

---

# NONPARAMETRIC CLUSTER ANALYSIS ON MULTIPLE OUTCOMES OF LONGITUDINAL DATA

Yang Lv<sup>1</sup>, Xiaolu Zhu<sup>2</sup>, Zhongyi Zhu<sup>3</sup> and Annie Qu<sup>4</sup>

<sup>1</sup> *School of Statistics, Capital University of Economics and Business, Beijing, China*

<sup>2</sup> *Amazon.com Inc., Seattle, Washington, U.S.*

<sup>3</sup> *Department of Statistics, School of Management, Fudan University, Shanghai, China*

<sup>4</sup> *Department of Statistics, University of Illinois at Urbana-Champaign*

## Supplementary Material

In this supplement, we provide simulation results under more settings, and give the technical proofs for Lemma 1-2, Theorem 1 and Corollary 1.

### S1 Other Simulations

In this section, we let the random error  $\boldsymbol{\varepsilon}_{im} = (\varepsilon_{i1m}, \dots, \varepsilon_{i10m})^T$  follow non-normal distributions such as the exponential distribution,  $t$ -distribution, or mixture distribution in the following Case 3–5:

Case 3 :  $\boldsymbol{\varepsilon}_{im} = \exp(\xi_{im}) - 1$ , where  $\xi_{im} \sim N(0, 0.25R)$ , and the correlation matrix  $R$  is the same as in Case 2.

Case 4 :  $\boldsymbol{\varepsilon}_{im} \sim t_3(0, 0.25R)$ , where the correlation  $R$  is the same as in Case 2.

Case 5 :  $\boldsymbol{\varepsilon}_{i1} \sim N(0, 0.25R)$  and  $\boldsymbol{\varepsilon}_{i2} \sim t_3(0, 0.04R)$ , where the correlation  $R$  is the same as in Case 2.

The simulation results based on 100 simulation runs are provided in the following Tables S1–S3, which show that the proposed method is robust against any parametric assumption and performs the best compared to other methods.

## S2 Proof of Lemma 1

Note that we can individually obtain  $\tilde{\boldsymbol{\beta}}_i$  and  $\tilde{b}_i$  as

$$\tilde{\boldsymbol{\beta}}_i = (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i \mathbf{y}_i, \quad (\text{S2.1})$$

$$\tilde{b}_i = (\mathbf{1}_{n_i}^T \boldsymbol{\Sigma}_i^{-1} \mathbf{1}_{n_i} + \lambda_2)^{-1} \mathbf{1}_{n_i}^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{y}_i - \mathbf{B}_i \tilde{\boldsymbol{\beta}}_i), \quad (\text{S2.2})$$

where  $\mathbf{W}_i = (\boldsymbol{\Sigma}_i + \frac{1}{\lambda_2} \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T)^{-1}$ . Then we have

$$\begin{aligned} \tilde{\mathbf{f}}_i &= \mathbf{B}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i \mathbf{y}_i \\ &= \mathbf{B}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i \mathbf{f}_i - \mathbf{B}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \lambda_1 \mathbf{D}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i)^{-1} \\ &\quad \mathbf{B}_i^T \mathbf{W}_i \mathbf{f}_i + \mathbf{B}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i (\mathbf{1}_{n_i} b_i + \boldsymbol{\varepsilon}_i). \end{aligned} \quad (\text{S2.3})$$

When  $\lambda_1 = 0$  in (S2.3), it refers to a regression spline estimator. More precisely, the regression spline estimator  $\tilde{\mathbf{f}}_i^{\text{reg}} = \mathbf{B}_i \tilde{\boldsymbol{\beta}}_i^{\text{reg}}$  is the minimizer of

$$(\mathbf{y}_i - \mathbf{B}_i \tilde{\boldsymbol{\beta}}_i^{\text{reg}})^T \mathbf{W}_i (\mathbf{y}_i - \mathbf{B}_i \tilde{\boldsymbol{\beta}}_i^{\text{reg}}), \quad (\text{S2.4})$$

thus  $E(\tilde{\mathbf{f}}_i^{\text{reg}}) = \mathbf{B}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i \mathbf{f}_i$ . Furthermore, we denote  $\mathbf{f}_i^s = \mathbf{B}_i \boldsymbol{\beta}_i$  as the best  $L_\infty$  approximation to  $\mathbf{f}_i$ , where  $\boldsymbol{\beta}_i \in \mathcal{R}^p$ . Thus, we obtain

$$\begin{aligned} \|\tilde{\mathbf{f}}_i - \mathbf{f}_i\|_{n_i}^2 &\leq \|E(\tilde{\mathbf{f}}_i^{\text{reg}}) - \mathbf{f}_i\|_{n_i}^2 + \|\mathbf{B}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \lambda_1 \mathbf{D}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i)^{-1} \\ &\quad \mathbf{B}_i^T \mathbf{W}_i \mathbf{f}_i\|_{n_i}^2 + \|\mathbf{B}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i (\mathbf{1}_{n_i} b_i + \boldsymbol{\varepsilon}_i)\|_{n_i}^2 \\ &= I_1 + I_2 + I_3, \end{aligned}$$

where the definitions of  $I_j, j = 1, 2, 3$ , should be apparent from the context. Under Assumption A1-A3, when  $k \rightarrow \infty$  and  $k^4 = o(n_0)$ , then from Theorem 1 in Zhu, Fung, and He (2008), it is straightforward to show that  $I_1 = O_p(h^{2r})$ .

Furthermore, let  $\mathbf{F}_i = \frac{1}{n_i} \mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i$ ,  $\mathbf{H}_i = \mathbf{F}_i + \frac{\lambda_1}{n_i} \mathbf{D}_i$ . From Zhu, Fung, and He (2008) Lemma A1, we have  $\|\mathbf{F}_i^{-1}\|_\infty = O(h^{-1})$ . Similarly to Lemma 6.2 in Cardot (2000), it

follows that  $\|\mathbf{D}_i\|_\infty = O(h^{1-2d})$ . Since  $\frac{\lambda_1 h^{-2d}}{n_0} = o(1)$ , then  $\|\mathbf{H}_i^{-1}\|_\infty = \|\mathbf{F}_i^{-1} - (\mathbf{I}_p + \mathbf{F}_i^{-1} \frac{\lambda_1}{n_i} \mathbf{D}_i)^{-1} \mathbf{F}_i^{-1} \frac{\lambda_1}{n_i} \mathbf{D}_i \mathbf{F}_i^{-1}\|_\infty = O(h^{-1})$ . We can write  $I_2 = \|\mathbf{B}_i(\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \lambda_1 \mathbf{D}_i (\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i (\mathbf{f}_i - \mathbf{f}_i^s + \mathbf{f}_i^s)\|_{n_i}^2 = \|\frac{\lambda_1}{n_i} \mathbf{B}_i \mathbf{H}_i^{-1} \mathbf{D}_i \mathbf{F}_i^{-1} \frac{1}{n_i} \mathbf{B}_i^T \mathbf{W}_i (\mathbf{f}_i - \mathbf{f}_i^s) + \frac{\lambda_1}{n_i} \mathbf{B}_i \mathbf{H}_i^{-1} \mathbf{D}_i \mathbf{F}_i^{-1} \mathbf{B}_i^T \mathbf{W}_i \mathbf{f}_i^s / n_i\|_{n_i}^2$ .

According to the proof of Theorem 1 in Chen and Wang (2011), we can show that

$$\begin{aligned} -\frac{\lambda_1}{n_i} \|\mathbf{B}_i \mathbf{H}_i^{-1} \mathbf{D}_i \mathbf{F}_i^{-1} \mathbf{B}_i^T \mathbf{W}_i \mathbf{f}_i^s / n_i\|_\infty &= O(\lambda_1 n_i^{-1} h^{-d}), \\ -\frac{\lambda_1}{n_i} \|\mathbf{B}_i \mathbf{H}_i^{-1} \mathbf{D}_i \mathbf{F}_i^{-1} \frac{1}{n_i} \mathbf{B}_i^T \mathbf{W}_i (\mathbf{f}_i - \mathbf{f}_i^s)\|_\infty &= o(\lambda_1 n_i^{-1} h^{r-2d}). \end{aligned}$$

Then  $I_2 = O_p(\frac{\lambda_1^2}{n_i^2} h^{-2d})$ . Next, it can be shown that

$$\begin{aligned} &E[\mathbf{B}_i(\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \mathbf{B}_i^T \mathbf{W}_i (\mathbf{1}_{n_i} b_i + \boldsymbol{\varepsilon}_i)]^T [\mathbf{B}_i(\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i + \lambda_1 \mathbf{D}_i)^{-1} \\ &\mathbf{B}_i^T \mathbf{W}_i (\mathbf{1}_{n_i} b_i + \boldsymbol{\varepsilon}_i)] \\ &\leq \lambda_{\max}(\mathbf{W}_i(\sigma_b^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T + \boldsymbol{\Sigma}_i^0)) \text{tr}\{\mathbf{B}_i^T \mathbf{W}_i \mathbf{B}_i n_i^{-1} \mathbf{H}_i^{-1} \mathbf{B}_i^T \mathbf{B}_i n_i^{-1} \mathbf{H}_i^{-1}\} \\ &\leq \lambda_{\max}(\mathbf{W}_i(\sigma_b^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T + \boldsymbol{\Sigma}_i^0)) \lambda_{\max}(\mathbf{F}_i \mathbf{H}_i^{-1}) \lambda_{\max}(\mathbf{B}_i^T \mathbf{B}_i) \text{tr}\{n_i^{-1} \mathbf{H}_i^{-1}\} \\ &= O(n_i^{-1} h^{-1}). \end{aligned}$$

The above inequalities are obtained from Zhou, Shen, and Wolfe (1998) Lemma 6.5. Since  $\lambda_{\max}(\mathbf{W}_i(\sigma_b^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T + \boldsymbol{\Sigma}_i^0)) < C$ , and for any  $x \in [0, 1]$ ,  $0 \leq \pi_{im}(x) \leq 1$ , thus  $\text{tr}\{\mathbf{B}_i^T \mathbf{B}_i\} \leq 1$  and  $\lambda_{\max}(\mathbf{B}_i^T \mathbf{B}_i) \leq 1$ . Then  $I_3 = O_p(n_i^{-1} h^{-1})$ . Consequently, we can show that for all  $i = 1, \dots, n$ ,

$$\begin{aligned} \|\tilde{\mathbf{f}}_i - \mathbf{f}_i\|_{n_i}^2 &\leq O_p(k^{-2r}) + O_p\left(\frac{\lambda_1^2}{n_i^2} k^{2d}\right) + O_p\left(\frac{k}{n_i}\right) \\ &\leq O_p(k^{-2r}) + O_p\left(\frac{\lambda_1^2}{n_0^2} k^{2d}\right) + O_p\left(\frac{k}{n_0}\right). \end{aligned}$$

Then, it follows that

$$\|\tilde{\mathbf{f}} - \mathbf{f}\|_N^2 = \frac{1}{N} \sum_{i=1}^n (\tilde{\mathbf{f}}_i - \mathbf{f}_i)^T (\tilde{\mathbf{f}}_i - \mathbf{f}_i) \leq O_p(k^{-2r}) + O_p\left(\frac{\lambda_1^2}{n_0^2} k^{2d}\right) + O_p\left(\frac{k}{n_0}\right).$$

### S3 Proof of Lemma 2

When the true group memberships  $\mathcal{G}_1, \dots, \mathcal{G}_G$  are known, the oracle approximation for  $\mathbf{f}$  conditional on the oracle estimate of random effects  $\tilde{\mathbf{b}}^{\text{or}}$  is denoted as  $\tilde{\mathbf{f}}^{\text{or}} = \mathbf{B}\tilde{\boldsymbol{\beta}}^{\text{or}}$ , where

$$\begin{aligned}\tilde{\boldsymbol{\beta}}_{(g)}^{\text{or}} &= \arg \min_{i \in \mathcal{G}_g} Q_{(g)}(\boldsymbol{\beta}_{(g)}, \mathbf{b}) \\ &= \arg \min_{i \in \mathcal{G}_g} \frac{1}{2} \sum_{i=1}^{|\mathcal{G}_g|} \{(\mathbf{y}_i - \mathbf{B}_i \boldsymbol{\beta}_{(g)} - \mathbf{1}_{\mathbf{n}_i} b_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{y}_i - \mathbf{B}_i \boldsymbol{\beta}_{(g)} - \mathbf{1}_{\mathbf{n}_i} b_i) + \\ &\quad \lambda_1 \boldsymbol{\beta}_{(g)}^T \mathbf{D}_i \boldsymbol{\beta}_{(g)} + \lambda_2 b_i^2\}.\end{aligned}$$

Similarly to the proof of Lemma 1, we can show under Assumptions A1 – A5 and  $\frac{\lambda_1 h^{-2d}}{N_0} = o(1)$  that

$$\|\tilde{\mathbf{f}}^{\text{or}} - \mathbf{f}\|_N^2 \leq O_p(k^{-2r}) + O_p\left(\frac{\lambda_1^2}{N_0^2} k^{2d}\right) + O_p\left(\frac{k}{N_0}\right).$$

### S4 Proof of Theorem 1

By the triangular inequality,  $\|\hat{\mathbf{f}} - \mathbf{f}\|_N^2 = \|\hat{\mathbf{f}} - \tilde{\mathbf{f}}^{\text{or}} + \tilde{\mathbf{f}}^{\text{or}} - \mathbf{f}\|_N^2 \leq \|\hat{\mathbf{f}} - \tilde{\mathbf{f}}^{\text{or}}\|_N^2 + \|\tilde{\mathbf{f}}^{\text{or}} - \mathbf{f}\|_N^2$ .

The proposed objective function is

$$\begin{aligned}H(\boldsymbol{\beta}, \mathbf{b}) &= \frac{1}{2}(\mathbf{Y} - \mathbf{B}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})^T \boldsymbol{\Sigma}^{-1}(\mathbf{Y} - \mathbf{B}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}) + \frac{1}{2} \lambda_1 \boldsymbol{\beta}^T \mathbf{D}_d \boldsymbol{\beta} + \frac{1}{2} \lambda_2 \|\mathbf{b}\|_2^2 + \\ &\quad \sum_{i,j \in \mathcal{L}} \rho(|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j|, \lambda_3).\end{aligned}\tag{S4.1}$$

We can obtain the estimate of random effects  $\mathbf{b}$  by minimizing (S4.1), then

$$\hat{\mathbf{b}} = (\mathbf{Z}^T \boldsymbol{\Sigma}^{-1} \mathbf{Z} + \lambda_2 \mathbf{I}_n)^{-1} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{B}\boldsymbol{\beta}).$$

Replacing  $\mathbf{b}$  in (S4.1) by  $\hat{\mathbf{b}}$ , we can obtain the profiled objective function

$$\begin{aligned}H_n(\boldsymbol{\beta}) &= \frac{1}{2}(\mathbf{Y} - \mathbf{B}\boldsymbol{\beta})^T \mathbf{W}(\mathbf{Y} - \mathbf{B}\boldsymbol{\beta}) + \frac{1}{2} \lambda_1 \boldsymbol{\beta}^T \mathbf{D}_d \boldsymbol{\beta} + \sum_{i,j \in \mathcal{L}} \rho(|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j|, \lambda_3) \\ &= Q_n(\boldsymbol{\beta}) + S_{\lambda_3}(\boldsymbol{\beta}),\end{aligned}\tag{S4.2}$$

where  $\mathbf{W} = (\boldsymbol{\Sigma} + \frac{1}{\lambda_2} \mathbf{Z}\mathbf{Z}^T)^{-1}$ ,  $\mathbf{W}_i = (\boldsymbol{\Sigma}_i + \frac{1}{\lambda_2} \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T)^{-1}$ ,  $Q_n(\boldsymbol{\beta}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{B}_i \boldsymbol{\beta}_i)^T \mathbf{W}_i (\mathbf{y}_i - \mathbf{B}_i \boldsymbol{\beta}_i) + \frac{1}{2} \sum_{i=1}^n \lambda_1 \boldsymbol{\beta}_i^T \mathbf{D}_i \boldsymbol{\beta}_i$  and  $S_{\lambda_3}(\boldsymbol{\beta}) = \sum_{i,j \in \mathcal{L}} \rho(|\boldsymbol{\beta}_i - \boldsymbol{\beta}_j|, \lambda_3)$ . Denote

$$\mathbf{D}_n^s = \partial^2 Q_n(\boldsymbol{\beta}) / \partial \boldsymbol{\beta} \partial \boldsymbol{\beta} = \mathbf{B}^T \mathbf{W} \mathbf{B} + \lambda_1 \mathbf{D}_d, \quad (\text{S4.3})$$

$$\begin{aligned} \mathbf{M}_n^s &= \text{Cov}(\partial Q_n(\boldsymbol{\beta}) / \partial \boldsymbol{\beta}) \\ &= \sum_{i=1}^n \mathbf{B}_i^T \mathbf{W}_i (\sigma_b^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T + \boldsymbol{\Sigma}_i^0) \mathbf{W}_i^T \mathbf{B}_i. \end{aligned} \quad (\text{S4.4})$$

Let  $\tau_n^s = \lambda_{\min}(\mathbf{D}_n^s (\mathbf{M}_n^s)^{-1} \mathbf{D}_n^s)$ , and  $\mathbf{B}\boldsymbol{\beta}^* = \mathbf{B}\tilde{\boldsymbol{\beta}}^{or} + (\tau_n^s)^{-\frac{1}{2}} \mathbf{u}$ , where  $\|\mathbf{u}\|_N = d_u$ .

Note that

$$\begin{aligned} (\tau_n^s)^{-1} &= \lambda_{\max}((\mathbf{D}_n^s)^{-1} (\mathbf{M}_n^s) (\mathbf{D}_n^s)^{-1}) \\ &\leq C_1 \lambda_{\max} \{ E[\mathbf{B}(\mathbf{B}^T \mathbf{W} \mathbf{B} + \lambda_1 \mathbf{D}_d)^{-1} \mathbf{B}^T \mathbf{W} \mathbf{Y} - E(\mathbf{B}\tilde{\boldsymbol{\beta}})]^T [\mathbf{B}(\mathbf{B}^T \mathbf{W} \mathbf{B} \\ &\quad + \lambda_1 \mathbf{D}_d)^{-1} \mathbf{B}^T \mathbf{W} \mathbf{Y} - E(\mathbf{B}\tilde{\boldsymbol{\beta}})] \} \\ &= C_1 E(\tilde{\mathbf{f}} - \mathbf{f}^s)^T (\tilde{\mathbf{f}} - \mathbf{f}^s) \\ &\leq C_1 N \|\tilde{\mathbf{f}} - \mathbf{f}\|_N^2 + C_1 N \|\mathbf{f} - \mathbf{f}^s\|_N^2, \end{aligned} \quad (\text{S4.5})$$

where the first inequality in (S4.5) can be derived from Lemma A.4 in Zhu and Qu (2018).

Thus, from Lemma 1,  $(\tau_n^s)^{-1} \leq A_1 (\frac{k}{n_0} + \frac{\lambda_1^2}{n_0^2} k^{2d} + k^{-2r})$  with  $A_1$  sufficiently large.

If  $N$  is sufficiently large, we have

$$\begin{aligned} \|\mathbf{B}(\tilde{\boldsymbol{\beta}}_{(g)}^{or} - \tilde{\boldsymbol{\beta}}_{(g')}^{or})\|_N &= \|\tilde{\mathbf{f}}_{(g)}^{or} - \tilde{\mathbf{f}}_{(g')}^{or}\|_N \\ &= \|\mathbf{f}_{(g)} - \mathbf{f}_{(g')} - (\mathbf{f}_{(g)} - \tilde{\mathbf{f}}_{(g)}^{or}) + (\mathbf{f}_{(g')} - \tilde{\mathbf{f}}_{(g')}^{or})\|_N \\ &\geq \|\mathbf{f}_{(g)} - \mathbf{f}_{(g')}\|_N - \|\mathbf{f}_{(g)} - \tilde{\mathbf{f}}_{(g)}^{or}\|_N - \|\mathbf{f}_{(g')} - \tilde{\mathbf{f}}_{(g')}^{or}\|_N \\ &\geq d_f. \end{aligned}$$

It is easy to prove that there exists a constant  $c$ , such that  $\|\tilde{\boldsymbol{\beta}}_{(g)}^{or} - \tilde{\boldsymbol{\beta}}_{(g')}^{or}\|_n \geq cd_f$ . From the definition of  $\rho_\tau(t, \lambda_3) = 0$ , when  $t = 0$  and  $\rho_\tau(t, \lambda_3) \geq 0$ , when  $t \neq 0$ , and since the minimum

distance  $d_f$  satisfies  $cd_f \geq \tau\lambda_3$ , we have  $S_{\lambda_3}(\tilde{\boldsymbol{\beta}}^{or}) = 0$ . By Taylor's expansion, we can obtain that

$$\begin{aligned}
L_n(u) &= H_n(\boldsymbol{\beta}^*) - H_n(\tilde{\boldsymbol{\beta}}^{or}) \\
&= Q_n(\boldsymbol{\beta}^*) - Q_n(\tilde{\boldsymbol{\beta}}^{or}) + S_{\lambda_3}(\boldsymbol{\beta}^*) \\
&= (\tau_n^s)^{-\frac{1}{2}} \dot{Q}_n^T(\tilde{\boldsymbol{\beta}}^{or}) \mathbf{u} + \frac{1}{2} (\tau_n^s)^{-1} \mathbf{u}^T \ddot{Q}_n(\tilde{\boldsymbol{\beta}}^{or}) \mathbf{u} + S_{\lambda_3}(\boldsymbol{\beta}^*) \\
&= (\tau_n^s)^{-\frac{1}{2}} \dot{Q}_n^T(\tilde{\boldsymbol{\beta}}^{or}) \mathbf{u} + \frac{1}{2} (\tau_n^s)^{-1} \mathbf{u}^T \mathbf{D}_n^s \mathbf{u} + S_{\lambda_3}(\boldsymbol{\beta}^*), \tag{S4.6}
\end{aligned}$$

where  $\dot{Q}_n(\tilde{\boldsymbol{\beta}}^{or}) = \frac{\partial}{\partial \boldsymbol{\beta}} Q(\boldsymbol{\beta})|_{\boldsymbol{\beta}=\tilde{\boldsymbol{\beta}}^{or}}$ ,  $\ddot{Q}_n(\tilde{\boldsymbol{\beta}}^{or}) = \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}} Q(\boldsymbol{\beta})|_{\boldsymbol{\beta}=\tilde{\boldsymbol{\beta}}^{or}}$ . Note that

$$\begin{aligned}
(\mathbf{M}_n^s)^{\frac{1}{2}} &= (\mathbf{D}_n^s)^{\frac{1}{2}} (\mathbf{D}_n^s)^{-\frac{1}{2}} (\mathbf{M}_n^s)^{\frac{1}{2}} (\mathbf{D}_n^s)^{-\frac{1}{2}} (\mathbf{D}_n^s)^{\frac{1}{2}} \\
&\leq (\mathbf{D}_n^s)^{\frac{1}{2}} \lambda_{max}((\mathbf{D}_n^s)^{-\frac{1}{2}} (\mathbf{M}_n^s)^{\frac{1}{2}} (\mathbf{D}_n^s)^{-\frac{1}{2}}) (\mathbf{D}_n^s)^{\frac{1}{2}} \\
&= \lambda_{min}((\mathbf{D}_n^s)^{\frac{1}{2}} (\mathbf{M}_n^s)^{-\frac{1}{2}} (\mathbf{D}_n^s)^{\frac{1}{2}})^{-1} \mathbf{D}_n^s \\
&= (\tau_n^s)^{-\frac{1}{2}} \mathbf{D}_n^s, \tag{S4.7}
\end{aligned}$$

and thus  $(\tau_n^s)^{-\frac{1}{2}} (\mathbf{M}_n^s)^{\frac{1}{2}} \leq (\tau_n^s)^{-1} \mathbf{D}_n^s$ . Consequently, if  $d_u$  is sufficiently large, then the second term in  $L_n(u)$  dominates the first term, which implies that, with probability tending to 1,  $L_n(u) > 0$  at  $\|\mathbf{u}\|_N = d_u$ . Hence we have

$$P\left\{ \inf_{\|\mathbf{u}\|_N = d_u} L_n(u) > 0 \right\} \rightarrow 1,$$

which entails that with probability tending to 1, there exists a local minimum of  $H_n(\boldsymbol{\beta})$  which lies in the ball  $\mathcal{B} = \{\boldsymbol{\beta} : \|\mathbf{B}(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}^{or})\|_N^2 = (\tau_n^s)^{-1} d_u^2 = A_2(\frac{k}{n_0} + \frac{\lambda_1^2}{n_0^2} k^{2d} + k^{-2r})\}$  with  $A_2$  sufficiently large. And  $\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} H_n(\boldsymbol{\beta})$ , then  $\|\mathbf{B}(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}^{or})\|_n^2 = O_p(\frac{k}{n_0} + \frac{\lambda_1^2}{n_0^2} k^{2d} + k^{-2r})$ . Thus by combining the result from Lemma 2, we can complete the proof.

## S5 Proof of Corollary 1

From Theorem 1, there exists a local minimizer  $\hat{\mathbf{f}} = \mathbf{B}\hat{\boldsymbol{\beta}}$ , where  $\hat{\boldsymbol{\beta}} \in \mathcal{B}$ . For any pair  $i, j$  such that  $G(i) = G(j)$ , we have  $\|\hat{\mathbf{f}}_i - \hat{\mathbf{f}}_j\|_N^2 = \|\hat{\mathbf{f}}_i - \mathbf{f}_i + \mathbf{f}_i - \mathbf{f}_j + \mathbf{f}_j - \hat{\mathbf{f}}_j\|_N^2 \leq 2 \max_i \|\hat{\mathbf{f}}_i - \mathbf{f}_i\|_N^2 + \|\mathbf{f}_i - \mathbf{f}_j\|_N^2 \leq O_p\left(\frac{k}{n_0} + \frac{\lambda_1^2}{n_0^2} k^{2d} + k^{-2r}\right) \rightarrow O_p(k^{-2r})$ , as  $n_0 \rightarrow \infty$ . Thus  $\hat{\mathbf{f}}_i$  and  $\hat{\mathbf{f}}_j$  will be in the same group with probability tending to 1. On the other side, for any pair  $i, j$  such that  $G(i) \neq G(j)$ , we have  $\|\hat{\mathbf{f}}_i - \hat{\mathbf{f}}_j\|_N^2 \geq \min \|\mathbf{f}_i - \mathbf{f}_j\|_N^2 - 2 \max_i \|\hat{\mathbf{f}}_i - \mathbf{f}_i\|_N^2 \geq d_f^2 - O_p(k^{-2r})$  as  $n_0 \rightarrow \infty$ , which indicates that  $\hat{\mathbf{f}}_i$  and  $\hat{\mathbf{f}}_j$  will be in different groups with probability tending to 1. Hence,  $P(\hat{\mathcal{G}} = \mathcal{G}) \rightarrow 1$ .

## References

- Cardot, H. (2000). Nonparametric estimation of smoothed principal components analysis of sampled noisy functions. *Journal of Nonparametric Statistics* 12, pp. 503-538.
- Chen, H. H. and Wang Y. J. (2011). A penalized spline approach to functional mixed effects model analysis. *Biometrics* 67, pp. 861-870.
- Zhou, S., Shen, X. and Wolfe, D. A. (1998). Local asymptotics for regression splines and confidence regions. *The Annals of Statistics* 26, pp. 1760-1782.
- Zhu, X. L. and Qu, A. (2018). Cluster analysis of longitudinal profiles with subgroups. *Electronic Journal of Statistics* 12, pp. 171-193.
- Zhu, Z. Y., Fung, W. K. and He, X. M. (2008). On the asymptotics of marginal regression splines with longitudinal data. *Biometrika* 95, pp. 907-917.



Table S1: Case3: Comparison results from the proposed nonparametric pairwise-grouping with three different working correlation structures (NPGGr-IN, NPGGr-AR(1), NPGGr-Ex), Gaussian Mixtures (bGM), K-means (bKmeans), SSClust, MixedEffects, and Kernel for balanced data.

	Methods	$\hat{K}$	Rand	aRand	Jaccard	AMSE
AR(1)	NPGGr-IN	3.02	0.9998	0.9995	0.9993	0.0420
	NPGGr-AR(1)	3.00	1.0000	1.0000	1.0000	0.0375
	NPGGr-Ex	3.00	1.0000	1.0000	1.0000	0.0380
	bGM	3.17	0.9941	0.9846	0.9819	0.0463
	bKmeans	3.00	0.9302	0.8606	0.8739	2.5818
	SSClust	7.98	0.8301	0.5504	0.4799	0.3549
	MixedEffects	4.68	0.9406	0.8555	0.8182	0.1199
	Kernel	4.24	0.9481	0.8743	0.8412	0.4984
Ex	NPGGr-IN	3.10	0.9990	0.9978	0.9970	0.0476
	NPGGr-AR(1)	3.00	1.0000	1.0000	1.0000	0.0369
	NPGGr-Ex	3.00	1.0000	1.0000	1.0000	0.0372
	bGM	3.06	0.9986	0.9968	0.9957	0.0420
	bKmeans	3.00	0.9360	0.8721	0.8842	2.3726
	SSClust	7.69	0.8297	0.5483	0.4786	0.2954
	MixedEffects	4.22	0.9524	0.8849	0.8544	0.0958
	Kernel	4.37	0.9433	0.8621	0.8264	0.5371

Table S2: Case4: Comparison results from the proposed nonparametric pairwise-grouping with three different working correlation structures (NPGGr-IN, NPGGr-AR(1), NPGGr-Ex), Gaussian Mixtures (bGM), K-means (bKmeans), SSClust, MixedEffects, and Kernel for balanced data.

	Methods	$\hat{K}$	Rand	aRand	Jaccard	AMSE
AR(1)	NPGGr-IN	3.70	0.9925	0.9826	0.9770	0.1601
	NPGGr-AR(1)	3.02	0.9998	0.9995	0.9993	0.0635
	NPGGr-Ex	3.03	0.9996	0.9992	0.9989	0.0679
	bGM	3.76	0.9788	0.9456	0.9351	0.1542
	bKmeans	3.00	0.9217	0.8438	0.8576	2.9920
	SSClust	6.42	0.8745	0.6782	0.6157	0.2549
	MixedEffects	5.18	0.9309	0.8323	0.7889	0.7447
	Kernel	5.29	0.9365	0.8454	0.8055	0.9188
Ex	NPGGr-IN	3.71	0.9925	0.9826	0.9769	0.1529
	NPGGr-AR(1)	3.06	0.9993	0.9984	0.9979	0.0800
	NPGGr-Ex	3.03	0.9996	0.9992	0.9989	0.0651
	bGM	3.57	0.9876	0.9685	0.9619	0.1383
	bKmeans	3.00	0.9260	0.8522	0.8654	2.8116
	SSClust	6.74	0.8600	0.6377	0.5713	0.3023
	MixedEffects	5.18	0.9258	0.8186	0.7730	0.6701
	Kernel	5.08	0.9353	0.8423	0.8019	0.9801

Table S3: Case5: Comparison results from the proposed nonparametric pairwise-grouping with three different working correlation structures (NPGGr-IN, NPGGr-AR(1), NPGGr-Ex), Gaussian Mixtures (bGM), K-means (bKmeans), SSClust, MixedEffects, and Kernel for balanced data.

	Methods	$\hat{K}$	Rand	aRand	Jaccard	AMSE
AR(1)	NPGGr-IN	3.00	1.0000	1.0000	1.0000	0.0329
	NPGGr-AR(1)	3.00	1.0000	1.0000	1.0000	0.0327
	NPGGr-Ex	3.00	1.0000	1.0000	1.0000	0.0328
	bGM	3.17	0.9965	0.9916	0.9892	0.0545
	bKmeans	3.00	0.9321	0.8639	0.8760	2.5680
	SSClust	8.40	0.8170	0.5093	0.4397	0.2825
	MixedEffects	4.68	0.9439	0.8632	0.8283	0.1128
	Kernel	3.94	0.9593	0.9015	0.8753	0.4525
Ex	NPGGr-IN	3.00	1.0000	1.0000	1.0000	0.0334
	NPGGr-AR(1)	3.00	1.0000	1.0000	1.0000	0.0333
	NPGGr-Ex	3.00	1.0000	1.0000	1.0000	0.0333
	bGM	3.13	0.9977	0.9946	0.9930	0.0510
	bKmeans	3.02	0.9284	0.8561	0.8682	2.6681
	SSClust	8.26	0.8082	0.4837	0.4127	0.2573
	MixedEffects	3.92	0.9769	0.9456	0.9294	0.0644
	Kernel	4.01	0.9582	0.8992	0.8720	0.4672