

SUPPLEMENTARY MATERIAL OF
“ESTIMATION AND VARIABLE SELECTION UNDER THE
FUNCTION-ON-SCALAR LINEAR MODEL
WITH COVARIATE MEASUREMENT ERROR”

Yifan Sun¹, Grace Y. Yi^{1,2,*}

¹*Department of Statistical and Actuarial Sciences, University of Western Ontario,*

London, Ontario, Canada N6A 5B7

²*Department of Computer Science, University of Western Ontario,*

London, Ontario, Canada N6A 5B7

Regularity conditions, technical proofs, and simulation results are displayed in Section S1, Section S2 and Section S3, respectively. In Section S4 we provide additional notes on the real data analysis. In Section S5, we extend the development in the main text to accommodate the generalized least squares loss function, with both theoretical properties and numerical studies reported.

S1 Regularity Conditions

For a square matrix \mathbf{A} , let $\rho_{\min}(\mathbf{A})$ and $\rho_{\max}(\mathbf{A})$ denote the minimum and maximum eigenvalues of \mathbf{A} , respectively. The following regularity condi-

*Corresponding author: Grace Y. Yi, gyi5@uwo.ca

tions are imposed:

A. Assumptions about covariates and responses:

In addition to assuming $\{\{\mathbf{X}_i, Y_i(t)\} : i = 1, 2, \dots, n\}$ to be i.i.d. for $t \in \mathcal{T}$, as in Section 2.1, we require that

A1. The dimension of \mathbf{X}_i , p , is fixed.

A2. $\mathbb{E}(X_{ij}^4) < \infty$ for $j = 1, \dots, p$, and $\rho_{\min}\{\mathbb{E}(\mathbf{X}_i \mathbf{X}_i^\top)\} > 0$.

A3. \mathcal{T} is a compact interval, as required in Section 2.1.

A4. \mathbf{X}_i is independent of $\{\varepsilon_i(t) : t \in \mathcal{T}\}$, as required in Section 2.1.

B. Assumption about the coefficient functions in model (2.1):

There exist a positive integer $q' \leq d$ and a constant $C_1 > 0$, together with a constant $\nu > 0$, such that the q' th order derivative of functions

$\beta_j(\cdot)$ satisfies

$$|\beta_j^{(q')}(t_1) - \beta_j^{(q')}(t_2)| \leq C_1 |t_1 - t_2|^\nu,$$

for all $t_1, t_2 \in \mathcal{T}$ and all $1 \leq j \leq p$. Define $q = q' + \nu$.

C. Assumptions about B-spline approximations:

C1. The number of knots satisfies

(i) $\frac{M^{2q}}{n} > C_2$ for some positive constant C_2 ;

(ii) $\frac{M}{n} \rightarrow 0$ as $n \rightarrow \infty$.

C2. For $j = 1, \dots, M$, let $h_j = \tau_j - \tau_{j-1}$. Assume that

$$\max_{1 \leq j \leq M} |h_{j+1} - h_j| = o(M^{-1}) \quad \text{and} \quad \max_{1 \leq j \leq M} h_j \leq C_3 \cdot \min_{1 \leq j \leq M} h_j$$

for some positive constant C_3 .

C3. There exists a positive constant $C_4 \geq 1$ such that

$$C_4^{-1}m/M \leq \rho_{\min}(\Phi^\top(\mathbf{t})\Phi(\mathbf{t})) \leq \rho_{\max}(\Phi^\top(\mathbf{t})\Phi(\mathbf{t})) \leq C_4m/M. \tag{S1.1}$$

These requirements imply the invertibility of $\Phi^\top(\mathbf{t})\Phi(\mathbf{t})$, yielding that $m \geq M + d$ since $\Phi(\mathbf{t})$ is an $m \times (M + d)$ matrix.

D. Assumptions about measurement error \mathbf{U}_i in model (3.7):

D1. For any $1 \leq j \leq p$, assume that $\mathbb{E}(U_{ij}^4) < \infty$.

D2. For $i = 1, \dots, n$, \mathbf{U}_i is independent of $\{Y_i(t) : t \in \mathcal{T}\}$ and \mathbf{X}_i , as required in Section 3.1.

E. Assumptions about the penalty function $P_\lambda(v)$ for $v \geq 0$:

E1. $P_\lambda(0) = 0$ and $P_\lambda(v) > 0$ for $v > 0$.

E2. The first derivative of $P_\lambda(v)$, denoted $P'_\lambda(v)$, exists on $(0, +\infty)$, and satisfies

(i) $P'_\lambda(v)$ is continuous and monotonically non-increasing;

(ii) $P'_\lambda(v) = 0$ for $v > a\lambda$ for some constant $a > 1$;

(iii) $\lim_{v \rightarrow 0^+} P'_\lambda(v) = \lambda$.

Assumption A1 states the setting we consider. Assumption A2 is imposed to avoid settings with the perfect covariate collinearity, which results in nonidentifiability issues; this assumption is widely used in the literature of function-on-scalar regression models (Wang et al., 2007; Parodi and Reimherr, 2018). Indeed, if $\mathbb{E}(\mathbf{X}_i \mathbf{X}_i^\top)$ is singular, there exists a non-zero vector $\boldsymbol{\gamma} \in \mathbb{R}^p$ such that $\boldsymbol{\gamma}^\top \mathbf{X}_i = 0$ almost surely. Then any function of the form $\boldsymbol{\beta}(t) + c \mathbf{f}_\boldsymbol{\gamma}(t)$ satisfies model (2.1), where $c \in \mathbb{R}$ is an arbitrary constant and $\mathbf{f}_\boldsymbol{\gamma}(t) \equiv \boldsymbol{\gamma}$.

Assumption B is commonly made for the function-on-scalar regression model, which allows the use of B-spline based methods to estimate the coefficient functions in the model. In fact, such $\beta_j(t)$ is smooth enough and can be approximated by some “best” function in the linear space spanned by all B-spline basis functions $\boldsymbol{\phi}(t)$, denoted $S_{d,M}$. To be specific, together with Condition C2, there always exist B-spline functions $\{\theta_j(t) \in S_{d,M} : j = 1, \dots, p, t \in \mathcal{T}\}$, such that

$$\sup_{1 \leq j \leq p} \|\theta_j - \beta_j\|_{L_\infty} \leq C'_1 M^{-q}, \quad (\text{S1.2})$$

for any positive integer M , where C'_1 is a positive constant functionally independent of M (De Boor, 1978, Theorem XII(6)), and thus yielding that $\|\theta_j - \beta_j\|_{L_\infty} \rightarrow 0$ as $M \rightarrow \infty$. For a B-spline function $\theta_j(t) \in S_{d,M}$

satisfying (S1.2), we write

$$\theta_j(t) = \boldsymbol{\phi}^\top(t) \mathbf{b}_{0,j}, \quad (\text{S1.3})$$

where $\mathbf{b}_{0,j} \in \mathbb{R}^{M+d}$. Let $\mathbf{b}_0 = (\mathbf{b}_{0,1}^\top, \dots, \mathbf{b}_{0,p}^\top)^\top$ denote the $(M+d)p \times 1$ vector. In order to establish the estimation consistency of $\hat{\beta}_j$, we therefore examine $\|\theta_j - \hat{\beta}_j\|_{L_\infty}$, and it suffices to find an upper bound of $\|\theta_j - \hat{\beta}_j\|_{L_\infty}$ which converges to 0 as $n \rightarrow \infty$, as discussed in the following derivation of Theorem 1.

Generally speaking, the number of B-spline inner knots should be divergent to derive the function consistency, as shown in (S1.2). However, a too large M results in an estimator with an overly large variance. Imposing condition C1 is to constrain the divergence rate of M . The imposition of Condition C(i) is to control the B-spline approximation error, which makes the terms involving the B-spline error (i.e., I_{22} and I_{25} in (S2.28) of the proof of Theorem 1) negligible, relative to other terms. Without this condition, the convergence rate for $\hat{\mathbf{b}}$ in Theorem 1 is not $O_p(\sqrt{M/n})$ but becomes $O_p(\sqrt{M/n} + M^{1/2-q})$. Condition C2 is identical to Assumption 3 of Zhou et al. (1998), which is easy to be met. For instance, the equally spaced knots satisfies Condition C2.

While constraint (S1.1) in Condition C3 may look somewhat stringent, they are readily satisfied for many settings. In fact, to satisfy (S1.1), it

suffices to require $\mathbf{a}^\top \{\Phi^\top(\mathbf{t})\Phi(\mathbf{t})/m\}\mathbf{a}$ to be of order $O(M^{-1})$ for any vector \mathbf{a} with $\|\mathbf{a}\| = 1$. To this end, let $s(t) = \mathbf{a}^\top \phi(t)$ for $t \in \mathcal{T}$. Then we write

$$\mathbf{a}^\top \{\Phi^\top(\mathbf{t})\Phi(\mathbf{t})/m\}\mathbf{a} = \int_{\mathcal{T}} s(t)^2 dF_m(t), \quad (\text{S1.4})$$

where $F_m(t) \triangleq \frac{1}{m} \sum_{k=1}^m \mathbb{I}(t_k \leq t)$ is the step function of observations $\mathbf{t} = \{t_1, \dots, t_m\}$. Suppose that there is a distribution function $F(t)$ with positive continuous density on the domain \mathcal{T} , which satisfies

$$\|F_m - F\|_{L_\infty} = o(M^{-1}). \quad (\text{S1.5})$$

On one hand, it can be proved that $\int_{\mathcal{T}} s(t)^2 dF(t) = O(M^{-1})$ (De Boor, 1978, Page 155). On the other hand, we have that $\int_{\mathcal{T}} s(t)^2 d(F_m - F)(t) = o(M^{-1})$ due to (S1.5). Therefore, the right-hand-side of equation (S1.4) is of order $O(M^{-1})$. The details of this proof can be found in Zhou et al. (1998, Lemma 6.1).

Assumption D1 is useful for technical derivations, and it is commonly made in the literature of measurement error models in a tacit manner, for example, by assuming the normality of \mathbf{U}_i in model (3.7) (e.g., Section 2.6 of Yi (2017)). Assumptions in E are widely used in the literature of variable selection. The well-known SCAD (Fan and Li, 2001) and MCP (Zhang, 2010) are two examples satisfying them.

Remark 1. Condition (S1.5) is easy to be met. For example, if the do-

main $\mathcal{T} \triangleq [0, 1]$ and the observations are equally spaced at $t_k = k/m$ for $k = 1, \dots, m$, then we take F as the uniform distribution function on $[0, 1]$, which satisfies $\|F_m - F\|_{L_\infty} = m^{-1}$. Hence (S1.5) holds under the assumption that $M/m \rightarrow 0$ as $m \rightarrow \infty$.

Remark 2. The results in Section 5 of the main text are established by treating the observation points $\mathbf{t} = \{t_1, \dots, t_m\}$ as deterministic or pre-specified. When the observation times are random variables, the results can still hold if proper modifications of the conditions are done; typically, the inequalities (S1.1) is now modified to hold with probability tending to 1. This new constraint can also be proved if $\|F_m - F\|_{L_\infty} = o_p(M^{-1})$, where the definitions of F_m and F are the same as before. In a special case where $\{t_1, \dots, t_m\}$ are independently sampled from a common distribution function, say, $F_0(t)$, then we have that $\|F_m - F_0\|_{L_\infty} = O_p(m^{-1/2})$ according to Massart (1990). Therefore, the condition $\|F_m - F\|_{L_\infty} = o_p(M^{-1})$ is satisfied if we set $F = F_0$ and assume $M^2/m \rightarrow 0$ as $m \rightarrow \infty$.

S2 Technical Details

S2.1 Preliminary Preparations

For a vector $\mathbf{a} = (a_1, \dots, a_k)^\top$, let $\|\mathbf{a}\|_\infty$ and $\|\mathbf{a}\|$ respectively denote the infinity norm and the Euclidean norm for \mathbf{a} , i.e., $\|\mathbf{a}\|_\infty = \max_{1 \leq j \leq k} |a_j|$ and $\|\mathbf{a}\| = \sqrt{\sum_{j=1}^k a_j^2}$. For a $p \times q$ matrix $\mathbf{A} = [a_{jk}]_{p \times q}$ with the (j, k) element $a_{jk} \in \mathbb{R}$, let $\|\mathbf{A}\|_F$ and $\|\mathbf{A}\|_2$ respectively represent the Frobenius norm and the matrix 2-norm (i.e., the spectral norm) for \mathbf{A} , i.e., $\|\mathbf{A}\|_F = \sqrt{\sum_{j=1}^p \sum_{k=1}^q a_{jk}^2}$ and $\|\mathbf{A}\|_2 = \sqrt{\rho_{\max}(\mathbf{A}^\top \mathbf{A})}$. An equivalent definition for the matrix 2-norm is given by

$$\|\mathbf{A}\|_2 = \sup_{\|\mathbf{w}\|=1} \|\mathbf{A}\mathbf{w}\|, \quad (\text{S2.6})$$

where \mathbf{w} is a vector with suitable dimension.

In the subsequent derivations, we frequently use the following properties, together with Lemmas 1 and 2.

Property 1: For matrices \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} having suitable dimensions,

$$(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top \text{ and } (\mathbf{A} \otimes \mathbf{C})(\mathbf{B} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}).$$

Property 2: For any matrices \mathbf{A} , \mathbf{B} and \mathbf{C} with suitable dimensions,

$$(\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{C}) = \text{vec}(\mathbf{ACB}).$$

Property 3: For any matrices \mathbf{A} and \mathbf{B} ,

$$\|\mathbf{A} \otimes \mathbf{B}\|_2 = \|\mathbf{A}\|_2 \|\mathbf{B}\|_2.$$

Property 4: For any matrices \mathbf{A} and \mathbf{B} with suitable dimensions, we have that

$$(a) \quad \|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_F \text{ and } \|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_2.$$

$$(b) \quad \|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F.$$

Property 5: For any symmetric positive semi-definite matrix \mathbf{G} ,

$$\rho_{\max}(\mathbf{G}^2) = \{\rho_{\max}(\mathbf{G})\}^2.$$

Property 6: For positive semi-definite matrices \mathbf{S} and \mathbf{T} ,

$$\rho_{\min}(\mathbf{S} \otimes \mathbf{T}) = \rho_{\min}(\mathbf{S})\rho_{\min}(\mathbf{T}).$$

Properties 1 and 2 can be found in Horn and Johnson (1991, Section 4.2) or proved by using the definition of the kronecker product. Properties 3 and 6 are directly resulted from Theorem 4.2.12 of Horn and Johnson (1991) and Properties 4 and 5 can be found in Horn and Johnson (2012, Section 5.6). We also present rigorous proofs for Properties 3-6 below.

Proof of Property 3, 4, 5 and 6:

For any rectangular matrices \mathbf{S}_1 and \mathbf{T}_1 , let the nonzero singular values of \mathbf{S}_1 and \mathbf{T}_1 be $\{\lambda_1, \dots, \lambda_s\}$ and $\{\mu_1, \dots, \mu_t\}$, respectively. According to Theorem 4.2.12 of Horn and Johnson (1991), all nonzero singular values for $\mathbf{S}_1 \otimes \mathbf{T}_1$ is given by $\{\lambda_j \mu_k : 1 \leq j \leq s, 1 \leq k \leq t\}$. Hence for any matrices \mathbf{A} and \mathbf{B} , we have that

$$\begin{aligned}\|\mathbf{A} \otimes \mathbf{B}\|_2^2 &= \rho_{\max}\{(\mathbf{A}^\top \mathbf{A}) \otimes (\mathbf{B}^\top \mathbf{B})\} \\ &= \rho_{\max}(\mathbf{A}^\top \mathbf{A}) \cdot \rho_{\max}(\mathbf{B}^\top \mathbf{B}) \\ &= \|\mathbf{A}\|_2^2 \|\mathbf{B}\|_2^2,\end{aligned}$$

where the first step is due to Property 1 and the definition of spectral norm, and the second step is because $\mathbf{A}^\top \mathbf{A}$ and $\mathbf{B}^\top \mathbf{B}$ are positive semi-definite, indicating that all eigenvalues are non-negative. Property 3 is then proved.

Similarly, for two positive semi-definite matrices \mathbf{S} and \mathbf{T} ,

$$\rho_{\min}(\mathbf{S} \otimes \mathbf{T}) = \rho_{\min}(\mathbf{S})\rho_{\min}(\mathbf{T}),$$

since their eigenvalues are all non-negative. Property 6 is proved.

Note that the Frobenius norm for matrix \mathbf{A} can be rewrite as $\|\mathbf{A}\|_F = \sqrt{\text{tr}(\mathbf{A}^\top \mathbf{A})}$, where $\text{tr}(\cdot)$ represents the trace operator. Property 4 (b) then directly follows since

$$\|\mathbf{A}\|_F^2 = \text{tr}(\mathbf{A}^\top \mathbf{A}) \geq \rho_{\max}(\mathbf{A}^\top \mathbf{A}) = \|\mathbf{A}\|_2^2.$$

For Property 4 (a), let $\mathbf{B} = [\mathbf{B}_1, \dots, \mathbf{B}_p]$, where \mathbf{B}_j for $j = 1, \dots, p$ are

column vectors of \mathbf{A} . Then we have that

$$\begin{aligned}
\|\mathbf{AB}\|_F^2 &= \text{tr}(\mathbf{B}^\top \mathbf{A}^\top \mathbf{AB}) \\
&= \sum_{j=1}^p \mathbf{B}_j^\top \mathbf{A}^\top \mathbf{AB}_j \\
&\leq \rho_{\max}(\mathbf{A}^\top \mathbf{A}) \sum_{j=1}^p \mathbf{B}_j^\top \mathbf{B}_j \\
&= \|\mathbf{A}\|_2^2 \|\mathbf{B}\|_F^2.
\end{aligned}$$

By the facts that $\|\mathbf{AB}\|_F$ also equals to $\text{tr}(\mathbf{ABB}^\top \mathbf{A}^\top)$ and $\rho_{\max}(\mathbf{BB}^\top) = \|\mathbf{B}\|_2^2$, we can prove the other inequality analogously.

Now we verify Property 5. Since \mathbf{G} is symmetric, there exists a real orthogonal matrix \mathbf{P} such that

$$\mathbf{G} = \mathbf{P}^\top \mathbf{\Lambda} \mathbf{P},$$

where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m)$ is the diagonal matrix with λ_j , $1 \leq j \leq m$, representing the j th eigenvalue of \mathbf{G} . Without loss of generality, we assume that $\lambda_1 = \rho_{\max}(\mathbf{G})$. It is then obvious that

$$\rho_{\max}(\mathbf{G}^2) = \rho_{\max}(\mathbf{P}^\top \mathbf{\Lambda}^2 \mathbf{P}) = \lambda_1^2 = \{\rho_{\max}(\mathbf{G})\}^2,$$

where the first step is due to that \mathbf{PP}^\top equals identity and the second step is because that all $\lambda_j \geq 0$, $j = 1, \dots, m$. \square

Lemma 1. *Let \mathbf{X} and \mathbf{Z} be two independent column vector variables with mean $\mathbf{0}$ and covariance matrices $\mathbf{\Sigma}_X$ and $\mathbf{\Sigma}_Z$ respectively. For any real-*

valued matrix \mathbf{C} with a suitable dimension, we have that

$$\text{Cov}\{(\mathbf{X} \otimes \mathbf{C})\mathbf{Z}\} = \Sigma_{\mathbf{X}} \otimes (\mathbf{C}\Sigma_{\mathbf{Z}}\mathbf{C}^{\top}).$$

Proof:

Let $\mathbf{X} = (X_1, \dots, X_p)^{\top}$ and let $\mathbf{W}_j = X_j\mathbf{C}\mathbf{Z}$ for $j = 1, \dots, p$. Then $(\mathbf{X} \otimes \mathbf{C})\mathbf{Z} = [\mathbf{W}_1^{\top}, \dots, \mathbf{W}_p^{\top}]^{\top}$. It is direct to verify that

$$\text{Var}(\mathbf{W}_j) = \mathbb{E}(X_j^2)\mathbf{C}\Sigma_{\mathbf{Z}}\mathbf{C}^{\top}$$

and

$$\text{Cov}(\mathbf{W}_j, \mathbf{W}_l) = \text{Cov}(X_j, X_l)\mathbf{C}\Sigma_{\mathbf{Z}}\mathbf{C}^{\top}.$$

It is then straightforward to combine these two results and obtain that

$$\text{Cov}([\mathbf{W}_1^{\top}, \dots, \mathbf{W}_p^{\top}]^{\top}) = \Sigma_{\mathbf{X}} \otimes (\mathbf{C}\Sigma_{\mathbf{Z}}\mathbf{C}^{\top}).$$

□

Lemma 2. *Let Σ be an $m \times m$ real-valued symmetric positive definite matrix and let Ψ be an $m \times d$ real-valued matrix with the full rank, where $d \leq m$. Then*

$$\rho_{\min}(\Psi^{\top}\Sigma\Psi) \geq \rho_{\min}(\Psi^{\top}\Psi)\rho_{\min}(\Sigma), \quad (\text{S2.7})$$

and

$$\rho_{\max}(\Psi^{\top}\Sigma\Psi) \leq \|\Psi^{\top}\Psi\|_2\rho_{\max}(\Sigma). \quad (\text{S2.8})$$

Proof:

Since Σ is real-valued and symmetric, there exists a real-valued orthogonal matrix \mathbf{P} such that

$$\Sigma = \mathbf{P}^\top \Lambda \mathbf{P},$$

where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ is the diagonal matrix with λ_j representing an eigenvalue of Σ for $1 \leq j \leq m$.

Define $\tilde{\Psi} = \mathbf{P}\Psi$ and write $\tilde{\Psi} = [\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_m]^\top$. Since $\mathbf{P}^\top \mathbf{P}$ equals the identity matrix, proving (S2.7) is equivalent to showing that

$$\rho_{\min}(\tilde{\Psi}^\top \Lambda \tilde{\Psi}) \geq \rho_{\min}(\tilde{\Psi}^\top \tilde{\Psi}) \rho_{\min}(\Lambda), \quad (\text{S2.9})$$

which is true because that

$$\begin{aligned} \rho_{\min}(\tilde{\Psi}^\top \Lambda \tilde{\Psi}) &= \rho_{\min} \left(\sum_{j=1}^m \lambda_j \boldsymbol{\psi}_j \boldsymbol{\psi}_j^\top \right) \\ &= \inf_{\|\mathbf{w}\|=1} \sum_{j=1}^m \lambda_j (\mathbf{w}^\top \boldsymbol{\psi}_j)^2 \\ &\geq \inf_{\|\mathbf{w}\|=1} \sum_{j=1}^m (\mathbf{w}^\top \boldsymbol{\psi}_j)^2 \rho_{\min}(\Lambda) \\ &= \rho_{\min}(\Lambda) \cdot \rho_{\min} \left(\sum_{j=1}^m \boldsymbol{\psi}_j \boldsymbol{\psi}_j^\top \right) \\ &= \rho_{\min}(\tilde{\Psi}^\top \tilde{\Psi}) \rho_{\min}(\Lambda), \end{aligned}$$

where the second and the fourth steps are due to the fact that $\rho_{\min}(\mathbf{Q}) = \inf_{\|\mathbf{w}\|=1} \mathbf{w}^\top \mathbf{Q} \mathbf{w}$ for any square matrix \mathbf{Q} . The proof of (S2.7) is then completed.

Expression (S2.8) can be proved using arguments analogous to those for (S2.7) and the identity $\rho_{\max}(\Psi^\top \Psi) = \|\Psi^\top \Psi\|_2$. \square

Lemma 3. *Consider the B-spline basis functions $\phi(t)$ in Section 2.2 and the linear space $S_{d,M}$ spanned by them given before (S1.2). For $t \in \mathcal{T}$, let $f(t) = \phi^\top(t)\mathbf{b} \in S_{d,M}$ denote a spline function over domain \mathcal{T} with coefficients $\mathbf{b} \in \mathbb{R}^{M+d}$. Assume Assumption C2. Then there exists positive constants E_1 and E_2 depending only on d such that*

$$E_1\sqrt{M}\|f\|_{L_2} \leq \|\mathbf{b}\| \leq E_2\sqrt{M}\|f\|_{L_2}.$$

Proof:

In Section 2.2, we use $\tau_1 < \tau_2 < \dots < \tau_{M+1}$ to denote knots on \mathcal{T} for the setup of B-spline, where τ_1 and τ_{M+1} represent the two endpoints of \mathcal{T} . Define $\tau_k = \tau_1$ for all $k = 1 - d, 1 - d + 1, \dots, 0$ and $\tau_k = \tau_{M+1}$ for all $k = M + 2, \dots, M + 2d + 1$. Let $\mathbf{b} = (b_1, \dots, b_{M+d})^\top$. According to equation (13) in Zhou et al. (1998), there exist positive constants \tilde{E}_1 and \tilde{E}_2 depending only on d , such that

$$\tilde{E}_1 \sum_{j=1}^{M+d} (\tau_{j+1} - \tau_{j-d}) b_j^2 \leq \|f\|_{L_2}^2 \leq \tilde{E}_2 \sum_{j=1}^{M+d} (\tau_{j+1} - \tau_{j-d}) b_j^2. \quad (\text{S2.10})$$

Due to Assumption C2 and the fact that $\min_{1 \leq j \leq M} h_j \leq 1/M$,

$$\max_{1 \leq j \leq M} h_j \leq C_3 \min_{1 \leq j \leq M} h_j \leq C_3/M.$$

Similarly, since $\max_{1 \leq j \leq M} h_j \geq 1/M$,

$$\min_{1 \leq j \leq M} h_j \geq C_3^{-1} \max_{1 \leq j \leq M} h_j \geq 1/(C_3 M).$$

Hence, we have that for any $j = 1, \dots, M$,

$$1/(C_3 M) \leq \min_{1 \leq j \leq M} h_j \leq \tau_{j+1} - \tau_{j-d} \leq (d+1) \max_{1 \leq j \leq M} h_j \leq C_3(d+1)/M. \quad (\text{S2.11})$$

Then by (S2.10) and (S2.11),

$$\frac{\tilde{E}_1}{C_3 M} \|\mathbf{b}\|^2 \leq \|f\|_{L_2}^2 \leq \frac{C_3(d+1)\tilde{E}_2}{M} \|\mathbf{b}\|^2,$$

which completes the proof. \square

S2.2 Proof of Theorem 1

The proof consists of two parts, each for proving (5.23) or (5.24).

Part 1: Prove (5.23).

Set $\delta_n = \sqrt{\frac{M}{n}}$. To show the existence of a local minimizer, say $\hat{\mathbf{b}}$, of $Q_n(\mathbf{b})$ satisfying $\|\hat{\mathbf{b}} - \mathbf{b}_0\| = O_p(\delta_n)$, it suffices to verify that, for any $\epsilon > 0$, there exists a positive constant K_ϵ such that

$$\mathbb{P} \left(\inf_{\|\mathbf{v}\|=K_\epsilon, \mathbf{v} \in \mathbb{R}^{p(M+d)}} Q_n(\mathbf{b}_0 + \delta_n \mathbf{v}) > Q_n(\mathbf{b}_0) \right) > 1 - \epsilon. \quad (\text{S2.12})$$

The proof of (S2.12) consists of the following steps, where for any given constant $K > 0$, we consider any \mathbf{v} satisfying $\|\mathbf{v}\| = K$.

Step 1: Identify a lower bound of $Q_n(\mathbf{b}_0 + \delta_n \mathbf{v}) - Q_n(\mathbf{b}_0)$.

For $i = 1, \dots, n$, let

$$\mathbf{Z}_i = \mathbf{X}_i^\top \otimes \Phi(\mathbf{t}) \text{ and } \mathbf{Z}_i^* = \mathbf{X}_i^{*\top} \otimes \Phi(\mathbf{t}). \quad (\text{S2.13})$$

By (3.11),

$$Q_n(\mathbf{b}) = \frac{1}{2} \sum_{i=1}^n \|\mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i^* \mathbf{b}\|^2 - \frac{n}{2} \mathbf{b}^\top [\Sigma \otimes \{\Phi^\top(\mathbf{t})\Phi(\mathbf{t})\}] \mathbf{b} + nm \sum_{j=1}^p P_\lambda(\|\mathbf{b}_j\|).$$

Therefore,

$$\begin{aligned}
Q_n(\mathbf{b}_0 + \delta_n \mathbf{v}) - Q_n(\mathbf{b}_0) &= \frac{1}{2} \sum_{i=1}^n \{ \|\mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i^*(\mathbf{b}_0 + \delta_n \mathbf{v})\|^2 - \|\mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i^* \mathbf{b}_0\|^2 \} \\
&\quad + nm \sum_{j=1}^p \{ P_\lambda(\|\mathbf{b}_{0,j} + \delta_n \mathbf{v}_j\|) - P_\lambda(\|\mathbf{b}_{0,j}\|) \} \\
&\quad - \frac{n}{2} (\mathbf{b}_0 + \delta_n \mathbf{v})^\top [\boldsymbol{\Sigma} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \}] (\mathbf{b}_0 + \delta_n \mathbf{v}) \\
&\quad + \frac{n}{2} \mathbf{b}_0^\top [\boldsymbol{\Sigma} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \}] \mathbf{b}_0 \\
&= \frac{\delta_n^2}{2} \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^{*\top} \mathbf{Z}_i^* \mathbf{v} - \delta_n \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^{*\top} (\mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i^* \mathbf{b}_0) \\
&\quad + nm \sum_{j=1}^p \{ P_\lambda(\|\mathbf{b}_{0,j} + \delta_n \mathbf{v}_j\|) - P_\lambda(\|\mathbf{b}_{0,j}\|) \} \\
&\quad - n \delta_n \mathbf{b}_0^\top [\boldsymbol{\Sigma} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \}] \mathbf{v} - \frac{n \delta_n^2}{2} \mathbf{v}^\top [\boldsymbol{\Sigma} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \}] \mathbf{v} \\
&\geq \frac{\delta_n^2}{2} \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^{*\top} \mathbf{Z}_i^* \mathbf{v} - \delta_n \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^{*\top} (\mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i^* \mathbf{b}_0) \\
&\quad + nm \sum_{j \notin J_0} \{ P_\lambda(\|\mathbf{b}_{0,j} + \delta_n \mathbf{v}_j\|) - P_\lambda(\|\mathbf{b}_{0,j}\|) \} \\
&\quad - n \delta_n \mathbf{b}_0^\top [\boldsymbol{\Sigma} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \}] \mathbf{v} - \frac{n \delta_n^2}{2} \mathbf{v}^\top [\boldsymbol{\Sigma} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \}] \mathbf{v} \\
&\triangleq I_1 + I_2 + I_3 + I_4 + I_5,
\end{aligned} \tag{S2.14}$$

where \mathbf{v}_j is the j th part of \mathbf{v} such that $\mathbf{v} = (\mathbf{v}_1^\top, \dots, \mathbf{v}_p^\top)^\top$. The inequality

is due to Assumption E1 and the fact that $\mathbf{b}_{0,j} = \mathbf{0}$ for all $j \in J_0$.

Step 2: Examine I_1 and I_5 in (S2.14).

By definition and model (3.7),

$$\mathbf{Z}_i^* = \mathbf{Z}_i + \mathbf{U}_i^\top \otimes \Phi(\mathbf{t}).$$

Hence the term I_1 in (S2.14) can be further split as

$$\begin{aligned} I_1 &= \frac{\delta_n^2}{2} \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^{*\top} \mathbf{Z}_i^* \mathbf{v} \\ &= \frac{\delta_n^2}{2} \sum_{i=1}^n \mathbf{v}^\top \mathbb{E}(\mathbf{Z}_i^\top \mathbf{Z}_i) \mathbf{v} + \delta_n^2 \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^\top \{\mathbf{U}_i^\top \otimes \Phi(\mathbf{t})\} \mathbf{v} \\ &\quad + \frac{\delta_n^2}{2} \sum_{i=1}^n \mathbf{v}^\top [(\mathbf{U}_i \mathbf{U}_i^\top) \otimes \{\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\}] \mathbf{v} \\ &\quad + \frac{\delta_n^2}{2} \sum_{i=1}^n \mathbf{v}^\top \{\mathbf{Z}_i^\top \mathbf{Z}_i - \mathbb{E}(\mathbf{Z}_i^\top \mathbf{Z}_i)\} \mathbf{v} \\ &\triangleq I_{1,1} + I_{1,2} + I_{1,3} + I_{1,4}. \end{aligned} \tag{S2.15}$$

Now we examine (S2.15) term by term.

For $I_{1,1}$ in (S2.15), we have that

$$\begin{aligned} I_{1,1} &= \frac{\delta_n^2}{2} \sum_{i=1}^n \mathbf{v}^\top \mathbb{E}(\mathbf{Z}_i^\top \mathbf{Z}_i) \mathbf{v} \\ &= \frac{n\delta_n^2}{2} \mathbf{v}^\top [\mathbb{E}(\mathbf{X}_i \mathbf{X}_i^\top) \otimes \{\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\}] \mathbf{v} \\ &\geq \frac{n\delta_n^2}{2} \|\mathbf{v}\|^2 \frac{m}{M} \rho_{\min} \{\mathbb{E}(\mathbf{X}_i \mathbf{X}_i^\top)\} \rho_{\min} \left\{ \frac{M}{m} \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \right\} \\ &\geq \delta_n^2 K^2 \frac{nm}{M} C_0 \end{aligned} \tag{S2.16}$$

for some positive constant C_0 , where the first equality is due to Property 1, the third step is due to Property 6 and the last inequality is due to Assumption A2 and Assumption C3 as well as $\|\mathbf{v}\| = K$.

For $I_{1,2}$ in (S2.15), we have that

$$\begin{aligned}
|I_{1,2}| &= \delta_n^2 \left| \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^\top \{ \mathbf{U}_i^\top \otimes \Phi(\mathbf{t}) \} \mathbf{v} \right| \\
&= \delta_n^2 \left| \mathbf{v}^\top \left[\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{v} \right| \\
&\leq \delta_n^2 \|\mathbf{v}\| \left\| \left[\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{v} \right\| \\
&\leq \delta_n^2 \|\mathbf{v}\|^2 \sup_{\|\mathbf{w}\|=1} \left\| \left[\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{w} \right\| \\
&= \delta_n^2 \|\mathbf{v}\|^2 \left\| \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right\|_2 \\
&= \delta_n^2 \|\mathbf{v}\|^2 \left\| \sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right\|_2 \|\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\|_2, \tag{S2.17}
\end{aligned}$$

where the second equality is due to (S2.13) and Property 1, the third step is due to the Cauchy-Schwarz inequality, the fifth step is due to the definition (S2.6), and the last equality is due to Property 3.

Note that

$$\left\| \sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right\|_2 \leq \left\| \sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right\|_F = \sqrt{\sum_{j,k=1}^p \left(\sum_{i=1}^n X_{ij} U_{ik} \right)^2} = O_p(\sqrt{n}) \tag{S2.18}$$

due to the central limit theorem, and by Assumption C3, we further know

that

$$\begin{aligned}
 \|\Phi^\top(\mathbf{t})\Phi(\mathbf{t})\|_2 &= \sqrt{\rho_{\max}\{\Phi^\top(\mathbf{t})\Phi(\mathbf{t})\Phi^\top(\mathbf{t})\Phi(\mathbf{t})\}} \\
 &= \rho_{\max}(\Phi^\top(\mathbf{t})\Phi(\mathbf{t})) \\
 &\leq C_4 m/M,
 \end{aligned} \tag{S2.19}$$

where the first step is due to the definition of the matrix 2-norm, and the second step is due to Property 5. Therefore, (S2.17) is bounded:

$$|I_{1,2}| = O_p(\delta_n^2 \sqrt{nm}/M) \cdot K^2 \tag{S2.20}$$

since $\|\mathbf{v}\| = K$.

Next, combining $I_{1,3}$ in (S2.15) and I_5 in (S2.14) gives that

$$I_{1,3} + I_5 = \frac{1}{2} \delta_n^2 \mathbf{v}^\top \left[\left\{ \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\} \otimes \{ \Phi^\top(\mathbf{t})\Phi(\mathbf{t}) \} \right] \mathbf{v}. \tag{S2.21}$$

For any $1 \leq j, k \leq p$,

$$\text{Var} \left\{ \sum_{i=1}^n (U_{ij} U_{ik} - \Sigma_{jk}) \right\} = n \text{Var}(U_{ij} U_{ik}) = O(n), \tag{S2.22}$$

due to Assumption D1, where Σ_{jk} is the (j, k) element of Σ . Therefore, by

Property 4 (b) and the definition of the Frobenius norm,

$$\begin{aligned}
 \left\| \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\|_2^2 &\leq \left\| \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\|_F^2 \\
 &= \sum_{j,k=1}^p \left\{ \sum_{i=1}^n (U_{ij} U_{ik} - \Sigma_{jk}) \right\}^2 \\
 &= O_p(n),
 \end{aligned}$$

where the last step comes from the Markov inequality and (S2.22), so

$$\left\| \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \boldsymbol{\Sigma}) \right\|_2 = O_p(\sqrt{n}). \quad (\text{S2.23})$$

Applying (S2.23) and (S2.19) to (S2.21) gives that

$$\begin{aligned} |I_{1,3} + I_5| &\leq \frac{1}{2} \delta_n^2 \|\mathbf{v}\|^2 \left\| \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \boldsymbol{\Sigma}) \right\|_2 \|\boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t})\|_2 \\ &= O_p(\delta_n^2 \sqrt{nm}/M) \cdot K^2, \end{aligned} \quad (\text{S2.24})$$

where the inequality is due to Property 3.

By (S2.13), (S2.19) and Properties 3 and 4,

$$\begin{aligned} |I_{1,4}| &= \frac{\delta_n^2}{2} \left| \sum_{i=1}^n \mathbf{v}^\top \{ \mathbf{Z}_i^\top \mathbf{Z}_i - \mathbb{E}(\mathbf{Z}_i^\top \mathbf{Z}_i) \} \mathbf{v} \right| \\ &= \frac{\delta_n^2}{2} \left| \mathbf{v}^\top \left[\left\{ \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top - n \mathbb{E}(\mathbf{X}_i \mathbf{X}_i^\top) \right\} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \} \right] \mathbf{v} \right| \\ &\leq \frac{n \delta_n^2}{2} \|\mathbf{v}\|^2 \left\| \frac{1}{n} \sum_{i=1}^n \{ \mathbf{X}_i \mathbf{X}_i^\top - \mathbb{E}(\mathbf{X}_i \mathbf{X}_i^\top) \} \right\|_F \|\boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t})\|_2 \\ &= O_p(\delta_n^2 \sqrt{nm}/M) \cdot K^2, \end{aligned} \quad (\text{S2.25})$$

where the last step comes from Assumption A2, the central limit theorem to bound the Frobenius norm term, and (S2.19).

Therefore, combining (S2.20), (S2.24) and (S2.25) gives that

$$|I_{1,2} + I_{1,3} + I_{1,4} + I_5| = O_p(\delta_n^2 \sqrt{nm}/M) \cdot K^2. \quad (\text{S2.26})$$

Step 3: Examine I_2 and I_4 in (S2.14).

To describe the difference between (2.3) and (S1.3), let $\Delta_j(t) \triangleq \beta_j(t) - \theta_j(t)$, which equals $\beta_j(t) - \mathbf{b}_{0,j}^\top \boldsymbol{\phi}(t)$, representing the B-spline approximation

error for $\beta_j(t)$. Then we write $\beta_j(t) = \mathbf{b}_{0,j}^\top \boldsymbol{\phi}(t) + \Delta_j(t)$. Inserting this into model (2.2) gives that

$$\begin{aligned}
 \mathbf{Y}_i(\mathbf{t}) &= \sum_{j=1}^p X_{ij} (\beta_j(t_1), \dots, \beta_j(t_m))^\top + \boldsymbol{\varepsilon}_i(\mathbf{t}) \\
 &= \sum_{j=1}^p X_{ij} \{ \boldsymbol{\Phi}(\mathbf{t}) \mathbf{b}_{0,j} + \Delta_j(\mathbf{t}) \} + \boldsymbol{\varepsilon}_i(\mathbf{t}) \\
 &= \boldsymbol{\Phi}(\mathbf{t}) [\mathbf{b}_{0,1}, \dots, \mathbf{b}_{0,p}] \mathbf{X}_i + \boldsymbol{\varepsilon}_i(\mathbf{t}) + \sum_{j=1}^p X_{ij} \Delta_j(\mathbf{t}) \\
 &= \{ \mathbf{X}_i^\top \otimes \boldsymbol{\Phi}(\mathbf{t}) \} \mathbf{b}_0 + \boldsymbol{\varepsilon}_i(\mathbf{t}) + \sum_{j=1}^p X_{ij} \Delta_j(\mathbf{t}), \tag{S2.27}
 \end{aligned}$$

where $\boldsymbol{\Delta}_j(\mathbf{t}) = (\Delta_j(t_1), \dots, \Delta_j(t_m))^\top$, and the last step is due to Property

2. Plugging (S2.27) into I_2 in (S2.14) yields that

$$\begin{aligned}
 I_2 &= -\delta_n \sum_{i=1}^n \mathbf{v}^\top \mathbf{Z}_i^{*\top} (\mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i^* \mathbf{b}_0) \\
 &= -\delta_n \mathbf{v}^\top \left[\sum_{i=1}^n (\mathbf{Z}_i^\top + \mathbf{U}_i \otimes \boldsymbol{\Phi}^\top(\mathbf{t})) \left\{ \boldsymbol{\varepsilon}_i(\mathbf{t}) + \sum_{j=1}^p X_{ij} \Delta_j(\mathbf{t}) - (\mathbf{U}_i^\top \otimes \boldsymbol{\Phi}(\mathbf{t})) \mathbf{b}_0 \right\} \right] \\
 &= -\delta_n \mathbf{v}^\top \sum_{i=1}^n \mathbf{Z}_i^\top \boldsymbol{\varepsilon}_i(\mathbf{t}) - \delta_n \mathbf{v}^\top \sum_{i=1}^n \mathbf{Z}_i^\top \left(\sum_{j=1}^p X_{ij} \Delta_j(\mathbf{t}) \right) \\
 &\quad + \delta_n \mathbf{v}^\top \sum_{i=1}^n \mathbf{Z}_i^\top (\mathbf{U}_i^\top \otimes \boldsymbol{\Phi}(\mathbf{t})) \mathbf{b}_0 - \delta_n \mathbf{v}^\top \sum_{i=1}^n (\mathbf{U}_i \otimes \boldsymbol{\Phi}^\top(\mathbf{t})) \boldsymbol{\varepsilon}_i(\mathbf{t}) \\
 &\quad - \delta_n \mathbf{v}^\top \sum_{i=1}^n (\mathbf{U}_i \otimes \boldsymbol{\Phi}^\top(\mathbf{t})) \left(\sum_{j=1}^p X_{ij} \Delta_j(\mathbf{t}) \right) + \delta_n \mathbf{v}^\top \sum_{i=1}^n (\mathbf{U}_i \otimes \boldsymbol{\Phi}^\top(\mathbf{t})) (\mathbf{U}_i^\top \otimes \boldsymbol{\Phi}(\mathbf{t})) \mathbf{b}_0 \\
 &\triangleq I_{2,1} + I_{2,2} + I_{2,3} + I_{2,4} + I_{2,5} + I_{2,6}.
 \end{aligned}$$

(S2.28)

Note that

$$\begin{aligned}
\left\| \sum_{i=1}^n \mathbf{Z}_i^\top \boldsymbol{\varepsilon}_i(\mathbf{t}) \right\| &= \left\| \sum_{i=1}^n (\mathbf{X}_i \otimes \boldsymbol{\Phi}^\top(\mathbf{t})) \boldsymbol{\varepsilon}_i(\mathbf{t}) \right\| \\
&= \left\| \text{vec} \left(\boldsymbol{\Phi}^\top(\mathbf{t}) \sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) \mathbf{X}_i^\top \right) \right\| \\
&= \left\| \boldsymbol{\Phi}^\top(\mathbf{t}) \sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) \mathbf{X}_i^\top \right\|_F \\
&\leq \|\boldsymbol{\Phi}(\mathbf{t})\|_2 \left\| \sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) \mathbf{X}_i^\top \right\|_F, \tag{S2.29}
\end{aligned}$$

where the second equality of is due to Property 2 and the last inequality is due to Property 4 (a).

On one hand,

$$\|\boldsymbol{\Phi}(\mathbf{t})\|_2 = \sqrt{\rho_{\max}(\boldsymbol{\Phi}^\top(\mathbf{t})\boldsymbol{\Phi}(\mathbf{t}))} = O(\sqrt{m/M}), \tag{S2.30}$$

by Assumption C3 and the definition of the matrix 2-norm. On the other hand, by the Chebyshev inequality, we have that

$$\begin{aligned}
\mathbb{P} \left(\left\| \sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) \mathbf{X}_i^\top \right\|_F \geq \sqrt{nm}R \right) &\leq \frac{1}{nmR^2} \mathbb{E} \left(\left\| \sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) \mathbf{X}_i^\top \right\|_F^2 \right) \\
&= \frac{n}{nmR^2} \sum_{j=1}^p \sum_{l=1}^m \mathbb{E}(X_{ij}^2) \mathbb{E}[\{\boldsymbol{\varepsilon}_i(t_l)\}^2] \\
&= O(1)/R^2 \rightarrow 0 \text{ as } R \rightarrow \infty,
\end{aligned}$$

where the second step is due to the independence assumption for $\boldsymbol{\varepsilon}_i(\mathbf{t})$ and \mathbf{X}_i , and the third step is due to Assumption A2 and the assumption that

$\sup_{t \in \mathcal{T}} \mathbb{E}[\{\varepsilon_i(t)\}^2] < \infty$. Hence,

$$\left\| \sum_{i=1}^n \varepsilon_i(\mathbf{t}) \mathbf{X}_i^\top \right\|_F = O_p(\sqrt{nm}). \quad (\text{S2.31})$$

Combining this rate and (S2.30) with (S2.29) gives that

$$\left\| \sum_{i=1}^n \mathbf{Z}_i^\top \varepsilon_i(\mathbf{t}) \right\| = O_p\left(\frac{m\sqrt{n}}{\sqrt{M}}\right),$$

which leads to

$$|I_{2,1}| \leq \delta_n \|\mathbf{v}\| \left\| \sum_{i=1}^n \mathbf{Z}_i^\top \varepsilon_i(\mathbf{t}) \right\| = O_p\left(\delta_n \frac{m\sqrt{n}}{\sqrt{M}}\right) \cdot K \quad (\text{S2.32})$$

Similarly, following the steps of (S2.31) and using the independence assumption for $\varepsilon_i(\mathbf{t})$ and \mathbf{U}_i and Assumption D1, we obtain that $I_{2,4}$ in (S2.28) can be bounded as

$$\begin{aligned} |I_{2,4}| &\leq \delta_n \|\mathbf{v}\| \left\| \sum_{i=1}^n (\mathbf{U}_i \otimes \Phi^\top(\mathbf{t})) \varepsilon_i(\mathbf{t}) \right\| \\ &\leq \delta_n \|\mathbf{v}\| \|\Phi(\mathbf{t})\|_2 \left\| \sum_{i=1}^n \varepsilon_i(\mathbf{t}) \mathbf{U}_i^\top \right\|_F \\ &= O_p\left(\delta_n \frac{m\sqrt{n}}{\sqrt{M}}\right) K. \end{aligned} \quad (\text{S2.33})$$

To examine $I_{2,2}$ in (S2.28), define the $m \times p$ matrix $\Delta(\mathbf{t}) = [\Delta_1(\mathbf{t}), \dots, \Delta_p(\mathbf{t})]$.

Then by (S2.13),

$$I_{2,2} = -\delta_n \mathbf{v}^\top \sum_{i=1}^n \{(\mathbf{X}_i \otimes \Phi^\top(\mathbf{t})) \Delta(\mathbf{t}) \mathbf{X}_i\}.$$

Note that by Property 2,

$$\begin{aligned}
\left\| \sum_{i=1}^n \{(\mathbf{X}_i \otimes \Phi^\top(\mathbf{t})) \Delta(\mathbf{t}) \mathbf{X}_i\} \right\| &= \left\| \text{vec} \left(\Phi^\top(\mathbf{t}) \Delta(\mathbf{t}) \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right) \right) \right\| \\
&= \left\| \Phi^\top(\mathbf{t}) \Delta(\mathbf{t}) \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right) \right\|_F \\
&\leq \|\Phi(\mathbf{t})\|_2 \times \|\Delta(\mathbf{t})\|_F \times \left\| \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right\|_2,
\end{aligned} \tag{S2.34}$$

where the last step is due to Property 4. By (S1.2), each element of $\Delta(\mathbf{t})$ is uniformly bounded by $C'_1 M^{-q}$, and hence,

$$\|\Delta(\mathbf{t})\|_F = O(\sqrt{m} M^{-q}). \tag{S2.35}$$

Furthermore, by the Cauchy-Schwarz inequality,

$$\begin{aligned}
\left\| \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right\|_2 &= \rho_{\max} \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right) \\
&= \sup_{\|\mathbf{w}\|=1} \mathbf{w}^\top \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right) \mathbf{w} \\
&= \sup_{\|\mathbf{w}\|=1} \left\{ \sum_{i=1}^n (\mathbf{w}^\top \mathbf{X}_i)^2 \right\} \\
&\leq \sum_{i=1}^n \|\mathbf{X}_i\|^2 \\
&= O_p(n),
\end{aligned} \tag{S2.36}$$

where the first step comes from Property 5, the second step is because that $\rho_{\max}(\mathbf{A}) = \sup_{\|\mathbf{w}\|=1} \mathbf{w}^\top \mathbf{A} \mathbf{w}$ for any square matrix \mathbf{A} , and the last step is due to Assumption A2. Hence, applying (S2.30), (S2.35), (S2.36) to (S2.34)

gives

$$\left\| \sum_{i=1}^n \{(\mathbf{X}_i \otimes \Phi^\top(\mathbf{t}))\Delta(\mathbf{t})\mathbf{X}_i\} \right\| = O_p(mM^{-q-1/2}n),$$

which leads to

$$|I_{2,2}| \leq \delta_n \|\mathbf{v}\| \left\| \sum_{i=1}^n \{(\mathbf{X}_i \otimes \Phi^\top(\mathbf{t}))\Delta(\mathbf{t})\mathbf{X}_i\} \right\| = O_p(\delta_n mM^{-q-1/2}n)K \quad (\text{S2.37})$$

since $\|\mathbf{v}\| = K$.

The term $I_{2,5}$ in (S2.28) can be treated analogously. By (S2.18), we can show that

$$\begin{aligned} |I_{2,5}| &\leq \delta_n \|\mathbf{v}\| \left\| \sum_{i=1}^n \{(\mathbf{U}_i \otimes \Phi^\top(\mathbf{t}))\Delta(\mathbf{t})\mathbf{X}_i\} \right\| \\ &\leq \delta_n \|\mathbf{v}\| \|\Phi(\mathbf{t})\|_2 \|\Delta(\mathbf{t})\|_F \left\| \sum_{i=1}^n \mathbf{U}_i \mathbf{X}_i^\top \right\|_2 \\ &= O_p(\delta_n mM^{-q-1/2}\sqrt{n})K. \end{aligned} \quad (\text{S2.38})$$

Now examine $I_{2,3}$ in (S2.28) using a similar strategy. First, define a $p \times (M + d)$ matrix $\tilde{\mathbf{U}}$ whose transpose is $[\mathbf{b}_{0,1}, \dots, \mathbf{b}_{0,p}]$. According to Assumption B, the true coefficient function $\beta_j(t)$ is uniformly bounded with respect of t , which implies that all elements of $\tilde{\mathbf{U}}$ is also uniformly bounded, i.e.,

$$\sup_{1 \leq j \leq p} \|\mathbf{b}_{0,j}\|_\infty = O(1). \quad (\text{S2.39})$$

Therefore,

$$\|\tilde{\mathbf{U}}\|_F = O(\sqrt{M}). \quad (\text{S2.40})$$

Then by (S2.13),

$$\begin{aligned}
|I_{2,3}| &= \delta_n \mathbf{v}^\top \left[\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{b}_0 \\
&\leq \delta_n \|\mathbf{v}\| \left\| \left[\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{b}_0 \right\| \\
&\leq \delta_n \|\mathbf{v}\| \left\| \text{vec} \left(\Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \tilde{\mathbf{U}} \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \right) \right\| \\
&= \delta_n \|\mathbf{v}\| \left\| \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \tilde{\mathbf{U}} \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right) \right\|_F \\
&\leq \delta_n \|\mathbf{v}\| \|\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\|_2 \|\tilde{\mathbf{U}}\|_F \left\| \sum_{i=1}^n \mathbf{X}_i \mathbf{U}_i^\top \right\|_2 \\
&= O_p(\delta_n m \sqrt{n} / \sqrt{M}) K, \tag{S2.41}
\end{aligned}$$

where the third step is due to Property 2, the fifth step is due to Property 4 and the last step is due to (S2.18), (S2.19), and (S2.40).

Now combining $I_{2,6}$ in (S2.28) and I_4 in (S2.14), we obtain that

$$I_{2,6} + I_4 = \delta_n \mathbf{v}^\top \left[\left\{ \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\} \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{b}_0,$$

which leads to

$$\begin{aligned}
 |I_{2,6} + I_4| &\leq \delta_n \|\mathbf{v}\| \left\| \left[\left\{ \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\} \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{b}_0 \right\| \\
 &\leq \delta_n \|\mathbf{v}\| \left\| \text{vec} \left[\{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \tilde{\mathbf{U}}^\top \left\{ \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\} \right] \right\| \\
 &= \delta_n \|\mathbf{v}\| \left\| \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \tilde{\mathbf{U}}^\top \left\{ \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\} \right\|_F \\
 &\leq \delta_n \|\mathbf{v}\| \left\| \sum_{i=1}^n (\mathbf{U}_i \mathbf{U}_i^\top - \Sigma) \right\|_2 \|\Phi(\mathbf{t})^\top \Phi(\mathbf{t})\|_2 \|\tilde{\mathbf{U}}\|_F \\
 &= O_p(\delta_n m \sqrt{n} / \sqrt{M}) K, \tag{S2.42}
 \end{aligned}$$

where the second and the fourth steps are due to Property 2 and Property 4 respectively, and the last step is due to (S2.23), (S2.19), and (S2.40).

Therefore, combining (S2.32), (S2.33), (S2.37), (S2.38), (S2.41) and (S2.42)

with (S2.28), we obtain that

$$\begin{aligned}
 |I_2 + I_4| &= O_p(\delta_n m M^{-q-1/2} n) K + O_p(\delta_n m \sqrt{n} / \sqrt{M}) K \\
 &= O_p(\delta_n m \sqrt{n} / \sqrt{M}) K, \tag{S2.43}
 \end{aligned}$$

where the last step is due to Assumption C1 (i).

Step 4: Examine I_3 in (S2.14).

By the chain rule,

$$\frac{dP_\lambda(\|\mathbf{w}\|)}{d\mathbf{w}} = P'_\lambda(\|\mathbf{w}\|) \frac{\mathbf{w}}{\|\mathbf{w}\|},$$

for any vector $\mathbf{w} \neq \mathbf{0}$. For $j \notin J_0$, the mean theorem gives that

$$P_\lambda(\|\mathbf{b}_{0,j} + \delta_n \mathbf{v}_j\|) - P_\lambda(\|\mathbf{b}_{0,j}\|) = \delta_n \mathbf{v}_j^\top P'_\lambda(\|\tilde{\mathbf{b}}_j\|) \frac{\tilde{\mathbf{b}}_j}{\|\tilde{\mathbf{b}}_j\|},$$

where $\tilde{\mathbf{b}}_j$ is a vector “between” $\mathbf{b}_{0,j}$ and $\mathbf{b}_{0,j} + \delta_n \mathbf{v}_j$. As $\|\mathbf{b}_{0,j}\|$ is bounded away from 0 for $j \notin J_0$, $\|\tilde{\mathbf{b}}_j\|$ is also bounded away from 0 for sufficiently large n . By Assumption E2 (ii), $P'_\lambda(t) = 0$ for any positive constant t as $\lambda \rightarrow 0$. Hence, when $n \rightarrow \infty$ and $\lambda \rightarrow 0$, we have that

$$P_\lambda(\|\mathbf{b}_{0,j} + \delta_n \mathbf{v}_j\|) - P_\lambda(\|\mathbf{b}_{0,j}\|) = 0,$$

which leads to

$$I_3 = 0. \tag{S2.44}$$

In conclusion, let $R \triangleq I_{1,2} + I_{1,3} + I_{1,4} + I_2 + I_3 + I_4 + I_5$. Then by (S2.26), (S2.43) and (S2.44), we obtain that

$$\begin{aligned} |R| &= O_p(\delta_n^2 \sqrt{n} m / M) K^2 + O_p(\delta_n m \sqrt{n} / \sqrt{M}) K \\ &= O_p(\delta_n m \sqrt{n} / \sqrt{M}) K, \end{aligned} \tag{S2.45}$$

where the second step is due to the definition of δ_n as $n \rightarrow \infty$.

According to (S2.14) and (S2.15), to prove (S2.12) it suffices to show that for any $\epsilon > 0$ there exists $K_\epsilon > 0$ such that

$$\mathbb{P} \left(\inf_{\|\mathbf{v}\|=K_\epsilon} I_{1,1} + R > 0 \right) > 1 - \epsilon. \tag{S2.46}$$

To this end, we prove a stronger version than (S2.46):

$$\mathbb{P} \left(\inf_{\|\mathbf{v}\|=K_\epsilon} I_{1,1} > \sup_{\|\mathbf{v}\|=K_\epsilon} |R| \right) > 1 - \epsilon. \tag{S2.47}$$

By the definition of δ_n , (S2.16) and (S2.45), we have that

$$I_{1,1} \geq m C_0 K^2 \quad \text{and} \quad |R| = O_p(m) \cdot K, \tag{S2.48}$$

for any $\|\mathbf{v}\| = K$. Then for any given $\epsilon > 0$, by (S2.48) we take K_ϵ such that $K_\epsilon C_0 > O(1)$. Then such K_ϵ makes (S2.47) hold. As a result, the assertion (S2.12) is then proved.

Part 2: Prove (5.24).

Given (5.23), it is straightforward that

$$\begin{aligned}
 \|\hat{\beta}_j - \beta_j\|_{L_\infty} &\leq \|\hat{\beta}_j - \theta_j\|_{L_\infty} + \|\theta_j - \beta_j\|_{L_\infty} \\
 &= \|(\hat{\mathbf{b}}_j - \mathbf{b}_{0,j})^\top \boldsymbol{\phi}\|_{L_\infty} + \|\theta_j - \beta_j\|_{L_\infty} \\
 &\leq \|\hat{\mathbf{b}}_j - \mathbf{b}_{0,j}\|_\infty + \|\theta_j - \beta_j\|_{L_\infty} \\
 &\leq \|\hat{\mathbf{b}}_j - \mathbf{b}_{0,j}\| + \|\theta_j - \beta_j\|_{L_\infty} \\
 &= O_p(\delta_n),
 \end{aligned}$$

where the first step is due to the triangle inequality, the second step comes from (3.13) and (S1.3), the third step is due to the fact that $\sum_{k=1}^{M+d} |\phi_k(t)| = 1$ for any $t \in \mathcal{T}$ (Farin et al., 2002, (6.6.7)), the fourth step is because of the relationship between the infinity norm and Euclidean norm for a vector, and the last step is due to Assumption C1(i), (S1.2) and (5.23).

According to Assumption A3, \mathcal{T} is bounded, and thus, its measure, denoted $|\mathcal{T}|$, is finite. Therefore, we have that for any $j = 1, \dots, p$,

$$\|\theta_j - \beta_j\|_{L_2} \leq \sqrt{|\mathcal{T}|} \cdot \|\theta_j - \beta_j\|_{L_\infty},$$

and thus, there exists a positive constant E_1 such that

$$\begin{aligned}
\|\hat{\beta}_j - \beta_j\|_{L_2} &\leq \|\hat{\beta}_j - \theta_j\|_{L_2} + \|\theta_j - \beta_j\|_{L_2} \\
&= \|(\hat{\mathbf{u}}_j - \mathbf{u}_{0,j})^\top \boldsymbol{\phi}\|_{L_2} + \sqrt{|\mathcal{T}|} \|\theta_j - \beta_j\|_{L_\infty} \\
&\leq E_1^{-1} M^{-1/2} \|\hat{\mathbf{u}}_j - \mathbf{u}_{0,j}\| + \sqrt{|\mathcal{T}|} \|\theta_j - \beta_j\|_{L_\infty} \\
&= O_p(\delta_n / \sqrt{M}),
\end{aligned}$$

where the third step follows from Lemma 3, and the last step is due to Assumption C1(i), (S1.2) and (5.23).

S2.3 Proof of Theorem 2

Proof of Theorem 2(i):

We show the result using proof by contradiction. Suppose that there exists some $k \in J_0$ such that $\hat{\mathbf{b}}_k \neq \mathbf{0}$. Then by (3.12), taking derivative of $Q_n(\mathbf{b})$ in (3.11) with respect to \mathbf{b}_k and evaluating it at $\hat{\mathbf{b}}$ yields that

$$\begin{aligned}
\mathbf{0} &= \frac{\partial Q_n(\hat{\mathbf{b}})}{\partial \mathbf{b}_k} \\
&= - \sum_{i=1}^n X_{ik}^* \boldsymbol{\Phi}^\top(\mathbf{t}) (\mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i^* \hat{\mathbf{b}}) \\
&\quad + nm P'_\lambda(\|\hat{\mathbf{b}}_k\|) \frac{\hat{\mathbf{b}}_k}{\|\hat{\mathbf{b}}_k\|} - n [\boldsymbol{\Sigma}_k \otimes \{\boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t})\}] \hat{\mathbf{b}} \\
&= - \sum_{i=1}^n (X_{ik} + U_{ik}) \boldsymbol{\Phi}^\top(\mathbf{t}) \left\{ \mathbf{Y}_i(\mathbf{t}) - \mathbf{Z}_i \hat{\mathbf{b}} - (\mathbf{U}_i^\top \otimes \boldsymbol{\Phi}(\mathbf{t})) \hat{\mathbf{b}} \right\} \\
&\quad + nm P'_\lambda(\|\hat{\mathbf{b}}_k\|) \frac{\hat{\mathbf{b}}_k}{\|\hat{\mathbf{b}}_k\|} - n [\boldsymbol{\Sigma}_k \otimes \{\boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t})\}] \hat{\mathbf{b}}, \tag{S2.49}
\end{aligned}$$

where $\boldsymbol{\Sigma}_k$ represents the k th row of $\boldsymbol{\Sigma}$, and we use (3.7) and (S2.13).

By (3.13) and (S2.13), $\mathbf{Z}_i \hat{\mathbf{b}} = \hat{\boldsymbol{\beta}}(\mathbf{t})^\top \mathbf{X}_i$, where the $p \times m$ matrix $\hat{\boldsymbol{\beta}}(\mathbf{t}) = [\hat{\boldsymbol{\beta}}(t_1), \dots, \hat{\boldsymbol{\beta}}(t_m)]$. Then further applying model (2.2) to the last expression in (S2.49) gives that

$$\begin{aligned}
 \mathbf{0} &= - \sum_{i=1}^n X_{ik} \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\varepsilon}_i(\mathbf{t}) - \sum_{i=1}^n X_{ik} \boldsymbol{\Phi}^\top(\mathbf{t}) (\boldsymbol{\beta}(\mathbf{t}) - \hat{\boldsymbol{\beta}}(\mathbf{t}))^\top \mathbf{X}_i \\
 &\quad + \sum_{i=1}^n X_{ik} \boldsymbol{\Phi}^\top(\mathbf{t}) (\mathbf{U}_i^\top \otimes \boldsymbol{\Phi}(\mathbf{t})) \hat{\mathbf{b}} - \sum_{i=1}^n U_{ik} \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\varepsilon}_i(\mathbf{t}) \\
 &\quad - \sum_{i=1}^n U_{ik} \boldsymbol{\Phi}^\top(\mathbf{t}) (\boldsymbol{\beta}(\mathbf{t}) - \hat{\boldsymbol{\beta}}(\mathbf{t}))^\top \mathbf{X}_i + \sum_{i=1}^n U_{ik} \boldsymbol{\Phi}^\top(\mathbf{t}) (\mathbf{U}_i^\top \otimes \boldsymbol{\Phi}(\mathbf{t})) \hat{\mathbf{b}} \\
 &\quad + nm P'_\lambda(\|\hat{\mathbf{b}}_k\|) \frac{\hat{\mathbf{b}}_k}{\|\hat{\mathbf{b}}_k\|} - n [\boldsymbol{\Sigma}_k \otimes \{\boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t})\}] \hat{\mathbf{b}} \\
 &\triangleq T_1 + T_2 + T_3 + T_4 + T_5 + T_6 + T_7 + T_8.
 \end{aligned} \tag{S2.50}$$

We now show that (S2.50) does not hold because the sum $\sum_{j=1}^8 T_j$ cannot equal zero. To this end, we examine the terms in (S2.50) individually with the following steps.

Step 1: Examination of T_7 in (S2.50):

According to Theorem 1, we have that $\|\hat{\mathbf{b}}_k\| = O_p(\delta_n)$, where $\delta_n = \sqrt{M/n}$ defined in Appendix A.1. By Assumption E2 (i) and the assumption that $\lambda/\delta_n \rightarrow \infty$, it is easily shown that $|P'_\lambda(\|\hat{\mathbf{b}}_k\|) - \lambda| \rightarrow 0$ as $n \rightarrow \infty$. Hence, for sufficiently large n , we claim that

$$\|T_7\| = nm \left| P'_\lambda(\|\hat{\mathbf{b}}_k\|) \right| > \frac{1}{2} nm \lambda. \tag{S2.51}$$

Step 2: Examination of $T_6 + T_8$ in (S2.50):

Simple calculations result in

$$\begin{aligned} T_6 + T_8 &= \sum_{i=1}^n (U_{ik} \mathbf{U}_i^\top) \otimes \{\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\} \hat{\mathbf{b}} - n [\Sigma_k \otimes \{\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\}] \hat{\mathbf{b}} \\ &= \left[\left\{ \sum_{i=1}^n (U_{ik} \mathbf{U}_i^\top - \Sigma_k) \right\} \otimes \{\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\} \right] \hat{\mathbf{b}}, \end{aligned}$$

where the first equality is due to Property 1.

For $1 \leq j \leq p$, we have that

$$\|\hat{\mathbf{b}}_j\|_\infty \leq \|\hat{\mathbf{b}}_j - \mathbf{b}_{0,j}\|_\infty + \|\mathbf{b}_{0,j}\|_\infty = O_p(1),$$

where we use Theorem 1, Condition C1, and (S2.39). Therefore,

$$\|\hat{\mathbf{b}}_j\| \leq \sqrt{M+d} \|\hat{\mathbf{b}}_j\|_\infty = O_p(\sqrt{M}),$$

which implies that

$$\|\hat{\mathbf{b}}\| = O_p(\sqrt{M}) \tag{S2.52}$$

by Condition C1 that the dimension p is finite. Hence, we arrive at

$$\|T_6 + T_8\| \leq \left\| \sum_{i=1}^n (U_{ik} \mathbf{U}_i^\top - \Sigma_k) \right\|_2 \|\Phi(\mathbf{t})^\top \Phi(\mathbf{t})\|_2 \|\hat{\mathbf{b}}\| = O_p(\sqrt{nm}/\sqrt{M}), \tag{S2.53}$$

where the first step is due to Property 3 and the second step is due to (S2.23), (S2.19) and (S2.52).

Step 3: Examination of T_1 and T_4 in (S2.50):

By Property 2, we can show that

$$T_1 = - \sum_{i=1}^n (X_{ik} \otimes \Phi^\top(\mathbf{t})) \boldsymbol{\varepsilon}_i(\mathbf{t}) = \text{vec} \left\{ \Phi^\top(\mathbf{t}) \left(\sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) X_{ik} \right) \right\},$$

which leads to

$$\|T_1\| = \left\| \Phi^\top(\mathbf{t}) \left(\sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) X_{ik} \right) \right\|_F \leq \|\Phi(\mathbf{t})\|_2 \left\| \sum_{i=1}^n \boldsymbol{\varepsilon}_i(\mathbf{t}) X_{ik} \right\| = O_p(\sqrt{nm}/\sqrt{M}), \quad (\text{S2.54})$$

where the inequality is due to Property 4 and the last step is due to (S2.30) and (S2.31).

Analogously, we can show that

$$\|T_4\| = \left\| - \sum_{i=1}^n U_{ik} \Phi^\top(\mathbf{t}) \boldsymbol{\varepsilon}_i(\mathbf{t}) \right\| = O_p(\sqrt{nm}/\sqrt{M}), \quad (\text{S2.55})$$

where the procedure is the same as (S2.54) except we replace X_{ik} in (S2.54) with U_{ik} .

Step 4: Examination of T_3 in (S2.50):

$$\begin{aligned} \|T_3\| &= \left\| \sum_{i=1}^n X_{ik} \Phi^\top(\mathbf{t}) (\mathbf{U}_i^\top \otimes \Phi(\mathbf{t})) \hat{\mathbf{b}} \right\| \\ &= \left\| \left[\left(\sum_{i=1}^n X_{ik} \mathbf{U}_i^\top \right) \otimes \{\Phi^\top(\mathbf{t}) \Phi(\mathbf{t})\} \right] \hat{\mathbf{b}} \right\| \\ &\leq \left\| \sum_{i=1}^n X_{ik} \mathbf{U}_i \right\|_2 \|\Phi(\mathbf{t})^\top \Phi(\mathbf{t})\|_2 \|\hat{\mathbf{b}}\| \\ &= O_p(\sqrt{nm}/\sqrt{M}), \end{aligned} \quad (\text{S2.56})$$

where the second and third steps are due to Property 1 and Property 3, respectively, and the last step is due to (S2.18), (S2.19) and (S2.52).

Step 5: Examination of T_2 and T_5 in (S2.50):

Write $\Theta = \beta(\mathbf{t}) - \hat{\beta}(\mathbf{t})$. Then by Theorem 1, $\|\Theta\|_F = O_p(\delta_n \sqrt{m})$. Then we have that

$$\begin{aligned} \|T_2\| &= \left\| \Phi^\top(\mathbf{t}) \Theta^\top \left(\sum_{i=1}^n X_{ik} \mathbf{X}_i \right) \right\| \\ &\leq \left\| \sum_{i=1}^n X_{ik} \mathbf{X}_i \right\|_2 \|\Phi(\mathbf{t})\|_2 \|\Theta\|_F \\ &= O_p(nm\delta_n/\sqrt{M}), \end{aligned} \tag{S2.57}$$

where the inequality is due to Property 4 and the last step is due to (S2.30) and (S2.36).

Similarly it can be proved that

$$\|T_5\| = \left\| \Phi^\top(\mathbf{t}) \Theta^\top \left(\sum_{i=1}^n U_{ik} \mathbf{X}_i \right) \right\| = O_p(\sqrt{n}m\delta_n/\sqrt{M}) \tag{S2.58}$$

where (S2.18) is used to replace (S2.36).

In conclusion, combining results (S2.53)-(S2.58) leads to

$$\|T_1 + T_2 + T_3 + T_4 + T_5 + T_6 + T_8\| = O_p(nm\delta_n/\sqrt{M}),$$

which, by (S2.51), is asymptotically ignorable relative to T_7 if $\sqrt{n}\lambda \rightarrow \infty$.

Therefore, with probability tending to 1, the sum $\sum_{j=1}^8 T_j$ in (S2.50) is dominated by T_7 , which can never be zero according to (S2.51). That

is, with probability tending to 1, (S2.50) is untrue, and thus, by proof of contradiction, there is no $k \in J_0$ such that $\hat{\mathbf{b}}_k \neq \mathbf{0}$.

Proof of Theorem 2(ii):

Corresponding to the active subvector \mathbf{X}_{I_i} defined before Theorem 1, let $\mathbf{X}_{I_i}^*$ denote its contaminated version, and let \mathbf{U}_{I_i} denote the corresponding subvector of \mathbf{U}_i in (3.7). For the cardinality $s = |J_0|$, we write $\mathbf{v} = (\mathbf{b}_1^\top, \dots, \mathbf{b}_s^\top)^\top$, $\mathbf{v}_0 = (\mathbf{b}_{0,1}^\top, \dots, \mathbf{b}_{0,s}^\top)^\top$ and $\hat{\mathbf{v}} = (\hat{\mathbf{b}}_1^\top, \dots, \hat{\mathbf{b}}_s^\top)^\top$. Concentrating on the subvector \mathbf{v} , let $\check{Q}_n(\mathbf{v}) = Q_n(\mathbf{v}, \mathbf{0})$, where $Q_n(\mathbf{b})$ is defined in (3.11), and $\mathbf{0}$ here is understood to have the dimension so that $(\mathbf{v}^\top, \mathbf{0}^\top)^\top$ is of dimension $p(M + d)$. Throughout the manuscript, we loosely use $\mathbf{0}$ to represent a zero vector or a zero matrix without differentiating them or saying its dimension.

Define

$$S(\mathbf{v}) = \partial \check{Q}_n(\mathbf{v}) / \partial \mathbf{v},$$

which, by (3.11), equals

$$\begin{aligned} S(\mathbf{v}) = & - \sum_{i=1}^n (\mathbf{X}_{I_i}^* \otimes \Phi^\top(\mathbf{t})) \left\{ \mathbf{Y}_i(\mathbf{t}) - \left(\mathbf{X}_{I_i}^{*\top} \otimes \Phi(\mathbf{t}) \right) \mathbf{v} \right\} \\ & - n \left[\Sigma_I \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \mathbf{v} + nm \mathbf{P}(\mathbf{v}) \mathbf{v}, \end{aligned} \quad (\text{S2.59})$$

where $\mathbf{P}(\mathbf{v}) \triangleq \text{diag}\{P'_\lambda(\|\mathbf{b}_1\|)/\|\mathbf{b}_1\|, \dots, P'_\lambda(\|\mathbf{b}_s\|)/\|\mathbf{b}_s\|\} \otimes \mathbf{I}_{M+d}$ and \mathbf{I}_{M+d} is the $(M + d)$ -dimensional identity matrix. Note that $S(\mathbf{v})$ is well-defined

if $\mathbf{b}_j \neq \mathbf{0}$ for all $j = 1, \dots, s$.

Because \mathbf{v}_0 is bounded away from $\mathbf{0}$ and $\hat{\mathbf{v}}$ converges in probability to \mathbf{v}_0 as $n \rightarrow \infty$, $\hat{\mathbf{v}}$ is also bounded in probability away from $\mathbf{0}$ as $n \rightarrow \infty$. Furthermore, by Theorem 1 and Theorem 2 (i), with probability tending to 1, $(\hat{\mathbf{v}}^\top, \mathbf{0}^\top)^\top$ is the minimizer of $Q_n(\mathbf{b})$, and thus, we have that

$$S(\hat{\mathbf{v}}) = \mathbf{0}. \tag{S2.60}$$

Combining (S2.59) and (S2.60) leads to

$$\begin{aligned} \mathbf{0} &= - \sum_{i=1}^n (\mathbf{X}_{I_i}^* \otimes \Phi^\top(\mathbf{t})) \left\{ \mathbf{Y}_i(\mathbf{t}) - \left(\mathbf{X}_{I_i}^{*\top} \otimes \Phi(\mathbf{t}) \right) \hat{\mathbf{v}} \right\} \\ &\quad - n \left[\Sigma_I \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \right] \hat{\mathbf{v}} + nm \mathbf{P}(\hat{\mathbf{v}}) \hat{\mathbf{v}}. \end{aligned} \tag{S2.61}$$

Since $\mathbf{b}_{0,j'} = \mathbf{0}$ for all $j' > s$, (S2.27) yields that

$$\mathbf{Y}_i(\mathbf{t}) = (\mathbf{X}_{I_i}^\top \otimes \Phi(\mathbf{t})) \mathbf{v}_0 + \boldsymbol{\varepsilon}_i(\mathbf{t}) + \sum_{j=1}^p X_{ij} \Delta_j(\mathbf{t}). \tag{S2.62}$$

Inserting (S2.62) into (S2.61) yields that

$$\begin{aligned}
 \mathbf{A}(\hat{\mathbf{v}} - \mathbf{v}_0) &= \frac{M}{nm} \sum_{i=1}^n (\mathbf{X}_{I_i}^* \otimes \Phi^\top(\mathbf{t})) \boldsymbol{\varepsilon}_i(\mathbf{t}) \\
 &\quad - \frac{M}{nm} \sum_{i=1}^n [(\mathbf{X}_{I_i}^* \mathbf{U}_{I_i}^\top - \boldsymbol{\Sigma}_I) \otimes \{\Phi^\top(\mathbf{t})\Phi(\mathbf{t})\}] \mathbf{v}_0 \\
 &\quad - \frac{M}{nm} \sum_{i=1}^n \left[(\mathbf{X}_{I_i}^* \mathbf{X}_{I_i}^{*\top} - \boldsymbol{\Sigma}_I - \boldsymbol{\Sigma}_{\mathbf{X}_I}) \otimes \{\Phi^\top(\mathbf{t})\Phi(\mathbf{t})\} \right] (\hat{\mathbf{v}} - \mathbf{v}_0) \\
 &\quad + \frac{M}{nm} \sum_{i=1}^n \left[(\mathbf{X}_{I_i}^* \otimes \Phi^\top(\mathbf{t})) \left\{ \sum_{j=1}^p X_{ij} \Delta_j(\mathbf{t}) \right\} \right] - M\mathbf{P}(\hat{\mathbf{v}})\hat{\mathbf{v}} \\
 &\triangleq I_1 + I_2 + I_3 + I_4 + I_5,
 \end{aligned} \tag{S2.63}$$

where $\mathbf{A} \triangleq \boldsymbol{\Sigma}_{\mathbf{X}_I} \otimes \left\{ \frac{M}{m} \Phi^\top(\mathbf{t})\Phi(\mathbf{t}) \right\}$ is a matrix whose eigenvalues are bounded and bounded away from 0 due to Assumption C3. Multiplying the inverse of \mathbf{A} on both sides of (S2.63), we can then investigate the limiting distribution of $\boldsymbol{\alpha}^\top(\hat{\mathbf{v}} - \mathbf{v}_0)$ for any vector $\boldsymbol{\alpha}$ satisfying $0 < \|\boldsymbol{\alpha}\| < \infty$. To this end, we examine (S2.63) term by term.

Step 1: Examination of I_1 and I_2 in (S2.63):

As defined in Section 5, $\boldsymbol{\Sigma}_{\mathbf{X}_I}$ is the covariance matrix of \mathbf{X}_{I_i} and $\boldsymbol{\Sigma}_I$ represents the covariance matrix of \mathbf{U}_{I_i} . Then the covariance matrix of I_1

in (S2.63) is computed by

$$\begin{aligned}
\text{Cov}(I_1) &= \frac{M^2}{nm^2} \text{Cov} \{ (\mathbf{X}_{I_i}^* \otimes \Phi^\top(\mathbf{t})) \boldsymbol{\varepsilon}_i(\mathbf{t}) \} \\
&= \frac{M^2}{nm^2} \text{Cov}(\mathbf{X}_{I_i}^*) \otimes \{ \Phi^\top(\mathbf{t}) \boldsymbol{\Sigma}_\varepsilon \Phi(\mathbf{t}) \} \\
&= \frac{M^2}{nm^2} (\boldsymbol{\Sigma}_{\mathbf{X}_I} + \boldsymbol{\Sigma}_I) \otimes \{ \Phi^\top(\mathbf{t}) \boldsymbol{\Sigma}_\varepsilon \Phi(\mathbf{t}) \},
\end{aligned}$$

where the second step is due to Lemma 1.

Let \mathbf{I}_s be the $s \times s$ identity matrix. Recall that $\mathbf{B}_{0s} = [\mathbf{b}_{0,1}, \dots, \mathbf{b}_{0,s}]$ is the $(M+d) \times s$ matrix formed by the blocks of \mathbf{b}_0 , defined in Section 5.

Then the covariance matrix of I_2 in (S2.63) is calculated as follows:

$$\begin{aligned}
\text{Cov}(I_2) &= \frac{M^2}{nm^2} \text{Cov} \left([(\mathbf{X}_{I_i}^* \mathbf{U}_{I_i}^\top - \boldsymbol{\Sigma}_I) \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \}] \mathbf{v}_0 \right) \\
&= \frac{M^2}{nm^2} \text{Cov} \left[\text{vec} \left\{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \left(\mathbf{U}_{I_i} \mathbf{X}_{I_i}^{*\top} \right) \right\} \right] \\
&= \frac{M^2}{nm^2} \text{Cov} \left([\mathbf{X}_{I_i}^* \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \}] \mathbf{U}_{I_i} \right) \\
&= \frac{M^2}{nm^2} \text{Cov} \left([\mathbf{X}_{I_i} \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \}] \mathbf{U}_{I_i} \right) \\
&\quad + \frac{M^2}{nm^2} \text{Cov} \left([\mathbf{U}_{I_i} \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \}] \mathbf{U}_{I_i} \right) \\
&= \frac{M^2}{nm^2} \boldsymbol{\Sigma}_{\mathbf{X}_I} \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \boldsymbol{\Sigma}_I \mathbf{B}_{0s}^\top \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \\
&\quad + \frac{M^2}{nm^2} \text{Cov} \left\{ \text{vec} \left(\Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \mathbf{U}_{I_i} \mathbf{U}_{I_i}^\top \right) \right\} \\
&= \frac{M^2}{nm^2} \boldsymbol{\Sigma}_{\mathbf{X}_I} \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \boldsymbol{\Sigma}_I \mathbf{B}_{0s}^\top \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \\
&\quad + \frac{M^2}{nm^2} \text{Cov} \left([\mathbf{I}_s \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \}] \text{vec} \left(\mathbf{U}_{I_i} \mathbf{U}_{I_i}^\top \right) \right) \\
&= \frac{M^2}{nm^2} \boldsymbol{\Sigma}_{\mathbf{X}_I} \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \boldsymbol{\Sigma}_I \mathbf{B}_{0s}^\top \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \} \\
&\quad + \frac{M^2}{nm^2} [\mathbf{I}_s \otimes \{ \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \mathbf{B}_{0s} \}] \text{Cov} \left\{ \text{vec} \left(\mathbf{U}_{I_i} \mathbf{U}_{I_i}^\top \right) \right\} [\mathbf{I}_s \otimes \{ \mathbf{B}_{0s}^\top \Phi^\top(\mathbf{t}) \Phi(\mathbf{t}) \}],
\end{aligned}$$

where the second, third and sixth steps are due to Property 2, the fourth step is because that \mathbf{X}_{I_i} has mean $\mathbf{0}$ and is independent from \mathbf{U}_{I_i} , and the fifth step is due to Lemma 1 and Property 2.

Define the $s^2 \times s^2$ matrix $\mathbf{\Gamma} = \text{Cov} \{ \text{vec} (\mathbf{U}_{I_i} \mathbf{U}_{I_i}^\top) \}$. Since $\boldsymbol{\varepsilon}_i(\mathbf{t})$ is independent of $\mathbf{X}_{I_i}^*$ and \mathbf{U}_{I_i} , I_1 and I_2 are uncorrelated. Thus,

$$\begin{aligned} \text{Cov}(I_1 + I_2) &= \frac{M^2}{nm^2} (\boldsymbol{\Sigma}_{\mathbf{X}_I} + \boldsymbol{\Sigma}_I) \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Sigma}_\varepsilon \boldsymbol{\Phi}(\mathbf{t}) \} \\ &\quad + \frac{M^2}{nm^2} \boldsymbol{\Sigma}_{\mathbf{X}_I} \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \mathbf{B}_{0s} \boldsymbol{\Sigma}_I \mathbf{B}_{0s}^\top \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \} \\ &\quad + \frac{M^2}{nm^2} [\mathbf{I}_s \otimes \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \mathbf{B}_{0s} \}] \mathbf{\Gamma} [\mathbf{I}_s \otimes \{ \mathbf{B}_{0s}^\top \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \}]. \end{aligned} \tag{S2.64}$$

By the assumption in Theorem 2 (ii) that $\rho_{\min}(\boldsymbol{\Sigma}_\varepsilon) \geq C_5 \xi_m$ with positive constants C_5 and $\{\xi_m : m = 1, 2, \dots\}$, $\boldsymbol{\Sigma}_\varepsilon$ is always positive definite. Furthermore, $\boldsymbol{\Sigma}_{\mathbf{X}_I}$ is also positive definite due to Assumption A2. Hence, $\boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha} \neq 0$ for any $\boldsymbol{\alpha} \neq \mathbf{0}$. Therefore, by the central limit theorem, we have that as $n \rightarrow \infty$,

$$\frac{\boldsymbol{\alpha}^\top \mathbf{A}^{-1} (I_1 + I_2)}{\sqrt{\boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha}}} \xrightarrow{d} N(0, 1). \tag{S2.65}$$

We next show that all the scaled remaining terms in (S2.63), $\boldsymbol{\alpha}^\top \mathbf{A}^{-1} I_k$ for $k = 3, 4, 5$, are asymptotically negligible relative to $\boldsymbol{\alpha}^\top \mathbf{A}^{-1} (I_1 + I_2)$.

Step 2: Examination of I_3 in (S2.63):

Following similar argument in (S2.23) ,

$$\begin{aligned}
& \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_{Ii}^* \mathbf{X}_{Ii}^{*\top} - \boldsymbol{\Sigma}_I - \boldsymbol{\Sigma}_{\mathbf{X}_I} \right\|_2 \\
\leq & \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_{Ii} \mathbf{X}_{Ii}^\top - \boldsymbol{\Sigma}_{\mathbf{X}_I} \right\|_2 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{U}_{Ii} \mathbf{U}_{Ii}^\top - \boldsymbol{\Sigma}_I \right\|_2 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_{Ii} \mathbf{U}_{Ii}^\top \right\|_2 \\
= & O_p(n^{-1/2})
\end{aligned}$$

since $\mathbb{E}(X_{ij}^4) < \infty$ and $\mathbb{E}(U_{ij}^4) < \infty$ for $j = 1, \dots, s$ by Assumptions A2 and D1, respectively. Hence, for I_3 in (S2.63), by (S2.19) and Theorem 1 we have that

$$\begin{aligned}
|\boldsymbol{\alpha}^\top \mathbf{A}^{-1} I_3| & \leq \|\mathbf{A}^{-1} \boldsymbol{\alpha}\| \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_{Ii}^* \mathbf{X}_{Ii}^{*\top} - \boldsymbol{\Sigma}_I - \boldsymbol{\Sigma}_{\mathbf{X}_I} \right\|_2 \\
& \quad \cdot \left\| \frac{M}{m} \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \right\|_2 \|\hat{\mathbf{v}} - \mathbf{v}_0\| \\
& = \|\mathbf{A}^{-1} \boldsymbol{\alpha}\| \cdot O_p(n^{-1/2}) \cdot O_p(\sqrt{M/n}). \tag{S2.66}
\end{aligned}$$

Meanwhile, we have assumed in Theorem 2 (ii) that $\rho_{\min}(\boldsymbol{\Sigma}_\varepsilon) \geq C_5 \xi_m$, and thus, by Assumption C3, we have that

$$\rho_{\min} \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Sigma}_\varepsilon \boldsymbol{\Phi}(\mathbf{t}) \} \geq \rho_{\min} \{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Phi}(\mathbf{t}) \} \rho_{\min}(\boldsymbol{\Sigma}_\varepsilon) \geq C'_5 \frac{m \xi_m}{M},$$

where C'_5 is a positive constant and the first inequality is due to Lemma 2.

Therefore, for some constant C_5'' ,

$$\begin{aligned}
 \boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha} &\geq \|\mathbf{A}^{-1} \boldsymbol{\alpha}\|^2 \rho_{\min}\{\text{Cov}(I_1)\} \\
 &= \frac{M^2}{nm^2} \|\mathbf{A}^{-1} \boldsymbol{\alpha}\|^2 \rho_{\min}(\boldsymbol{\Sigma}_I + \boldsymbol{\Sigma}_{\mathbf{X}_I}) \rho_{\min}\{\boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Sigma}_\varepsilon \boldsymbol{\Phi}(\mathbf{t})\} \\
 &\geq C_5'' \frac{M}{nm} \xi_m \|\mathbf{A}^{-1} \boldsymbol{\alpha}\|^2,
 \end{aligned} \tag{S2.67}$$

where the equality is due to Property 6, and the last step is due to the facts that $\boldsymbol{\Sigma}_{\mathbf{X}_I}$ is positive definite and $\boldsymbol{\Sigma}_I$ is semi-positive definite, which implies that the minimum eigenvalue for $\boldsymbol{\Sigma}_{\mathbf{X}_I} + \boldsymbol{\Sigma}_I$ is bounded away from 0. As a result of (S2.66) and (S2.67),

$$|\boldsymbol{\alpha}^\top \mathbf{A}^{-1} I_3| = o_p\left(\sqrt{\boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha}}\right),$$

since $\frac{n\xi_m}{m} \rightarrow \infty$ as required in Theorem 2 (ii).

Step 3: Examination of I_4 in (S2.63):

Recall that $\boldsymbol{\Delta}(\mathbf{t}) \triangleq [\boldsymbol{\Delta}_1(\mathbf{t}), \dots, \boldsymbol{\Delta}_p(\mathbf{t})]$. By (S2.30), (S2.35) and (S2.36),

we have that

$$\begin{aligned}
|\boldsymbol{\alpha}^\top \mathbf{A}^{-1} I_4| &\leq \frac{M}{nm} \left\| \sum_{i=1}^n \left[(\mathbf{X}_{I_i}^* \otimes \boldsymbol{\Phi}^\top(\mathbf{t})) \left\{ \sum_{j=1}^p X_{ij} \boldsymbol{\Delta}_j(\mathbf{t}) \right\} \right] \right\| \|\mathbf{A}^{-1} \boldsymbol{\alpha}\| \\
&= \frac{M}{nm} \left\| \text{vec} \left\{ \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Delta}(\mathbf{t}) \left(\sum_{i=1}^n \mathbf{X}_{I_i}^* \mathbf{X}_i^\top \right) \right\} \right\| \|\mathbf{A}^{-1} \boldsymbol{\alpha}\| \\
&= \frac{M}{nm} \left\| \boldsymbol{\Phi}^\top(\mathbf{t}) \boldsymbol{\Delta}(\mathbf{t}) \left(\sum_{i=1}^n \mathbf{X}_{I_i}^* \mathbf{X}_i^\top \right) \right\|_F \|\mathbf{A}^{-1} \boldsymbol{\alpha}\| \\
&\leq \frac{M}{nm} \|\mathbf{A}^{-1} \boldsymbol{\alpha}\| \|\boldsymbol{\Phi}(\mathbf{t})\|_2 \|\boldsymbol{\Delta}(\mathbf{t})\|_F \left\| \sum_{i=1}^n \mathbf{X}_{I_i}^* \mathbf{X}_i^\top \right\|_2 \\
&= \frac{M}{nm} \|\mathbf{A}^{-1} \boldsymbol{\alpha}\| \cdot O_p \left(\sqrt{\frac{m}{M}} \right) \cdot O \left(\frac{\sqrt{m}}{M^q} \right) \cdot O_p(n),
\end{aligned}$$

where the second step is due to Property 2, and the fourth is due to Property

4. Thus, we have that

$$|\boldsymbol{\alpha}^\top \mathbf{A}^{-1} I_4| = o_p \left(\sqrt{\boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha}} \right),$$

due to the assumption $\frac{M^{2q} \xi_m}{nm} \rightarrow \infty$ as required in Theorem 2 (ii) and (S2.67).

Step 4: Examination of I_5 in (S2.63):

Because \mathbf{v}_0 is bounded away from $\mathbf{0}$ and $\hat{\mathbf{v}}$ converges in probability to \mathbf{v}_0 as $n \rightarrow \infty$, $\hat{\mathbf{v}}$ is bounded away from $\mathbf{0}$ in probability as $n \rightarrow \infty$. Therefore, we have that $\boldsymbol{\alpha}^\top \mathbf{A}^{-1} I_5 = 0$ with probability tending to 1 since $P'_\lambda(t) = 0$ for any fixed $t > 0$ and sufficiently small λ by Assumption E2 (ii).

In conclusion,

$$\frac{\boldsymbol{\alpha}^\top \mathbf{A}^{-1} (I_3 + I_4 + I_5)}{\sqrt{\boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha}}} = o_p(1).$$

Combining this with result (S2.65) and equation (S2.63) leads to

$$\frac{\boldsymbol{\alpha}^\top(\hat{\boldsymbol{v}} - \boldsymbol{v}_0)}{\sqrt{\boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha}}} \xrightarrow{d} N(0, 1) \quad \text{as } n \rightarrow \infty. \quad (\text{S2.68})$$

We next give an explicit expression for the asymptotic variance of $\boldsymbol{\alpha}^\top(\hat{\boldsymbol{v}} - \boldsymbol{v}_0)$. By the definition of \mathbf{A} and inserting the expression (S2.64) into the denominator of (S2.68), we have that

$$\begin{aligned} & \boldsymbol{\alpha}^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha} \\ &= \frac{1}{n} \boldsymbol{\alpha}^\top \left(\left\{ \boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1} (\boldsymbol{\Sigma}_{\mathbf{X}_I} + \boldsymbol{\Sigma}_I) \boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1} \right\} \right. \\ & \quad \otimes \left[\left\{ \boldsymbol{\Phi}^\top(\boldsymbol{t}) \boldsymbol{\Phi}(\boldsymbol{t}) \right\}^{-1} \left\{ \boldsymbol{\Phi}^\top(\boldsymbol{t}) \boldsymbol{\Sigma}_\varepsilon \boldsymbol{\Phi}(\boldsymbol{t}) \right\} \left\{ \boldsymbol{\Phi}^\top(\boldsymbol{t}) \boldsymbol{\Phi}(\boldsymbol{t}) \right\}^{-1} \right] \\ & \quad \left. + \boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1} \otimes (\mathbf{B}_{0s} \boldsymbol{\Sigma}_I \mathbf{B}_{0s}^\top) + (\boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1} \otimes \mathbf{B}_{0s}) \boldsymbol{\Gamma} (\boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1} \otimes \mathbf{B}_{0s}^\top) \right) \boldsymbol{\alpha}, \quad (\text{S2.69}) \end{aligned}$$

where we use Property 1 and the fact that $(\mathbf{A}_1 \otimes \mathbf{B}_1)^{-1} = \mathbf{A}_1^{-1} \otimes \mathbf{B}_1^{-1}$ for any nonsingular matrices \mathbf{A}_1 and \mathbf{B}_1 (Horn and Johnson, 1991, Section 4.2).

For $j = 1, \dots, s$, let $\boldsymbol{\alpha}_j = \mathbf{e}_j \otimes \boldsymbol{\phi}(t)$, where \mathbf{e}_j is the s -dimensional vector with the j th element being 1 and the rest being 0. Then $\|\boldsymbol{\alpha}_j\| = \|\boldsymbol{\phi}(t)\|$ is bounded away from 0 for any t . Indeed, for any fixed $t \in \mathcal{T}$, at most $d + 1$ consecutive $\phi_j(t)$ are positive while the remaining are all zero due to the construction of the B-spline basis. Without loss of generality, suppose that $\phi_j(t) > 0$ for $j = 1, \dots, d + 1$, and $\phi_j(t) = 0$ otherwise. Let $\boldsymbol{\alpha} = (a_1, \dots, a_{M+d})$ with $a_j = 1$ for $j = 1, \dots, d + 1$ and $a_j = 0$ otherwise.

By Cauchy-Schwartz's inequality, we have that

$$(d+1)\|\boldsymbol{\phi}(t)\|^2 = \|\mathbf{a}\|^2\|\boldsymbol{\phi}(t)\|^2 \geq \left(\sum_{j=1}^{M+d} a_j \phi_j(t)\right)^2 = \left(\sum_{j=1}^{M+d} \phi_j(t)\right)^2 = 1,$$

where the last equality follows from Farin et al. (2002, (6.6.7)). Hence,

$$\|\boldsymbol{\alpha}_j\| \geq \sqrt{1/(d+1)}.$$

With $\boldsymbol{\alpha} = \boldsymbol{\alpha}_j$, the equation (S2.69) can be simplified as

$$\begin{aligned} & \boldsymbol{\alpha}_j^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha}_j \\ &= \frac{1}{n} \left((\Omega_{jj} + \tilde{\Omega}_{jj}) \boldsymbol{\phi}^\top(t) \{\boldsymbol{\Phi}^\top(t) \boldsymbol{\Phi}(t)\}^{-1} \{\boldsymbol{\Phi}^\top(t) \boldsymbol{\Sigma}_\epsilon \boldsymbol{\Phi}(t)\} \{\boldsymbol{\Phi}^\top(t) \boldsymbol{\Phi}(t)\}^{-1} \boldsymbol{\phi}(t) \right. \\ & \quad \left. + \Omega_{jj} \boldsymbol{\phi}^\top(t) \mathbf{B}_{0s} \boldsymbol{\Sigma}_I \mathbf{B}_{0s}^\top \boldsymbol{\phi}(t) + [\boldsymbol{\Omega}_j^\top \otimes \{\boldsymbol{\phi}^\top(t) \mathbf{B}_{0s}\}] \boldsymbol{\Gamma} [\boldsymbol{\Omega}_j \otimes \{\mathbf{B}_{0s}^\top \boldsymbol{\phi}(t)\}] \right) \\ & \triangleq \sigma_j^2(t), \end{aligned}$$

where Ω_{jj} and $\boldsymbol{\Omega}_j$ are the (j, j) element and the j th column of $\boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1}$, respectively; and $\tilde{\Omega}_{jj}$ is the (j, j) element of $\boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1} \boldsymbol{\Sigma}_I \boldsymbol{\Sigma}_{\mathbf{X}_I}^{-1}$. Consequently, we have that

$$\frac{\boldsymbol{\alpha}_j^\top (\hat{\mathbf{v}} - \mathbf{v}_0)}{\sqrt{\sigma_j^2(t)}} \xrightarrow{d} N(0, 1) \quad \text{as } n \rightarrow \infty. \quad (\text{S2.70})$$

Finally, we derive the limiting distribution of $\hat{\beta}_j(t)$ for $j = 1, \dots, s$. By definition,

$$\begin{aligned} \hat{\beta}_j(t) - \beta_j(t) &= \boldsymbol{\phi}^\top(t) \hat{\mathbf{b}}_j - \beta_j(t) \\ &= \boldsymbol{\phi}^\top(t) (\hat{\mathbf{b}}_j - \mathbf{b}_{0,j}) + \boldsymbol{\phi}^\top(t) \mathbf{b}_{0,j} - \beta_j(t) \\ &= \boldsymbol{\alpha}_j^\top (\hat{\mathbf{v}} - \mathbf{v}_0) + \boldsymbol{\phi}^\top(t) \mathbf{b}_{0,j} - \beta_j(t). \end{aligned} \quad (\text{S2.71})$$

By (S1.2), $\boldsymbol{\phi}^\top(t)\mathbf{b}_{0,j} - \beta_j(t) = O(M^{-q})$, which, by (S2.67), leads to

$$\boldsymbol{\phi}^\top(t)\mathbf{b}_{0,j} - \beta_j(t) = o_p\left(\sqrt{\boldsymbol{\alpha}_j^\top \mathbf{A}^{-1} \text{Cov}(I_1 + I_2) \mathbf{A}^{-1} \boldsymbol{\alpha}_j}\right) \quad (\text{S2.72})$$

since $\frac{M^{2q}\xi_m}{nm} \rightarrow \infty$ as required in Theorem 2 (ii). Therefore, combining

(S2.70), (S2.71) and (S2.72) yields that

$$\frac{\hat{\beta}_j(t) - \beta_j(t)}{\sqrt{\sigma_j^2(t)}} \xrightarrow{d} N(0,1) \quad \text{as } n \rightarrow \infty.$$

S3 Simulation Studies

In this section, we conduct simulation studies to evaluate the finite-sample performance of the proposed method and also demonstrate the deleterious effects of ignoring measurement error in inferential procedures.

S3.1 Simulation Design

For $i = 1, \dots, n$, we first generate covariate vector \mathbf{X}_i independently from a centered multivariate normal distribution with covariance matrix whose (j, j') element is $0.5^{|j-j'|}$ for $1 \leq j, j' \leq p$, and then generate response $Y_i(t)$ from model (2.1), where the random error process $\varepsilon_i(t)$ is generated from a centered Gaussian process with covariance function $\Sigma_\varepsilon(t, t') = \sigma_\varepsilon^2 \rho_\varepsilon^{15|t-t'|^2}$ for $t, t' \in \mathcal{T}$. Here σ_ε^2 , which is set to be 0.1, 0.25 or 0.5, representing the variance of $\varepsilon_i(t)$ for $t \in \mathcal{T}$, and ρ_ε is set to be 0.1 or 0.5 to reflect varying

auto-correlation. The common observation points are equally spaced on the domain $\mathcal{T} = [0, 1]$ and the number of observations, m , is set to be 30. Let the number of covariates $p = 6$ with the first three to be active. The functional coefficients are given by $\beta_1(t) = 2t^2$, $\beta_2(t) = \sqrt{2} \cos(3\pi t/2 + \pi/2)$, $\beta_3(t) = \sqrt{2} \sin(\pi t/2) + 3\sqrt{2} \sin(3\pi t/2)$, and $\beta_4(t) = \beta_5(t) = \beta_6(t) = 0$ for $t \in \mathcal{T}$.

The measurement error \mathbf{U}_i in model (3.7) follows a centered multivariate normal distribution with element (j, j') in the covariance matrix Σ to be $\Sigma_{j,j'} = \sigma_{\mathbf{U}}^2 \rho_{\mathbf{U}}^{|j-j'|}$ for $1 \leq j, j' \leq p$, where $\sigma_{\mathbf{U}}^2$ is the variance of each element of \mathbf{U}_i and $\rho_{\mathbf{U}}$ reflects the auto-correlation of the elements in \mathbf{U}_i . To reflect different degrees of measurement error, we set $\sigma_{\mathbf{U}}^2 = 0.1, 0.2$ or 0.4 , and $\rho_{\mathbf{U}} = 0.1$ or 0.5 . Thus, the resulting reliability ratio for each variable, called a marginal reliability ratio, is $\text{Var}(X_{ij}) / \text{Var}(X_{ij}^*) = 1 / (1 + \sigma_{\mathbf{U}}^2) = 0.91, 0.83$ or 0.71 , respectively, where we consider a common marginal reliability ratio for all the covariates.

We consider the sample size $n = 50, 100, 200$, or 500 and simulate $N \triangleq 300$ datasets for each parameter configuration. When implementing the proposed method, we determine tuning parameters M and λ by minimizing CV (4.21) or BIC (4.22) where we set the grid to be $\{M_k = 5k : k = 1, 2, 3\} \times \{\lambda_k = 10^{-1.5+2k/9} : k = 0, 1, \dots, 9\}$, and take the threshold

parameter τ in (4.15) to be 10^{-4} . The impact on the choice of τ is studied in Section S3.5 of the supplementary material. The SCAD penalty (Fan and Li, 2001) is used as the penalty function.

S3.2 Analysis Methods and Performance Metrics

We analyze simulated data $\{\{\mathbf{Y}_i(\mathbf{t}), \mathbf{X}_i^*\} : i = 1, \dots, n; t \in \mathcal{T}\}$ using two methods. The first method, called the naive method, disregards the difference between \mathbf{X}_i^* and \mathbf{X}_i , which basically optimizes the objective function (3.11) with $\boldsymbol{\Sigma} = \mathbf{0}$. The second method, called the proposed method, accommodates the measurement error effects by implementing (3.12). As a comparison, we analyze the precisely measured values $\{\{\mathbf{Y}_i(\mathbf{t}), \mathbf{X}_i\} : i = 1, \dots, n; t \in \mathcal{T}\}$ by implementing (3.12) with $L_n(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}^*)$ replaced by $\tilde{L}_n(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X})$ in (2.6), and call this procedure the oracle method. In implementing those methods, we consider both the CV and BIC.

We evaluate the performance of the resulting estimator using three metrics, the mean integrated squared error (MISE):

$$\text{MISE} = \frac{1}{N} \sum_{r=1}^N \sum_{j=1}^p \int_{\mathcal{T}} \left\{ \hat{\beta}_j^{[r]}(t) - \beta_j(t) \right\}^2 dt,$$

the *integrated squared bias* (ISB):

$$\text{ISB} = \sum_{j=1}^p \int_{\mathcal{T}} \left\{ \frac{1}{N} \sum_{r=1}^N \hat{\beta}_j^{[r]}(t) - \beta_j(t) \right\}^2 dt,$$

and the integrated variance (IV): $IV = \text{MISE} - \text{ISB}$, where $\hat{\beta}_j^{[r]}(t)$ is the estimated function of $\beta_j(t)$ for $j = 1, \dots, p$ obtained from applying each of the aforementioned methods to the r th simulated dataset.

Further, regarding active covariates as “positive” and inactive covariates as “negative”, we evaluate the finite sample performance of selection by the *positive selection ratio* (PSR), defined as the proportion of truly discovered active covariates among all the active covariates, and the *negative selection ratio* (NSR), defined as the proportion of truly discovered inactive covariates among all inactive covariates.

S3.3 Simulation Results - Impact of the Sample Size

First, we evaluate how the performance of the three methods may vary with the sample size, where we set $\sigma_U^2 = 0.1$, $\rho_U = 0.1$, $\sigma_\varepsilon^2 = 0.25$, and $\rho_\varepsilon = 0.1$, and we report in Table S3.1 the results corresponding to CV and BIC (in parentheses). Clearly, the naive method incurs notable biases in estimation, reflected by considerably larger values of ISB and IV than those of the oracle method. On the contrary, the proposed method greatly reduces the bias of the naive method, with the produced ISB values nearly identical to those output by the oracle method. While the proposed method yields larger IV values than the naive method, the overall metric MISE, combining both ISB

and IV, displays a lot smaller values for the proposed method than the naive method, especially when the sample size becomes large. The fact that the proposed method produces a larger IV value than the naive method reflects the price paid to reduce the bias in estimating the coefficient functions, and that a smaller value in the overall metric MISE for the proposed method than the naive method demonstrates the overall benefit of correcting for the measurement error effects in the inferential procedure.

In terms of variable selection, the impact of ignoring the measurement error effects may not be as profound as that for estimation; the naive method yields acceptable PSR values, especially when the sample size is not small. However, NSR values produced by the naive method noticeably deviate from those of the oracle method. As expected, the proposed method significantly improves the performance of the naive method, with the produced variable selection results fairly close to those of the oracle method.

Overall, the naive method yields unreliable results in terms of both estimation and variable selection, and biased results can exacerbate as the sample size increases. In contrast, the proposed method produces satisfactory results, which are fairly comparable to those output by the oracle method; and those results are further improved by increasing the sample

size.

While using CV and BIC yields different results, they reveal the same trend for the dependence on the sample size. In terms of estimation, using CV always results in smaller ISB and IV, and thus MISE, than BIC does, suggesting that using CV produces more accurate function estimates than using BIC. Regarding variable selection, however, using BIC seems preferred to using CV, because using BIC yields higher NSR than using CV.

S3.4 Simulation Results - Impact of the Measurement Error Degree

Here, we evaluate how different degrees of measurement error may affect the performance of the proposed method as well as the naive method. We consider the case with $n = 200$, $\sigma_\varepsilon^2 = 0.25$, and $\rho_\varepsilon = 0.1$. To reflect different magnitudes in measurement error, we consider $\rho_U = 0.1$ or 0.5 and $\sigma_U^2 = 0.2$ or 0.4 , and report in Table S3.2 the results obtained from the three methods using both CV and BIC, with the results corresponding to BIC presented in parentheses.

Again, we observe the same patterns as in Section S3.3, except for cases with $\sigma_U^2 = 0.4$ and BIC used. The proposed method outperforms the naive method in producing smaller values of MISE and larger values of NSR, while

Table S3.1: Simulated results in Section S3.3: Evaluation of the impact of the sample size n under setting with $\sigma_{\mathcal{U}}^2 = 0.1$, $\rho_{\mathcal{U}} = 0.1$, $\sigma_{\varepsilon}^2 = 0.25$, and $\rho_{\varepsilon} = 0.1$. The entries record the results obtained from using CV and BIC, with those corresponding to BIC included in parentheses.

n	Method	MISE	ISB	IV	PSR	NSR
50	naive	0.390(0.411)	0.249(0.261)	0.141(0.150)	0.998(0.999)	0.689(0.746)
	proposed	0.229(0.270)	0.004(0.009)	0.225(0.261)	1.000(0.996)	0.819(0.896)
	oracle	0.024(0.024)	0.000(0.000)	0.024(0.024)	1.000(1.000)	0.932(0.986)
100	naive	0.334(0.335)	0.263(0.264)	0.071(0.071)	1.000(1.000)	0.730(0.754)
	proposed	0.094(0.127)	0.000(0.007)	0.094(0.120)	1.000(1.000)	0.892(0.941)
	oracle	0.012(0.011)	0.000(0.000)	0.012(0.011)	1.000(1.000)	0.918(0.997)
200	naive	0.277(0.277)	0.239(0.239)	0.038(0.038)	1.000(1.000)	0.758(0.784)
	proposed	0.047(0.063)	0.001(0.003)	0.046(0.060)	1.000(1.000)	0.911(0.934)
	oracle	0.006(0.005)	0.000(0.000)	0.006(0.005)	1.000(1.000)	0.933(0.997)
500	naive	0.280(0.276)	0.263(0.260)	0.017(0.016)	1.000(1.000)	0.562(0.653)
	proposed	0.016(0.024)	0.000(0.001)	0.016(0.023)	1.000(1.000)	0.918(0.952)
	oracle	0.002(0.002)	0.000(0.000)	0.002(0.002)	1.000(1.000)	0.987(0.999)

both methods produce values of PSR close to 1. As the measurement error variance $\sigma_{\mathcal{U}}^2$ increases from 0.2 to 0.4, or equivalently, the marginal reliability ratio decreases from 0.83 to 0.71, both the naive and proposed methods

produce increased MISE values yet decreased NSR values. Interestingly, with a given σ_U^2 , when the autocorrelation ρ_U increases, an opposite pattern is observed for both methods. The impact of increasing ρ_U on the naive method seems more noticeable than on the proposed method.

Regarding the use of CV or BIC, again, as observed in Table S3.1, using CV tends to produce estimates with smaller MISE and NSR than using BIC. However, different from the patterns in Table S3.1, applying BIC to the proposed method does not yield satisfactory PSR values for the case with $\rho_U = 0.1$ and $\sigma_U^2 = 0.4$, as the application of CV does.

Additionally, we evaluate how the performance of the proposed method may be influenced by other factors, including the association strength in the response model and the magnitudes of ρ_ε and σ_ε^2 . Further, we assess the impact of different treatments of Σ on the performance of the proposed method. In particular, we conduct simulation studies for the case where Σ is estimated from repeated surrogate measurements of the covariates. We also consider the case where Σ is misspecified and conduct simulation studies accordingly. The details of these additional numerical studies are deferred to Sections S3.6- S3.8 of the supplementary material.

In summary, the simulation studies demonstrate that the naive method ignoring the measurement error effects produces biased results, and that

Table S3.2: Simulated results for Section S3.4: Evaluation of the impact of different measurement error degrees under the setting with $n = 200$, $\sigma_\varepsilon^2 = 0.25$, and $\rho_\varepsilon = 0.1$. The entries record the results obtained from using CV and BIC, with those corresponding to BIC included in parentheses.

(ρ_U, σ_U^2)	Method	MISE	ISB	IV	PSR	NSR
	oracle	0.006(0.006)	0.000(0.000)	0.006(0.006)	1.000(1.000)	0.918(0.997)
(0.1, 0.2)	naive	0.815(0.810)	0.752(0.750)	0.063(0.061)	1.000(1.000)	0.522(0.594)
	proposed	0.110(0.449)	0.004(0.306)	0.106(0.143)	1.000(0.999)	0.878(0.993)
(0.1, 0.4)	naive	1.984(1.976)	1.898(1.895)	0.085(0.081)	1.000(1.000)	0.359(0.448)
	proposed	0.362(2.241)	0.017(1.814)	0.345(0.427)	0.997(0.698)	0.776(0.999)
(0.5, 0.2)	naive	0.360(0.362)	0.319(0.320)	0.041(0.042)	1.000(1.000)	0.892(0.919)
	proposed	0.077(0.132)	0.001(0.026)	0.076(0.106)	1.000(1.000)	0.912(0.940)
(0.5, 0.4)	naive	1.003(1.007)	0.942(0.945)	0.061(0.062)	1.000(1.000)	0.822(0.857)
	proposed	0.207(0.720)	0.003(0.392)	0.204(0.328)	1.000(0.954)	0.827(0.983)

the proposed method greatly improves its performance and produces satisfactory results under a variety of settings. Using CV with the proposed method seems to be preferred over the use of BIC, especially for the case of parameter estimation where the reliability ratio is very low and the sample size is not big enough; in terms of inactive variable detection, the use of BIC may outperform the use of CV.

S3.5 Simulation Results - Effects of τ

As mentioned in Section 4 of the main text, the matrix $\mathbf{X}^{*\top}\mathbf{X}^* - n\boldsymbol{\Sigma}$ may not be positive definite in some finite sample cases, and thus, (4.15) is essential since the unique Cholesky decomposition can only be done for positive definite matrices. To see how the nature of $\mathbf{X}^{*\top}\mathbf{X}^* - n\boldsymbol{\Sigma}$ may depend on the sample size n , we consider the data generated as in Section S3.1 and report the frequencies that $\mathbf{X}^{*\top}\mathbf{X}^* - n\boldsymbol{\Sigma}$ is not positive definite in Table S3.3. It shows that the frequency increases with an increasing σ_U^2 and decreases with an increasing sample size n . Sufficiently large sample size almost always guarantees the matrix $\mathbf{X}^{*\top}\mathbf{X}^* - n\boldsymbol{\Sigma}$ to be positive definite .

Further simulation studies are conducted to study the impact of the choice of the threshold parameter τ in (4.15). As opposed to the mean integrated squared error (MISE) considered in Section S3.2, we calculate the standard deviation of the integrated squared error $\sum_{j=1}^p \int_{\mathcal{T}} \left\{ \hat{\beta}_j^{[r]}(t) - \beta_j(t) \right\}^2 dt$ for $r = 1, \dots, N$, and let SDISE represents it, where $N = 300$ is the number of replicates and $\hat{\beta}_j^{[r]}(t)$ is the estimated function of $\beta_j(t)$ for $j = 1, \dots, p$ by applying the proposed method to the r th simulated dataset. We display values of MISE, SDISE (in parentheses), mean of PSR, and NSR in Table S3.4 under the same setting in Section S3.1. As shown in the table, the performance of the proposed method does not seem to be very sensitive to

Table S3.3: Frequencies of non-positive definiteness of matrix $\mathbf{X}^{*\top}\mathbf{X}^* - n\boldsymbol{\Sigma}$, with $\rho_U=0.1$.

n	σ_U^2	marginal reliability ratio	Frequency
50	0.2	0.833	0.000
	0.4	0.714	0.063
	0.6	0.625	0.257
100	0.2	0.833	0.000
	0.4	0.714	0.000
	0.6	0.625	0.020
200	0.2	0.833	0.000
	0.4	0.714	0.000
	0.6	0.625	0.000

the choice of τ for a wide range of values. However, it is risky to choose relatively large τ for an extremely small sample size. For larger σ_U^2 (i.e., smaller marginal reliability ratio), it is safer to choose a smaller τ to obtain accurate estimates. The results based on the BIC tuning parameter selection strategy are similar and thus omitted.

S3. SIMULATION STUDIES

Table S3.4: Simulation results obtained from using CV under the setting with $\rho_U = 0.1$, $\sigma_\varepsilon^2 = 0.25$ and $\rho_\varepsilon = 0.1$.

σ_U^2	n	τ	MISE (SDISE)	PSR	NSR
0.4	50	0.0001	2.532 (2.726)	0.778	0.931
		0.001	2.532 (2.726)	0.778	0.931
		0.01	3240 (36227)	0.691	0.907
	100	0.0001	0.999 (1.181)	0.973	0.811
		0.001	0.999 (1.182)	0.973	0.811
		0.01	0.999 (1.183)	0.973	0.811
	200	0.0001	0.396 (0.570)	1.000	0.813
		0.001	0.396 (0.570)	1.000	0.813
		0.01	0.396 (0.570)	1.000	0.813
0.6	50	0.0001	4.156 (3.533)	0.574	0.948
		0.001	4.156 (3.533)	0.574	0.948
		0.01	25381 (135641)	0.480	0.810
	100	0.0001	2.080 (2.379)	0.827	0.911
		0.001	2.080 (2.380)	0.827	0.911
		0.01	2.652 (4.346)	0.803	0.901
	200	0.0001	0.896 (1.233)	0.981	0.797
		0.001	0.896 (1.233)	0.981	0.797
		0.01	0.896 (1.233)	0.981	0.797

S3.6 Simulation Results - Effect of ρ_ε and σ_ε^2

We evaluate how different association strength in the response model may influence the performance of the proposed estimator, and report the results in Table S3.5 for the case with $n = 100$, $\sigma_{\mathcal{U}}^2 = 0.1$ and $\rho_{\mathcal{U}} = 0.1$, using CV or BIC (results are presented in parentheses), where we consider $\rho_\varepsilon = 0.1$ or 0.5 and $\sigma_\varepsilon^2 = 0.1$ or 0.5 . As expected, the oracle method always produces the best results among the three methods for all settings. Compared to the naive method, the proposed method produces consistently smaller values of MISE and larger values of NSR, yet both methods give satisfactory values of PSR. For all the three methods, an increase in σ_ε^2 leads to an increase in MISE, but has relatively little effect on changing values of PSR and NSR. The impact of ρ_ε does not appear to be substantial. Under the setting with a small value of $\sigma_{\mathcal{U}}^2$ (i.e., $\sigma_{\mathcal{U}}^2 = 0.1$), the performance of all the three methods appear stable, regardless of whether CV or BIC is used, though using CV tends to give smaller values for MISE and NSR than BIC does in all settings.

S3. SIMULATION STUDIES

Table S3.5: Simulation results for Section S3.6: Evaluation of the impact of different association strengths in the response model, where $n = 100$, $\sigma_U^2 = 0.1$, and $\rho_U = 0.1$. The entries record the results obtained from using CV and BIC, with the results corresponding to BIC included in parentheses.

$(\rho_\varepsilon, \sigma_\varepsilon^2)$	Method	MISE	ISB	IV	PSR	NSR
(0.1, 0.1)	naive	0.325 (0.328)	0.260 (0.264)	0.065 (0.064)	1.000 (1.000)	0.761 (0.788)
	proposed	0.089 (0.112)	0.000 (0.004)	0.088 (0.108)	1.000 (1.000)	0.908 (0.927)
	oracle	0.005 (0.005)	0.000 (0.000)	0.005 (0.005)	1.000 (1.000)	0.943 (0.999)
(0.1, 0.5)	naive	0.343 (0.347)	0.262 (0.264)	0.081 (0.083)	1.000 (1.000)	0.733 (0.791)
	proposed	0.106 (0.137)	0.000 (0.007)	0.106 (0.130)	1.000 (1.000)	0.894 (0.951)
	oracle	0.024 (0.023)	0.000 (0.000)	0.023 (0.023)	1.000 (1.000)	0.910 (0.989)
(0.5, 0.1)	naive	0.329 (0.330)	0.264 (0.265)	0.065 (0.065)	1.000 (1.000)	0.767 (0.779)
	proposed	0.088 (0.119)	0.000 (0.006)	0.088 (0.112)	1.000 (1.000)	0.911 (0.932)
	oracle	0.005 (0.005)	0.000 (0.000)	0.005 (0.005)	1.000 (1.000)	0.947 (0.991)
(0.5, 0.5)	naive	0.345 (0.352)	0.262 (0.267)	0.083 (0.085)	1.000 (1.000)	0.749 (0.791)
	proposed	0.111 (0.141)	0.000 (0.008)	0.111 (0.133)	1.000 (1.000)	0.881 (0.939)
	oracle	0.025 (0.025)	0.000 (0.000)	0.025 (0.024)	1.000 (1.000)	0.928 (0.982)

S3.7 Simulation Study with the Measurement Error Covariance

Matrix Σ Estimated

In this subsection, we investigate the performance of the proposed method for the case where Σ is unknown, but replicate surrogate measurement are available for certain study subjects, allowing for its estimation. In particular, we assume that half of n samples have additional K repeated surrogate measurements. For $k = 1, \dots, K + 1$ and $i = 1, \dots, n/2$, with n assumed to be even, let $\mathbf{X}_{i[k]}^*$ denote the k th surrogate measurement of \mathbf{X}_i . Then using the method of moments, we estimate Σ by

$$\frac{2}{n} \sum_{i=1}^{n/2} \left\{ \frac{1}{K} \sum_{k=1}^{K+1} (\mathbf{X}_{i[k]}^* - \bar{\mathbf{X}}_i^*) (\mathbf{X}_{i[k]}^* - \bar{\mathbf{X}}_i^*)^\top \right\}, \quad (\text{S3.1})$$

where $\bar{\mathbf{X}}_i^* \triangleq \frac{1}{K+1} \sum_{k=1}^{K+1} \mathbf{X}_{i[k]}^*$. See equation (4.3) of Carroll et al. (2006) for a more general formula applicable to varying numbers of replicates for different subjects.

We generate $\{\{\mathbf{X}_i, Y_i(t)\} : t \in \mathcal{T}; i = 1, \dots, n\}$ and $\{\mathbf{X}_i^* : i = \frac{n}{2} + 1, \dots, n\}$ using the same manner of Section S3.1, and $\{\mathbf{X}_{i[k]}^* : k = 1, \dots, K + 1\}$ are generated independently from model (3.7) in the main text for $i = 1, \dots, \frac{n}{2}$, where $\sigma_U = 0.1$, $\rho_U = 0.1$, $\sigma_\varepsilon^2 = 0.25$ and $\rho_\varepsilon = 0.1$, together with $n = 50$ or 200 . First, we apply (S3.1) to the repeated surrogate measurements $\{\mathbf{X}_{i[k]}^* : k = 1, \dots, K + 1; i = 1, \dots, \frac{n}{2}\}$ to obtain an estimate

$\hat{\Sigma}$ of Σ . Next, having Σ replaced by $\hat{\Sigma}$, we apply the proposed method to the data $\{\{\mathbf{X}_i^*, Y_i(t)\} : t \in \mathcal{T}, i = 1, \dots, n\}$, with \mathbf{X}_i^* taken as $\mathbf{X}_{i[1]}^*$ for $i = 1, \dots, \frac{n}{2}$. where CV is used for the tuning parameter selection as in Section S3.1.

We consider $K = 0, 1$ or 2 , where we let $K = 0$ represent the case for which no replicated surrogate measurements are available and Σ is taken as known when applying the proposed method. The results are reported in Table S3.6 in the same manner as for Table S3.1.

As expected, the performance of the proposed method improves as K increases due to the better estimation of Σ with a larger K . When $n = 200$, the proposed method with estimated Σ from $K = 4$ performs analogously to that for the case with $K = 0$.

S3.8 Effects of the Misspecified Measurement Error Covariance Matrix

In this subsection, we assess the sensitivity of the proposed method to the misspecification of the measurement error covariance matrix Σ . In Section S3.1, the covariance matrix Σ used to simulate data is specified as the matrix, denoted as Σ_s , with the (j, j') element set to be $\sigma_U^2 \rho_U^{|j-j'|}$ for $1 \leq j, j' \leq p$, where we consider $\sigma_U^2 = \rho_U = 0.1$. When implementing the

Table S3.6: Simulation results for Section S3.7: Evaluation of the impact of the number of additional repeated measurements K using CV, where $\sigma_{\mathbf{U}} = 0.1$, $\rho_{\mathbf{U}} = 0.1$, $\sigma_{\varepsilon}^2 = 0.25$, $\rho_{\varepsilon} = 0.1$, and $n = 50$ or 200 .

		MISE	ISB	IV	PSR	NSR
$n=50$	$K = 0$	0.229	0.004	0.225	1.000	0.819
	$K = 1$	0.340	0.003	0.337	0.997	0.783
	$K = 4$	0.239	0.004	0.235	0.998	0.823
$n=200$	$K = 0$	0.047	0.001	0.046	1.000	0.911
	$K = 1$	0.054	0.001	0.053	1.000	0.878
	$K = 4$	0.047	0.001	0.046	1.000	0.903

proposed method, the covariance matrix Σ for the error term \mathbf{U}_i is mistaken as a working matrix, denoted as Σ_w , whose (j, j') element is set as $\tilde{\sigma}_{\mathbf{U}}^2 \tilde{\rho}_{\mathbf{U}}^{|j-j'|}$, where in case 1, we take $\tilde{\sigma}_{\mathbf{U}}^2 = \gamma_1 \sigma_{\mathbf{U}}^2$ and $\tilde{\rho}_{\mathbf{U}} = \rho_{\mathbf{U}}$ for a positive constant γ_1 ; and in case 2, we let $\tilde{\sigma}_{\mathbf{U}}^2 = \sigma_{\mathbf{U}}^2$ and $\tilde{\rho}_{\mathbf{U}} = \gamma_2 \rho_{\mathbf{U}}$ for a constant γ_2 . When $\gamma_1 = 1$ or $\gamma_2 = 1$, the working matrix Σ_w is identical to the true matrix Σ_s used to generate data.

Tables S3.7 and S3.8 display the results for cases 1 and 2, respectively, in the same manner as for Table S3.1, where we take $\gamma_1 \in \{0.5, 0.8, 1.0, 1.2, 1.5, 2.0\}$, $\gamma_2 \in \{-1.0, -0.5, 0, 0.5, 1.0, 2.0\}$, and the results for the naive method are displayed in parentheses. Clearly, the values of MISE and ISB minimize

when true Σ is used in both cases. The values of IV tend to increase as γ_1 increases or γ_2 decreases; they decrease as n increases. The values of PSR are always 1 for all settings. When n is large, the best NSR is obtained when true Σ is used. However, best NSR values may be reached for $\gamma_1 > 1$ or $\gamma_2 < 1$ if n is small. Similar performance to the true Σ -based method is achievable when the discrepancy between the working matrix Σ_w and Σ_s is moderate. However, a large discrepancy, such as $\gamma_1 = 2$, can lead to the proposed method performing even worse than the naive method, as shown in Table S3.7.

Table S3.7: Simulation results for Section S3.8: Misspecification of Σ for Case 1 using CV, where $\sigma_\varepsilon^2 = 0.25$, $\rho_\varepsilon = 0.1$, and $n = 100$ or 500 . The results for the naive method is included in parentheses.

		$\gamma_1 = 0.5$	0.8	1	1.2	1.5	2
$n = 100$	MISE	0.148	0.096	0.094 (0.334)	0.126	0.278	1.003
	ISB	0.074	0.012	0.000 (0.263)	0.018	0.127	0.646
	IV	0.074	0.084	0.094 (0.071)	0.108	0.151	0.357
	PSR	1.000	1.000	1.000 (1.000)	1.000	1.000	1.000
	NSR	0.861	0.891	0.892 (0.730)	0.916	0.903	0.837
$n = 500$	MISE	0.086	0.025	0.016 (0.280)	0.037	0.151	0.742
	ISB	0.071	0.011	0.000 (0.263)	0.017	0.119	0.666
	IV	0.015	0.014	0.016 (0.017)	0.020	0.032	0.076
	PSR	1.000	1.000	1.000 (1.000)	1.000	1.000	1.000
	NSR	0.786	0.904	0.918 (0.562)	0.913	0.823	0.611

Table S3.8: Simulation results for Section S3.8: Misspecification of Σ for Case 2 using CV, where $\sigma_\varepsilon^2 = 0.25$, $\rho_\varepsilon = 0.1$, and $n = 100$ or 500 . The results for the naive method is included in parentheses.

		$\gamma_2 = -1$	-0.5	0	0.5	1	2
$n = 100$	MISE	0.166	0.133	0.112	0.098	0.094 (0.334)	0.097
	ISB	0.042	0.022	0.009	0.002	0.000 (0.263)	0.006
	IV	0.124	0.111	0.103	0.096	0.094 (0.071)	0.091
	PSR	1.000	1.000	1.000	1.000	1.000 (1.000)	1.000
	NSR	0.900	0.912	0.911	0.923	0.892 (0.730)	0.874
$n = 500$	MISE	0.076	0.047	0.030	0.020	0.016 (0.280)	0.021
	ISB	0.044	0.022	0.009	0.002	0.000 (0.263)	0.005
	IV	0.031	0.025	0.021	0.018	0.016 (0.017)	0.016
	PSR	1.000	1.000	1.000	1.000	1.000 (1.000)	1.000
	NSR	0.768	0.851	0.889	0.916	0.918 (0.562)	0.881

S4 Analysis Results Additional to Section 6

This section records additional results for Section 6, which reports some results for $a = c = 1$ only. First, for $a = c = 1$, we display in Figure S4.1 the estimates of $\beta_j(t)$ corresponding to each of the twelve covariates obtained by using either CV or BIC to choose a suitable tuning parameter value. Estimates of $X_{i,9}$ and $X_{i,10}$ are quite different with the use of CV or BIC, for which the corresponding covariates $X_{i,9}$ and $X_{i,10}$ are excluded by CV but not by BIC. Additionally, all the covariates excluded by BIC are also excluded by CV.

Next, we report sensitivity in estimating $\beta_1(t)$, the coefficient function corresponding to the BMI Z-score ($X_{i,1}$), when different degrees of measurement error are present in the data, and display the estimates of $\beta_1(t)$ in Figure S4.2, where $a = 0, 1, 2$ and 3 and $c = 0, 5, 10$ and 15 are considered. Here $a = 0$ corresponds to the naive method which ignores the differences between \mathbf{X}_i and \mathbf{X}_i^* . It seems that the value of a has more noticeable impact on estimation results than that of c . All other estimated functions share analogous trends and the variable selection results are consistent among all the chosen values for a and c , and are not reported here.

S4. ANALYSIS RESULTS ADDITIONAL TO SECTION ??

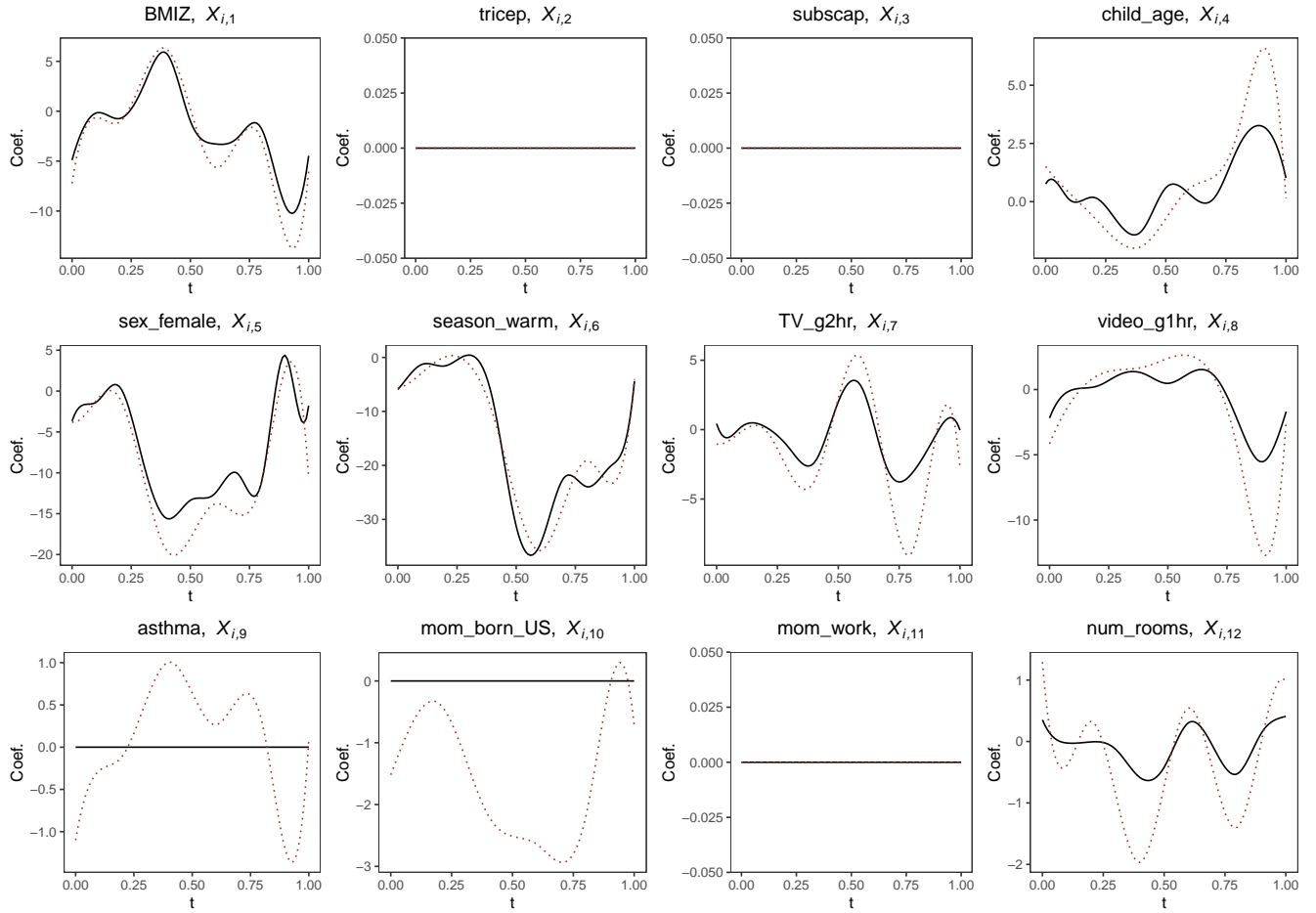


Figure S4.1: The estimates of $\beta_j(t)$ corresponding to each of the twelve covariates obtained by using CV (solid) or BIC (dotted) to choose a tuning parameter value with $a = c = 1$.

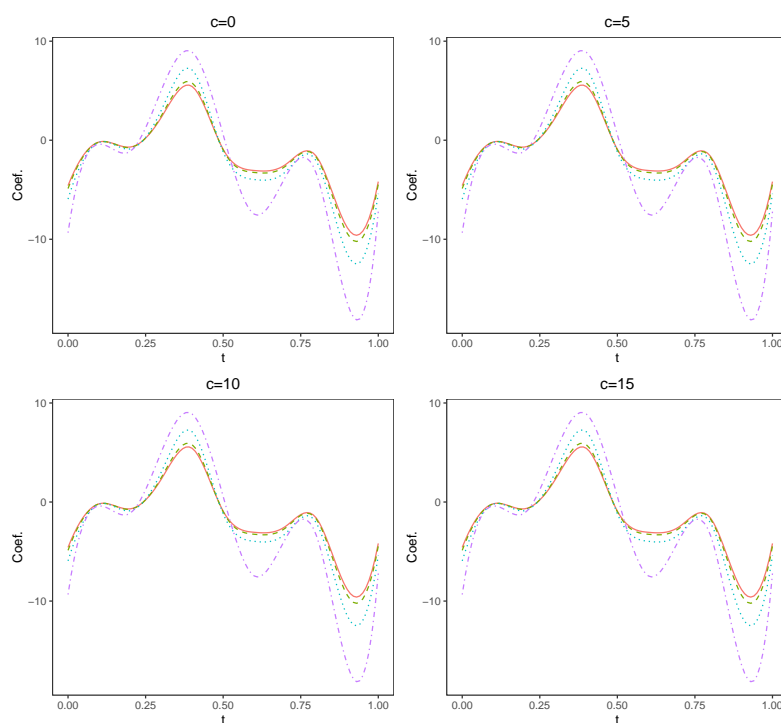


Figure S4.2: BMI Z-score coefficient function estimates for different values of a and c : the estimates corresponding to $a = 0, 1, 2$ and 3 are expressed by curves marked with solid red, dash green, dotted blue and dash-dotted purple, respectively.

S5 Extension to Accommodating the Generalized Least Squares Loss Function

S5.1 Method and Properties

The simple least squares loss function (2.6) focuses on expressing the differences between the responses $\mathbf{Y}(\mathbf{t})$ and their approximate mean $(\mathbf{X} \otimes \Phi(\mathbf{t}))\mathbf{b}$, without accounting for the correlation within the error process $\boldsymbol{\varepsilon}(\mathbf{t})$ over \mathbf{t} . In this section, we extend the method in the main text to accommodate generalized least squares loss functions. Let \mathbf{G} denote a symmetric positive-definite $m \times m$ matrix. Define

$$\tilde{L}_n^G(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}) \triangleq \frac{1}{2} \sum_{i=1}^n \|\mathbf{Y}_i(\mathbf{t}) - (\mathbf{X}_i^\top \otimes \Phi(\mathbf{t}))\mathbf{b}\|_G^2, \quad (\text{S5.2})$$

where $\|\mathbf{a}\|_G^2$ represents $\mathbf{a}^\top \mathbf{G} \mathbf{a}$ for any m -dimensional vector \mathbf{a} . The new loss function (S5.2) differs from (2.6) in the inclusion of a weight matrix \mathbf{G} .

Under the model (3.7), (3.9) can be modified as

$$\mathbb{E} \left\{ \tilde{L}_n^G(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}^*) \right\} = \mathbb{E} \left\{ \tilde{L}_n^G(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}) \right\} + \frac{n}{2} \mathbf{b}^\top [\boldsymbol{\Sigma} \otimes \{\Phi^\top(\mathbf{t})\mathbf{G}\Phi(\mathbf{t})\}] \mathbf{b},$$

which motivates the generalized version of (3.10):

$$L_n^G(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}^*) \triangleq \tilde{L}_n^G(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}^*) - \frac{n}{2} \mathbf{b}^\top [\boldsymbol{\Sigma} \otimes \{\Phi^\top(\mathbf{t})\mathbf{G}\Phi(\mathbf{t})\}] \mathbf{b}. \quad (\text{S5.3})$$

Define

$$Q_n^G(\mathbf{b}) = L_n^G(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}^*) + nm \sum_{j=1}^p P_\lambda(\|\mathbf{b}_j\|),$$

which is identical to (3.11) except replacing the loss function $L_n(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}^*)$ with $L_n^G(\mathbf{b}; \mathbf{Y}(\mathbf{t}), \mathbf{X}^*)$.

Replacing $Q_n(\mathbf{b})$ in (3.12) with $Q_n^G(\mathbf{b})$, we let $\hat{\mathbf{b}}^G$ denote the resulting estimators, and let $\hat{\beta}_j^G$ for $j = 1, \dots, p$ denote the estimators determined by (3.13) with $\hat{\mathbf{b}}_j$ replaced by the j th element of $\hat{\mathbf{b}}^G$. Analogous to Theorems 1 and 2 in the main text, we establish asymptotic properties for $\hat{\mathbf{b}}^G$ and $\hat{\beta}_j^G$ with $j = 1, \dots, p$.

Corollary 1. *Assume that the weight matrix \mathbf{G} satisfies*

$$C_4'^{-1} \zeta_m' \leq \rho_{\min}(\mathbf{G}) \leq \rho_{\max}(\mathbf{G}) \leq C_4' \zeta_m, \quad (\text{S5.4})$$

where C_4' is a positive constant, ζ_m and ζ_m' are sequences of m satisfying $\zeta_m/\zeta_m' = o(\sqrt{n/M})$. Then under the conditions of Theorem 1, there exists a local minimizer $\hat{\mathbf{b}}^G$ of $Q_n^G(\mathbf{b})$ such that

$$\|\hat{\mathbf{b}}^G - \mathbf{b}_0\| = O_p(\zeta_m \sqrt{M/n/\zeta_m'}),$$

and hence, for $j = 1, \dots, p$,

$$\|\hat{\beta}_j^G - \beta_j\|_{L_\infty} = O_p(\zeta_m \sqrt{M/n/\zeta_m'}) \text{ and } \|\hat{\beta}_j^G - \beta_j\|_{L_2} = O_p(\zeta_m \sqrt{1/n/\zeta_m'}).$$

Remark 3. In Corollary 1, ζ_m is allowed to be divergent and ζ_m' is allowed to converge to 0 as $m \rightarrow \infty$. For example, if we choose \mathbf{G} to be the inverse

of Σ_ε , then ζ_m can be ξ_m^{-1} , as discussed in Remark 2, and thus, can be divergent for some settings of $\varepsilon_i(t)$. Also, the largest eigenvalue of Σ_ε may be divergent as $m \rightarrow \infty$, and thus, the minimum eigenvalue of the inverse of Σ_ε , \mathbf{G} , converges to 0, showing that $\zeta'_m \rightarrow 0$. Corollary 1 implies that incorporating \mathbf{G} may result in a lower convergence rate than Theorem 1.

The following corollary encompasses Theorem 2 as a special case by setting $\zeta'_m = \zeta_m = 1$ and \mathbf{G} to be the identity matrix.

Corollary 2. *Suppose the conditions in Corollary 1 hold. Assume further that $\lambda\sqrt{n/M} \cdot \zeta'_m/\zeta_m \rightarrow \infty$ and $\lambda\sqrt{n} \cdot \zeta'_m/\zeta_m^2 \rightarrow \infty$ as $n \rightarrow \infty$.*

(i) *Then with probability tending to 1, the minimizer $\hat{\mathbf{b}}$ satisfies $\hat{\mathbf{b}}_j = \mathbf{0}$ for all $j \in J_0$, and thus, $\hat{J}_0 = J_0$.*

(ii) *Assume $\rho_{\min}(\Sigma_\varepsilon) \geq C_5 \xi_m$, where C_5 is a positive constant and $\{\xi_m : m = 1, 2, \dots\}$ is a sequence of constants satisfying $\frac{n\xi_m \zeta_m^{14}}{m\zeta_m^4} \rightarrow \infty$ and $\frac{M^{2q} \xi_m \zeta_m'^2}{nm\zeta_m^2} \rightarrow \infty$ as $n \rightarrow \infty$ and $m \rightarrow \infty$. Then for any $t \in \mathcal{T}$ and $j \notin J_0$, we have that*

$$\frac{\hat{\beta}_j(t) - \beta_j(t)}{\sqrt{\sigma_j^2(t)}} \xrightarrow{d} N(0, 1) \quad \text{as } n \rightarrow \infty,$$

where

$$\begin{aligned} & \sigma_j^2(t) \\ &= \frac{1}{n} \left((\Omega_{jj} + \tilde{\Omega}_{jj}) \boldsymbol{\phi}^\top(t) \{\boldsymbol{\Phi}^\top(t) \mathbf{G} \boldsymbol{\Phi}(t)\}^{-1} \{\boldsymbol{\Phi}^\top(t) \mathbf{G} \boldsymbol{\Sigma}_\varepsilon \mathbf{G} \boldsymbol{\Phi}(t)\} \{\boldsymbol{\Phi}^\top(t) \mathbf{G} \boldsymbol{\Phi}(t)\}^{-1} \boldsymbol{\phi}(t) \right. \\ & \quad \left. + \Omega_{jj} \boldsymbol{\phi}^\top(t) \mathbf{B}_{0s} \boldsymbol{\Sigma}_I \mathbf{B}_{0s}^\top \boldsymbol{\phi}(t) + [\boldsymbol{\Omega}_j^\top \otimes \{\boldsymbol{\phi}^\top(t) \mathbf{B}_{0s}\}] \boldsymbol{\Gamma} [\boldsymbol{\Omega}_j \otimes \{\mathbf{B}_{0s}^\top \boldsymbol{\phi}(t)\}] \right) \end{aligned}$$

with Ω_{jj} , $\boldsymbol{\Omega}_j$ and $\tilde{\Omega}_{jj}$ given in Theorem 2.

Comparing the variance $\sigma_j^2(t)$ in Corollary 2 and the counterpart in Theorem 2, the only difference is the inclusion of \mathbf{G} in the first term. Specifically, if we set $\mathbf{G} = \boldsymbol{\Sigma}_\varepsilon^{-1}$, the first term can be simplified as

$$\frac{1}{n} (\Omega_{jj} + \tilde{\Omega}_{jj}) \boldsymbol{\phi}^\top(t) \{\boldsymbol{\Phi}^\top(t) \boldsymbol{\Sigma}_\varepsilon^{-1} \boldsymbol{\Phi}(t)\}^{-1} \boldsymbol{\phi}(t).$$

The proofs of Corollaries 1 and 2 are similar to those for Theorems 1 and 2 with slight modifications, and thus, are omitted.

S5.2 Implementation

The algorithm of minimizing $Q_n^G(\mathbf{b})$ is analogous to Section 4 in the main text, except the inclusion of the matrix \mathbf{G} . To be specific, let $\mathbf{W}^G \triangleq (\mathbf{X}^{*\top} \mathbf{X}^* - n\boldsymbol{\Sigma}) \otimes \{\boldsymbol{\Phi}^\top(t) \mathbf{G} \boldsymbol{\Phi}(t)\}$ be the counterpart of \mathbf{W} in (4.14). We then obtain $\overline{\mathbf{W}}^G$ by (4.15) and the Cholesky decomposition $\overline{\mathbf{W}}^G = (\mathbf{V}^G)^\top \mathbf{V}^G$.

Define $\tilde{\mathbf{b}}^G \triangleq (\overline{\mathbf{W}}^G)^{-1} (\mathbf{X}^* \otimes \boldsymbol{\Phi}(t))^\top \mathbb{G} \mathbf{Y}(t)$, where $\mathbb{G} \triangleq \text{diag}(\mathbf{G}, \dots, \mathbf{G})$

with n blocks of \mathbf{G} . We then consider minimizing the objective function

$$\begin{aligned}\tilde{Q}_n^G(\mathbf{b}) &= \frac{1}{2nm} \left\| \mathbf{V}^G \tilde{\mathbf{b}}^G - \mathbf{V}^G \mathbf{b} \right\|^2 + \sum_{j=1}^p P_\lambda(\|\mathbf{b}_j\|) \\ &= \frac{1}{2(M+d)p} \left\| \sqrt{\frac{(M+d)p}{nm}} \mathbf{V}^G \tilde{\mathbf{b}}^G - \sqrt{\frac{(M+d)p}{nm}} \mathbf{V}^G \mathbf{b} \right\|^2 + \sum_{j=1}^p P_\lambda(\|\mathbf{b}_j\|).\end{aligned}$$

The remaining procedure is identical to the discussion for (4.19) in the main text, except that the term $\mathbf{V}^G \tilde{\mathbf{b}}^G$ is now computed by solving the equation

$$(\mathbf{V}^G)^\top \tilde{\mathbf{Y}} = (\mathbf{X}^* \otimes \Phi(\mathbf{t}))^\top \mathbb{G} \mathbf{Y}(\mathbf{t})$$

for unknown $\tilde{\mathbf{Y}}$.

For the tuning parameters selection, the CV method is similar to (4.21) except the ordinary $\|\cdot\|$ term is replaced by $\|\cdot\|_G$ and the $\Phi^\top(\mathbf{t})\Phi(\mathbf{t})$ is now $\Phi^\top(\mathbf{t})\mathbf{G}\Phi(\mathbf{t})$. Analogously, BIC is the same as (4.22) with $\|\cdot\|$ replaced by $\|\cdot\|_G$, where $\|\mathbf{a}\|_G^2$ represents $\mathbf{a}^\top \mathbb{G} \mathbf{a}$ for any mn -dimensional vector \mathbf{a} .

S5.3 A Simulation Study

In the literature of generalized least squares (GLS) methods, the matrix \mathbf{G} is commonly set as the inverse matrix Σ_ε^{-1} . In this subsection, we compare the performance of the original method based on least squares (LS) which minimizes (3.11), to the extended approach with a generalized least squares function, where we set \mathbf{G} to be Σ_ε^{-1} or $\hat{\Sigma}_\varepsilon^{-1}$, and let GLS0 and GLS refer to those methods, respectively. Here $\hat{\Sigma}_\varepsilon^{-1}$ is simply obtained using the

residuals, formed as follows: the residual for unit i is given by

$$\hat{e}_i \triangleq \mathbf{Y}_i(\mathbf{t}) - (\mathbf{X}_i^{*\top} \otimes \Phi(\mathbf{t}))\hat{\mathbf{b}} \text{ for } i = 1, \dots, n,$$

where $\hat{\mathbf{b}}$ denotes the estimate obtained from the LS method (3.12). To account for the use of the contaminated \mathbf{X}_i^* , we estimate Σ_ε by adding a debias term as follows:

$$\hat{\Sigma}_\varepsilon = \frac{1}{n} \sum_{i=1}^n \hat{e}_i \hat{e}_i^\top - \Phi(\mathbf{t}) \hat{\mathbf{B}} \Sigma \hat{\mathbf{B}}^\top \Phi^\top(\mathbf{t}), \quad (\text{S5.5})$$

where $\hat{\mathbf{B}}$ is the $(M+d) \times p$ matrix version of $\hat{\mathbf{b}}$ satisfying $\text{vec}(\hat{\mathbf{B}}) = \hat{\mathbf{b}}$. We can then compute the inverse $\hat{\Sigma}_\varepsilon^{-1}$ for GLS. If there is no measurement error, the first part of $\hat{\Sigma}_\varepsilon$ in (S5.5) is often considered in multivariate responses models. See, for example, Sofer et al. (2014) and the reference therein. In functional response models, Chen et al. (2016) proposed a similar procedure.

The setting is analogous to those in Section S3.1 except for the random error. Indeed, direct calculation of Σ_ε^{-1} based on the definition, i.e., $\Sigma_\varepsilon(t, t') = \sigma_\varepsilon^2 \rho_\varepsilon^{15|t-t'|^2}$ for $t, t' \in \mathcal{T}$, is unstable or infeasible since those $\varepsilon_i(t_1), \dots, \varepsilon_i(t_m)$ are highly correlated even when ρ_ε is very small. Alternatively, we consider enhancing the diagonal elements of $\hat{\Sigma}_\varepsilon$ by defining $\Sigma_\varepsilon(t, t') = \sigma_\varepsilon^2 \rho_\varepsilon^{15|t-t'|^2} + \sigma_\varepsilon'^2 \mathbf{1}(t = t')$ for $t, t' \in \mathcal{T}$ now, where σ_ε' is a positive constant.

Repeating the experiments for three methods 300 times, Table S5.1 re-

S5. EXTENSION TO ACCOMMODATING THE GENERALIZED LEAST
SQUARES LOSS FUNCTION

ports MISE and SDISE (displayed in parentheses), the definition of which is given in Section S3.5) and means of PSR and NSR under different sample sizes n using the CV method. Clearly, GLS0 is consistently the best, though it only slightly outperforms LS. When n is relatively small, GLS can perform worse than LS because of inaccurate estimation of Σ_ε .

Table S5.1: Simulation results for comparing three methods under the setting with $\sigma_U^2 = 0.1$, $\rho_U = 0.1$, $\sigma_\varepsilon^2 = \sigma_{\varepsilon'}^2 = 0.25$, and $\rho_\varepsilon = 0.1$

n	Method	MISE (SDISE)	PSR	NSR
100	LS	0.097 (0.090)	1.000	0.903
	GLS	0.124 (0.142)	0.998	0.848
	GLS0	0.093 (0.088)	1.000	0.913
200	LS	0.047 (0.041)	1.000	0.894
	GLS	0.046 (0.036)	1.000	0.896
	GLS0	0.047 (0.040)	1.000	0.912
500	LS	0.017 (0.016)	1.000	0.922
	GLS	0.016 (0.014)	1.000	0.911
	GLS0	0.015 (0.012)	1.000	0.953

References

- Carroll, R. J., D. Ruppert, L. A. Stefanski, and C. M. Crainiceanu (2006). *Measurement Error in Nonlinear Models: A Modern Perspective*. Chapman and Hall/CRC.
- Chen, Y., J. Goldsmith, and R. T. Ogden (2016). Variable selection in function-on-scalar regression. *Stat* 5(1), 88–101.
- De Boor, C. (1978). *A Practical Guide to Splines*. Springer-Verlag New York.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association* 96(456), 1348–1360.
- Farin, G., J. Hoschek, and M. S. Kim (2002). *Handbook of Computer Aided Geometric Design*. Elsevier.
- Horn, R. A. and C. R. Johnson (1991). *Topics in Matrix Analysis*. Cambridge University Press.
- Horn, R. A. and C. R. Johnson (2012). *Matrix Analysis*. Cambridge University Press.
- Massart, P. (1990). The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *The Annals of Probability*, 1269–1283.
- Parodi, A. and M. Reimherr (2018). Simultaneous variable selection and smoothing for high-dimensional function-on-scalar regression. *Electronic Journal of Statistics* 12(2), 4602–4639.
- Sofer, T., L. Dicker, and X. Lin (2014). Variable selection for high dimensional multivariate outcomes. *Statistica Sinica* 24(4), 1633.

REFERENCES

- Wang, L., G. Chen, and H. Li (2007). Group scad regression analysis for microarray time course gene expression data. *Bioinformatics* 23(12), 1486–1494.
- Yi, G. Y. (2017). *Statistical Analysis with Measurement Error or Misclassification: Strategy, Method and Application*. Springer.
- Zhang, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics* 38(2), 894–942.
- Zhou, S., X. Shen, and D. A. Wolfe (1998). Local asymptotics for regression splines and confidence regions. *The Annals of Statistics* 26(5), 1760–1782.