2:** ADMM for Updating  $\mathbf{A}$ .
 

---

**Input:** Data set  $\{\mathbf{y}_i, \mathbf{X}_i, t_i\}_{i=1}^n$ ; Fixed  $\tilde{\mathbf{S}}^{(t+1)}$ ,  $\mathbf{B}^{(t)}$ , and  $\mathbf{C}^{(t)}$ ; Primal and dual feasibility  $\epsilon^{\text{pri}}$  and  $\epsilon^{\text{dual}}$ ; Parameters of varying penalty strategy  $\mu$  and  $\eta$ ; Random initial  $\mathbf{A} \in \mathbb{R}^{p \times R_1}$ .

**Output:**  $\tilde{\mathbf{A}}^{(t+1)}$  and  $\mathbf{\Gamma}^{(t+1)}$  (as for the alternative of  $\tilde{\mathbf{A}}^{(t+1)}\tilde{\mathbf{S}}_{(1)}^{(t+1)}$ ).

**for**  $k = 0, 1, \dots$  **do**

- 1. Update  $\mathbf{A}_{(k+1)}$  using (S.5).
- 2. Update  $\mathbf{\Gamma}_{(k+1)}$  through by rows through (S.6).
- 3. Update  $\boldsymbol{\nu}_{(k+1)}$  through the third line of (S.4).
- 4. Calculate the primal and dual residuals  $r_{(k+1)}$  and  $s_{(k+1)}$  as in (S.7).

**if**  $r_{(k+1)} \leq \epsilon^{\text{pri}}$  and  $s_{(k+1)} \leq \epsilon^{\text{dual}}$  **then**

| Output  $\mathbf{A}_{(k+1)}$  as  $\tilde{\mathbf{A}}^{(t+1)}$  and  $\mathbf{\Gamma}_{(k+1)}$  as  $\mathbf{\Gamma}^{(t+1)}$ .

**else**

| Update  $\rho$  using the varying penalty strategy (S.8) and back to step 1 with  $k \leftarrow k + 1$ .

**end**

**end**

---

where  $\mathbf{U}_{\mathbf{B},i}^{(t)} = [\{\mathbf{C}^{(t)} \otimes (\mathbf{x}_i^\top \mathbf{A}^{(t+1)})\}(\mathbf{S}_{(2)}^{(t+1)})^\top] \otimes \mathbf{b}^\top(t_i)$  and  $\text{Mat}_{K,R_2}(\cdot)$  transforms a vector with length  $K \cdot R_2$  into a matrix with  $K \times R_2$  dimension by its columns. We denote  $\text{grad}_{H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})}$  as the Riemannian gradient of  $H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})$ , and define a projection operator  $\mathcal{P}_{\mathbf{B}}(\boldsymbol{\xi})$  be

$$\mathcal{P}_{\mathbf{B}}(\boldsymbol{\xi}) = (\mathbf{I} - \mathbf{B}\mathbf{B}^\top)\boldsymbol{\xi} + \frac{1}{2}\mathbf{B}(\mathbf{B}^\top\boldsymbol{\xi} - \boldsymbol{\xi}^\top\mathbf{B}),$$

where  $\boldsymbol{\xi}$  is a generic matrix with dimension  $K \times R_2$ . According to Absil et al. (2009), the Riemannian gradient  $\text{grad}_{H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})}$  at  $\mathbf{B}$  is simply the projection of  $\nabla H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})$  onto the tangent space, that is,

$$\text{grad}_{H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})} = \mathcal{P}_{\mathbf{B}}\{\nabla H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B})\}. \quad (\text{S.9})$$

After getting the Riemannian gradient, we perform a line-search on the negative Riemannian gradient and retract the update back onto the Stiefel manifold. One of the commonly used retractions of Stiefel manifold is (Absil et al., 2009)

$$R_{\mathbf{B}}(\boldsymbol{\xi}) = \text{qf}(\mathbf{B} + \boldsymbol{\xi}),$$

where  $\text{qf}(\cdot)$  denotes the Q factor of the QR decomposition for a matrix. Therefore, we update

$$\mathbf{B}_{(k+1)} = \text{qf}\{\mathbf{B}_{(k)} - \tau_{\mathbf{B}_{(k)}} \text{grad } H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}_{(k)})\}, \quad (\text{S.10})$$

where  $\tau_{\mathbf{B}_{(k)}}$  is the Armijo step size. We select  $\tau_{\mathbf{B}_{(k)}}$  such that, for some fixed  $\epsilon_{\mathbf{B}} \in (0, 1)$ , the following inequality holds

$$\begin{aligned} H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}_{(k+1)}) &\leq H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}_{(k)}) \\ &\quad - \epsilon_{\mathbf{B}} \tau_{\mathbf{B}_{(k)}} \|\text{grad } H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}_{(k)})\|_F^2. \end{aligned} \quad (\text{S.11})$$

The above discussion leads to Algorithm S.3 as follows.

---

**Algorithm S.3:** Stiefel Manifold Optimization for Updating  $\mathbf{B}$ .

---

**Input** : Data set  $\{\mathbf{y}_i, \mathbf{X}_i, t_i\}_{i=1}^n$ ; Fixed  $\mathbf{S}^{(t+1)}$ ,  $\mathbf{A}^{(t+1)}$ , and  $\mathbf{C}^{(t)}$ ; Random initial  $\mathbf{B}_{(0)} \in \mathbb{R}^{K \times R_2}$ ;

$\epsilon_{\mathbf{B}} > 0$  and  $\epsilon > 0$ .

**Output:**  $\mathbf{B}^{(t+1)}$

**for**  $k$  from 0, 1, ... **do**

1. Compute the Riemannian gradient at  $\mathbf{B}_{(k)}$  as in (S.9).
2. Obtain  $\mathbf{B}_{(k+1)}$  as in (S.10) with step size  $\tau_{\mathbf{B}_{(k)}}$ , where  $\tau_{\mathbf{B}_{(k)}}$  is chosen according to (S.11).
3. Calculate the descent of the objective function

$$d_{\mathbf{B}} = H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}_{(k)}) - H_{\mathbf{S}^{(t+1)}, \mathbf{A}^{(t+1)}, \mathbf{C}^{(t)}}(\mathbf{B}_{(k+1)}).$$

If  $d_{\mathbf{B}} \leq \epsilon$ , then output  $\mathbf{B}_{(k+1)}$  as  $\mathbf{B}^{(t+1)}$ . If not, return to step 1.

**end**

---

## S.2 Proof of Theorem 1

To present the proof, we use  $C$  with or without subscripts to represent generic positive constants that may change values from line to line. Recall that  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^\top \in \mathbb{R}^{n \times pK}$ , where  $\mathbf{z}_i = \mathbf{b}(t_i) \otimes \mathbf{x}_i$ . For any  $\mathcal{I} \subset \{1, \dots, p\}$ , define  $\mathbf{Z}_{\mathcal{I}} \in \mathbb{R}^{|\mathcal{I}| \times pK}$  such that it has the same columns as  $\mathbf{Z}$  for those predictors within the index  $\mathcal{I}$  and zero columns otherwise. Denote  $\mathcal{P}_{\mathcal{I}}$  as the projection matrix onto the column space of  $\mathbf{Z}_{\mathcal{I}}$ .

With the notation of  $\mathbf{z}_i$ , we can write

$$\{\mathbf{G} \bar{\times}_2 \mathbf{b}(t_i)\}^\top \mathbf{x}_i = (\mathbf{A} \otimes \mathbf{B}) \mathbf{S}_{(3)}^\top \mathbf{C}^\top \mathbf{z}_i, \quad i = 1, \dots, n,$$

in (3.5) of the main paper by using the matricization operator of a tensor (Kolda and Bader, 2009). Let  $\mathbf{D} = (\mathbf{A} \otimes \mathbf{B}) \mathbf{S}_{(3)}^\top \mathbf{C}^\top \in \mathbb{R}^{pK \times q}$ . It can be seen that  $\mathbf{D} = \mathbf{G}_{(3)}^\top$ , where  $\mathbf{G}_{(3)}$  is the mode-3 matricization of  $\mathbf{G}$ . Let  $\mathbf{D}_j \in \mathbb{R}^{K \times q}$  be the collection of rows of  $\mathbf{D}$  associated with the predictor  $j$ .

The optimization problem (3.6) of the main paper can then be rewritten as

$$\arg \min_{\mathbf{D}} \|\mathbf{Y} - \mathbf{Z}\mathbf{D}\|_F^2 + \lambda \sum_{j=1}^p \|\mathbf{D}_j\|_F. \quad (\text{S.12})$$

Respectively denoting  $\mathbf{D}_0 = \mathbf{G}_{0,(3)}^\top$  and  $\widehat{\mathbf{D}}$  as the solution of (S.12) (which is equal to  $\widehat{\mathbf{G}}_{(3)}^\top$ ), we have  $\Delta_{\mathbf{G}}^2 = \|\mathbf{Z}(\mathbf{D} - \mathbf{D}_0)\|_F^2$  and  $\Delta_{\widehat{\mathbf{G}}}^2 = \|\mathbf{Z}(\widehat{\mathbf{D}} - \mathbf{D}_0)\|_F^2$ . Recall that  $\mathcal{J}(\mathbf{G}) = \{j : \mathbf{D}_j \neq \mathbf{0}\}$  and  $\mathcal{J}(\widehat{\mathbf{G}}) = \{j : \widehat{\mathbf{D}}_j \neq \mathbf{0}\}$ . Denote  $\tilde{\mathcal{J}} = \mathcal{J}(\mathbf{G}) \cup \mathcal{J}(\widehat{\mathbf{G}})$  for simplicity.

By definition, we can obtain

$$\|\mathbf{Y} - \mathbf{Z}\widehat{\mathbf{D}}\|_F^2 + \lambda \sum_{j \in \mathcal{J}(\widehat{\mathbf{G}})} \|\widehat{\mathbf{D}}_j\|_F \leq \|\mathbf{Y} - \mathbf{Z}\mathbf{D}\|_F^2 + \lambda \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j\|_F,$$

which implies

$$\begin{aligned} \Delta_{\widehat{\mathbf{G}}}^2 + \lambda \sum_{j \in \mathcal{J}(\mathbf{G})^c} \|\widehat{\mathbf{D}}_j - \mathbf{D}_j\|_F &\leq \Delta_{\mathbf{G}}^2 + 2\langle \mathbf{E}, \mathbf{Z}_{\tilde{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}) \rangle \\ &\quad + 2\langle \mathbf{R}, \mathbf{Z}_{\tilde{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}) \rangle + \lambda \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \widehat{\mathbf{D}}_j\|_F. \end{aligned} \quad (\text{S.13})$$

Denote  $d_1(\cdot)$  as the largest singular value of a matrix. Since  $\text{rank}(\mathbf{D}) \leq R_3$  and  $\text{rank}(\mathbf{S}_{(3)}) \leq R_1 R_2$ , we have

$$\text{rank}\{\mathbf{Z}_{\tilde{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D})\} \leq 2R = 2 \min\{R_3, R_1 R_2\}.$$



It yields

$$\begin{aligned}
 2\langle \mathbf{E}, \mathbf{Z}_{\bar{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}) \rangle &= 2\langle \mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}, \mathbf{Z}_{\bar{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}) \rangle \\
 &\leq 2d_1(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E})\|\mathbf{Z}_{\bar{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D})\|_* \\
 &\leq 2d_1(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E})\sqrt{2R}\|\mathbf{Z}_{\bar{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D})\|_F \\
 &\leq 2d_1(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E})\sqrt{2R}(\Delta_{\widehat{\mathbf{G}}} + \Delta_{\mathbf{G}}) \\
 &\leq 16Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) + \frac{1}{4}(\Delta_{\widehat{\mathbf{G}}}^2 + \Delta_{\mathbf{G}}^2),
 \end{aligned} \tag{S.14}$$

where  $\mathcal{P}_{\bar{\mathcal{J}}}$  is the projection matrix of the space spanned by  $\mathbf{Z}_{\bar{\mathcal{J}}}$  and  $\|\cdot\|_*$  denotes the nuclear norm. The second, third, fourth, and fifth lines of (S.14) are due to the Cauchy-Schwarz inequality, the relationship between the nuclear norm and Frobenius norm, the triangle inequality, and the inequality of arithmetic and geometric means, respectively. It follows from (S.13) and (S.14),

$$\begin{aligned}
 \frac{3}{4}\Delta_{\widehat{\mathbf{G}}}^2 + \lambda \sum_{j \in \mathcal{J}(\mathbf{G})^c} \|\widehat{\mathbf{D}}_j - \mathbf{D}_j\|_F &\leq \frac{5}{4}\Delta_{\widehat{\mathbf{G}}}^2 + 16Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) \\
 &\quad + 2\langle \mathbf{R}, \mathbf{Z}_{\bar{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}) \rangle + \lambda \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \widehat{\mathbf{D}}_j\|_F.
 \end{aligned}$$

Now, there are two cases as follows.

**Case 1.** If

$$\lambda \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \widehat{\mathbf{D}}_j\|_F \geq \frac{5}{4}\Delta_{\widehat{\mathbf{G}}}^2 + 16Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) + 2\langle \mathbf{R}, \mathbf{Z}_{\bar{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}) \rangle, \tag{S.15}$$

it straightforwardly implies

$$\sum_{j \in \mathcal{J}(\mathbf{G})^c} \|\widehat{\mathbf{D}}_j - \mathbf{D}_j\|_F \leq 2 \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \widehat{\mathbf{D}}_j\|_F.$$

By Condition  $\mathcal{M}(\mathcal{J}(\mathbf{G}), \delta_{\mathcal{J}(\mathbf{G})})$ , we have

$$\begin{aligned}
 \lambda \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \widehat{\mathbf{D}}_j\|_F &\leq \frac{4}{n\delta_{\mathcal{J}(\mathbf{G})}}\lambda^2|\mathcal{J}(\mathbf{G})| + \frac{n\delta_{\mathcal{J}(\mathbf{G})}}{16} \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \widehat{\mathbf{D}}_j\|_F^2 \\
 &\leq \frac{4}{n\delta_{\mathcal{J}(\mathbf{G})}}\lambda^2|\mathcal{J}(\mathbf{G})| + \frac{1}{16}\|\mathbf{Z}\widehat{\mathbf{D}} - \mathbf{Z}\mathbf{D}\|_F^2 \\
 &\leq \frac{4}{n\delta_{\mathcal{J}(\mathbf{G})}}\lambda^2|\mathcal{J}(\mathbf{G})| + \frac{1}{8}(\Delta_{\widehat{\mathbf{G}}}^2 + \Delta_{\mathbf{G}}^2).
 \end{aligned} \tag{S.16}$$

Using (S.15), (S.16), and the Cauchy-Schwarz inequality, we further obtain

$$\frac{3}{4}\Delta_{\widehat{\mathbf{G}}}^2 \leq 2\lambda \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \widehat{\mathbf{D}}_j\|_F \leq \frac{8}{n\delta_{\mathcal{J}(\mathbf{G})}}\lambda^2|\mathcal{J}(\mathbf{G})| + \frac{1}{4}(\Delta_{\widehat{\mathbf{G}}}^2 + \Delta_{\mathbf{G}}^2).$$

Thus,

$$\Delta_{\hat{\mathbf{G}}}^2 \leq 2\Delta_{\mathbf{G}}^2 + \frac{16\lambda^2|\mathcal{J}(\mathbf{G})|}{n\delta_{\mathcal{J}(\mathbf{G})}}. \quad (\text{S.17})$$

For notational simplicity, we let

$$\gamma = 24n\lambda_{\max}(\boldsymbol{\Sigma})K\sigma^2\{1 + \log(p)\}. \quad (\text{S.18})$$

Taking  $\lambda^2 = 768\gamma R_3 R$  in (S.17) will then lead to

$$\Delta_{\hat{\mathbf{G}}}^2 \leq C_1\Delta_{\mathbf{G}}^2 + C_2\|\mathbf{R}\|_F^2 + \frac{C_3 R_3 R \lambda_{\max}(\boldsymbol{\Sigma}) K \sigma^2 \{1 + \log(p)\} |\mathcal{J}(\mathbf{G})|}{\delta_{\mathcal{J}(\mathbf{G})}}. \quad (\text{S.19})$$

**Case 2.** If

$$\lambda \sum_{j \in \mathcal{J}(\mathbf{G})} \|\mathbf{D}_j - \hat{\mathbf{D}}_j\|_F < \frac{5}{4}\Delta_{\mathbf{G}}^2 + 16Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) + 2\langle \mathbf{R}, \mathbf{Z}_{\bar{\mathcal{J}}}(\hat{\mathbf{D}} - \mathbf{D}) \rangle,$$

using the Cauchy-Schwarz inequality, we have

$$\begin{aligned} 3\Delta_{\hat{\mathbf{G}}}^2 &\leq 10\Delta_{\mathbf{G}}^2 + 128Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) + 16\langle \mathbf{R}, \mathbf{Z}_{\bar{\mathcal{J}}}(\hat{\mathbf{D}} - \mathbf{D}) \rangle \\ &\leq 10\Delta_{\mathbf{G}}^2 + 128Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) + 16\|\mathbf{R}\|_F(\Delta_{\hat{\mathbf{G}}} + \Delta_{\mathbf{G}}). \end{aligned} \quad (\text{S.20})$$

The quadratic form of inequality (S.20) implies that

$$\Delta_{\hat{\mathbf{G}}} \leq \frac{16\|\mathbf{R}\|_F + \sqrt{\zeta_2}}{6}, \quad (\text{S.21})$$

where

$$\zeta_2 = 256\|\mathbf{R}\|_F^2 + 12\{10\Delta_{\mathbf{G}}^2 + 128Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) + 16\|\mathbf{R}\|_F\Delta_{\mathbf{G}}\}.$$

Plugging (S.21) into (S.20) shows that

$$9\Delta_{\hat{\mathbf{G}}}^2 \leq 304\|\mathbf{R}\|_F^2 + 108\Delta_{\mathbf{G}}^2 + 768Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}). \quad (\text{S.22})$$

According to Lemma 8 of Bunea et al. (2012), we have

$$\mathbb{P}(d_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) - 6K\sigma^2\{|\tilde{\mathcal{J}}| + |\tilde{\mathcal{J}}|\log(ep/|\tilde{\mathcal{J}}|)\} - 6\sigma^2q > \tilde{t}) \leq \frac{16\sigma^2 \exp(-q/2)}{\tilde{t}}$$

for  $\tilde{t} > 0$ . Taking  $\tilde{t} = 6K\sigma^2 \log(p)$ , we then have

$$\begin{aligned} d_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) &\leq \tilde{t} + 6K\sigma^2\{|\tilde{\mathcal{J}}| + |\tilde{\mathcal{J}}|\log(ep/|\tilde{\mathcal{J}}|)\} + 6\sigma^2q \\ &\leq 12K\sigma^2\{1 + \log(p)\}\{|\mathcal{J}(\mathbf{G})| + |\mathcal{J}(\hat{\mathbf{G}})|\} + 6\sigma^2q \end{aligned} \quad (\text{S.23})$$

with probability at least

$$1 - \frac{8 \exp(-q/2)}{3K \log(p)}. \quad (\text{S.24})$$

Next, observe that when  $\widehat{\mathbf{C}}$  of the solution of (3.6) of the main paper is given, the other components  $(\widehat{\mathbf{S}}, \widehat{\mathbf{A}}, \widehat{\mathbf{B}})$  can be obtained through solving

$$\arg \min_{\mathbf{S}, \mathbf{A}, \mathbf{B}} \|\mathbf{Y}\widehat{\mathbf{C}} - \mathbf{Z}(\mathbf{A} \otimes \mathbf{B})\mathbf{S}_{(3)}^\top\|_F^2 + \lambda \sum_{j=1}^p \|\mathbf{a}_j^\top \mathbf{S}_{(1)}\|_F,$$

where  $\mathbf{a}_j^\top$  is the  $j$ -th row of  $\mathbf{A}$ . By a QR decomposition, we can represent  $\mathbf{A}\mathbf{S}_{(1)} = \mathbf{U}\mathbf{V}^\top$  such that  $\mathbf{U} \in \mathbb{R}^{p \times R_1}$ ,  $\mathbf{V} \in \mathbb{R}^{R_2 R_3 \times R_1}$  and  $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_{R_1}$ . Thus, the solution of (3.6) of the main paper, which is in the form of  $(\mathbf{S}, \mathbf{A}, \mathbf{B}, \mathbf{C})$ , can be reparametrized into the form of  $(\mathbf{U}, \mathbf{V}, \mathbf{B}, \mathbf{C})$ . In other words, we can find a  $(\widehat{\mathbf{V}}, \widehat{\mathbf{U}})$  such that  $\widehat{\mathbf{G}}_{(1)} = \widehat{\mathbf{A}}\widehat{\mathbf{S}}_{(1)}(\widehat{\mathbf{C}} \otimes \widehat{\mathbf{B}})^\top = \widehat{\mathbf{U}}\widehat{\mathbf{V}}^\top(\widehat{\mathbf{C}} \otimes \widehat{\mathbf{B}})^\top$ . Let

$$\widetilde{\mathbf{Z}}_i = [\mathbf{x}_i^\top \otimes \{\mathbf{I}_{R_3} \otimes \mathbf{b}^\top(t_i)\widehat{\mathbf{B}}\}] (\mathbf{I}_p \otimes \widehat{\mathbf{V}}) \in \mathbb{R}^{R_3 \times p R_1}, \quad i = 1, \dots, n. \quad (\text{S.25})$$

Therefore,  $\widehat{\mathbf{U}}$  can be obtained from solving

$$\widehat{\mathbf{U}} = \arg \min_{\mathbf{U}} \|\mathbf{Y}\widehat{\mathbf{C}} - (\widetilde{\mathbf{Z}}_1 \mathbf{u}, \dots, \widetilde{\mathbf{Z}}_n \mathbf{u})^\top\|_F^2 + \lambda \sum_{j=1}^p \|\mathbf{u}_j\|_2, \quad (\text{S.26})$$

where  $\mathbf{u} = \text{vec}(\mathbf{U}^\top) \in \mathbb{R}^{p R_1}$  and  $\mathbf{u}_j$  is the  $j$ -th row of  $\mathbf{U}$ . Since  $\mathbf{Y}\widehat{\mathbf{C}} \in \mathbb{R}^{n \times R_3}$ , we write the  $i$ -th row of  $\mathbf{Y}\widehat{\mathbf{C}}$  as  $\gamma_i^\top$ ,  $i = 1, \dots, n$ , that is,

$$\mathbf{Y}\widehat{\mathbf{C}} = (\gamma_1, \dots, \gamma_n)^\top.$$

Using this notation, (S.26) can be written as

$$\widehat{\mathbf{U}} = \arg \min_{\mathbf{U}} \sum_{i=1}^n \|\gamma_i - \widetilde{\mathbf{Z}}_i \mathbf{u}\|_2^2 + \lambda \sum_{j=1}^p \|\mathbf{u}_j\|_2.$$

Analogously, let  $\widehat{\mathbf{u}} = \text{vec}(\widehat{\mathbf{U}})$  and  $\widehat{\mathbf{u}}_j$  be the  $j$ -th row of  $\widehat{\mathbf{U}}$ . To solve  $\widehat{\mathbf{u}}_j$  from the above displayed objective function, we have

$$\sum_{i=1}^n \widetilde{\mathbf{Z}}_{ij}^\top (\gamma_i - \widetilde{\mathbf{Z}}_i \widehat{\mathbf{u}}) = \lambda \frac{\widehat{\mathbf{u}}_j}{\|\widehat{\mathbf{u}}_j\|_2} \in \mathbb{R}^{R_1},$$

where  $\widetilde{\mathbf{Z}}_{ij} \in \mathbb{R}^{R_3 \times R_1}$  is the submatrix of  $\widetilde{\mathbf{Z}}_i$  defined in (S.25), associated the  $j$ -th predictor,  $j \in \mathcal{J}(\widehat{\mathbf{G}})$ . For simplicity, denote

$$\widetilde{\mathbf{Z}}_\# = (\widetilde{\mathbf{Z}}_1^\top, \dots, \widetilde{\mathbf{Z}}_n^\top)^\top = (\widetilde{\mathbf{Z}}_{\#,1}, \dots, \widetilde{\mathbf{Z}}_{\#,p}) \in \mathbb{R}^{n R_3 \times p R_1},$$

where  $\tilde{\mathbf{Z}}_{\# , j} = (\tilde{\mathbf{Z}}_{1j}^\top, \dots, \tilde{\mathbf{Z}}_{nj}^\top)^\top \in \mathbb{R}^{nR_3 \times R_1}$ . We then have

$$\left\| \tilde{\mathbf{Z}}_{\# , j}^\top [(\boldsymbol{\gamma}_1^\top, \dots, \boldsymbol{\gamma}_n^\top)^\top - \{(\tilde{\mathbf{Z}}_1 \hat{\mathbf{u}})^\top, \dots, (\tilde{\mathbf{Z}}_n \hat{\mathbf{u}})^\top\}^\top] \right\|_2^2 = \lambda^2, \quad j \in \mathcal{J}(\hat{\mathbf{G}}).$$

Thus,

$$\begin{aligned} \lambda^2 |\mathcal{J}(\hat{\mathbf{G}})| &= \sum_{j \in \mathcal{J}(\hat{\mathbf{G}})} \left\| \tilde{\mathbf{Z}}_{\# , j}^\top [(\boldsymbol{\gamma}_1^\top, \dots, \boldsymbol{\gamma}_n^\top)^\top - \{(\tilde{\mathbf{Z}}_1 \hat{\mathbf{u}})^\top, \dots, (\tilde{\mathbf{Z}}_n \hat{\mathbf{u}})^\top\}^\top] \right\|_2^2 \\ &\leq 2\lambda_{\max}(\tilde{\mathbf{Z}}_{\#} \tilde{\mathbf{Z}}_{\#}^\top) \{ \|\mathbf{Z}(\mathbf{A}_0 \otimes \mathbf{B}_0) \mathbf{S}_{0,(3)}^\top \mathbf{C}_0^\top \hat{\mathbf{C}} - \mathbf{Z}(\hat{\mathbf{A}} \otimes \hat{\mathbf{B}}) \hat{\mathbf{S}}_{(3)}^\top \hat{\mathbf{C}}^\top \hat{\mathbf{C}}\|_F^2 + \|\mathcal{P}_{\tilde{\mathcal{J}}} \mathbf{E} \hat{\mathbf{C}}\|_F^2 + \|\mathbf{R} \hat{\mathbf{C}}\|_F^2 \} \\ &\leq 2\lambda_{\max}(\tilde{\mathbf{Z}}_{\#} \tilde{\mathbf{Z}}_{\#}^\top) \{ R_3 \Delta_{\hat{\mathbf{G}}}^2 + R_3 d_1^2(\mathcal{P}_{\tilde{\mathcal{J}}} \mathbf{E}) + R_3 \|\mathbf{R}\|_F^2 \}. \end{aligned} \quad (\text{S.27})$$

Note that the eigenvalues of

$$\tilde{\mathbf{Z}}_{\#}^\top \tilde{\mathbf{Z}}_{\#} = \sum_{i=1}^n \tilde{\mathbf{Z}}_i^\top \tilde{\mathbf{Z}}_i = \sum_{i=1}^n (\mathbf{I}_p \otimes \hat{\mathbf{V}}^\top) \{ \mathbf{x}^\top \otimes (\mathbf{I}_{R_3} \otimes \mathbf{b}^\top \hat{\mathbf{B}}) \}^\top \{ \mathbf{x}^\top \otimes (\mathbf{I}_{R_3} \otimes \mathbf{b}^\top \hat{\mathbf{B}}) \} (\mathbf{I}_p \otimes \hat{\mathbf{V}})$$

are the same as those of

$$\sum_{i=1}^n \{ \mathbf{x}^\top \otimes (\mathbf{I}_{R_3} \otimes \mathbf{b}^\top \hat{\mathbf{B}}) \}^\top \{ \mathbf{x}^\top \otimes (\mathbf{I}_{R_3} \otimes \mathbf{b}^\top \hat{\mathbf{B}}) \}. \quad (\text{S.28})$$

Similarly, it can be shown that (S.28) has the same eigenvalues as of

$$\sum_{i=1}^n \{ \mathbf{x} \otimes \hat{\mathbf{B}}^\top \mathbf{b} \} \{ \mathbf{x}^\top \otimes \mathbf{b}^\top \hat{\mathbf{B}} \}.$$

By (S.27) and Lemma S.1 , we obtain

$$\lambda^2 |\mathcal{J}(\hat{\mathbf{G}})| \leq 2n\lambda_{\max}(\boldsymbol{\Sigma}) \{ R_3 \Delta_{\hat{\mathbf{G}}}^2 + R_3 d_1^2(\mathcal{P}_{\tilde{\mathcal{J}}} \mathbf{E}) + R_3 \|\mathbf{R}\|_F^2 \}. \quad (\text{S.29})$$

For simplicity, let

$$\alpha = 6\sigma^2 q + 16\sigma^2 \exp(-q/2) + 12K\sigma^2 |\mathcal{J}(\mathbf{G})| \{1 + \log(p)\}. \quad (\text{S.30})$$

Recall that  $\gamma = 24n\lambda_{\max}(\boldsymbol{\Sigma})K\sigma^2 \{1 + \log(p)\}$ . It follows from (S.23) and (S.29) that

$$d_1^2(\mathcal{P}_{\tilde{\mathcal{J}}} \mathbf{E}) \leq \alpha + \frac{\gamma R_3}{\lambda^2} (\Delta_{\hat{\mathbf{G}}}^2 + d_1^2(\mathcal{P}_{\tilde{\mathcal{J}}} \mathbf{E}) + \|\mathbf{R}\|_F^2),$$

with probability at least (S.24), which yields

$$\left(1 - \frac{\gamma R_3}{\lambda^2}\right) d_1^2(\mathcal{P}_{\tilde{\mathcal{J}}} \mathbf{E}) \leq \alpha + \frac{\gamma R_3}{\lambda^2} (\Delta_{\hat{\mathbf{G}}}^2 + \|\mathbf{R}\|_F^2),$$

that is,

$$d_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) \leq \frac{\alpha}{1 - \gamma R_3/\lambda^2} + \frac{\gamma R_3/\lambda^2}{1 - \gamma R_3/\lambda^2} (\Delta_{\mathbf{G}}^2 + \|\mathbf{R}\|_F^2), \quad (\text{S.31})$$

Plugging (S.31) into (S.22) yields

$$\begin{aligned} & \frac{9 - 9\gamma R_3/\lambda^2 - 768\gamma R_3 R/\lambda^2}{1 - \gamma R_3/\lambda^2} (\Delta_{\mathbf{G}}^2) \\ & \leq 108(\Delta_{\mathbf{G}}^2) + \frac{768R\alpha}{1 - \gamma R_3/\lambda^2} + \left( \frac{768\gamma R_3 R/\lambda^2}{1 - \gamma R_3/\lambda^2} + 304 \right) \|\mathbf{R}\|_F^2. \end{aligned} \quad (\text{S.32})$$

Taking  $\lambda^2 = 768\gamma R_3 R$  in (S.32) yields

$$\Delta_{\mathbf{G}}^2 \leq C_1 \Delta_{\mathbf{G}}^2 + C_2 q R \sigma^2 + C_3 R K \sigma^2 |\mathcal{J}(\mathbf{G})| \{1 + \log(p)\} + C_4 \|\mathbf{R}\|_F^2, \quad (\text{S.33})$$

with probability at least (S.24). By definitions, it is obviously to see  $\lambda_{\max}(\mathbf{\Sigma}) \geq \delta_{\mathcal{J}(\mathbf{G})}$ . With (5.3) of the main paper, (S.19), and (S.33), it shows that

$$\Delta_{\mathbf{G}}^2 \leq C_1 \Delta_{\mathbf{G}}^2 + C_2 q R \sigma^2 + C_3 \frac{R_3 R K |\mathcal{J}(\mathbf{G})| \lambda_{\max}(\mathbf{\Sigma}) \sigma^2 \log(p)}{\delta_{\mathcal{J}(\mathbf{G})}} + C_4 \frac{nsq}{K^{2\tau}},$$

with probability at least (S.24), which finishes the proof.

## S.3 Proof of Corollaries

### S.3.1 Proof of Corollary 1

By definitions, (5.3), (5.5), and (5.7) of the main paper, we have

$$\begin{aligned} \Delta_{\hat{\mathbf{F}}}^2 &= \sum_{i=1}^n \|\hat{\mathbf{F}}(t_i)^\top \mathbf{x}_i - \mathbf{F}_0(t_i)^\top \mathbf{x}_i\|^2 \\ &= \sum_{i=1}^n \|\hat{\mathbf{F}}(t_i)^\top \mathbf{x}_i - \{\mathbf{G}_0 \bar{\times}_2 \mathbf{b}(t_i)\}^\top \mathbf{x}_i + \{\mathbf{G}_0 \bar{\times}_2 \mathbf{b}(t_i)\}^\top \mathbf{x}_i - \mathbf{F}_0(t_i)^\top \mathbf{x}_i\|^2 \\ &\leq 2 \sum_{i=1}^n \left\{ \|\hat{\mathbf{F}}(t_i)^\top \mathbf{x}_i - \{\mathbf{G}_0 \bar{\times}_2 \mathbf{b}(t_i)\}^\top \mathbf{x}_i\|^2 + \|\{\mathbf{G}_0 \bar{\times}_2 \mathbf{b}(t_i)\}^\top \mathbf{x}_i - \mathbf{F}_0(t_i)^\top \mathbf{x}_i\|^2 \right\} \\ &\leq 2S_5 \Delta_{\mathbf{G}}^2 + 2S_6 q R \sigma^2 + 2S_7 \frac{R_3 R K |\mathcal{J}(\mathbf{G})| \lambda_{\max}(\mathbf{\Sigma}) \sigma^2 \log(p)}{\delta_{\mathcal{J}(\mathbf{G})}} + (2S_8 + 2S_3) \frac{nsq}{K^{2\tau}}, \end{aligned} \quad (\text{S.34})$$

with probability at least (S.24).

### S.3.2 Proof of Corollary 2

Theorem 1 of the main paper implies that when  $\Sigma$  satisfies Condition  $\mathcal{M}(\mathcal{J}(\mathbf{G}_0), \delta_{\mathcal{J}(\mathbf{G}_0)})$ ,

$$\Delta_{\widehat{\mathbf{G}}}^2 \leq S_6 q R \sigma^2 + S_7 \frac{R_3 R K s \lambda_{\max}(\Sigma) \sigma^2 \log(p)}{\delta_{\mathcal{J}(\mathbf{G}_0)}} + S_8 \frac{nsq}{K^{2\tau}}, \quad (\text{S.35})$$

with probability at least (S.24). Similarly, applying the arguments as in the proof of Corollary 1, it can be shown the convergence rate of  $\Delta_{\widehat{\mathbf{F}}}^2$  has the upper bound

$$O_p \left( q R \sigma^2 + \frac{R_3 R K s \lambda_{\max}(\Sigma) \sigma^2 \log(p)}{\delta_{\mathcal{J}(\mathbf{G}_0)}} + \frac{nsq}{K^{2\tau}} \right),$$

by specifying  $\mathbf{G} = \mathbf{G}_0$  in (S.34) under the assumption that  $\Sigma$  satisfies Condition  $\mathcal{M}(\mathcal{J}(\mathbf{G}_0), \delta_{\mathcal{J}(\mathbf{G}_0)})$ . In the above bound, let

$$K \asymp \left\{ \frac{n \delta_{\mathcal{J}(\mathbf{G}_0)} q}{R_3 R \lambda_{\max}(\Sigma) \log(p)} \right\}^{1/(2\tau+1)}, \quad (\text{S.36})$$

we then have

$$\Delta_{\widehat{\mathbf{F}}}^2 = O_p \left( q R + \left\{ \frac{R_3 R \lambda_{\max}(\Sigma) \log(p)}{\delta_{\mathcal{J}(\mathbf{G}_0)}} \right\}^{2\tau/(2\tau+1)} s(nq)^{1/(2\tau+1)} \right),$$

which finishes the proof of (5.9) of the main paper.

Next, we will show the relevant predictors will be selected consistently. Recall that  $\mathbf{D}_0 = \mathbf{G}_{0,(3)}^\top$  and  $\widehat{\mathbf{D}} = \widehat{\mathbf{G}}_{(3)}^\top$ .

By (5.1) of the main paper and the orthonormality of  $\mathbf{b}(t)$ , we observe that for  $j \in \mathcal{J}(\mathbf{G}_0)$ ,

$$\begin{aligned} \sum_{l=1}^q \|f_{0,jl}\|_2^2 &\leq \sum_{l=1}^q \left( \left\| \sum_{k=1}^K G_{0,jkl} b_k \right\|_2 + \left\| \sum_{k=1}^K G_{0,jkl} b_k - f_{0,jl} \right\|_2 \right)^2 \\ &\leq 2 \|\mathbf{D}_{0,j}\|_F^2 + \frac{Cq}{K^{2\tau}}, \end{aligned}$$

where  $\mathbf{D}_{0,j} \in \mathbb{R}^{K \times q}$  is the submatrix of  $\mathbf{D}_0$  associated with the predictor  $j$ . Now, with  $K$  chosen to be (S.36), and the assumption that  $\sum_{l=1}^q \|f_{0,jl}\|_2^2 \geq S_{11}$ , it follows that when (5.10) of the main paper is satisfied, we have

$$\|\mathbf{D}_{0,j}\|_F > C, \quad (\text{S.37})$$

for some positive constant  $C$ . On the other hand, the triangular inequality implies

$$\|\mathbf{D}_{0,j}\|_F \leq \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F + \|\widehat{\mathbf{D}}_j\|_F. \quad (\text{S.38})$$

Thus, if we can show as  $n$  tends to infinity

$$\mathbb{P}\{\|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F < C/2, j \in \mathcal{J}(\mathbf{G}_0)\} \rightarrow 1, \quad (\text{S.39})$$

then following (S.37)–(S.39), we obtain

$$\mathbb{P}\{\|\widehat{\mathbf{D}}_j\|_F > C/2, j \in \mathcal{J}(\mathbf{G}_0)\} \rightarrow 1.$$

Using the orthonormality of  $\mathbf{b}(t)$ , the above result will complete the proof of this corollary.

It remains to prove (S.39). Recall that

$$\begin{aligned} \frac{3}{4}\Delta_{\mathbf{G}}^2 + \lambda \sum_{j \in \mathcal{J}(\mathbf{G}_0)^c} \|\widehat{\mathbf{D}}_j - \mathbf{D}_{0,j}\|_F &\leq 16Rd_1^2(\mathcal{P}_{\widetilde{\mathcal{J}}}\mathbf{E}) \\ &\quad + 2\langle \mathbf{R}, \mathbf{Z}_{\widetilde{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}_0) \rangle + \lambda \sum_{j \in \mathcal{J}(\mathbf{G}_0)} \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F, \end{aligned}$$

where  $\widetilde{\mathcal{J}} = \mathcal{J}(\mathbf{G}_0) \cup \mathcal{J}(\widehat{\mathbf{G}})$ , and there are two cases as follows.

**Case 1.** If

$$\lambda \sum_{j \in \mathcal{J}(\mathbf{G}_0)} \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F \geq 16Rd_1^2(\mathcal{P}_{\widetilde{\mathcal{J}}}\mathbf{E}) + 2\langle \mathbf{R}, \mathbf{Z}_{\widetilde{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}_0) \rangle,$$

it straightforwardly implies

$$\sum_{j \in \mathcal{J}(\mathbf{G}_0)^c} \|\widehat{\mathbf{D}}_j - \mathbf{D}_{0,j}\|_F \leq 2 \sum_{j \in \mathcal{J}(\mathbf{G}_0)} \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F.$$

By Condition  $\mathcal{M}(\mathcal{J}(\mathbf{G}_0), \delta_{\mathcal{J}(\mathbf{G}_0)})$ , we have

$$\begin{aligned} \delta_{\mathcal{J}(\mathbf{G}_0)} \sum_{j \in \mathcal{J}(\mathbf{G}_0)} \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F^2 &\leq \text{tr}\{(\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j)^\top \boldsymbol{\Sigma} (\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j)\} \\ &\leq \frac{\Delta_{\mathbf{G}}^2}{n}. \end{aligned} \tag{S.40}$$

Using (S.19), (S.36), and (5.3) of the main paper, we have

$$\begin{aligned} \Delta_{\mathbf{G}}^2 &\leq C_1 \frac{nsq}{K^{2\tau}} + \frac{C_2 R_3 R K s \lambda_{\max}(\boldsymbol{\Sigma}) \sigma^2 \{1 + \log(p)\}}{\delta_{\mathcal{J}(\mathbf{G}_0)}} \\ &\leq C_3 \left\{ \frac{R_3 R \lambda_{\max}(\boldsymbol{\Sigma}) \log(p)}{\delta_{\mathcal{J}(\mathbf{G}_0)}} \right\}^{2\tau/(2\tau+1)} s(nq)^{1/(2\tau+1)}. \end{aligned} \tag{S.41}$$

With (5.10), we prove (S.39) by plugging (S.41) in (S.40).

**Case 2.** If

$$\lambda \sum_{j \in \mathcal{J}(\mathbf{G}_0)} \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F < 16Rd_1^2(\mathcal{P}_{\widetilde{\mathcal{J}}}\mathbf{E}) + 2\langle \mathbf{R}, \mathbf{Z}_{\widetilde{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}_0) \rangle, \tag{S.42}$$

we will use the following arguments to prove (S.39). By (S.31), we have with probability at least (S.24),

$$\begin{aligned} Rd_1^2(\mathcal{P}_{\widetilde{\mathcal{J}}}\mathbf{E}) &\leq \frac{\alpha R}{1 - 1/(768R)} + \frac{1/768}{1 - 1/(768R)} (\Delta_{\mathbf{G}}^2 + \|\mathbf{R}\|_F^2) \\ &\leq C_1 \alpha R + C_2 (\Delta_{\mathbf{G}}^2 + \|\mathbf{R}\|_F^2), \end{aligned} \tag{S.43}$$

where  $\alpha$  is defined as in (S.30). It follows from (S.42), (5.3) of the main paper, (S.35), and (S.43) that

$$\begin{aligned} \lambda \sum_{j \in \mathcal{J}(\mathbf{G}_0)} \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F &< 16Rd_1^2(\mathcal{P}_{\bar{\mathcal{J}}}\mathbf{E}) + 2\langle \mathbf{R}, \mathbf{Z}_{\bar{\mathcal{J}}}(\widehat{\mathbf{D}} - \mathbf{D}_0) \rangle \\ &\leq C_1 Rq + C_2 \frac{nsq}{K^{2\tau}} + C_3 \frac{R_3 R K s \lambda_{\max}(\boldsymbol{\Sigma}) \log(p)}{\delta_{\mathcal{J}(\mathbf{G}_0)}}, \end{aligned}$$

with probability at least (S.24). Using (5.6) of the main paper and (S.36), it implies

$$\begin{aligned} \lambda^2 &= S_4 R_3 R n \lambda_{\max}(\boldsymbol{\Sigma}) K \sigma^2 \{1 + \log(p)\} \\ &\geq S_4 R_3 R n \lambda_{\max}(\boldsymbol{\Sigma}) K \sigma^2 \log(p) \\ &\asymp C \{R_3 R \lambda_{\max}(\boldsymbol{\Sigma}) \log(p)\}^{2\tau/(2\tau+1)} n^{2(\tau+1)/(2\tau+1)} \{q \delta_{\mathcal{J}(\mathbf{G}_0)}\}^{1/(2\tau+1)}. \end{aligned}$$

We then have

$$\begin{aligned} &\sum_{j \in \mathcal{J}(\mathbf{G}_0)} \|\mathbf{D}_{0,j} - \widehat{\mathbf{D}}_j\|_F \\ &\leq C_1 q^{(4\tau+1)/(4\tau+2)} R^{(\tau+1)/(2\tau+1)} \{R_3 \lambda_{\max}(\boldsymbol{\Sigma}) \log(p)\}^{-\tau/(2\tau+1)} n^{-(\tau+1)/(2\tau+1)} \delta_{\mathcal{J}(\mathbf{G}_0)}^{-1/(4\tau+2)} \\ &\quad + C_2 s \{R R_3 \lambda_{\max}(\boldsymbol{\Sigma}) \log(p)\}^{\tau/(2\tau+1)} n^{-\tau/(2\tau+1)} q^{1/(4\tau+2)} \delta_{\mathcal{J}(\mathbf{G}_0)}^{-(4\tau+1)/(4\tau+2)}, \end{aligned} \tag{S.44}$$

with probability at least (S.24). Using (5.10) of the main paper, (S.36), and (S.44), we have (S.39), which completes the proof.

## S.4 Random Design

This section is organized as follows. We first provide two auxiliary lemmas that are used as preliminary results (Lemmas S.1 and S.2). Lemma S.3 show that the Condition  $\mathcal{M}(\mathcal{J}, \delta_{\mathcal{J}})$  can be satisfied with high probability (tending to one) when  $\mathbf{x}$  and  $t$  are random, and the sample size  $n$  is large enough compared with  $|\mathcal{J}|^2 q^2 K^2 + |\mathcal{J}|^2 q K \log p$ . A special application of this result is for  $\mathcal{J}(\mathbf{G}_0)$  with  $|\mathcal{J}(\mathbf{G}_0)| = s$ . The term  $\lambda_{\max}(\boldsymbol{\Sigma})/\delta_{\mathcal{J}}$  appeared in Theorem 1 of the main paper is determined in Lemmas S.4 and S.5 for the cases when  $p$  and  $q$  are diverging with the sample size  $n$  and fixed constants, respectively.

For notational simplicity, we use  $C$  with or without subscripts to represent generic constants that may change values from line to line.



**Lemma S.1.** Let  $\mathbf{z}_{i,\mathcal{I}} = \mathbf{x}_{i,\mathcal{I}} \otimes \mathbf{b}(t_i)$ , where  $\mathcal{I} \subset \{1, \dots, p\}$ ,  $i = 1, \dots, n$ . If the following eigenvalue bounds

$$L_{\mathcal{I}} \leq \lambda_{\min} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}_{i,\mathcal{I}} \mathbf{z}_{i,\mathcal{I}}^{\top} \right) \leq \lambda_{\max} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}_{i,\mathcal{I}} \mathbf{z}_{i,\mathcal{I}}^{\top} \right) \leq U_{\mathcal{I}}$$

hold, then

$$L_{\mathcal{I}} \leq \lambda_{\min} \left( \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{Z}}_{i,\mathcal{I}} \bar{\mathbf{Z}}_{i,\mathcal{I}}^{\top} \right) \leq \lambda_{\max} \left( \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{Z}}_{i,\mathcal{I}} \bar{\mathbf{Z}}_{i,\mathcal{I}}^{\top} \right) \leq U_{\mathcal{I}},$$

where  $\bar{\mathbf{Z}}_{i,\mathcal{I}} = \mathbf{x}_{i,\mathcal{I}} \otimes \mathbf{B}^{\top} \mathbf{b}(t_i)$ , for any  $\mathbf{B}$  satisfying  $\mathbf{B}^{\top} \mathbf{B} = \mathbf{I}_{R_2}$ .

*Proof.* By definition, for all  $w_{jk} \in \mathbb{R}$ , we have

$$\left\{ \sum_{j \in \mathcal{I}} \sum_{k=1}^K \frac{1}{n} \sum_{i=1}^n x_{ij} b_k(t_i) w_{jk} \right\}^2 \leq U_{\mathcal{I}} \sum_{j \in \mathcal{I}} \sum_{k=1}^K w_{jk}^2.$$

Denote  $\mathbf{B} = (\mathbf{B}_1, \dots, \mathbf{B}_K)^{\top}$ . We then have

$$\begin{aligned} \left\{ \sum_{j \in \mathcal{I}} \sum_{r_2=1}^{R_2} \frac{1}{n} \sum_{i=1}^n x_{ij} \sum_{k=1}^K B_{kr_2} b_k(t_i) \bar{w}_{jr_2} \right\}^2 &= \left\{ \sum_{j \in \mathcal{I}} \sum_{k=1}^K \frac{1}{n} \sum_{i=1}^n x_{ij} b_k(t_i) \sum_{r_2=1}^{R_2} B_{kr_2} \bar{w}_{jr_2} \right\}^2 \\ &\leq U_{\mathcal{I}} \sum_{j \in \mathcal{I}} \sum_{k=1}^K (\bar{w}_j^{\top} \mathbf{B}_k \mathbf{B}_k^{\top} \bar{w}_j) \\ &= U_{\mathcal{I}} \sum_{j \in \mathcal{I}} (\bar{w}_j^{\top} \mathbf{B}^{\top} \mathbf{B} \bar{w}_j) \\ &= U_{\mathcal{I}} \sum_{j \in \mathcal{I}} \sum_{r_2=1}^{R_2} \bar{w}_{jr_2}^2, \end{aligned}$$

where  $\bar{w}_{jr_2} \in \mathbb{R}$  and  $\bar{\mathbf{w}}_j = (\bar{w}_{j1}, \dots, \bar{w}_{jR_2})^{\top}$ . Therefore we finish the proof of the upper bound. The proof of lower bound can be obtained similarly.  $\square$

**Lemma S.2.** Suppose the conditional density of  $t$  given  $\mathbf{x}$ , denoted as  $f_{t|\mathbf{x}}(t)$ , is bounded away from 0 and infinity by a constant, that is, there exists  $C > 0$

$$1/C < f_{t|\mathbf{x}}(t) < C, \quad t \in [0, 1].$$

If  $\|\mathbf{x}\|_{\psi_2} \leq \kappa$ , then

$$\|\mathbf{x} \otimes \mathbf{b}(t)\|_{\psi_2} \leq C\kappa\sqrt{K} \tag{S.45}$$

and

$$\lambda_{\max} [\mathbb{E}\{\mathbf{x} \otimes \mathbf{b}(t)\}\{\mathbf{x}^{\top} \otimes \mathbf{b}^{\top}(t)\}] \leq C\kappa^2. \tag{S.46}$$

Further, if  $\lambda_{\min}\{\mathbb{E}(\mathbf{x}\mathbf{x}^\top)\} \geq C$ , then

$$\lambda_{\min}[\mathbb{E}\{\mathbf{x} \otimes \mathbf{b}(t)\}\{\mathbf{x}^\top \otimes \mathbf{b}^\top(t)\}] \geq C. \quad (\text{S.47})$$

*Proof.* Denote  $\{\tilde{b}_k(t)\}_{k=1}^K$  as the ordinary B-spline basis with the same knots and order as  $\{b_k(t)\}$ . By equivalence, there exists a nonnegative matrix  $\Psi \in \mathbb{R}^{K \times K}$  such that

$$\mathbf{b}(t) = \Psi \tilde{\mathbf{b}}(t), \quad (\text{S.48})$$

where  $\tilde{\mathbf{b}}(t) = (\tilde{b}_1(t), \dots, \tilde{b}_K(t))^\top$ . It follows that

$$\mathbf{I}_K = \Psi \left\{ \int \tilde{\mathbf{b}}(t) \tilde{\mathbf{b}}^\top(t) dt \right\} \Psi^\top$$

and

$$\Psi^{-1}(\Psi^\top)^{-1} = \int \tilde{\mathbf{b}}(t) \tilde{\mathbf{b}}^\top(t) dt.$$

By the property of B-spline (see, e.g., De Boor, 1973, 1976) and Assumption 3, for  $1 \leq q \leq +\infty$  and  $\mathbf{v} \in \mathbb{R}^K$ , we have

$$C_\zeta \|\mathbf{v}\|_q \leq h_n^{-1/q} \left\| \sum_{k=1}^K v_k \tilde{b}_k \right\|_q \leq C \|\mathbf{v}\|_q, \quad (\text{S.49})$$

where  $C_\zeta$  and  $C$  are two positive constants and  $C_\zeta$  depends on the order of B-spline  $\zeta$ . Thus,

$$C_1 \frac{1}{K} \leq \lambda_{\min}\{\Psi^{-1}(\Psi^\top)^{-1}\} \leq \lambda_{\max}\{\Psi^{-1}(\Psi^\top)^{-1}\} \leq C_2 \frac{1}{K},$$

which yields

$$C_2 K \leq \lambda_{\min}\{\Psi \Psi^\top\} \leq \lambda_{\max}\{\Psi \Psi^\top\} \leq C_1 K. \quad (\text{S.50})$$

Recall

$$\|\mathbf{x} \otimes \mathbf{b}(t)\|_{\psi_2} = \sup_{\|\mathbf{w}\|_2 \leq 1} \|\langle \mathbf{x} \otimes \mathbf{b}(t), \mathbf{w} \rangle\|_{\psi_2}.$$

For all  $\mathbf{w} \in \mathbb{R}^{pK}$  with  $\|\mathbf{w}\|_2 \leq 1$ , we then have

$$\begin{aligned} \|\langle \mathbf{x} \otimes \mathbf{b}(t), \mathbf{w} \rangle\|_{L^q}^q &= \mathbb{E}_{\mathbf{x}} \int \left| \sum_{j,k} x_j b_k(t) w_{jk} \right|^q f_{t|\mathbf{x}}(t) dt \\ &= \mathbb{E}_{\mathbf{x}} \int \left| \sum_{j,k} x_j \left( \sum_{\tilde{k}} \Psi_{k,\tilde{k}} \tilde{b}_{\tilde{k}}(t) \right) w_{jk} \right|^q f_{t|\mathbf{x}}(t) dt \end{aligned}$$

$$\begin{aligned}
 &\leq C_1 \mathbb{E}_{\mathbf{x}} \int \left| \sum_{\tilde{k}} \left( \sum_{k,j} \Psi_{k,\tilde{k}} x_j w_{jk} \right) \tilde{b}_{\tilde{k}}(t) \right|^q dt \\
 &\leq C^q \frac{1}{K} \mathbb{E}_{\mathbf{x}} \left\| \left\{ \left( \sum_{k,j} \Psi_{k,1} x_j w_{jk} \right), \dots, \left( \sum_{k,j} \Psi_{k,K} x_j w_{jk} \right) \right\} \right\|_q^q \\
 &\leq C^q \frac{1}{K} \sum_{\tilde{k}=1}^K \mathbb{E}_{\mathbf{x}} \left| \sum_{k,j} \Psi_{k,\tilde{k}} x_j w_{jk} \right|^q \\
 &\leq C^q \frac{1}{K} \sum_{\tilde{k}=1}^K \mathbb{E}_{\mathbf{x}} \left| \sum_j x_j \left( \sum_k \Psi_{k,\tilde{k}} w_{jk} \right) \right|^q \\
 &\leq C^q \frac{1}{K} \sum_{k=1}^K \|\langle \mathbf{x}, \{\Psi_{:,k}^{\top}(\mathbf{w}_1, \dots, \mathbf{w}_p)\}^{\top} \rangle, \|_{L_q}^q \\
 &\leq C^q \frac{1}{K} \sum_{k=1}^K (\sqrt{q} \|\mathbf{x}\|_{\psi_2} \|\Psi_{:,k}^{\top}(\mathbf{w}_1, \dots, \mathbf{w}_p)\|_2)^q \\
 &\leq C^q \frac{1}{K} q^{q/2} \kappa^q \sum_{k=1}^K \left( \|\Psi_{:,k}^{\top}(\mathbf{w}_1, \dots, \mathbf{w}_p)\|_2 \right)^{q/2} \\
 &\leq C^q K^{q/2-1} q^{q/2} \kappa^q, \tag{S.51}
 \end{aligned}$$

where  $\mathbf{w}_j = (w_{j1}, \dots, w_{jK})^{\top}$ ,  $w_{jk}$  is the  $\{(j-1)K+k\}$ -th element of  $\mathbf{w}$  and  $\Psi_{:,k} = (\Psi_{1,k}, \dots, \Psi_{K,k})^{\top}$ . Therefore, we obtain

$$\|\langle \mathbf{x} \otimes \mathbf{b}(t), \mathbf{w} \rangle\|_{L_q} \leq C \kappa \sqrt{q} K^{1/2-1/q} \leq C \kappa \sqrt{qK},$$

which yields the result of (S.45). Taking  $q = 2$  in (S.51) will lead to the result of (S.46). Similarly,

$$\begin{aligned}
 \|\langle \mathbf{x} \otimes \mathbf{b}(t), \mathbf{w} \rangle\|_{L_2}^2 &= \mathbb{E}_{\mathbf{x}} \int \left| \sum_{\tilde{k}} \left( \sum_{k,j} x_j \Psi_{k,\tilde{k}} w_{jk} \right) \tilde{b}_{\tilde{k}}(t) \right|^2 f_{t|\mathbf{x}}(t) dt \\
 &\geq C \frac{1}{K} \mathbb{E}_{\mathbf{x}} \left\| \left\{ \left( \sum_{k,j} \Psi_{k,1} x_j w_{jk} \right), \dots, \left( \sum_{k,j} \Psi_{k,K} x_j w_{jk} \right) \right\} \right\|_2^2 \\
 &\geq C \frac{1}{K} \sum_{\tilde{k}=1}^K \mathbb{E}_{\mathbf{x}} \left| \sum_j x_j \left( \sum_k \Psi_{k,\tilde{k}} w_{jk} \right) \right|^2 \\
 &\geq C \frac{1}{K} \sum_{\tilde{k}=1}^K \|\Psi_{:,k}^{\top}(\mathbf{w}_1, \dots, \mathbf{w}_p)\|_2^2 \\
 &\geq C \|\mathbf{w}\|_2^2,
 \end{aligned}$$

which proves (S.47). □

Before presenting the next lemmas, we first introduce  $\gamma_{\alpha}$ -functionals defined in (Talagrand, 2005; Banerjee et al., 2015). In Lemmas S.3–S.5,  $\gamma_2$  will be used in the proofs.

**Definition S.1** ( $\gamma_\alpha$ -functionals). Consider a metric space  $(T, d)$  and for a finite set  $\mathcal{A} \subset T$ , let  $|\mathcal{A}|$  denote its cardinality. An admissible sequence is an increasing sequence of subsets  $\{\mathcal{A}_n, n \geq 0\}$  of  $T$ , such that  $|\mathcal{A}_0| = 1$  and for  $n \geq 1$ ,  $|\mathcal{A}_n| = 2^{2^n}$ . Given  $\alpha > 0$ , we define the  $\gamma_\alpha$ -functional as

$$\gamma_\alpha(T, d) = \inf \sup_{t \in T} \sum_{n=0}^{\infty} \text{Diam}\{A_n(t)\},$$

where  $A_n(t)$  is the unique element of  $\mathcal{A}_n$  that contains  $t$ ,  $\text{Diam}\{A_n(t)\}$  is the diameter of  $A_n$  according to  $d$ , and the infimum is over all admissible sequences of  $T$ .

The next lemma shows that Condition  $\mathcal{M}(\mathcal{J}, \delta_{\mathcal{J}})$  is satisfied with high probability under random models.

**Lemma S.3.** *If the conditional density of  $t$  given  $\mathbf{x}$  is bounded by a constant,  $\|\mathbf{x}\|_{\psi_2} \leq \kappa$  and  $\lambda_{\min}\{\mathbb{E}(\mathbf{x}\mathbf{x}^\top)\} \geq C_1$ , then Condition  $\mathcal{M}(\mathcal{J}, \delta_{\mathcal{J}})$  holds for some  $\delta_{\mathcal{J}} > 0$  with high probability when  $n \geq C_3(|\mathcal{J}|^2 q^2 K^2 + |\mathcal{J}|^2 q K \log p)$ , where  $C_3$  depends on  $C_1$  and  $\delta_{\mathcal{J}}$ .*

*Proof.* Recall  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^\top$ , where  $\mathbf{z}_i = \mathbf{x}_i \otimes \mathbf{b}(t_i)$ , which are i.i.d. copies of  $\mathbf{z} = \mathbf{x} \otimes \mathbf{b}(t)$ . Let  $\Sigma_{\mathbf{z}} = \mathbb{E}(\mathbf{z}\mathbf{z}^\top)$ .

Consider the class of functions:

$$\mathcal{F} = \left[ f_{\mathbf{M}}\{\mathbf{x} \otimes \mathbf{b}(t)\} = \frac{1}{\sqrt{\text{tr}(\mathbf{M}^\top \Sigma_{\mathbf{z}} \mathbf{M})}} \sqrt{\text{tr}[\mathbf{M}^\top \{\mathbf{x} \otimes \mathbf{b}(t)\} \{\mathbf{x} \otimes \mathbf{b}(t)\}^\top \mathbf{M}]} : \text{vec}(\mathbf{M}) \in \mathcal{A} \right],$$

where

$$\mathcal{A} = \left\{ \text{vec}(\mathbf{M}) : \|\mathbf{M}\|_F^2 \leq 1 \text{ and } 2 \sum_{j \in \mathcal{J}} \|\mathbf{M}_j\|_F \geq \sum_{j \in \mathcal{J}^c} \|\mathbf{M}_j\|_F \right\}.$$

For  $\mathbf{M} \in \mathcal{A}$ , we have

$$\sum_{j=1}^p \|\mathbf{M}_j\|_F \leq 3 \sum_{j \in \mathcal{J}} \|\mathbf{M}_j\|_F \leq 3|\mathcal{J}| \|\mathbf{M}\|_F = 3|\mathcal{J}|.$$

Define

$$\bar{\mathcal{A}} = \left\{ \text{vec}(\mathbf{M}) : \sum_{j=1}^p \|\mathbf{M}_j\|_F \leq 3|\mathcal{J}| \text{ and } 2 \sum_{j \in \mathcal{J}} \|\mathbf{M}_j\|_F \geq \sum_{j \in \mathcal{J}^c} \|\mathbf{M}_j\|_F \right\}.$$

We then have

$$\mathcal{A} \subset \bar{\mathcal{A}}. \tag{S.52}$$

Note that

$$\|f_{\mathbf{M}}\|_{L_2}^2 = \frac{1}{\text{tr}(\mathbf{M}^\top \Sigma_{\mathbf{z}} \mathbf{M})} \mathbb{E}(\text{tr}[\mathbf{M}^\top \{\mathbf{x} \otimes \mathbf{b}(t)\} \{\mathbf{x} \otimes \mathbf{b}(t)\}^\top \mathbf{M}]) = 1.$$

We then have  $\mathcal{F} \subset \mathcal{S}_{L_2} := \{f : \|f\|_{L_2} = 1\}$ . By definitions, we have

$$\|f_M\|_{\psi_2} = \|\|\tilde{\mathbf{z}}^\top \widetilde{\mathbf{M}}\|_2\|_{\psi_2} \leq \sum_{l=1}^q \|\tilde{\mathbf{z}}^\top \widetilde{\mathbf{M}}_l\|_{\psi_2} \leq \|\tilde{\mathbf{z}}\|_{\psi_2} \sum_{l=1}^q \|\widetilde{\mathbf{M}}_l\|_2 \leq \sqrt{q} \|\tilde{\mathbf{z}}\|_{\psi_2},$$

where  $\tilde{\mathbf{z}}^\top = \mathbf{z}^\top \boldsymbol{\Sigma}^{-1/2}$ ,  $\widetilde{\mathbf{M}} = \boldsymbol{\Sigma}^{1/2} \mathbf{M}$ , and  $\widetilde{\mathbf{M}}_l$  is the  $l$ -th column of  $\widetilde{\mathbf{M}}$ . Further, noting that

$$\langle \tilde{\mathbf{z}}, \mathbf{w} \rangle = \langle \mathbf{z}^\top \boldsymbol{\Sigma}_z^{-1/2}, \mathbf{w}^\top \rangle = \langle \mathbf{z}, \boldsymbol{\Sigma}_z^{-1/2} \mathbf{w} \rangle,$$

and

$$\|\boldsymbol{\Sigma}_z^{-1/2} \mathbf{w}\|_2 \leq C \|\mathbf{w}\|_2,$$

it follows from Lemma S.2 that

$$\|\tilde{\mathbf{z}}\|_{\psi_2} \leq C\sqrt{q} \|\mathbf{z}\|_{\psi_2} \leq C\sqrt{qK}\kappa.$$

By Theorem 2.1.1 of Talagrand (2005), we have

$$\gamma_2(\mathcal{F} \cap \mathcal{S}_{L_2}, \|\cdot\|_{\psi_2}) \leq C\sqrt{qK}\kappa \gamma_2(\mathcal{F} \cap \mathcal{S}_{L_2}, \|\cdot\|_{L_2}) \leq C\sqrt{qK}\kappa w(\mathcal{A}),$$

where  $\gamma_2$  is the  $\gamma_2$ -functional defined as in Definition S.1.

Using Theorem 10 of Banerjee et al. (2015), we chose

$$\theta = C\sqrt{qK}\kappa^2 \frac{w(\mathcal{A})}{\sqrt{n}} \geq C\kappa \frac{\gamma_2(\mathcal{F} \cap \mathcal{S}_{L_2}, \|\cdot\|_{\psi_2})}{\sqrt{n}}$$

to satisfy the lower bound in equation (125) of Banerjee et al. (2015). As a result,

$$\sup_{\mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \frac{1}{\text{tr}(\mathbf{M}^\top \boldsymbol{\Sigma}_z \mathbf{M})} \text{tr}[\mathbf{M}^\top \{\mathbf{x}_i \otimes \mathbf{b}(t_i)\} \{\mathbf{x}_i \otimes \mathbf{b}(t_i)\}^\top \mathbf{M}] - 1 \right| \leq C\sqrt{qK}\kappa^2 \frac{w(\mathcal{A})}{\sqrt{n}},$$

with probability at least

$$1 - \exp\{-CqKw^2(\mathcal{A})\}.$$

Applying equation (53) of Banerjee et al. (2015) and (S.52), we have

$$w(\mathcal{A}) \leq C|\mathcal{J}| \sqrt{Kq + \log p}.$$

The definition of  $\mathcal{A}$  and Lemma S.2 together show that, for  $\text{vec}(\mathbf{M}) \in \mathcal{A}$ ,

$$\text{tr}(\mathbf{M}^\top \boldsymbol{\Sigma}_z \mathbf{M}) \geq C\|\mathbf{M}\|_F^2 \geq C \sum_{j \in |\mathcal{J}|} \|\mathbf{M}_j\|_F^2.$$

When  $n \geq C_3(|\mathcal{J}|^2 q^2 K^2 + |\mathcal{J}|^2 q K \log p)$ , we have

$$\frac{1}{n} \sum_{i=1}^n \text{tr}[\mathbf{M}^\top \{\mathbf{x}_i \otimes \mathbf{b}(t_i)\} \{\mathbf{x}_i \otimes \mathbf{b}(t_i)\}^\top \mathbf{M}] \geq C_4 \text{tr}(\mathbf{M}^\top \boldsymbol{\Sigma}_z \mathbf{M}) \geq C_2 \sum_{j \in |\mathcal{J}|} \|\mathbf{M}_j\|_F^2$$

with high probability. Therefore, when  $n \geq C_3(|\mathcal{J}|^2 q^2 K^2 + |\mathcal{J}|^2 q K \log p)$ , Condition  $\mathcal{M}(\mathcal{J}, \delta_{\mathcal{J}})$  is satisfied for some  $\delta_{\mathcal{J}} = C_2 > 0$  with high probability.  $\square$

When  $q$  and  $p$  are allowed to be growing with the sample size  $n$ , the following lemma shows that  $\lambda_{\max}(\boldsymbol{\Sigma})$  is bounded by  $O(K)$  with probability tending to one. Therefore, together with Lemma S.3, it implies that the term  $\lambda_{\max}(\boldsymbol{\Sigma})/\delta_{\mathcal{J}}$  in Theorem 1 of the main paper is upper bounded by  $O(K)$  with high probability.

**Lemma S.4.** Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$ ,  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^\top$ , where  $\mathbf{z}_i = \mathbf{x}_i \otimes \mathbf{b}(t_i)$ ,  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$ . If

$$\lambda_{\max}\left(\frac{\mathbf{X}^\top \mathbf{X}}{n}\right) \leq U_x < \infty, \quad (\text{S.53})$$

then

$$\lambda_{\max}\left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n}\right) \leq C_1 K U_x. \quad (\text{S.54})$$

*Proof.* For any  $\mathbf{v} = (v_1, \dots, v_p)^\top \in \mathbb{R}^p$ , we have

$$\frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p x_{ij} v_j \right)^2 \leq U_x \sum_{j=1}^p v_j^2 \quad (\text{S.55})$$

due to the definition of  $\lambda_{\max}(\cdot)$ . It can be shown that for all  $\tilde{\mathbf{v}} = (\tilde{v}_1, \dots, \tilde{v}_p)^\top \in \mathbb{R}^p$ , we have

$$\frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p |x_{ij} \tilde{v}_j| \right)^2 \leq U_x \sum_{j=1}^p \tilde{v}_j^2.$$

To see this, if  $x_{ij} \leq 0$ , we take  $v_j = -|\tilde{v}_j|$ , and if  $x_{ij} > 0$ , we take  $v_j = |\tilde{v}_j|$ . The aforementioned inequality is thus obtained by (S.55). Next, consider  $\mathbf{w} = (\mathbf{w}_1^\top, \dots, \mathbf{w}_p^\top)^\top \in \mathbb{R}^{pK}$ , where  $\mathbf{w}_j = (w_{j1}, \dots, w_{jK})^\top \in \mathbb{R}^K$ . We now borrow some notations used in the proof of Lemma S.2. Recall that  $\{\tilde{b}_k(t)\}_{k=1}^K$  is the B-spline basis with the same knots and order as  $\{b_k(t)\}$ , and there exists  $\boldsymbol{\Psi} = (\Psi_{k, \tilde{k}})_{K \times K}$  such that (S.48) holds. By definitions and (S.50), we have

$$\begin{aligned} \frac{\mathbf{w}^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{w}}{n} &= \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{j=1}^p x_{ij} \sum_{k=1}^K b_k(t) w_{jk} \right\}^2 \\ &\leq \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{j=1}^p \left| \sum_{k=1}^K b_k(t) w_{jk} \right| |x_{ij}| \right\}^2 \\ &\leq U_x \sum_{j=1}^p \left( \sum_{k=1}^K b_k(t) w_{jk} \right)^2 \\ &= C U_x K \|\mathbf{w}\|_2^2, \end{aligned} \quad (\text{S.56})$$

which finishes the proof.  $\square$

When  $q$  and  $p$  are fixed (that is, not growing with the sample size  $n$ ), Lemma S.5 shows that with probability tending to one,  $\lambda_{\max}(\boldsymbol{\Sigma})/\delta_{\mathcal{J}}$  is upper bounded by a constant.

**Lemma S.5.** *If the conditional density of  $t$  given  $\mathbf{x}$  is bounded by a constant,  $\|\mathbf{x}\|_{\psi_2} \leq \kappa$  and  $\lambda_{\min}\{\mathbb{E}(\mathbf{x}\mathbf{x}^\top)\} \geq C_1$ , then Condition  $\mathcal{M}(\mathcal{J}, \delta_{\mathcal{J}})$  holds for some constant  $\delta_{\mathcal{J}} > 0$  and  $\lambda_{\max}(\boldsymbol{\Sigma})$  is upper bounded by a constant, with high probability when  $n \geq CK^2$ .*

*Proof.* The following proof is similar to that of Lemma S.3. Recall  $\mathbf{Z} = (z_1, \dots, z_n)^\top$ , where  $z_i = \mathbf{x}_i \otimes \mathbf{b}(t_i)$ , which are i.i.d. copies of  $\mathbf{z} = \mathbf{x} \otimes \mathbf{b}(t)$ . Let  $\boldsymbol{\Sigma}_z = \mathbb{E}(\mathbf{z}\mathbf{z}^\top)$ . Consider the class of functions:

$$\mathcal{F}_w = \left[ f_w\{\mathbf{x} \otimes \mathbf{b}(t)\} = \frac{1}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}_z \mathbf{w}}} \sqrt{\mathbf{w}^\top \{\mathbf{x} \otimes \mathbf{b}(t)\} \{\mathbf{x} \otimes \mathbf{b}(t)\}^\top \mathbf{w}} : \mathbf{w} \in \mathcal{A}_w \right],$$

where

$$\mathcal{A}_w = \{\mathbf{w} : \|\mathbf{w}\|_F^2 \leq 1\}.$$

It is trivial to see  $\|f_w\|_{L_2}^2 = 1$ , which yields  $\mathcal{F}_w \subset \mathcal{S}_{L_2} := \{f : \|f\|_{L_2} = 1\}$ . By definitions and Lemma S.2, we have

$$\|f_w\|_{\psi_2} = \|\|\mathbf{z}^\top \mathbf{w}\|_2\|_{\psi_2} \leq \|\|\mathbf{z}\|_{\psi_2}\|_{\psi_2} \|\mathbf{w}\|_2 \leq C\sqrt{K}\kappa.$$

By Theorem 2.1.1 of Talagrand (2005), we have

$$\gamma_2(\mathcal{F}_w \cap \mathcal{S}_{L_2}, \|\cdot\|_{\psi_2}) \leq C\sqrt{K}\kappa \gamma_2(\mathcal{F}_w \cap \mathcal{S}_{L_2}, \|\cdot\|_{L_2}) \leq C\sqrt{K}\kappa w(\mathcal{A}_w),$$

where  $\gamma_2$  is the  $\gamma_2$ -functional defined as in Definition S.1.

Using Theorem 10 of Banerjee et al. (2015), we chose

$$\theta = C\sqrt{K}\kappa^2 \frac{w(\mathcal{A}_w)}{\sqrt{n}} \geq C\kappa \frac{\gamma_2(\mathcal{F}_w \cap \mathcal{S}_{L_2}, \|\cdot\|_{\psi_2})}{\sqrt{n}}$$

to satisfy the lower bound in equation (125) of Banerjee et al. (2015). As a result,

$$\sup_{\mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \frac{1}{\mathbf{w}^\top \boldsymbol{\Sigma}_z \mathbf{w}} [\mathbf{w}^\top \{\mathbf{x}_i \otimes \mathbf{b}(t_i)\} \{\mathbf{x}_i \otimes \mathbf{b}(t_i)\}^\top \mathbf{w}] - 1 \right| \leq C\sqrt{K}\kappa^2 \frac{w(\mathcal{A}_w)}{\sqrt{n}}, \quad (\text{S.57})$$

with probability at least

$$1 - \exp\{-CKw^2(\mathcal{A}_w)\}.$$

By the covering number argument, we have

$$N(\epsilon, \mathcal{A}_w, l_2) \leq (C_1/\epsilon)^{C_2 K}.$$

It follows from the Dudley's integral entropy bound (see, e.g., Theorem 3.1 of Koltchinskii, 2011) that the Gaussian width satisfies

$$w(\mathcal{A}_w) \leq C\sqrt{K}. \tag{S.58}$$

If  $n \geq C\kappa^4 K^2$  for some constant  $C$ , by (S.46), (S.47), (S.57), and (S.58), we have

$$C_1 \leq \lambda_{\min}\left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n}\right) \leq \lambda_{\max}\left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n}\right) \leq C_2 \kappa^2$$

with high probability, which finishes the proof. □

## S.5 Additional Numerical Results

### S.5.1 Additional Experiments of Synthetic Data

We extend our simulation setting to larger numbers of response variables to show the trend of the performance of the proposed method when  $q$  increases. To match the measure of estimation error used in our theoretical investigation (see (5.4) and (5.5) in Section 5 of the main paper), we calculate the integrated squared error (ISE) in this additional study, that is,

$$\text{ISE} = \sum_{j=1}^p \sum_{l=1}^q \int_0^1 \{\widehat{f}_{jl}(t) - f_{jl}(t)\}^2 dt,$$

where  $\widehat{f}_{jl}(t)$ 's are the estimated coefficient functions of various methods. More precisely, we calculate the ISE for the scenarios of  $q = 15, 30, \text{ and } 50$  under our simulation setting of  $n = 200$  or  $400$ ,  $p = 51$ , and  $\text{SNR} = 20$ . For each scenario, 50 replicates are generated and the proposed all-mode reduction model is trained, where the tuning parameters are determined as in Section 4.6 of the main paper. The results are shown in Figure S.1 as a line chart, with a bar denoted the standard errors of each scenario. Figure S.1 shows that the ISE has an rising trend as the number of responses increases, which is consistent with our main theorem.



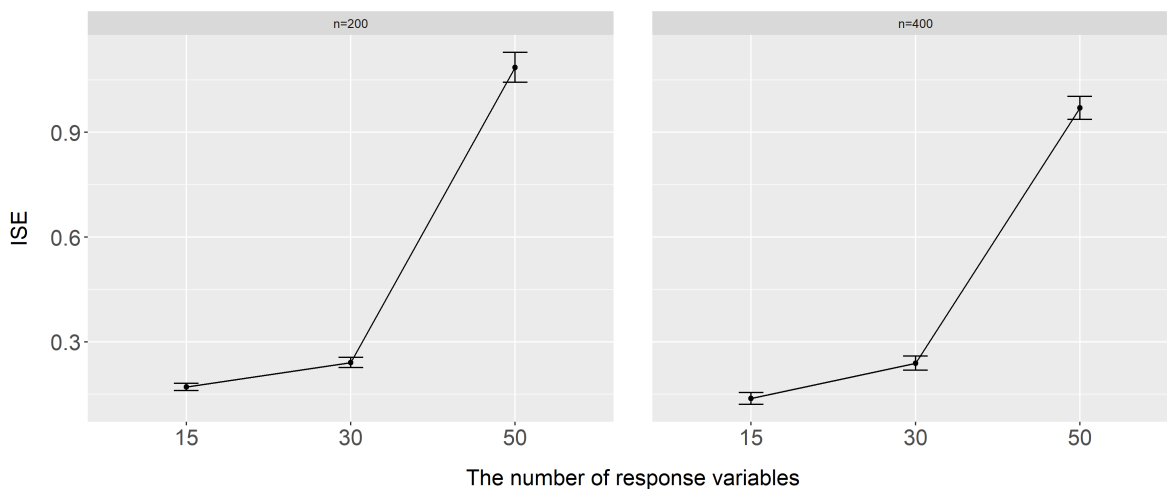


Figure S.1: The ISE for different numbers of the response variables.

### S.5.2 Additional Results of Real Data

In this subsection, we plot the estimated coefficient functions of one relevant SNP, rs9321440, identified by the proposed all-mode reduction method. Figure S.2 depicts the fitted coefficient functions and their average based on 50 replicates of random splitting. From the estimation of varying coefficient functions, we observe that the SNP rs9321440 may have different varying effects on the phenotypes of *height*, *bi-deltoid girth*, *right arm girth-upper third*, *hip girth*, and *thigh girth* given distinct body weights. To be specific, for both *height* and *bi-deltoid girth*, the estimated coefficients show negative patterns, and the effects decrease first but then increase with respect to the increase of body weights. This similarity in patterns could be explained by the high correlations between these two phenotypes (Chalmers et al., 2021), where both are mainly due to the skeleton of a human. For *right arm girth-upper third*, *hip girth*, and *thigh girth*, these measurements are characterized by the body fat and highly correlated (Freedman and Rimm, 1989). Their corresponding estimated coefficients have similar patterns and show positive effects. As for the phenotype of *waist girth*, the effect of this SNP may not vary with body weights significantly.

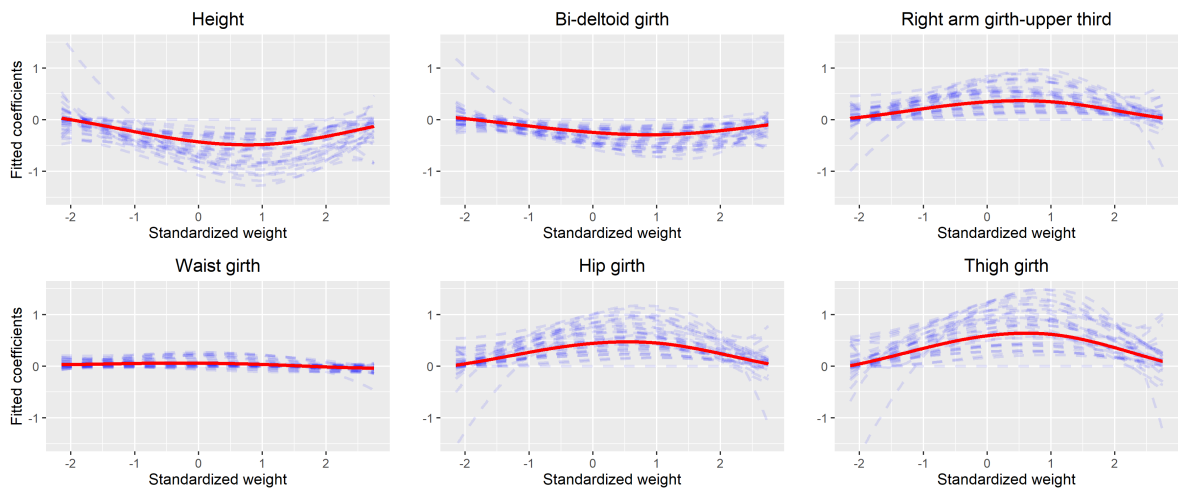


Figure S.2: Fitted coefficient functions of the biologically confirmed SNP rs9321440. In each panel, the blue dashed lines and the red solid line are the fitted functions and their average based on 50 replicates of random splitting, respectively.

## References

- Absil, P.-A., R. Mahony, and R. Sepulchre (2009). *Optimization algorithms on matrix manifolds*. Princeton, New Jersey: Princeton University Press.
- Banerjee, A., S. Chen, F. Fazayeli, and V. Sivakumar (2015). Estimation with norm regularization. *arXiv preprint arXiv:1505.02294*.
- Boyd, S., N. Parikh, E. Chu, B. Peleato, and J. Eckstein (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning* 3(1), 1–122.
- Bunea, F., Y. She, and M. H. Wegkamp (2012). Joint variable and rank selection for parsimonious estimation of high-dimensional matrices. *The Annals of Statistics* 40(5), 2359–2388.
- Chalmers, P. N., S. R. Lindsay, W. Smith, J. Kawakami, R. Hill, R. Z. Tashjian, and J. D. Keener (2021). Infraspinatus and deltoid length and patient height: implications for lateralization and distalization in reverse total shoulder arthroplasty. *Journal of Shoulder and Elbow Surgery* 30(4), 712–719.
- De Boor, C. (1973). The quasi-interpolant as a tool in elementary polynomial spline theory. *Approximation Theory*, 269–276.
- De Boor, C. (1976). Splines as linear combinations of b-splines. a survey. Technical report, Wisconsin Univ Madison Mathematics Research Center.
- Freedman, D. S. and A. A. Rimm (1989). The relation of body fat distribution, as assessed by six girth measurements, to diabetes mellitus in women. *American Journal of Public Health* 79(6), 715–720.
- Gabay, D. and B. Mercier (1976). A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications* 2(1), 17–40.
- Hunter, D. R. and R. Li (2005). Variable selection using mm algorithms. *The Annals of Statistics* 33(4), 1617–1642.

Kolda, T. G. and B. W. Bader (2009). Tensor decompositions and applications. *SIAM Review* 51(3), 455–500.

Koltchinskii, V. (2011). *Oracle inequalities in empirical risk minimization and sparse recovery problems*. New York: Springer Science & Business Media.

Talagrand, M. (2005). *The generic chaining*. Berlin: Springer.

Zhu, Y. (2017). An augmented admm algorithm with application to the generalized lasso problem. *Journal of Computational and Graphical Statistics* 26(1), 195–204.