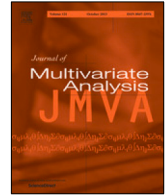




Contents lists available at ScienceDirect

## Journal of Multivariate Analysis

journal homepage: [www.elsevier.com/locate/jmva](http://www.elsevier.com/locate/jmva)

## Semiparametric penalized quadratic inference functions for longitudinal data in ultra-high dimensions

Brittany Green<sup>a</sup>, Heng Lian<sup>b</sup>, Yan Yu<sup>c,\*</sup>, Tianhai Zu<sup>d</sup><sup>a</sup> Department of Information Systems, Analytics, and Operations, University of Louisville, Louisville, KY, USA<sup>b</sup> Department of Mathematics, City University of Hong Kong, Hong Kong<sup>c</sup> Department of Operations, Business Analytics, & Information Systems, University of Cincinnati, Cincinnati, OH, USA<sup>d</sup> Department of Management Science and Statistics, University of Texas at San Antonio, San Antonio, TX, USA

## ARTICLE INFO

## Article history:

Received 8 January 2022

Received in revised form 3 March 2023

Accepted 5 March 2023

Available online 23 March 2023

## AMS 2020 subject classifications:

primary 62G08

secondary 62G20

## Keywords:

Longitudinal data

Model selection

Multivariate correlated response

Partially linear model

Single-index model

## ABSTRACT

In many biomedical and health studies, multivariate data arise from repeated measurements on a sample of subjects over time. In order to analyze such longitudinal data, we need to consider the correlations from the same subject, and it is inappropriate to use a simple multivariate model assuming independence structure. Motivated by a large scale longitudinal public health study that requires longitudinal data analysis with correlated multivariate discrete responses from repeated measurements and very high dimensional covariates, we adopt a flexible semiparametric approach for simultaneous variable selection and estimation without the requirement of specifying the full likelihood. Specifically, we propose generalized partially linear single-index models using penalized quadratic inference functions for longitudinal data in ultra-high dimensions. A key feature is that we allow the number of single-index covariates in the nonparametric term to diverge and even to be in ultra-high dimensions. The penalized quadratic inference functions easily incorporate within-subject correlation and pursue efficient estimation, and the single-index models can incorporate nonlinearity and some interactions while avoiding the curse of dimensionality. In this challenging setting, we contribute both an efficient algorithm and new asymptotic theory for our proposed approach for diverging and even ultra-dimensional covariates and a multivariate correlated response in longitudinal data. We apply our method to investigate diabetes status within a continuing longitudinal public health study with very high-dimensional genetic variables and phenotype variables.

© 2023 Elsevier Inc. All rights reserved.

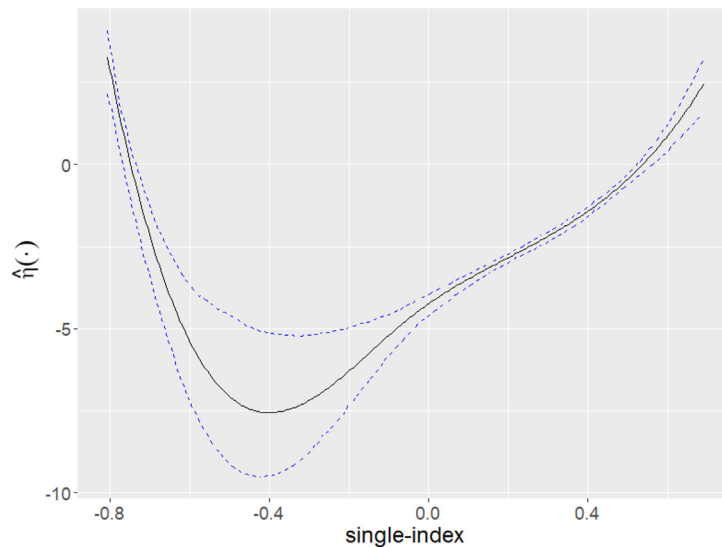
## 1. Introduction

Numerous large scale public health research studies are longitudinal, where participants have repeated measurements taken over time. To analyze such kind of multivariate data that exhibit clear correlation among within-subject responses, we need to incorporate the correlation structure rather than assuming a simple multivariate regression model under independence.

One main goal of analyzing these types of longitudinal studies is to identify the genetic and phenotype factors related to a disease to provide insight into more effective treatment and disease prevention strategies. For example, researchers have discovered risk factors linked to various disease mechanisms (e.g., Meigs et al. [23]) using the Framingham data,

\* Corresponding author.

E-mail address: [yan.yu@uc.edu](mailto:yan.yu@uc.edu) (Y. Yu).



**Fig. 1.** Curve estimates for real data application to diabetes analysis. We apply penalized quadratic inference function approach to generalized partially linear single-index model for longitudinal data in high dimensions. Here the response variable of interest is diabetes status, a binary correlated discrete response, and all continuous phenotype and genotype variables are embedded in the single-index term. Blue dot curves are the 95% confidence intervals.

an ongoing large scale multi-generational health study [6]. Within these types of large-scale longitudinal studies, while the true correlation among participants is usually difficult to uncover, incorporating within participant dependence can lead to more efficient estimation. Moreover, the disease of interest measured over time is sometimes a correlated discrete measure, such as diabetes status. This non-normal correlated response creates a challenge in specifying the joint likelihood for longitudinal data.

In addition to the longitudinal nature of large scale public health studies, more recently, the genotype of participants is also collected. These genetic factors are very high dimensional as demonstrated in the Framingham data which collects a complex array of genetic data, including tens of thousands of single nucleotide polymorphisms (SNPs) from each participant. Notably, previous research has linked some genetic factors to disease. For instance, genomic studies have helped identify mechanisms of hypertension and diabetes [39]. In these very high dimensional settings, variable selection is imperative to identify the important risk factors, since usually only a few covariates relate to the response and including non-important variables lessens estimation efficiency and impedes inference. In addition, not only do these types of studies have high dimensional genetic data, recent research has found that genetic data interacts with phenotype variables. For example, Taylor et al. [31] show that the genetic effects of hypertension are altered under phenotype factors such as BMI.

A motivating example of this paper to exemplify this complexity focuses on identifying factors related to diabetes status, a correlated discrete response, from the offspring cohort of the longitudinal Framingham data. Diabetes affects millions of people worldwide, and identifying genetic and phenotype factors that relate to diabetes can help inform preventative measures and further uncover biological measures of diabetes [10]. In particular, one research question of interest is to determine which SNPs in high dimensions and phenotype factors relate to diabetes status among participants over time. While a traditional approach to this problem employs a linear model, due to the complexity of gene expression and interactions with phenotype factors, inflexible models with parametric assumptions may not account for the potential nonlinearity and synergy between genetic and phenotype data in high dimensional longitudinal data. To demonstrate, Fig. 1 shows a clear nonlinear relationship under our proposed model between diabetes and the combination of SNPs and phenotype factors in the Framingham data.

To balance interpretability and flexibility for accurate estimation for this high dimensional set of risk factors, we consider a flexible semiparametric approach for repeated observations. Specifically, we adopt generalized partially linear single-index models (GPLSIM) [4] for longitudinal data in ultra-high dimensions. Generalized partially linear single-index models achieve dimension reduction by reducing the high-dimensional predictors to a univariate index within a flexible function. Moreover, single-index models can capture some interactions among covariates as opposed to additive models. This is advantageous since SNPs and phenotypic risk factors do not relate to disease in isolation: the compound impact of the genotype–phenotype interaction has been shown to outperform the impact of using the conventional risk factors in isolation (e.g., Franks [10], Taylor et al. [31]).

Moreover, diabetes status, the outcome of interest measured during multiple waves of the Framingham data, is a correlated discrete response. This poses a challenge as the full joint likelihood can be intractable for correlated discrete

data. To tackle this challenge, we employ the penalized quadratic inference function (QIF) to account for within-subject correlation, perform model selection, and seek efficient estimation for diverging and potentially ultra-high dimensional longitudinal data. Previous research employing penalized generalized estimating equations (GEE) for diverging and ultra-high dimensional longitudinal data for linear and semiparametric models includes Wang et al. [38] and Green et al. [12]. However, generalized estimating equations are known to be less efficient and overfit the model compared to the quadratic inference function approach when the working correlation matrix is misspecified (e.g., Cho and Qu [5], Qu et al. [27]). In addition, the quadratic inference function has applicability in various model setups and shows promising results (e.g., Qu and Li [26], Wang et al. [36,37]).

As a result of the multiple benefits of the quadratic inference function, a handful of works have proposed research employing partially linear single-index models using the quadratic inference function for longitudinal data (e.g., Bai et al. [2] and Lai et al. [18]). However, these works assume the dimension of both the single-index and partially linear covariates is fixed. Also in fixed finite dimensions, Ma et al. [21] and Li et al. [19] incorporated variable selection for the partially linear single-index model employing the quadratic inference function for continuous longitudinal responses with the identity link function. In these studies, the real-data applications considered a fixed low-to-moderate dimensional setup: the application in Ma et al. [21] included a total of 11 covariates, and the application in Li et al. [19] included 13 covariates. In contrast, the Framingham data analyzed using our approach have 878 participants but involves over 500,000 SNP covariates and phenotype variables even within the nonparametric portion.

As opposed to the previous works, our approach allows the number of covariates in both the nonparametric and linear components of the generalized partially linear single-index model to diverge and even in ultra-high dimensions. We also allow the number of important covariates to diverge. This is especially pertinent for our motivating example, since there are more than 500,000 genetic SNP variables. In particular, allowing flexible modeling and determining the sparse set of important covariates can lead to more accurate estimation. However, as a result of incorporating diverging covariates, we encounter added challenges in both computation and theory when we allow ultra-high dimensional data within the nonlinear, unknown, flexible function with potentially diverging support of the single index. We establish asymptotic theory for ultra-dimensional covariates for model selection and estimation including the oracle property, which is much more challenging to establish than fixed-dimensional theory. In addition, while previous approaches exist for estimating the coefficients of the generalized partially linear single-index model, efficient estimation becomes even more difficult when introducing ultra-high dimensional correlated data. This is because of the potentially high dimensional covariates in a nonlinear unknown function estimated nonparametrically together with the non-convex smoothly clipped absolute deviation (SCAD) penalty function as in Fan and Li [7], all within a longitudinal framework. Therefore, to select the sparse set of important covariates and estimate the corresponding coefficients, we provide a computationally efficient iterative algorithm. This approach implements strategic approximations to reduce the computational burden and increase the effectiveness of the algorithm, even for discrete correlated responses.

## 2. Quadratic inference function for semiparametric longitudinal data analysis

### 2.1. Model

For each subject  $i$  with  $i \in \{1, \dots, n\}$ , assume correlation among its observations over time  $t \in \{1, \dots, T_i\}$ , but independence from other subjects. We observe, for subject  $i$ , a correlated multivariate response vector  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{it}, \dots, Y_{iT_i})$  over repeated measurements over time. We also observe at each time  $t$ , the  $p_n \times 1$  dimensional single-index covariate vector  $\mathbf{X}_{it} = (X_{it,1}, \dots, X_{it,p_n})^\top$  and the  $q_n \times 1$  dimensional linear covariate vector  $\mathbf{Z}_{it} = (Z_{it,1}, \dots, Z_{it,q_n})^\top$ . Notably, we allow both the number of single-index and partially linear covariates,  $p_n$  and  $q_n$ , to diverge and potentially be in the exponential order. Also, the number of observations  $T_i$  can be different for each subject  $i$  for imbalanced observations.

We consider generalized partially linear single-index models for longitudinal data in ultra-high dimensions to allow flexibility while avoiding the curse-of-dimensionality, that is,

$$E(Y_{it} | \mathbf{X}_{it}, \mathbf{Z}_{it}) = \mu_{it} = g^{-1}(\eta(\mathbf{X}_{it}^\top \boldsymbol{\beta}_0) + \mathbf{Z}_{it}^\top \boldsymbol{\gamma}_0), \quad i \in \{1, \dots, n\}; t \in \{1, \dots, T_i\}. \quad (1)$$

Here  $\eta(\cdot)$  is an unknown flexible function estimated non-parametrically; and  $g(\cdot)$  is a link function. Note that the responses can be modeled from a general exponential family including special cases such as Gaussian or Binomial for either continuous or discrete responses. The coefficient vector for the single-index covariates is  $\boldsymbol{\beta}_0 = (\beta_{01}, \dots, \beta_{0p_n})^\top$  and the coefficient vector for the linear covariates is  $\boldsymbol{\gamma}_0$  with  $\boldsymbol{\gamma}_0 = (\gamma_{01}, \dots, \gamma_{0q_n})^\top$ . We assume a sparse set of important covariates, which is common in modern statistics literature [14]. In particular, there is a nonzero subset of  $p_{sn}$  coefficients from the total  $p_n$  coefficients and a subset of  $q_{sn}$  nonzero coefficients from the total  $q_n$  linear coefficients, where the rest of the coefficients are zero. For model identifiability, we assume  $\|\boldsymbol{\beta}_0\| = 1$  with the first component positive [40].

We estimate the flexible univariate function of the conditional mean  $\eta(\cdot)$  non-parametrically using polynomial splines. We first assume the support of the single-index  $\mathbf{X}_{it}^\top \boldsymbol{\beta}_0$  is  $[a, b]$ . We note that the length of this support can be diverging due to the potentially ultra-high dimensional covariates, thus practically, we use the support of  $\mathbf{X}_{it}^\top \boldsymbol{\beta}_0$  based on a given  $\boldsymbol{\beta}$ . We then divide this support based on  $H'$  interior knots to create subintervals  $[c_k, c_{k+1}]$ , where  $k \in \{0, \dots, H'\}$  is determined by the partition,  $a = c_0 < c_1 < \dots < c_{H'} < c_{H'+1} = b$ . Given we approximate  $\eta(\cdot)$  with a degree  $s \geq 2$  polynomial over each interval, a polynomial spline of order  $s$  is a  $s - 1$  degree polynomial on each interval and globally

$s - 2$  times differentiable [29]. We consider the nonparametric estimation of the unknown flexible function  $\eta(\cdot)$  with the single-index  $u_{it} = \mathbf{X}_{it}^\top \boldsymbol{\beta}_0$  as a linear combination of B-spline bases  $\eta(u_{it}) \approx \mathbf{G}^\top(u_{it})\boldsymbol{\theta}$ , and  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_H)^\top$  are the basis coefficients of size  $H \equiv H_n = 1 + s + H'$ . Accordingly, since  $\eta(\mathbf{X}_{it}^\top \boldsymbol{\beta})$  is estimated by  $\mathbf{G}^\top(\mathbf{X}_{it}^\top \boldsymbol{\beta})\boldsymbol{\theta}$ , the conditional mean  $\mu_{it} = g^{-1}(\eta(\mathbf{X}_{it}^\top \boldsymbol{\beta}_0) + \mathbf{Z}_{it}^\top \boldsymbol{\gamma}_0)$  now becomes  $g^{-1}(\mathbf{G}^\top(\mathbf{X}_{it}^\top \boldsymbol{\beta}_0)\boldsymbol{\theta}_0 + \mathbf{Z}_{it}^\top \boldsymbol{\gamma}_0)$ . Therefore, in the remainder of the paper, we pursue estimation of the column coefficient vector  $\boldsymbol{\alpha}_0 = (\boldsymbol{\theta}_0^\top \boldsymbol{\beta}_0^\top \boldsymbol{\gamma}_0^\top)^\top$ , the spline coefficients for approximating the nonparametric function  $\eta$ , the single-index coefficients, and the partially linear coefficients, respectively.

## 2.2. Quadratic inference function and estimation

To incorporate the correlation among a participant's observations in longitudinal data, the quadratic inference function (QIF) approach may be used to estimate the spline basis, single-index, and partially linear coefficient vector,  $\boldsymbol{\alpha}$ , for the given covariates. The QIF from Qu et al. [27] substitutes a linear combination of basis matrices in place of the inverse of the working correlation matrix,  $\mathbf{R}^{-1} \approx a_1 \mathbf{M}_1 + \dots + a_m \mathbf{M}_m$ . Here  $\mathbf{M}_1, \dots, \mathbf{M}_m$  are predetermined, known symmetric matrices and  $\mathbf{a} = (a_1, \dots, a_m)^\top$  are constant coefficients. For most of the common correlation structures, available linear combinations of basis matrices exist to approximate  $\mathbf{R}^{-1}$  as further discussed in Section 5.2.

To estimate the coefficient vector  $\boldsymbol{\alpha}_0$ , Qu et al. [27] extend generalized estimating equations from Zeger and Liang [42] by adopting the generalized method of moments estimator from Hansen [13] to minimize the following quadratic inference function

$$Q_n(\boldsymbol{\alpha}) = \mathbf{g}_n^\top \mathbf{W}_n^{-1} \mathbf{g}_n, \quad (2)$$

where the extended score vector is

$$\mathbf{g}_n = \frac{1}{n} \sum_i \mathbf{g}_i(\boldsymbol{\alpha}) = \frac{1}{n} \begin{Bmatrix} \sum_i \mathbf{V}_i^\top(\boldsymbol{\alpha}) \mathbf{A}_i(\boldsymbol{\alpha})^{1/2} \mathbf{M}_1 \mathbf{A}_i(\boldsymbol{\alpha})^{-1/2} (\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \\ \sum_i \mathbf{V}_i^\top(\boldsymbol{\alpha}) \mathbf{A}_i(\boldsymbol{\alpha})^{1/2} \mathbf{M}_2 \mathbf{A}_i(\boldsymbol{\alpha})^{-1/2} (\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \\ \vdots \\ \sum_i \mathbf{V}_i^\top(\boldsymbol{\alpha}) \mathbf{A}_i(\boldsymbol{\alpha})^{1/2} \mathbf{M}_m \mathbf{A}_i(\boldsymbol{\alpha})^{-1/2} (\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \end{Bmatrix},$$

and  $\mathbf{W}_n = (1/n) \sum_i \mathbf{g}_i(\boldsymbol{\alpha}) \mathbf{g}_i(\boldsymbol{\alpha})^\top$ . The resulting estimates are determined as  $\hat{\boldsymbol{\alpha}} = \arg \min_{\boldsymbol{\alpha}} Q_n(\boldsymbol{\alpha})$ , and  $Q_n(\boldsymbol{\alpha})$  is known as the quadratic inference function [27]. For  $\mathbf{g}_n$  in Eq. (2), we introduce matrix notation for subject  $i$ , where  $\mathbf{Y}_i$  denotes the  $T_i \times 1$  response vector, and  $\boldsymbol{\mu}_i(\boldsymbol{\alpha}) = g^{-1}(\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta})\boldsymbol{\theta} + \mathbf{Z}_i \boldsymbol{\gamma})$  is the spline approximated marginal mean, where  $\mathbf{X}_i$  is the  $T_i \times p_n$  dimensional single-index covariate matrix,  $\mathbf{G}$  is the corresponding  $T_i \times H$  spline basis, and  $\mathbf{Z}_i$  is the  $T_i \times q_n$  dimensional linear covariate matrix. Further,  $\mathbf{A}_i(\boldsymbol{\alpha}) = \text{diag}\{\sigma_{i1}^2(\boldsymbol{\alpha}), \dots, \sigma_{iT_i}^2(\boldsymbol{\alpha})\}$  is a diagonal matrix of the variance of  $\mathbf{Y}_i$ . Here  $\sigma_{it}^2(\boldsymbol{\alpha}) = \phi \dot{\mu}(h_{it})$  is the spline approximated marginal variance per subject  $i$  and observation  $t$ , where the systematic component  $h_{it} = \mathbf{G}^\top(\mathbf{X}_{it}^\top \boldsymbol{\beta})\boldsymbol{\theta} + \mathbf{Z}_{it}^\top \boldsymbol{\gamma}$ ,  $\dot{\mu}(\cdot)$  is the first derivative with respect to  $h_{it}$ , and the scaling constant  $\phi = 1$  as in Wang et al. [36].  $\mathbf{V}_i(\boldsymbol{\alpha})$  is defined as  $\mathbf{V}_i(\boldsymbol{\alpha}) = (\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}), \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta})\boldsymbol{\theta} \mathbf{X}_i \mathbf{J}(\boldsymbol{\beta}), \mathbf{Z}_i))$ , where  $\dot{\mathbf{G}}(\cdot)$  is the first derivative of the spline basis. We define the Jacobian matrix to be  $\mathbf{J}(\boldsymbol{\beta}) = \partial \boldsymbol{\beta} / \partial \boldsymbol{\beta}^{(-1)} = ((-\boldsymbol{\beta}^{(-1)} / (1 - \|\boldsymbol{\beta}^{(-1)}\|^2)^{1/2})^\top, I_{(p-1) \times (p-1)})^\top$ , where we reparameterize  $\boldsymbol{\beta}$  to be a function of  $\boldsymbol{\beta}^{(-1)} = (\boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_p)^\top$  using the “delete-one-component” method for identifiability as in Yu et al. [41] and Yu and Ruppert [40].

Notably, the quadratic inference function does not need to estimate the coefficients  $\mathbf{a} = (a_1, \dots, a_m)^\top$ . This may be especially beneficial in a semiparametric high dimensional setting with likely even more nuisance parameters than a moderate dimensional setting [5,37]. Moreover, Qu et al. [27] show that the resulting estimators from minimizing this quadratic inference function are the most efficient estimators given the same class of estimating functions, which includes generalized estimating equations. This is beneficial, because a misspecified working correlation structure may lead to efficiency loss, but in practice, the true correlation structure is oftentimes unknown. In addition, the QIF only requires the first two moments of the response distribution alleviating the difficulty of specifying the full joint likelihood for correlated discrete responses [27].

## 3. Penalized quadratic inference function for ultra-high dimensional data

Variable selection is essential, since over selecting variables can adversely affect estimation efficiency and inference for a sparse set of important variables. Nonetheless, it is challenging to identify a sparse set of important variables in high dimensions [8]. In our motivating example, there are more than 500,000 SNPs and the correlated discrete response creates a further challenge in specifying the joint likelihood. Thus, for concurrent variable selection and estimation with diverging and even potentially ultra-high dimensional covariates along with correlated discrete responses, we adopt the penalized quadratic inference function for the generalized partially linear single-index model in (1). We define our penalized quadratic inference function to minimize

$$Q_p(\boldsymbol{\alpha}) = Q_n(\boldsymbol{\alpha}) + \sum_{j=1}^{p_n} q_{\lambda_p}(|\beta_j|) + \sum_{k=1}^{q_n} q_{\lambda_q}(|\gamma_k|). \quad (3)$$

Here  $Q_n(\alpha)$  refers to the QIF equation (2),  $q_{\lambda_p}(|\beta_j|)$  with  $j \in \{1, \dots, p_n\}$  is the penalty function for each single-index coefficient with corresponding tuning parameter  $\lambda_p$ . Similarly,  $q_{\lambda_q}(|\gamma_k|)$  with  $k \in \{1, \dots, q_n\}$  is the penalty function for each linear coefficient with corresponding tuning parameter  $\lambda_q$ .

While other penalty functions are available, we employ the smoothly clipped absolute deviation (SCAD) penalty. The first derivative of the SCAD penalty function is defined as  $\dot{q}_\lambda(\zeta) = \lambda\{I(\zeta \leq \lambda) + ((a\lambda - \zeta)_+ / (a-1)\lambda)I(\zeta > \lambda)\}$ , for  $q_\lambda(0) = 0$  and  $a > 2$  for a given regularization parameter  $\lambda$ . As suggested in Fan and Li [7], we set  $a = 3.7$ .

#### 4. Asymptotic properties

In our proposed semiparametric approach for longitudinal data, the total number of covariates can be ultra-high dimensional in both the nonparametric and partially linear portions. Additionally, the true important covariates,  $p_{sn}$  and  $q_{sn}$ , can be diverging. Particularly for longitudinal data, few existing approaches allow diverging or even ultra-high dimensional variables with nonlinearity for discrete responses. In this challenging setting, we establish important theoretical properties for the estimators of both the partially linear and, more importantly, the nonparametric single-index components, not only in moderately high dimensions but even in ultra-high dimensions.

We first establish asymptotic theory, namely, convergence rate and asymptotic normality, in the oracle case (i.e., when the exact true covariates are given ahead of time). In this case, we use subscript (s) to denote the oracle. That is, we let  $\mathbf{X}_{(s)j}$  be the true  $T_i \times p_{sn}$  dimensional single-index covariate matrix, and  $\mathbf{Z}_{(s)ji}$  be the true  $T_i \times q_{sn}$  dimensional linear covariate matrix. The true non-zero  $p_{sn}$ -dimensional single-index parameter vector is  $\beta_{0(s)} = \{\beta_{01}, \dots, \beta_{0p_{sn}}\}^\top$ , and the true non-zero  $q_{sn}$ -dimensional partially linear vector is  $\gamma_{0(s)} = \{\gamma_{01}, \dots, \gamma_{0q_{sn}}\}^\top$ . The true non-zero spline parameters corresponding to the single-index  $\mathbf{X}_{(s)j}\beta_{0(s)}$  are  $\theta_{0(s)}$ . Then for the oracle estimators, we let  $\hat{\alpha}_{(s)}$  refer to the oracle estimator for the true important parameters  $\alpha_{0(s)} = (\theta_{0(s)} \beta_{0(s)} \gamma_{0(s)})$ , and  $\hat{\zeta}_{(s)}$  refer to the estimate for  $\zeta_{0(s)} = (\beta_{0(s)}^{(-1)} \gamma_{0(s)})$ . Also, we let  $\mathbf{V}_{0(s)ji} = (\mathbf{G}(\mathbf{X}_{(s)j}\beta_{0(s)}), \text{diag}\{\dot{\eta}(\mathbf{X}_{(s)j}\beta_{0(s)})\} \mathbf{X}_{(s)j} \mathbf{J}(\beta_{0(s)}), \mathbf{Z}_{(s)ji})$  with  $\mathbf{J}(\beta_{0(s)}) = \partial \beta_{0(s)} / \partial \beta_{0(s)}^{(-1)}$ , and we define the true conditional mean of the response as  $\mu_{0(s)ji} = \mathbf{g}^{-1}(\eta(\mathbf{X}_{(s)j}\beta_{0(s)}) + \mathbf{Z}_{(s)ji}\gamma_{0(s)})$ .

We use  $\dot{f}, \ddot{f}$  to denote the first and second derivative of functions. In the theoretical study, we assume that  $T_i \equiv T$  for simplicity, and we assume  $T$  is fixed. For two  $T$ -dimensional vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,  $\mathbf{a} \odot \mathbf{b}$  denotes the Hadamard product taken component-wise resulting in another  $T$ -dimensional vector. For a matrix  $\mathbf{A}$  with  $T$  rows,  $\mathbf{A} \odot \mathbf{a}$  is the matrix of the same size as  $\mathbf{A}$  resulting from applying the Hadamard product to each column of  $\mathbf{A}$ . We assume in (A2) below that  $\eta$  is smooth. Under that assumption, there exists a vector  $\theta_0$  such that  $\|\mathbf{G}(\cdot)\theta_0 - \eta(\cdot)\|_\infty \leq CH_n^{-d}$ .

We rely on the following assumptions.

(A1)  $\sup_{1 \leq i \leq n, 1 \leq t \leq T} \|\mathbf{X}_{(s)it}\| = O_p(\sqrt{p_{sn}})$ ,  $\sup_{1 \leq i \leq n, 1 \leq t \leq T} \|\mathbf{Z}_{(s)it}\| = O_p(\sqrt{q_{sn}})$ , and  $(\mathbf{Y}_i, \mathbf{X}_i, \mathbf{Z}_i)$  are i.i.d.

(A2)  $\eta \in \mathcal{H}^d(M)$  for some  $d \geq 2$  and a constant  $M > 0$ , where  $\mathcal{H}^d$  contains all functions  $\eta$  such that  $|\eta^{(d_1)}(x) - \eta^{(d_1)}(y)| \leq M|x - y|^{d-d_1}$ ,  $d_1$  is the largest integer strictly smaller than  $d$  and  $\eta^{(d_1)}$  is the  $d_1$ -th order derivative of  $\eta$ . We also assume  $\eta$  is a bounded function.

(A3)  $E[\|\mathbf{Y}_i - \mu_{0(s)ji}\|^{2+\delta}] < \infty$  for some  $\delta > 0$ .

(A4) There exist positive constants  $c_1$  and  $c_2$  such that  $c_1 \leq \lambda_{\min}(n^{-1} \sum_i \mathbf{V}_{0(s)ji}^\top \mathbf{V}_{0(s)ji}) \leq \lambda_{\max}(n^{-1} \sum_i \mathbf{V}_{0(s)ji}^\top \mathbf{V}_{0(s)ji}) \leq c_2$ , where  $\lambda_{\min}$  and  $\lambda_{\max}$  denote the smallest and largest eigenvalues of a matrix, respectively.

(A5) On the set  $\{\alpha_{(s)} : \|\alpha_{(s)} - \alpha_{0(s)}\| \leq Cr_n\}$ , where  $C$  is a positive constant,  $\dot{\mu}_i(\alpha)$ , is uniformly bounded away from 0 and  $\infty$ , and  $\ddot{\mu}_i(\alpha_{0(s)})$  is uniformly bounded.

(A6)  $\mathbf{M}_1 = \mathbf{I}$ ,  $\mathbf{M}_2, \dots, \mathbf{M}_m$  are linear independent non-negative definite matrices with bounded eigenvalues.

**Remark 1.** (A1) is trivially satisfied if the predictors are bounded random variables. (A2) assumes that the nonparametric function  $\eta$  is smooth so that it can be approximated well by a spline function. In particular, under (A2), we can find  $\theta_0$  such that  $\sup_u |\mathbf{G}^\top(u)\theta_0 - \eta(u)| \leq CH_n^{-d}$  for some constant  $C > 0$  [29]. For (A4), we note that  $\mathbf{V}_{0(s)ji}$  consists of three components. For  $\mathbf{G}(\mathbf{X}_{(s)j}\beta_{0(s)})$ , using the same argument as in Lemma 3 of [16], if the density of  $\mathbf{X}_{(s)j}\beta_{0(s)}$  is bounded away from zero and infinity,  $\sum_i \mathbf{G}^\top(\mathbf{X}_{(s)j}\beta_{0(s)})\mathbf{G}(\mathbf{X}_{(s)j}\beta_{0(s)})/n$  has eigenvalues bounded away from zero and infinity with probability approaching one. For the other two components, as long as the dimension  $p_{(s)n}$  and  $q_{(s)n}$  do not diverge too fast, it is reasonable to assume they also have eigenvalues bounded away from zero and infinity. In order for the assumption to hold for  $\sum_i \mathbf{V}_{0(s)ji}^\top \mathbf{V}_{0(s)ji}/n$ , we thus require the three components are not strongly correlated.

Further, for use in the proof of asymptotic normality, we define the following projection. Let

$\mathcal{M}_t = \{g : Eg^2(\mathbf{X}_t^\top \beta_0) < \infty\}$ . For any random vector  $\mathbf{a} \in \mathbf{R}^T$  that is a function of  $(\mathbf{X}_i, \mathbf{Z}_i)$ , we define  $E_{\mathcal{M}}[\mathbf{a}] = \mathbf{g}(\mathbf{X}_i \beta_0) = (g_1(\mathbf{X}_i^\top \beta_0), \dots, g_T(\mathbf{X}_i^\top \beta_0))^\top$ , where  $\mathbf{g} = (g_1, \dots, g_T)^\top$  is the minimizer of

$$E[(\mathbf{a} - \mathbf{g}(\mathbf{X}_i \beta_0))^\top \Omega (\mathbf{a} - \mathbf{g}(\mathbf{X}_i \beta_0))] \quad (4)$$

over  $g_j \in \mathcal{M}_j$ ,  $j \in \{1, \dots, T\}$ , where

$$\Omega = \mathbf{F}_i (E[\mathbf{F}_i^\top \mathbf{R} \mathbf{F}_i])^{-1} \mathbf{F}_i^\top,$$



$\mathbf{F}_i = (\mathbf{A}_{0i}^{1/2} \mathbf{M}_1 \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i}, \dots, \mathbf{A}_{0i}^{1/2} \mathbf{M}_m \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i}), \mathbf{R} = \mathbf{A}_{0i}^{-1/2} E[\boldsymbol{\epsilon}_i \boldsymbol{\epsilon}_i^\top] \mathbf{A}_{0i}^{-1/2}$ . Here  $\mathbf{R}$  is the true correlation matrix of the error term.

Interpreting  $\Omega$  as a weight matrix (in the non-longitudinal case it is just a  $1 \times 1$  weight), the above can indeed be regarded as a projection. Using projections in the proof of asymptotic normality in the parametric part is an important technique in various semiparametric models [43]. This definition of projection can be extended to the case when  $\mathbf{a}$  is a random matrix with  $T$  columns such that the projection is obtained row by row.

Using the defined projection, we write  $E_{\mathcal{M}}[\mathbf{Z}_{(s)i}] = \{f_{tj}(\mathbf{X}_{(s)i} \boldsymbol{\beta}_{0(s)})\}_{1 \leq t \leq T, 1 \leq j \leq q_{sn}}$  and  $E_{\mathcal{M}}[\text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{X}_i] = \{m_{tj}(\mathbf{X}_i \boldsymbol{\beta}_0)\}_{1 \leq t \leq T, 1 \leq j \leq p_{sn}}$ . Define

$$\mathbf{U}_{0(s)i}^\top = \left( \mathbf{J}^\top(\boldsymbol{\beta}_{0(s)}) \mathbf{X}_{(s)i}^\top \text{diag}(\dot{\eta}(\mathbf{X}_{(s)i} \boldsymbol{\beta}_{0(s)})) - E_{\mathcal{M}}[\mathbf{J}^\top(\boldsymbol{\beta}_{0(s)}) \mathbf{X}_{(s)i}^\top \text{diag}(\dot{\eta}(\mathbf{X}_{(s)i} \boldsymbol{\beta}_{0(s)}))] \right) / (\mathbf{Z}_{(s)i}^\top - E_{\mathcal{M}}[\mathbf{Z}_{(s)i}^\top])$$

(A7)  $f_{tj}, m_{tj} \in \mathcal{H}^{d'}$  for some  $d' \geq 1$ .

(A8) There exist positive constants  $c_3$  and  $c_4$  such that  $c_3 \leq \lambda_{\min}(n^{-1} \sum_i \mathbf{U}_{0(s)i}^\top \mathbf{U}_{0(s)i}) \leq$

$\lambda_{\max}(n^{-1} \sum_i \mathbf{U}_{0(s)i}^\top \mathbf{U}_{0(s)i}) \leq c_4$ .

(A9)  $\boldsymbol{\epsilon}_i = (\epsilon_{i1}, \dots, \epsilon_{iT})^\top$ , where  $\epsilon_{it} = Y_{it} - \mu_{it}$ , is a sub-exponential random vector.

**Theorem 1** (Convergence Rate of Oracle Estimator). Under assumptions (A1)–(A6) and that

$(H_n^3 + H_n^2 p_{sn} + q_{sn}) r_n^2 \rightarrow 0$ , we have

$$\|\hat{\boldsymbol{\alpha}}_{(s)} - \boldsymbol{\alpha}_{0(s)}\| = O_p(r_n),$$

where  $r_n = \sqrt{(H_n + p_{sn} + q_{sn})/n} + H_n^{-d}$ .

**Theorem 2** (Asymptotic Normality of Oracle Estimator). Under assumptions (A1)–(A8) and that

$n(H_n^3 + H_n^2 p_{sn} + q_{sn})^{1/2} r_n^3 \rightarrow 0$ , then for any unit vector  $\mathbf{a} \in \mathbf{R}^{p_{sn} + q_{sn} - 1}$ ,

$$\sqrt{n} \mathbf{a}^\top \boldsymbol{\Sigma}_{(s)}^{-1/2} (\hat{\boldsymbol{\xi}}_{(s)} - \boldsymbol{\xi}_{0(s)}) \xrightarrow{d} N(0, 1),$$

where  $\boldsymbol{\Sigma}_{(s)} = E[\mathbf{U}_{0(s)}^\top \mathbf{F}_{0(s)}] (E[\mathbf{F}_{0(s)}^\top \mathbf{R} \mathbf{F}_{0(s)}])^{-1} E[\mathbf{F}_{0(s)}^\top \mathbf{U}_{0(s)}]$ .

The assumptions  $(H_n^3 + H_n^2 p_{sn} + q_{sn}) r_n^2 \rightarrow 0$  and  $n(H_n^3 + H_n^2 p_{sn} + q_{sn})^{1/2} r_n^2 \rightarrow 0$  implicitly put constraints on the allowed divergence rate for  $p_{sn}, q_{sn}$ . For example, if setting  $H_n = n^{\frac{2d-1}{2d+1}}$ , which balances the term  $\sqrt{H_n/n}$  and  $H_n^{-d}$  in the definition  $r_n$ , then  $(H_n^3 + H_n^2 p_{sn} + q_{sn}) r_n^2 \rightarrow 0$  implies  $p_{sn} = o_p(n^{\frac{2d-1}{2(2d+1)}})$  and  $q_{sn} = o_p(n^{\frac{1}{2}})$ , while  $n(H_n^3 + H_n^2 p_{sn} + q_{sn})^{1/2} r_n^2 \rightarrow 0$  implies  $p_{sn} = o_p(n^{\frac{2d-1}{4(2d+1)}})$  and  $q_{sn} = o_p(n^{\frac{1}{4}})$ .

Next, we establish asymptotic properties when the single-index covariates and the partially linear covariates are ultra-high dimensional for our semiparametric penalized quadratic inference function estimators.

**Theorem 3** (Asymptotic Normality of PQIF Estimator). Under the same assumptions for Theorem 2 and (A9) and

$\left( \sqrt{\frac{(H_n^3 + H_n^2 p_{sn} + q_{sn}) \ln^2 p_n}{n}} + \sqrt{H_n + p_{sn} + q_{sn}} \right) r_n \ll \lambda_p \ll \min_{j \leq p_{sn}} |\beta_{0j}|, \sqrt{H_n + p_{sn} + q_{sn}} \left( 1 + \sqrt{\frac{\ln^2 q_n}{n}} \right) r_n \ll \lambda_q \ll \min_{j \leq q_{sn}} |\gamma_{0j}|$ , there is an  $r_n$ -consistent local minimizer  $\hat{\boldsymbol{\alpha}} = (\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}})$  such that for any unit vector  $\mathbf{a}$ ,

(i)

$$\sqrt{n} \mathbf{a}^\top \boldsymbol{\Sigma}_{(s)}^{-1/2} (\hat{\boldsymbol{\xi}}_{(s)} - \boldsymbol{\xi}_{0(s)}) \xrightarrow{d} N(0, 1),$$

where  $\boldsymbol{\Sigma}_{(s)} = E[\mathbf{U}_{0(s)}^\top \mathbf{F}_{0(s)}] (E[\mathbf{F}_{0(s)}^\top \mathbf{R} \mathbf{F}_{0(s)}])^{-1} E[\mathbf{F}_{0(s)}^\top \mathbf{U}_{0(s)}]$ .

(ii)  $\hat{\beta}_{p_{sn}+1} = \dots = \hat{\beta}_{p_n} = \hat{\gamma}_{q_{sn}+1} = \dots = \hat{\gamma}_{q_n} = 0$  with probability approaching one.

The assumptions  $\left( \sqrt{\frac{(H_n^3 + H_n^2 p_{sn} + q_{sn}) \ln^2 p_n}{n}} + \sqrt{H_n + p_{sn} + q_{sn}} \right) r_n \ll \lambda_p \ll \min_{j \leq p_{sn}} |\beta_{0j}|$  and  $\sqrt{H_n + p_{sn} + q_{sn}} \left( 1 + \sqrt{\frac{\ln^2 q_n}{n}} \right) r_n \ll \lambda_q \ll \min_{j \leq q_{sn}} |\gamma_{0j}|$  implicitly put constraints on the allowed rates of  $p_n$  and  $q_n$ , which also depend on

the divergence rates of  $p_{sn}$  and  $q_{sn}$ . For example, if  $H_n = n^{\frac{1}{2d+1}}$ ,  $p_{sn}$  and  $q_{sn}$  are fixed, and  $\min_{j \leq p_{sn}} |\beta_{0j}|$  and  $\min_{j \leq q_{sn}} |\gamma_{0j}|$  are bounded away from zero, these assumptions would impose the constraint  $\ln p_n = o_p(n^{\frac{2d-1}{2d+1}})$  and  $\ln q_n = o_p(n^{\frac{2d}{2d+1}})$ .

One key contribution we make is to establish the above challenging theoretical properties for ultra-high dimensional correlated response data. Here we allow not only ultra-high dimensional partially linear covariates but importantly also the single-index covariates within the nonlinear flexible unknown function to be ultra-high dimensional. We select

important variables that are potentially diverging in the generalized partially linear single-index model framework with PQIF functions. We relegate detailed technical proofs along with four lemmas to Section 9. We hope similar techniques can be adopted for other highly nonlinear ultra-high dimensional modeling.

## 5. Algorithm and detailed implementation

### 5.1. Algorithm

In high dimensions, computational challenges arise when estimating the spline, single-index, and partially linear coefficients of the penalized quadratic inference function (3). These computational difficulties result from the non-convex SCAD penalty function and the potentially ultra-high dimensional variables in the nonparametric and linear components of the generalized partially linear single-index model. By adopting an iterative algorithm for estimation, we prevent the prohibitive computational cost of estimating all parameters in one step. Moreover, we employ two strategic approximations: a linear approximation of the unknown flexible function and a local quadratic approximation of the non-convex SCAD penalty.

First, we address the computational burden due to the potentially ultra-high dimensional covariates embedded inside the likely nonlinear, unknown function estimated nonparametrically. This nonlinear optimization over a high-dimensional space likely demands a considerable amount of computation. To ease this burden, we apply a linear approximation of  $\eta(\cdot)$  to convert this penalized nonlinear problem into a penalized linear problem. Then we can take advantage of existing linear algorithms, which are more computationally expedient than nonlinear estimation, particularly in high dimensions. Specifically, we apply the first order Taylor series approximation of  $\eta(\mathbf{X}_i\boldsymbol{\beta})$  at the point  $\mathbf{X}_i\boldsymbol{\beta}_0$  as  $\eta(\mathbf{X}_i\boldsymbol{\beta}_0) \cong \eta(\mathbf{X}_i\boldsymbol{\beta}_0) + \text{diag}(\dot{\eta}(\mathbf{X}_i\boldsymbol{\beta}_0))\mathbf{X}_i(\boldsymbol{\beta} - \boldsymbol{\beta}_0)$ , where  $\dot{\eta}(\cdot)$  is the first derivative of the function  $\eta(\cdot)$ . Notably, the parameters to be estimated  $\boldsymbol{\beta}$  are now outside of the unknown, flexible function. We can now use the linear penalized quadratic inference function method as in Cho and Qu [5] to estimate the parameters.

Still, while this problem is now linear, minimization of the linear penalized quadratic inference function proves challenging because of the non-convex SCAD penalty. Following Cho and Qu [5], we approximate the linear penalized quadratic inference function by implementing a local quadratic approximation of the SCAD penalty and then minimize this approximated quadratic function using the Newton–Raphson algorithm. In particular, on iteration  $k$  of the iterative algorithm, we let the column coefficient vector  $\boldsymbol{\zeta}^{(k)} = (\boldsymbol{\beta}^{(k)} \boldsymbol{\gamma}^{(k)})$ . Then as in Fan and Li [7], the local quadratic approximation of the SCAD penalty function for iteration  $k$  and coefficient  $\zeta_j^{(k)}$  is  $q_\lambda(|\zeta_j^{(k)}|) \cong q_\lambda(|\zeta_j^{(k)}|) + (1/2)(\dot{q}_\lambda(|\zeta_j^{(k)}|)/|\zeta_j^{(k)}|)(\zeta_j^{(k)2} - \zeta_j^{(k)})^2$  with  $\zeta_j \approx \zeta_j^{(k)}$  and  $\zeta_j^{(k)} \neq 0$ . See Cho and Qu [5] for the full estimation algorithm for the linear penalized quadratic inference function problem.

Altogether, our two-step iterative algorithm consists of, in step one, estimating the spline coefficients  $\boldsymbol{\theta}$  through a quadratic inference function in conjunction with a polynomial B-spline basis. In this step, the single-index and linear coefficient estimates  $\boldsymbol{\beta}, \boldsymbol{\gamma}$  are given from the previous iteration or as initial values during the first iteration of the algorithm. In the second step, given the estimated spline coefficients  $\hat{\boldsymbol{\theta}}$ , we estimate the linear and single-index coefficients  $\boldsymbol{\beta}, \boldsymbol{\gamma}$  via the linear approximation of  $\eta(\cdot)$ , which converts the nonlinear optimization to a linear penalized quadratic inference function method from Cho and Qu [5]. In this step, the conditional mean with  $\eta(\cdot)$  estimated by polynomial splines  $g^{-1}(\mathbf{G}(\mathbf{X}_i\boldsymbol{\beta})\boldsymbol{\theta} + \mathbf{Z}_i\boldsymbol{\gamma})$  becomes  $g^{-1}(\mathbf{G}(\mathbf{X}_i\boldsymbol{\beta})\hat{\boldsymbol{\theta}} + \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i\boldsymbol{\beta})\hat{\boldsymbol{\theta}})\mathbf{X}_i(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \mathbf{Z}_i\boldsymbol{\gamma})$  upon implementation of the linear approximation of  $\eta(\mathbf{X}_i\boldsymbol{\beta})$  at the point  $\mathbf{X}_i\hat{\boldsymbol{\beta}}$ . We continue this process of first estimating the spline coefficients and then estimating the linear and single-index coefficients until convergence.

A detailed description of the iterative algorithm is the following:

Step 0: Initialize  $\hat{\boldsymbol{\zeta}}^{(0)} = \begin{pmatrix} \hat{\boldsymbol{\beta}}^{(0)} \\ \hat{\boldsymbol{\gamma}}^{(0)} \end{pmatrix}$ . See Section 5.2 for details on obtaining initial values under various scenarios.

Step 1: Given  $\hat{\boldsymbol{\zeta}}^{(k-1)} = \begin{pmatrix} \hat{\boldsymbol{\beta}}^{(k-1)} \\ \hat{\boldsymbol{\gamma}}^{(k-1)} \end{pmatrix}$ , estimate the spline coefficients,  $\hat{\boldsymbol{\theta}}^{(k-1)}$ , by minimizing the quadratic inference

function  $\mathbf{g}_n^\top \mathbf{W}_n^{-1} \mathbf{g}_n$ , where  $\mathbf{g}_n = (1/n) \sum_i \mathbf{g}_i(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}^{(k-1)}, \hat{\boldsymbol{\gamma}}^{(k-1)})$ .

Step 2: Given the estimated spline coefficients  $\hat{\boldsymbol{\theta}}^{(k-1)}$ , the  $k_{th}$  penalized estimator of  $\hat{\boldsymbol{\zeta}}^{(k)} = (\hat{\boldsymbol{\beta}}^{(k)} \hat{\boldsymbol{\gamma}}^{(k)})$  is determined by minimizing the penalized quadratic inference function

$$Q_n(\hat{\boldsymbol{\theta}}^{(k-1)}, \boldsymbol{\beta}, \boldsymbol{\gamma}) + \sum_{j=1}^{p_n} q_{\lambda_p}(|\beta_j|) + \sum_{k=1}^{q_n} q_{\lambda_q}(|\gamma_k|), \quad (5)$$

where  $Q_n(\hat{\boldsymbol{\theta}}^{(k-1)}, \boldsymbol{\beta}, \boldsymbol{\gamma}) = \mathbf{g}_n^\top \mathbf{W}_n^{-1} \mathbf{g}_n$  and  $\mathbf{g}_n = (1/n) \sum_i \mathbf{g}_i(\hat{\boldsymbol{\theta}}^{(k-1)}, \boldsymbol{\beta}, \boldsymbol{\gamma})$ . We also assume the identifiability constraints  $\|\boldsymbol{\beta}\| = 1$  and  $\beta_1 > 0$ .

Using the approach from Cho and Qu [5], we estimate the single-index and partially linear parameters  $\boldsymbol{\beta}, \boldsymbol{\gamma}$  by taking a quadratic approximation of the penalized quadratic inference function with the SCAD penalty and applying the Newton–Raphson algorithm. For more details on implementation, refer to Cho and Qu [5]. We then repeat steps 1 and 2 until convergence. In practice, we observe that the algorithm converges in around 3 steps.

## 5.2. Practical implementation information for algorithm

In this section, we detail practical implementation aspects contributing to the performance of our approach and provide suggestions for desirable performance. The aspects we must decide on include an appropriate linear combination of basis matrices to use for the quadratic inference function, the degree, number, and placement of knots for the B-spline basis for univariate smoothing, and the penalty parameter for variable selection. Further implementation decisions include determining initial values and selecting a screening method to prohibit extensive computational burden when incorporating ultra-high dimensional covariates. We provide additional studies in the Supplementary Material.

### Choice of Basis Matrices for Quadratic Inference Function

For the quadratic inference function, one must choose a linear combination of basis matrices  $\mathbf{M}_i$  to approximate the inverse of the working correlation structure, as in (2) and (3). As detailed in Qu et al. [27], for the exchangeable correlation structure,  $\mathbf{R}^{-1}$  can be approximated by  $a_1\mathbf{I} + a_2\mathbf{M}_2$ , where  $\mathbf{M}_2$  has 1 on the off-diagonal and 0s on the diagonal. Here  $a_1 = -(m-2)\rho + 1/l_1$  and  $a_2 = \rho/l_1$  with  $l_1 = (m-1)\rho^2 - (n-2)\rho - 1$ . For the AR(1) case,  $\mathbf{R}^{-1}$  can be approximated by  $a_1\mathbf{I} + a_2\mathbf{M}_2 + a_3\mathbf{M}_3$ , where  $\mathbf{M}_2$  has 0 everywhere except the two main off-diagonals have 1 s, and  $\mathbf{M}_3$  has 0 everywhere except at (1, 1) and (T, T), which are 1 s where T is the largest time point.  $a_1 = (1 + \rho^2)/l_2$ ,  $a_2 = -\rho/l_2$ , and  $a_3 = -\rho^2/l_2$  with  $l_2 = 1 - \rho^2$ . When no previous information exists about the possible correlation structure, one may apply the approach from Zhou and Qu [44] to consistently select the correlation structure. This approach employs informative basis matrices in a sufficient class of the true structure. For further details on more working correlation structures, we refer the reader to [26,27,44].

### Initial Values and Screening

In ultra-high dimensions, the computational time of most approaches tends to be prohibitive for practical usage [8]. Due to this, screening is commonly used in relevant literature to rapidly reduce the covariate dimension in practice (e.g., Cai et al. [3] and Fang et al. [9]). We echo Fan and Lv [8] and integrate sure independence screening to first efficiently reduce dimensionality to benefit step 2 of our iterative algorithm, which is the linear penalized quadratic inference function.

Alternatively, in moderately high dimensions, the linear penalized quadratic inference function can estimate parameters in a satisfactory amount of time for practical usage without employing a screening approach. Relevant literature often uses initial values with good properties in moderately high dimensions. For example, Wang et al. [35] incorporate initial values from linear generalized estimating equations. One may also use sure independence screening estimates at the original dimension as initial values [8]. For our approach in moderately high dimensions, we may use linear quadratic inference function estimates as initial values.

### Tuning Parameters for Penalization

Common to many optimization problems, tuning parameter selection is essential to achieve desirable variable selection performance. One must find a suitable value for the tuning parameter of the penalty function  $\lambda_p$  for the single-index parameters and  $\lambda_q$  for the linear parameters of the model. Since the role of the quadratic inference function is analogous to negative twice the log likelihood, we can employ the high dimensional Bayesian information criterion (HBIC) for tuning parameter selection as in Wang et al. [34]. The model selection criterion is

$$HBIC(\lambda_p, \lambda_q) = Q_n(\hat{\theta}, \hat{\beta}_{\lambda_p}, \hat{\gamma}_{\lambda_q}) + d_\lambda \frac{\ln(n)}{2n} C_n, \quad (6)$$

where  $d_\lambda$  is the number of important nonzero coefficients from both the single-index and partially linear portions of the fitted model.  $C_n$  is considered as  $\ln(\ln(p_n + q_n))$  and is based on empirical evidence as in the literature Wang et al. [34]. We identify the minimum HBIC by searching through a grid of separate  $\lambda_p$  and  $\lambda_q$  values, which affords fewer modeling restrictions.

### Spline Smoothing Tuning Parameters

Regarding the spline basis, one must select the number, degree, and placement of knots. Selection criteria for the number of interior knots for the B-spline basis are determined by various approaches. As described in Ma et al. [21], a BIC type criterion focusing on consistency can be used to select the number of knots, or when efficiency is of interest, AIC and cross-validation approaches can be implemented to select the number of knots [15]. In Ruppert and Carroll [28] and Yu and Ruppert [40], the number of knots implemented depends on the characteristics and shape of the function to be estimated. In particular, monotonicity and discontinuity dictate the number of knots that are needed in an empirical analysis.

In terms of knot placement, knots are usually equally spaced, or knots are placed at the quantiles of the support of the estimated single-index. We note that during the iterative algorithm the estimated index values and consequently knot placement may change in practice. Moreover, a knot must be near any discontinuities for best performance. In our simulation studies, we find equally spaced knots with 2 interior knots perform well. Huang et al. [16] claim that when the number of knots does not greatly influence performance, one can fix the number of knots. Regarding the choice of degree, in practice, usually quadratic and cubic splines are implemented [21,41].



**Table 1**

Summary of variable selection results for the generalized partially linear single-index model for the binary response example. The total numbers of covariates are  $p_n + q_n = 500$  and  $p_n + q_n = 5000$  with  $p_n = q_n$  calculated over 200 simulations. “True%” is the percentage of times the true important variables are selected over the replications. “TN” is the average of the true negatives over the replications, and “FN” is the average of the false negatives over the replications. “MSEp” and “MSEq” are the average mean squared errors for the single-index parameters and partially linear parameters over all simulation replications.

	Structure	Single-index covariates				Partially linear covariates			
		True%	TNs	FNs	MSEp	True%	TNs	FNs	MSEq
$p_n + q_n = 500$	Independence	86	246.81	0	0.0240	100	247	0	0.0180
	AR(1)	92	246.92	0	0.0120	100	247	0	0.0160
	Exchangeable	99	246.98	0	0.0080	100	247	0	0.0130
$p_n + q_n = 5000$	Independence	80	2496.91	0.12	0.0586	93	2496.93	0	0.0216
	AR(1)	79	2496.94	0.15	0.0658	95	2496.95	0	0.0195
	Exchangeable	97	2497.00	0.03	0.0179	98	2496.98	0	0.0156

## 6. Simulation studies

We demonstrate the estimation and variable selection performance of our generalized partially linear single-index model using the penalized quadratic inference function through simulation studies. We mainly focus our investigation on a correlated binary response example in very high dimensions. We also examine our approach under imbalanced data and a complex correlation structure. We present additional settings and continuous response examples in the Supplementary Material. For both the single-index and linear covariates, “True%” is the percentage of simulation replications that select only the relevant variables, “TN” represents the true negatives, the number of true zero coefficients that the model sets to zero, and “FN” represents the false negatives, the number of true nonzero coefficients that are set falsely to zero. For the estimation results, “MSEp” is the average mean squared error of the single-index coefficient estimates  $\|\hat{\beta} - \beta_0\|^2$  and “MSEq” is the same for the partially linear coefficients  $\|\hat{\gamma} - \gamma_0\|^2$ .

### Binary Response Example

We consider the correlated binary responses from the marginal mean

$$\ln\left(\frac{p_{it}}{1-p_{it}}\right) = \sin\left(\frac{(\mathbf{X}_{it}^\top \beta - a)\pi}{b-a}\right) + \mathbf{Z}_{it}^\top \gamma.$$

There are  $n = 400$  subjects with  $T_i = T = 10$  time points for all subjects, and  $p_n + q_n = 500$  and  $p_n + q_n = 5000$  with  $p_n = q_n$ . The coefficient vector is  $\beta_0 = (1, 1, 1, 0, \dots, 0)^\top / \sqrt{3}$  and  $\gamma_0 = (1, 1, 1, 0, \dots, 0)^\top$ , and  $a$  and  $b$  are constants equal to  $\sqrt{3}/2 - 1.645/\sqrt{12}$  and  $\sqrt{3}/2 + 1.645/\sqrt{12}$  respectively. The covariates  $\mathbf{X}_{it}$  are sampled from an independent uniform distribution with minimum at 0 and maximum at 0.5, and the covariates  $\mathbf{Z}_{it}$  are sampled from an independent normal distribution with mean at 0 and standard deviation of 0.5. We use the method described in Macke et al. [22] to simulate correlated binary data with exchangeable correlation structure with  $\rho = 0.2$ . We evaluate the performance using independence, AR(1), and exchangeable working correlation matrices.

Table 1 reports the estimation and variable selection results, showing that our GPLSIM PQIF model selects the true important variables with a high percentage over 200 simulation replications. On average, the model also correctly accounts for the number of true negatives of the covariates. In particular, when  $p_n + q_n = 500$ , the number of true negatives is very close to the true amount of 247 for the single-index covariates, and the number of true negatives is close to the true amount of 247 for the linear covariates in all scenarios. Similarly, in the  $p_n + q_n = 5000$  case, the number of true negatives is close to the true amount of 2497 for the single-index case, and the number of true negatives is also close to the true amount of 2497 for the linear covariates. Moreover, the number of true important single-index or linear variables under selected is low, since the number of false negatives in the  $p_n + q_n = 500$  case is 0, and it is below 5% in all instances of the  $p_n + q_n = 5000$  case. Here the model selection performance under various correlation structures is quite similar. However, the performance is better in almost all cases under the true exchangeable correlation structure.

In terms of estimation, the MSEs in Table 1 are small for all parameter estimates. Further, Table 2 indicates that all parameter estimates are close to the true parameter values. As is the case with variable selection for this example, the parameter estimation performance under various correlation structures is quite similar. However, it is better in almost all cases under the true exchangeable correlation structure.

### Complex Correlation Structure

We further investigate a more complex correlation structure: a mixture of an AR(1) correlation structure and an exchangeable correlation structure with  $\rho = 0.5$  and  $p_n = 250$  and  $q_n = 250$ . In particular, the correlation structure simulated is  $\text{Corr}(Y_{it}, Y_{ik}) = (\rho^{|t-k|}/2 + \rho/2)$  with  $t, k \in \{1, \dots, T\}$  when  $t \neq k$ ; and  $\text{Corr}(Y_{it}, Y_{ik}) = 1$  when  $t = k$  for each subject  $i$ . Table 3 indicates that our proposed model, labeled as “GPLSIM SCAD”, yields good selection results for an underlying complex correlation structure even when the working correlation matrix is misspecified. This is consistent with

**Table 2**

Summary of parameter estimates for the generalized partially linear single-index model for the binary response example. The total numbers of covariates are  $p_n + q_n = 500$  and  $p_n + q_n = 5000$  with  $p_n = q_n$ . The sample mean, bias, and standard error are calculated over 200 simulations for single-index and partially linear parameter estimates.

	par.	Independence			AR(1)			Exchangeable		
		mean	bias	se	mean	bias	se	mean	bias	se
$p_n + q_n = 500$	$\beta_1$	0.5671	-0.0103	0.0590	0.5739	-0.0035	0.0544	0.5757	-0.0017	0.0484
	$\beta_2$	0.5659	-0.0114	0.0630	0.5731	-0.0042	0.0519	0.577	-0.0004	0.0504
	$\beta_3$	0.5780	0.0006	0.0602	0.5746	-0.0027	0.0557	0.5722	-0.0052	0.0520
	$\gamma_1$	1.0024	0.0024	0.0778	1.0156	0.0156	0.0734	1.0171	0.0171	0.0625
	$\gamma_2$	0.9997	-0.0003	0.0772	1.0123	0.0123	0.0703	1.0168	0.0168	0.0654
	$\gamma_3$	0.9938	-0.0062	0.0812	1.0098	0.0098	0.0709	1.0108	0.0108	0.0659
$p_n + q_n = 5000$	$\beta_1$	0.5855	0.0081	0.0746	0.5626	-0.0148	0.1359	0.5778	0.0005	0.0551
	$\beta_2$	0.5241	-0.0533	0.1694	0.5374	-0.0400	0.1729	0.5556	-0.0217	0.1098
	$\beta_3$	0.5718	-0.0056	0.1216	0.5751	-0.0022	0.1189	0.5831	0.0057	0.0497
	$\gamma_1$	1.0022	0.0022	0.0833	1.0144	0.0144	0.0798	1.0165	0.0165	0.0761
	$\gamma_2$	1.0125	0.0125	0.0786	1.0190	0.0190	0.0717	1.0228	0.0228	0.0625
	$\gamma_3$	1.0017	0.0017	0.0811	1.0052	0.0052	0.0808	1.0091	0.0091	0.0681

**Table 3**

Summary of estimation and variable selection results for the generalized partially linear single-index model for the binary response example with complex correlation with different penalties for  $p_n = 250$  and  $q_n = 250$ . The true correlation structure is  $\text{Corr}(Y_{it}, Y_{ik}) = (\rho^{|t-k|}/2 + \rho/2)$  with  $t, k = 1, \dots, T$  when  $t \neq k$ ,  $\rho = 0.5$ , and  $\text{Corr}(Y_{it}, Y_{ik}) = 1$  when  $t = k$  for each subject  $i$ . The “GPLSIM SCAD” and “GPLSIM LASSO” refer to our proposed generalized partially linear single-index model using the penalized quadratic inference function with SCAD penalty and LASSO penalty respectively. Similarly, The “Linear PQIF SCAD” and “Linear PQIF LASSO” refer to linear penalized quadratic inference function with SCAD penalty and LASSO penalty respectively. The remaining settings are the same as in the binary response example. “TN” is the average of the true negatives over the replications, and “FN” is the average of the false negatives over the replications. “MSEp” and “MSEq” are the average mean squared errors for the single-index parameters and partially linear parameters over all simulation replications.

Model	Structure	Single-index covariates				Partially linear covariates			
		True%	TNs	FNs	MSEp	True%	TNs	FNs	MSEq
GPLSIM SCAD	Independence	93	246.93	0	0.0169	100	247	0	0.0143
	AR(1)	88	246.88	0	0.0142	99	246.99	0	0.0106
	Exchangeable	96	246.96	0	0.0127	100	247	0	0.0123
GPLSIM LASSO	Independence	87	246.74	0.13	0.0929	86	246.85	0	0.1298
	AR(1)	79	246.33	0.18	0.1805	100	247	0	0.0976
	Exchangeable	89	246.71	0.09	0.0868	99	247	0	0.0840
Linear PQIF SCAD	Independence	8	245	0	0.5370	78	246.85	0	0.0177
	AR(1)	32	246.10	0.78	0.5488	90	247	0	0.0100
	Exchangeable	28	246.41	0.98	0.5246	88	246.80	0	0.0138
Linear PQIF LASSO	Independence	2	246.09	1.37	0.7846	29	245.81	0	0.0652
	AR(1)	22	245.98	0.42	0.3213	52	246.29	0	0.0784
	Exchangeable	55	246.68	0.34	0.2427	50	246.30	0	0.0654

similar literature, e.g., Qu et al. [27], that QIF can yield a consistent estimator under a misspecified working correlation matrix.

We further examine the estimation and variable selection results with a popular LASSO penalty [32] in comparison to SCAD penalty. Table 3 indicates when using the LASSO penalty for all correlation structures, the GPLSIM PQIF model tends to over-select variables that are not important in the single-index portion. In addition, the parameter estimates for the SCAD penalty are closer to the true parameter values compared to those using LASSO penalty. This is similarly observed for the linear PQIF model, where for all correlation structures the LASSO penalty appears to perform worse in two areas: non-important covariates are selected in the single-index portion and all parameter estimates are farther from the true values.

#### Unequal $T_i$ per Subject

Next, we examine unequal  $T_i$  per subject for imbalanced data. We simulate 400 subjects with 10 observations per subject with  $p_n + q_n = 500$ , and 400 subjects with 5 observations per subject for  $p_n + q_n = 500$ . The remaining set is the same as in the binary response example with a simulated exchangeable correlation structure and  $\rho = 0.5$ . The results in Table 4 indicate that the estimated parameters are close to the true parameters. The variable selection results have a high true percentage, with a range of 97%–100% for the single-index parameters and near 100% for the linear parameters. Our algorithm uses the high-dimensional linear penalized QIF from Cho and Qu [5] and linear QIF from Qu et al. [27] as a base, thus incorporating imbalanced data is natural.

**Table 4**

Summary of estimation and variable selection results for the binary response example with unequal  $T_i$  with  $p_n = 250$  and  $q_n = 250$ . This analysis with unequal numbers of observations per subject  $T_i$  has  $n = n_1 + n_2$  subjects. In particular,  $n_1 = 400$  subjects with  $T = 10$  observations per subject and  $n_2 = 400$  subjects with  $T = 5$  observations per subject. The remaining settings are the same as in the binary response example. “True%” is the percentage of times the true important variables are selected over the replications. “TN” is the average of the true negatives over the replications, and “FN” is the average of the false negatives over the replications.

Structure	Single-index covariates				Partially linear covariates			
	True%	TNs	FNs	MSEp	True%	TNs	FNs	MSEq
Independence	97	246.96	0	0.0175	100	247	0	0.0174
AR(1)	100	247	0	0.0089	100	247	0	0.0142
Exchangeable	100	247	0	0.0092	100	247	0	0.0168

## 7. An application to diabetes analysis

Diabetes is a widespread disease associated with various health complications including but not limited to an increased risk of stroke and vision loss [33]. Due to the known impact of both phenotype and high dimensional genotype variables on diabetes [10], identifying important genetic and phenotype risk factors may allow early diagnosis and more effective treatment. To investigate this relationship between risk factors and diabetes status using our proposed model, we use the ongoing Framingham Heart Study data [6]. This study is a continuing longitudinal study of cardiovascular disease, and researchers have also used this data to investigate various diseases such as diabetes (e.g., Meigs et al. [23]). For more in-depth details on the phenotype and genotype variables we used, please visit the Framingham study page at [https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000007.v32.p13](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000007.v32.p13).

To investigate factors related to diabetes status, we use a subset of 878 participants among the offspring cohort of the Framingham data measured over four exams. To avoid a large gap between exam 1 and exam 2, we use a subset that covers four waves of exams to ensure the diabetes indicators and diabetes-related quantitative traits are comparable. According to the Framingham Data dictionary, the participants' diabetes status variable, the correlated binary response  $Y_i$ , is derived by an algorithm considering blood glucose test results, treatment status, and other information per the protocol of `vr_diab_ex09_1_1002s`. We include ten phenotype covariates that are potentially linked to diabetes from previous research: age, systolic blood pressure (SBP), total cholesterol (TC), smoking status (SMK), cigarette per day (CPD), body mass index (BMI), weight (WGT), Ventricular rate (VENT\_RT), Triglyceride (TRIG), and HDL cholesterol. We also include high dimensional genetic data from participants as covariates in our analysis. Participants were genotyped on the Affymetrix 500K creating 500,568 single nucleotide polymorphisms (SNPs). After removing SNPs that are vastly missing or possessing zero variance, quality control filters yielded 53,722 SNPs for modeling. For computational efficiency, we also follow the common practice and employ a combined screening method as in [11,17].

We apply our proposed generalized partially linear single-index model with PQIF approach to this subset of Framingham data to investigate the relationship between the longitudinal binary response, diabetes status, and genetic and phenotype risk factors. The logit link function is used. The only binary phenotype risk factor is smoking status, which naturally goes into the partially linear term when fitting the model while the rest of the phenotype variables and SNPs are embedded in the single-index term. Specifically, for participant  $i$  over 4 waves,  $\ln(p_{it}/(1-p_{it})) = \eta(\mathbf{X}_{it}^T \boldsymbol{\beta}) + Z_{it}\gamma$ , where all phenotype and genotype variables except smoking status enter the single-index term, and smoking status ( $Z_{it}$ ) enters partially linearly. We find the best penalty parameter using HBIC, and we use quadratic degree with 3 equally spaced knots. See Section 5.2 for more information on tuning parameters and setup.

Fig. 1 indicates a clear nonlinear relationship between the flexible function and the single-index made up of genetic and phenotype risk factors that were selected as important using our approach. A linear model is likely to misspecify the relationship. For comparison, we also report the linear SCAD penalized quadratic inference function along with our proposed generalized partially linear single-index model using the SCAD penalized quadratic inference function. For the linear model, the same covariates are included. We use the BIQIF from Cho and Qu [5] to determine the best performing penalty parameter for the linear model.

Table 5 reports the 3-fold cross-validation area under the curve (AUC) and the out-of-sample model AUC with a 70/30 training and testing split of the data. Under the exchangeable structure, our proposed GPLSIM model clearly outperforms the linear penalized QIF model in terms of higher AUC for both cross-validation and out-of-sample testing. Under the correlation structure of AR(1) and independence, although with less margin, our proposed partially linear single-index model using penalized QIF still consistently outperforms the linear penalized QIF model.

Table 5 also reports the number of phenotype variables selected and the number of genotype variables selected by the model. Under the exchangeable structure, four phenotype variables have been selected. They are SBP, WGT, TRIG, and HDL. All of them have been identified as risk factors for diabetes in the previous literature and have been used as key components in predictive models of incidence of diabetes mellitus such as [20,30]. Among the large amount of genotype variables, the proposed model selected six SNPs, three of them have been confirmed by literature: rs5018648, rs10946398 and rs4506565 (Diabetes Genetics Initiative (2007); [1,24,25]).

**Table 5**

Summary of results for real data application to diabetes analysis. “CV AUC” and “OOS AUC” are cross-validation model area under the curve with 3-fold cross-validation and out-of-sample model area under the curve using 70/30 training testing split data set. “PQIF-GPLSIM” refers to our proposed generalized partially linear single-index model using penalized QIF, and “PQIF-linear” refers to the linear penalized QIF model. The numbers of selected phenotype and genotype variables are reported with full data.

		OOS AUC	CV AUC	# of phenotype selected	# of gene selected
PQIF-GPLSIM	Exchangeable	0.809	0.819	4	6
	AR(1)	0.801	0.794	5	6
	Independence	0.802	0.799	4	6
PQIF-Linear	Exchangeable	0.786	0.787	5	4
	AR(1)	0.773	0.781	3	6
	Independence	0.781	0.783	4	5

## 8. Conclusion

In public health, researchers are increasingly conducting large-scale longitudinal studies to investigate the relationship between disease and genetic and phenotype factors. These studies can provide insight into more effective treatment and disease prevention strategies. To incorporate correlation with a discrete response and to account for the complexity and synergy among genetic factors and phenotype variables, we propose an approach via penalized quadratic inference functions for generalized partially linear single-index models. Specifically, the quadratic inference functions can yield efficient estimation when the working correlation structure is misspecified, and the generalized partially linear single-index models are flexible and can incorporate some interactions. We allow genetic factors that are diverging and even in the exponential order with the number of participants not only in the linear portion of the model, but importantly also in the nonlinear portion estimated non-parametrically. We establish theoretical results such as asymptotic normality and the oracle property in ultra-high dimensions. Moreover, we develop an efficient estimation algorithm for computational expediency. We employ our approach to investigate diabetes status for an ongoing longitudinal public health study with genetic factors in very high dimensions.

## 9. Technical proofs

We provide detailed proofs and supporting lemmas for the asymptotic properties of estimators in ultra-high dimensions for our proposed generalized partially linear single-index model using the penalized quadratic inference function. In Section 9.1, we first prove [Theorem 1](#), the convergence rate under the oracle setting. In Section 9.2, we provide the proof of [Theorem 2](#), asymptotic normality, in the oracle setting. In Section 9.3, we prove [Theorem 3](#) which determines asymptotic properties for the penalized quadratic inference function estimator when both the single-index and linear variables can diverge and even be in ultra-high dimensions. In both the oracle and the ultra-high dimensional settings, the important variables can diverge for both the partially linear portion and the single-index portion. For simplicity of notation, we drop all  $n$  subscript in this section.

### 9.1. Proof of convergence rate for oracle estimators

Let  $\mathbf{V}_i(\alpha) = (\mathbf{G}(\mathbf{X}_i\beta), \text{diag}\{\dot{\mathbf{G}}(\mathbf{X}_i\beta)\theta\}\mathbf{X}_i\mathbf{J}(\beta), \mathbf{Z}_i)$ ,  $\mathbf{V}_{0i} = (\mathbf{G}(\mathbf{X}_i\beta_0), \text{diag}\{\dot{\mathbf{G}}(\mathbf{X}_i\beta_0)\}\mathbf{X}_i\mathbf{J}(\beta_0), \mathbf{Z}_i)$ ,  $\mathbf{K}_{i\ell}(\alpha) = \mathbf{V}_i^T(\alpha)\mathbf{A}_i^{1/2}(\alpha)\mathbf{M}_\ell\mathbf{A}_i^{-1/2}(\alpha)$ ,  $\mathbf{K}_{0i\ell} = \mathbf{V}_{0i}^T\mathbf{A}_{0i}^{1/2}\mathbf{M}_\ell\mathbf{A}_{0i}^{-1/2}$ ,  $\mathbf{K}_i(\alpha) = (\mathbf{K}_{i1}^T(\alpha), \dots, \mathbf{K}_{im}^T(\alpha))^T$ , and  $\mathbf{K}_{0i} = (\mathbf{K}_{0i1}^T, \dots, \mathbf{K}_{0im}^T)^T$ . Define  $\mathbf{g}_i(\alpha) = \mathbf{K}_i(\alpha)(\mathbf{Y}_i - \mu_i(\alpha))$ ,  $\mathbf{g}_{0i}(\alpha) = \mathbf{K}_{0i}(\mathbf{Y}_i - \mu_i(\alpha))$ ,  $\mathbf{g}_{0i} = \mathbf{K}_{0i}\epsilon_i$  with  $\epsilon_i = \mathbf{Y}_i - \mu_{0i}$ ,  $\bar{\mathbf{g}}_n(\alpha) = \frac{1}{n} \sum_i \mathbf{g}_i(\alpha)$ ,  $\bar{\mathbf{g}}_{0n}(\alpha) = \frac{1}{n} \sum_i \mathbf{g}_{0i}(\alpha)$ ,  $\bar{\mathbf{g}}_{0n} = \frac{1}{n} \sum_i \mathbf{g}_{0i}$ ,  $\mathbf{W}_n(\alpha) = \frac{1}{n} \sum_i \mathbf{g}_{0i}(\alpha)\mathbf{g}_{0i}^T(\alpha)$ ,  $\mathbf{W}_{0n} = \frac{1}{n} \sum_i \mathbf{g}_{0i}\mathbf{g}_{0i}^T$ ,  $\mathbf{Q}_n(\alpha) = \bar{\mathbf{g}}_n(\alpha)^T\mathbf{W}_n^{-1}(\alpha)\bar{\mathbf{g}}_n(\alpha)$ , and  $\mathbf{Q}_{0n}(\alpha) = \bar{\mathbf{g}}_{0n}(\alpha)^T\mathbf{W}_{0n}^{-1}\bar{\mathbf{g}}_{0n}(\alpha)$ .

Let  $r_n = \sqrt{\frac{H+p_s+q_s}{n}} + H^{-d}$ . To get the convergence rate, we need to show that for  $L$  sufficiently large, with probability approaching one as  $n \rightarrow \infty$ ,

$$\inf_{\|\alpha - \alpha_0\| = Lr_n} Q_n(\alpha) - Q_n(\alpha_0) \geq Cr_n^2. \quad (7)$$

The above will be implied by (8) and (9),

$$\sup_{\|\alpha - \alpha_0\| \leq Cr_n} |Q_n(\alpha) - Q_{0n}(\alpha)| = o_p(r_n^2), \quad (8)$$

and for  $L$  sufficiently large,

$$\inf_{\|\alpha - \alpha_0\| = Lr_n} Q_{0n}(\alpha) - Q_{0n}(\alpha_0) \geq CL^2r_n^2. \quad (9)$$

As a first step we prove (8).

$$\begin{aligned} Q_n(\alpha) - Q_{0n}(\alpha) &= \bar{\mathbf{g}}_n(\alpha)^\top \mathbf{W}_n^{-1}(\alpha) \bar{\mathbf{g}}_n(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha)^\top \mathbf{W}_{0n}^{-1} \bar{\mathbf{g}}_{0n}(\alpha) \\ &= (\bar{\mathbf{g}}_n(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha))^\top \mathbf{W}_n^{-1}(\alpha) \bar{\mathbf{g}}_n(\alpha) + \bar{\mathbf{g}}_{0n}(\alpha)^\top (\mathbf{W}_n^{-1}(\alpha) - \mathbf{W}_{0n}^{-1}) \bar{\mathbf{g}}_n(\alpha) + \bar{\mathbf{g}}_{0n}(\alpha)^\top \mathbf{W}_{0n}^{-1} (\bar{\mathbf{g}}_n(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha)) \\ &= O_p \left( \sqrt{\frac{H^3 + H^2 p_s + q_s}{n}} r_n^2 + \sqrt{H^3 + H^2 p_s + q_s} r_n^3 \right) = o_p(r_n^2), \end{aligned} \quad (10)$$

using (a), (c), (d), and (e) from Lemma 1 results below.

In Step 2, we prove (9).

$$Q_{0n}(\alpha) - Q_{0n}(\alpha_0) = (\bar{\mathbf{g}}_{0n}(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha_0))^\top \mathbf{W}_{0n} (\bar{\mathbf{g}}_{0n}(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha_0)) + 2\bar{\mathbf{g}}_{0n}(\alpha_0)^\top \mathbf{W}_{0n} (\bar{\mathbf{g}}_{0n}(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha_0)) \geq CL^2 r_n^2,$$

using (a), (b), and (e) from Lemma 1 results below.

**Lemma 1.** We establish the following properties:

- (a)  $\sup_{\|\alpha - \alpha_0\| \leq Cr_n} \|\bar{\mathbf{g}}_{0n}(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha_0)\| = O_p(r_n).$
- (b)  $\inf_{\|\alpha - \alpha_0\| = Lr_n} \|\bar{\mathbf{g}}_{0n}(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha_0)\| \geq CLr_n.$
- (c)  $\sup_{\|\alpha - \alpha_0\| \leq Cr_n} \|\bar{\mathbf{g}}_n(\alpha) - \bar{\mathbf{g}}_{0n}(\alpha)\| = O_p \left( \sqrt{\frac{H^3 + H^2 p_s + q_s}{n}} r_n \right).$
- (d)  $\sup_{\|\alpha - \alpha_0\| \leq Cr_n} \|\mathbf{W}_n^{-1}(\alpha) - \mathbf{W}_{0n}^{-1}\| = O_p \left( \sqrt{H^3 + H^2 p_s + q_s} r_n \right).$
- (e)  $\|\bar{\mathbf{g}}_{0n}(\alpha_0)\| = O_p(r_n).$

**Proof of Lemma 1.** We first prove (e), where

$$\sum_i \mathbf{K}_{0i}(\mathbf{Y}_i - \mu_i(\alpha_0)) = \sum_i \mathbf{K}_{0i} \epsilon_i + \sum_i \mathbf{K}_{0i}(\mu_{0i} - \mu_i(\alpha_0)) = O_p \left( \sqrt{n(H + p_s + q_s)} + nH^{-d} \right) = O_p(nr_n).$$

Here we used that for any unit vector  $\mathbf{a}$  of appropriate dimension,

$$\|\mathbf{a}^\top \sum_i \mathbf{K}_i(\mu_{0i} - \mu_i(\alpha_0))\| \leq \left( \sum_i |\mathbf{a}^\top \mathbf{K}_i \mathbf{K}_i^\top \mathbf{a}|^2 \right)^{1/2} \left( \sum_i \|\mu_{0i} - \mu_i(\alpha_0)\|^2 \right)^{1/2} = O_p(nH^{-d}).$$

For (c), the proof is based on repeated application of Taylor's expansion, but the diverging dimension makes it quite messy and therefore hard to keep track of the higher order terms. Let  $\delta_{it} = \mathbf{G}^\top(\mathbf{X}_{it}^\top \beta_0) \theta_0 - \eta(\mathbf{X}_{it}^\top \beta_0)$  and  $\delta_i = (\delta_{i1}, \dots, \delta_{iT})^\top$ . Then employing Taylor's expansion,

$$\begin{aligned} h_{it}(\alpha) - h_{0it} &= \mathbf{G}^\top(\mathbf{X}_{it}^\top \beta) \theta - \mathbf{G}^\top(\mathbf{X}_{it}^\top \beta_0) \theta_0 + \mathbf{Z}_{it}^\top(\gamma - \gamma_0) + \delta_{it} \\ &= \mathbf{G}^\top(\mathbf{X}_{it}^\top \beta_0) (\theta - \theta_0) + \dot{\mathbf{G}}^\top(\mathbf{X}_{it}^\top \beta_0) \theta_0 \mathbf{X}_{it}^\top (\beta - \beta_0) + \mathbf{Z}_{it}^\top(\gamma - \gamma_0) \\ &\quad + \dot{\mathbf{G}}^\top(\mathbf{X}_{it}^\top \beta^*) (\theta - \theta_0) \mathbf{X}_{it}^\top (\beta - \beta_0) + \ddot{\mathbf{G}}^\top(\mathbf{X}_{it}^\top \beta^*) \theta_0 (\mathbf{X}_{it}^\top (\beta - \beta_0))^2 + \delta_{it} \\ &= (\mathbf{G}^\top(\mathbf{X}_{it}^\top \beta_0), \dot{\eta}(\mathbf{X}_{it}^\top \beta_0) \mathbf{X}_{it}^\top, \mathbf{Z}_{it}^\top) (\alpha - \alpha_0) + \delta_{it} + O_p \left( \left( \sqrt{H^3 p_s} + p_s \right) r_n^2 \right), \end{aligned}$$

where  $\beta^*$  lies between  $\beta_0$  and  $\beta$  in the following quantities, where superscript  $*$  always indicates such values that arise from Taylor's expansion, or

$$\mathbf{h}_i(\alpha) - \mathbf{h}_{0i} = (\mathbf{G}(\mathbf{X}_i \beta_0), \dot{\eta}(\mathbf{X}_i \beta_0) \mathbf{X}_i, \mathbf{Z}_i) (\alpha - \alpha_0) + \delta_i + O_p \left( \left( \sqrt{H^3 p_s} + p_s \right) r_n^2 \right) = O_p \left( \sqrt{H + p_s + q_s} r_n \right). \quad (11)$$

We also have

$$\sum_i \|\mathbf{h}_i(\alpha) - \mathbf{h}_i(\alpha_0)\|^2 = O_p(nr_n^2 + nH^2 \min\{H, p_s\} r_n^4).$$

Similarly, we have by Taylor's expansion

$$\begin{aligned} \mathbf{A}_i^{1/2}(\alpha) - \mathbf{A}_{0i}^{1/2} &= \frac{1}{2} \mathbf{A}_{0i}^{-1/2} \text{diag}(\ddot{\mu}_{0i}) \text{diag}(\mathbf{h}_i(\alpha) - \mathbf{h}_{0i}) + O_p(\|\mathbf{h}_i(\alpha) - \mathbf{h}_{0i}\|^2) \\ &= \frac{1}{2} \mathbf{A}_{0i}^{-1/2} \text{diag}(\ddot{\mu}_{0i}) \text{diag} \{ (\mathbf{G}(\mathbf{X}_i \beta_0), \dot{\eta}(\mathbf{X}_i \beta_0) \mathbf{X}_i, \mathbf{Z}_i) (\alpha - \alpha_0) + \delta_i \} \\ &\quad + O_p \left( \left( \sqrt{H^3 p_s} + p_s \right) r_n^2 + (H + p_s + q_s) r_n^2 \right), \end{aligned}$$



$$\sum_i \left\| \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) - \mathbf{A}_{0i}^{1/2} \right\|^2 = \sum_i \left\| \frac{1}{2} \mathbf{A}_{0i}^{-1/2} \text{diag}(\ddot{\boldsymbol{\mu}}_{0i}) \text{diag} \{ (\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0), \dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{X}_i, \mathbf{Z}_i) (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) + \boldsymbol{\delta}_i \} \right\|^2 \\ + O_p(nH^2 \min\{H, p_s\} r_n^4) = O_p(nr_n^2 + nH^2 \min\{H, p_s\} r_n^4),$$

$$\mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{A}_{0i}^{-1/2} = -\frac{1}{2} \mathbf{A}_{0i}^{-3/2} \text{diag}(\ddot{\boldsymbol{\mu}}_{0i}) \text{diag}(\mathbf{h}_i(\boldsymbol{\alpha}) - \mathbf{h}_{0i}) + O_p(\|\mathbf{h}_i(\boldsymbol{\alpha}) - \mathbf{h}_{0i}\|^2) \\ = -\frac{1}{2} \mathbf{A}_{0i}^{-3/2} \text{diag}(\ddot{\boldsymbol{\mu}}_{0i}) \text{diag} \{ (\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0), \dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{X}_i, \mathbf{Z}_i) (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) \} \\ + O_p\left(\left(\sqrt{H^3 p_s} + p_s\right) r_n^2 + (H + p_s + q_s) r_n^2\right),$$

$$\sum_i \left\| \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{A}_{0i}^{-1/2} \right\|^2 = O_p(nr_n^2 + nH^2 \min\{H, p_s\} r_n^4) \\ \mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}) - \mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0) = \dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}_0) \odot (\mathbf{X}_i (\boldsymbol{\beta} - \boldsymbol{\beta}_0)) + O_p(H^{5/2} p_s r_n^2),$$

and

$$\sum_i \left\| \mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}) - \mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0) \right\|^2 = \sum_i \left\| \dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}_0) \odot (\mathbf{X}_i (\boldsymbol{\beta} - \boldsymbol{\beta}_0)) \right\|^2 + O_p(nH^5 p_s r_n^4) = O_p(nH^3 r_n^2 + nH^5 p_s r_n^4).$$

Using the identity

$$A_1 B_1 C_1 - A_0 B_0 C_0 = (A_1 - A_0) B_0 C_0 + A_0 (B_1 - B_0) C_0 + A_0 B_0 (C_1 - C_0) + A_0 (B_1 - B_0) (C_1 - C_0) \\ + (A_1 - A_0) B_0 (C_1 - C_0) + (A_1 - A_0) (B_1 - B_0) C_0 + (A_1 - A_0) (B_1 - B_0) (C_1 - C_0),$$

we have

$$\mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}) \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} = \dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}_0) \odot (\mathbf{X}_i (\boldsymbol{\beta} - \boldsymbol{\beta}_0)) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \\ + \frac{1}{2} \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{A}_{0i}^{-1/2} \text{diag}(\ddot{\boldsymbol{\mu}}_{0i}) \text{diag} \{ (\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0), \dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{X}_i, \mathbf{Z}_i) (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) + \boldsymbol{\delta}_i \} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \\ - \frac{1}{2} \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-3/2} \text{diag}(\ddot{\boldsymbol{\mu}}_{0i}) \text{diag} \{ (\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0), \dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{X}_i, \mathbf{Z}_i) (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) + \boldsymbol{\delta}_i \} \\ + O_p\left(\sqrt{H^3 p_s (H + p_s + q_s)} r_n^2 + \sqrt{H} (H + p_s + q_s) r_n^2\right),$$

and

$$\sum_i \left\| \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}) \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \right\|^2 = O_p(nH^3 r_n^2).$$

Using  $\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha}) = \boldsymbol{\epsilon}_i + (\boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha}))$ ,  $\boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha}) = O_p(\sqrt{H + p_s + q_s} r_n)$ , and  $\sum_i \|\boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha})\|^2 = O_p(nr_n^2)$ , we obtain

$$\sum_i \left( \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}) \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \right) (\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \\ = \sum_i \left( \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}) \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{G}^\top(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \right) (\boldsymbol{\epsilon}_i + \boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) = O_p(\sqrt{nH^3} r_n).$$

Other components of  $\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i}(\boldsymbol{\alpha})$  can be similarly dealt with. More specifically, we have

$$\mathbf{J}^\top(\boldsymbol{\beta}) \mathbf{X}_i^\top \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}) \boldsymbol{\theta}) - \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) = \frac{\partial \mathbf{J}^\top(\boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \\ + \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\ddot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \odot (\mathbf{X}_i (\boldsymbol{\beta} - \boldsymbol{\beta}_0))) \\ + \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}_0) (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \dot{\boldsymbol{\delta}}_i) + O_p\left(\sqrt{p_s (H^3 + p_s)} r_n^2\right) \\ = O_p\left(\left(\sqrt{p_s (H^3 p_s + p_s)} r_n\right)\right),$$

where  $\dot{\boldsymbol{\delta}}_i = \dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}_0) \boldsymbol{\theta}_0 - \dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)$ .

$$\mathbf{J}^\top(\boldsymbol{\beta}) \mathbf{X}_i^\top \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}) \boldsymbol{\theta}) \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{R}^{-1} \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \\ = \frac{\partial \mathbf{J}^\top(\boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} + \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\ddot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \odot (\mathbf{X}_i (\boldsymbol{\beta} - \boldsymbol{\beta}_0))) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2}$$

$$\begin{aligned}
& + \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}_0)(\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \dot{\boldsymbol{\delta}}_i) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} + \frac{1}{2} \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{A}_{0i}^{-1/2} \text{diag}(\dot{\boldsymbol{\mu}}_{0i}) \\
& \text{diag}\{(\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0), \dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{X}_i, \mathbf{Z}_i)(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) + \boldsymbol{\delta}_i\} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} - \frac{1}{2} \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-3/2} \text{diag}(\dot{\boldsymbol{\mu}}_{0i}) \\
& \text{diag}\{(\mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0), \dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{X}_i, \mathbf{Z}_i)(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) + \boldsymbol{\delta}_i\} + O_p\left(\sqrt{(H^3 p_s + p_s^2)(H + p_s + q_s) r_n^2} + \sqrt{p_s}(H + p_s + q_s) r_n^2\right) \\
& = O_p\left(\sqrt{H^3 p_s + p_s^2} r_n + \sqrt{p_s}(H + p_s + q_s) r_n\right),
\end{aligned}$$

and

$$\sum_i \left\| \mathbf{J}^\top(\boldsymbol{\beta}) \mathbf{X}_i^\top \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}) \boldsymbol{\theta}) \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{R}^{-1} \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \right\|^2 = O_p(n H^2 p_s r_n^2),$$

which in turn implies, using  $\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha}) = \boldsymbol{\epsilon}_i + (\boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha}))$ ,

$$\begin{aligned}
& \sum_i \left( \mathbf{J}^\top(\boldsymbol{\beta}) \mathbf{X}_i^\top \text{diag}(\dot{\mathbf{G}}(\mathbf{X}_i \boldsymbol{\beta}) \boldsymbol{\theta}) \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{J}^\top(\boldsymbol{\beta}_0) \mathbf{X}_i^\top \text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \right) (\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \\
& = O_p(n H^2 p_s r_n^2).
\end{aligned}$$

We can similarly show

$$\sum_i \left( \mathbf{Z}_i^\top \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{Z}_i^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \right) (\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) = O_p(n q_s r_n^2),$$

since  $\sum_i \left\| \mathbf{Z}_i^\top \mathbf{A}_i^{1/2}(\boldsymbol{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\boldsymbol{\alpha}) - \mathbf{Z}_i^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} \right\|^2 = O_p(n q_s r_n^2)$ .

Thus, we finally get

$$\sup_{\|\boldsymbol{\alpha} - \boldsymbol{\alpha}_0\| \leq C r_n} \left\| \sum_i \mathbf{g}_i(\boldsymbol{\alpha}) - \sum_i \mathbf{g}_{0i}(\boldsymbol{\alpha}) \right\| = O_p\left(\sqrt{n(H^3 + H^2 p_s + q_s) r_n}\right).$$

To prove (a) and (b), for any unit vector  $\mathbf{a}$ ,

$$\begin{aligned}
\mathbf{a}^\top \sum_i (\mathbf{g}_{0i}(\boldsymbol{\alpha}) - \mathbf{g}_{0i}(\boldsymbol{\alpha}_0)) & = \mathbf{a}^\top \sum_i \mathbf{K}_{0i} (\boldsymbol{\mu}_i(\boldsymbol{\alpha}_0) - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \\
& \leq \left( \sum_i |\mathbf{a}^\top \mathbf{K}_{0i} \mathbf{K}_{0i} \mathbf{a}| \right)^{1/2} \left( \sum_i \|\boldsymbol{\mu}_i(\boldsymbol{\alpha}_0) - \boldsymbol{\mu}_i(\boldsymbol{\alpha})\|^2 \right)^{1/2} = O_p(n r_n).
\end{aligned}$$

And for the lower bound, we similarly have for  $\mathbf{a} = \boldsymbol{\alpha}_0 - \boldsymbol{\alpha}$ , and  $\ell \in \{1, \dots, m\}$ ,

$$\begin{aligned}
\mathbf{a}^\top \sum_i (\mathbf{g}_{0i\ell}(\boldsymbol{\alpha}) - \mathbf{g}_{0i\ell}(\boldsymbol{\alpha}_0)) & = \mathbf{a}^\top \sum_i \mathbf{K}_{0i\ell} (\boldsymbol{\mu}_i(\boldsymbol{\alpha}_0) - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \\
& = \sum_i \mathbf{a}^\top \mathbf{K}_{0i\ell} \text{diag}(\dot{\boldsymbol{\mu}}_{0i}) \mathbf{V}_{0i}(\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}) \\
& \quad + \sum_i \mathbf{a}^\top \mathbf{K}_{0i\ell} (\boldsymbol{\mu}_i(\boldsymbol{\alpha}_0) - \boldsymbol{\mu}_i(\boldsymbol{\alpha}) - \text{diag}(\dot{\boldsymbol{\mu}}_{0i}) \mathbf{V}_{0i}(\boldsymbol{\alpha}_0 - \boldsymbol{\alpha})) \\
& = \sum_i (\boldsymbol{\alpha}_0 - \boldsymbol{\alpha})^\top \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i}(\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}) + O_p(n H^2 (H + p_s) r_n^4).
\end{aligned}$$

Thus,

$$\begin{aligned}
\|\boldsymbol{\alpha} - \boldsymbol{\alpha}_0\| \left\| \sum_i (\mathbf{g}_{0i\ell}(\boldsymbol{\alpha}) - \mathbf{g}_{0i\ell}(\boldsymbol{\alpha}_0)) \right\| & \geq \sum_i (\boldsymbol{\alpha}_0 - \boldsymbol{\alpha})^\top \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i}(\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}) - O_p(n H^2 (H + p_s) r_n^4) \\
& = C n \|\boldsymbol{\alpha} - \boldsymbol{\alpha}_0\|^2,
\end{aligned}$$

which implies (b).

To prove (d), start with

$$\begin{aligned}
\mathbf{W}_n(\boldsymbol{\alpha}) - \mathbf{W}_{0n} & = \frac{1}{n} \sum_i \mathbf{g}_i(\boldsymbol{\alpha}) \mathbf{g}_i^\top(\boldsymbol{\alpha}) - \frac{1}{n} \sum_i \mathbf{g}_{0i} \mathbf{g}_{0i}^\top \\
& = \frac{1}{n} \sum_i (\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i}) \mathbf{g}_{0i}^\top + \frac{1}{n} \sum_i \mathbf{g}_{0i} (\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i})^\top + \frac{1}{n} \sum_i (\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i}) (\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i})^\top;
\end{aligned}$$

$$\begin{aligned}
\left\| \sum_i (\mathbf{g}_{0i}(\boldsymbol{\alpha}) - \mathbf{g}_{0i}) \mathbf{g}_{0i}^\top \right\|^2 &= \left\| \sum_i \mathbf{K}_{0i} (\boldsymbol{\mu}_i(\boldsymbol{\alpha}) - \boldsymbol{\mu}_{0i}) \boldsymbol{\epsilon}_i^\top \mathbf{K}_{0i}^\top \right\|^2 = O_p \left( \sum_i \|\mathbf{K}_{0i}^\top \mathbf{K}_{0i} (\boldsymbol{\mu}_i(\boldsymbol{\alpha}) - \boldsymbol{\mu}_{0i})\|^2 \right) \\
&= O_p(n(H + p_s + q_s)^2 r_n^2); \\
&\leq \left\| \sum_i (\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i}(\boldsymbol{\alpha})) \mathbf{g}_{0i}^\top \right\|^2 \leq \left\| \sum_i (\mathbf{K}_i(\boldsymbol{\alpha}) - \mathbf{K}_{0i}) (\boldsymbol{\epsilon}_i + \boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) \boldsymbol{\epsilon}_i^\top \mathbf{K}_{0i}^\top \right\|^2 \\
&= O_p \left( \left\| \sum_i (\mathbf{K}_i(\boldsymbol{\alpha}) - \mathbf{K}_{0i}) \mathbf{K}_{0i}^\top \right\|^2 + \sum_i \|\mathbf{K}_{0i}^\top (\mathbf{K}_i(\boldsymbol{\alpha}) - \mathbf{K}_{0i}) (\boldsymbol{\mu}_i(\boldsymbol{\alpha}) - \boldsymbol{\mu}_{0i})\|^2 \right) \\
&= O_p(n^2(H^3 + H^2 p_s + q_s) r_n^2),
\end{aligned}$$

using that  $\sum_i \|\mathbf{K}_i(\boldsymbol{\alpha}) - \mathbf{K}_{0i}\|^2 = O_p(n(H^3 + H^2 p_s + q_s) r_n^2)$ . Then

$$\begin{aligned}
\sum_i (\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i}) (\mathbf{g}_i(\boldsymbol{\alpha}) - \mathbf{g}_{0i})^\top &= \sum_i ((\mathbf{K}_i(\boldsymbol{\alpha}) - \mathbf{K}_{0i}) (\boldsymbol{\epsilon}_i + \boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha})) + \mathbf{K}_{0i} (\boldsymbol{\mu}_{0i} - \boldsymbol{\mu}_i(\boldsymbol{\alpha})))^\otimes \\
&= O_p \left( \sum_i \|\mathbf{K}_i(\boldsymbol{\alpha}) - \mathbf{K}_{0i}\|^2 + \sum_i \|\mathbf{K}_i(\boldsymbol{\alpha}) - \mathbf{K}_{0i}\|^2 \|\boldsymbol{\mu}_i(\boldsymbol{\alpha}) - \boldsymbol{\mu}_{0i}\|^2 \right. \\
&\quad \left. + \sum_i \|\boldsymbol{\mu}_i(\boldsymbol{\alpha}) - \boldsymbol{\mu}_{0i}\|^2 \right) \\
&= O_p(n(H^3 + H^2 p_s + q_s) r_n^2).
\end{aligned}$$

Therefore,

$$\begin{aligned}
\|\mathbf{W}_n(\boldsymbol{\alpha}) - \mathbf{W}_{0n}\| &= O_p \left( \frac{H + p_s + q_s}{\sqrt{n}} r_n + \sqrt{H^3 + H^2 p_s + q_s} r_n + (H^3 + H^2 p_s + q_s) r_n^2 \right) \\
&= O_p(\sqrt{H^3 + H^2 p_s + q_s} r_n).
\end{aligned}$$

□

## 9.2. Proof of asymptotic normality for oracle estimators

Let  $\mathbf{N}_{0i} = (\text{diag}(\dot{\eta}(\mathbf{X}_i \boldsymbol{\beta}_0)) \mathbf{X}_i \mathbf{J}(\boldsymbol{\beta}_0), \mathbf{Z}_i)$  and define

$$\begin{aligned}
\mathbf{P} &= \arg \min_{\mathbf{Q}} (\mathbf{N} - \mathbf{GQ})^\top \begin{pmatrix} \mathbf{A}_{01}^{1/2} \mathbf{M}_1 \mathbf{A}_{01}^{1/2} \mathbf{V}_{01} & \cdots & \mathbf{A}_{01}^{1/2} \mathbf{M}_m \mathbf{A}_{01}^{1/2} \mathbf{V}_{01} \\ \vdots & \ddots & \vdots \\ \mathbf{A}_{0n}^{1/2} \mathbf{M}_1 \mathbf{A}_{0n}^{1/2} \mathbf{V}_{0n} & \cdots & \mathbf{A}_{0n}^{1/2} \mathbf{M}_m \mathbf{A}_{0n}^{1/2} \mathbf{V}_{0n} \end{pmatrix} \mathbf{W}_{0n}^{-1} \\
&\times \begin{pmatrix} \mathbf{V}_{01}^\top \mathbf{A}_{01}^{1/2} \mathbf{M}_1 \mathbf{A}_{01}^{1/2} & \cdots & \mathbf{V}_{0n}^\top \mathbf{A}_{0n}^{1/2} \mathbf{M}_1 \mathbf{A}_{0n}^{1/2} \\ \vdots & \ddots & \vdots \\ \mathbf{V}_{01}^\top \mathbf{A}_{01}^{1/2} \mathbf{M}_m \mathbf{A}_{01}^{1/2} & \cdots & \mathbf{V}_{0n}^\top \mathbf{A}_{0n}^{1/2} \mathbf{M}_m \mathbf{A}_{0n}^{1/2} \end{pmatrix} (\mathbf{N} - \mathbf{GQ}).
\end{aligned} \tag{12}$$

Note (12) is the empirical version of (4). We write  $\boldsymbol{\zeta} = (\boldsymbol{\beta}^{(-1)\top}, \boldsymbol{\gamma}^\top)^\top$  for the parameters of the parametric portion. Using the reparametrization  $\boldsymbol{\theta}^* = \boldsymbol{\theta} + \mathbf{P}\boldsymbol{\zeta}$ , there is a 1-1 mapping between  $(\boldsymbol{\theta}^*, \boldsymbol{\zeta})$  and  $(\boldsymbol{\theta}, \boldsymbol{\zeta})$ . Thus the problem of minimizing  $Q_n(\boldsymbol{\theta}, \boldsymbol{\zeta})$  over  $(\boldsymbol{\theta}, \boldsymbol{\zeta})$  is equivalent to minimizing over  $(\boldsymbol{\theta}^*, \boldsymbol{\zeta})$ . We will show in Lemma 2 that  $\|\mathbf{P}\|_{op}$  is bounded despite its diverging dimension. This means that an  $r_n$ -consistent estimator  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\zeta}})$  is equivalent to an  $r_n$ -consistent estimator  $(\widehat{\boldsymbol{\theta}}^*, \widehat{\boldsymbol{\zeta}})$ . In the following, we always regard the parameters as  $(\boldsymbol{\theta}^*, \boldsymbol{\zeta})$ , and we simply write  $Q_n(\boldsymbol{\theta}^*, \boldsymbol{\zeta})$  for the QIF objective we are minimizing when using such a reparametrization and do the same for other quantities that depend on the parameters  $\boldsymbol{\alpha} = (\boldsymbol{\theta}, \boldsymbol{\zeta})$ . Fixing  $\boldsymbol{\theta}^*$  at  $\widehat{\boldsymbol{\theta}}^* = \widehat{\boldsymbol{\theta}} + \mathbf{P}\widehat{\boldsymbol{\zeta}}$ , then obviously  $\widehat{\boldsymbol{\zeta}}$  minimizes  $Q_n(\widehat{\boldsymbol{\theta}}^*, \boldsymbol{\zeta})$ .

Let  $\mathbf{U}_i = \mathbf{N}_{0i} - \mathbf{G}(\mathbf{X}_i \boldsymbol{\beta}_0) \mathbf{P}$ , where this can be interpreted as orthogonalized predictors for the parametric part. Define  $Q_{0n}(\boldsymbol{\zeta}) = \bar{\mathbf{g}}_{0n}(\boldsymbol{\zeta})^\top \mathbf{W}_{0n} \bar{\mathbf{g}}_{0n}(\boldsymbol{\zeta})$ , with  $\bar{\mathbf{g}}_{0n}(\boldsymbol{\zeta}) = \frac{1}{n} \sum_i \mathbf{g}_{0i}(\boldsymbol{\zeta})$ ,  $\mathbf{g}_{0i}(\boldsymbol{\zeta}) = (\mathbf{g}_{0i1}^\top(\boldsymbol{\zeta}), \dots, \mathbf{g}_{0im}^\top(\boldsymbol{\zeta}))^\top$ ,  $\mathbf{g}_{0i\ell}(\boldsymbol{\zeta}) = \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} (\boldsymbol{\epsilon}_i - \mathbf{A}_{0i} \mathbf{G}(\mathbf{X}_i^\top \boldsymbol{\beta}) (\boldsymbol{\theta}^* - \boldsymbol{\theta}_0^*) - \mathbf{A}_{0i} \mathbf{U}_i (\boldsymbol{\zeta} - \boldsymbol{\zeta}_0))$ .

Let  $\tilde{\zeta}$  be the minimizer of  $Q_{0n}(\zeta)$ . We first establish the asymptotic normality of  $\tilde{\zeta}$ . Obviously,  $Q_{0n}(\zeta)$  is a quadratic function of  $\zeta$  with a close-form minimizer

$$\tilde{\zeta} = \zeta_0 + (\mathbf{S}_{0n}^\top \mathbf{W}_{0n}^{-1} \mathbf{S}_{0n})^{-1} \mathbf{S}_{0n}^\top \mathbf{W}_{0n}^{-1} \frac{1}{n} \sum_i \mathbf{K}_{0i} \left( \epsilon_i - \mathbf{A}_{0i} \mathbf{G}(\mathbf{X}_i^\top \beta_0) (\hat{\theta}^* - \theta_0^*) \right),$$

where  $\mathbf{S}_{0n} = \frac{1}{n} \sum_i \mathbf{K}_{0i} \mathbf{A}_{0i} \mathbf{U}_i$ . To establish the asymptotic normality of  $\tilde{\zeta}$ , we need to show

$$\mathbf{S}_{0n}^\top \mathbf{W}_{0n}^{-1} \frac{1}{n} \sum_i \mathbf{K}_{0i} \mathbf{A}_{0i} \mathbf{G}(\mathbf{X}_i^\top \beta_0) (\hat{\theta}^* - \theta_0^*) = o_p(n^{-1/2}).$$

We note that  $\hat{\theta}^* - \theta_0^*$  has a nonparametric rate and a naive bound would fail to show the  $o_p(n^{-1/2})$  rate above. However, it turns out the above is exactly zero due to the definition of  $\mathbf{P}$ . In fact, the first order condition of the optimization problem (12) is just

$$\mathbf{S}_{0n}^\top \mathbf{W}_{0n}^{-1} \frac{1}{n} \sum_i \mathbf{K}_{0i} \mathbf{A}_{0i} \mathbf{G}(\mathbf{X}_i^\top \beta_0) = \mathbf{0}.$$

To show  $\hat{\zeta}$  has the same asymptotic distribution as  $\tilde{\zeta}$ , we need to establish

$$\sup_{\|\zeta - \zeta_0\| \leq Cn} |Q_n(\hat{\theta}^*, \zeta) - Q_{0n}(\zeta)| = o_p(1/n), \quad (13)$$

and

$$Q_{0n}(\zeta) - Q_{0n}(\tilde{\zeta}) \geq C \|\zeta - \tilde{\zeta}\|^2. \quad (14)$$

Indeed, if (13) and (14) hold, we will have for any  $\epsilon > 0$ ,

$$\inf_{\|\zeta - \tilde{\zeta}\| = \epsilon/\sqrt{n}} Q_n(\hat{\theta}^*, \zeta) - Q_n(\hat{\theta}^*, \tilde{\zeta}) \geq C \|\zeta - \tilde{\zeta}\|^2 - o_p(1/n) > 0.$$

Since  $\hat{\zeta}$  minimizes  $Q_n(\hat{\theta}^*, \zeta)$ , the above implies  $\|\hat{\zeta} - \tilde{\zeta}\|_\infty \leq \|\hat{\zeta} - \tilde{\zeta}\| = o_p(n^{-1/2})$ , and thus  $\hat{\zeta}$  has the same asymptotic distribution as  $\tilde{\zeta}$ , which finishes the proof.

Note that (13) is already shown in (10), where the more stringent assumption for Theorem 2 makes the rate  $o_p(1/n)$  instead of  $o_p(r_n^2)$ , and (14) is shown in Lemma 3.

**Lemma 2.** The operator norm (largest singular value) of  $\mathbf{P}$  defined in (12) is bounded.

**Proof of Lemma 2.** We have the closed-form

$$\mathbf{P} = \left\{ \left( \frac{1}{n} \sum_i \mathbf{G}(\mathbf{X}_i \beta_0)^\top \mathbf{A}_{0i} \mathbf{K}_{0i}^\top \right) \mathbf{W}_{0n}^{-1} \left( \frac{1}{n} \sum_i \mathbf{K}_{0i} \mathbf{A}_{0i} \mathbf{G}(\mathbf{X}_i \beta_0) \right) \right\}^{-1} \left( \frac{1}{n} \sum_i \mathbf{G}(\mathbf{X}_i \beta_0)^\top \mathbf{A}_{0i} \mathbf{K}_{0i}^\top \right) \mathbf{W}_{0n}^{-1} \left( \frac{1}{n} \sum_i \mathbf{K}_{0i} \mathbf{A}_{0i} \mathbf{N}_{0i} \right).$$

The sample averages that appear above (note  $\mathbf{W}_{0n}$  is also a sample average) are obviously converging to their population counterparts, and we thus only consider the population quantities.

First,  $E[\mathbf{G}(\mathbf{X}_i \beta_0)^\top \mathbf{A}_{0i} \mathbf{K}_{0i}^\top]$  and  $E[\mathbf{K}_{0i} \mathbf{A}_{0i} \mathbf{N}_{0i}]$  are both submatrices of  $E[\mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i}]$ , and thus their operator norms are bounded. Second, the quantity that we take the inverse of is a principal submatrix

$$\sum_\ell E[\mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i}] \mathbf{W}_{0n}^{-1} E[\mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i}],$$

whose eigenvalues are bounded away from zero and infinity, and thus

$$\left\{ \left( \frac{1}{n} \sum_i \mathbf{G}(\mathbf{X}_i \beta_0)^\top \mathbf{A}_{0i} \mathbf{K}_{0i}^\top \right) \mathbf{W}_{0n}^{-1} \left( \frac{1}{n} \sum_i \mathbf{K}_{0i} \mathbf{A}_{0i} \mathbf{G}(\mathbf{X}_i \beta_0) \right) \right\}^{-1}$$

also has bounded eigenvalues.  $\square$

The next lemma proves (14).

**Lemma 3.**  $Q_{0n}(\zeta) - Q_{0n}(\tilde{\zeta}) \geq C \|\zeta - \tilde{\zeta}\|^2$ .

**Proof of Lemma 3.** Using that  $Q_{0n}(\zeta)$  is a quadratic form with minimizer  $\tilde{\zeta}$ , we have

$$Q_{0n}(\zeta) - Q_{0n}(\tilde{\zeta}) \geq \lambda_{\min}(\mathbf{D}) \|\zeta - \tilde{\zeta}\|^2,$$

where  $\lambda_{\min}(\mathbf{D})$  denotes the minimum eigenvalue of the matrix

$$\mathbf{D} = \left( \frac{1}{n} \sum_i \begin{pmatrix} \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_1 \mathbf{A}_{0i}^{1/2} \mathbf{U}_i \\ \vdots \\ \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{U}_i \end{pmatrix} \right)^\top \mathbf{W}_{0n}^{-1} \left( \frac{1}{n} \sum_i \begin{pmatrix} \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_1 \mathbf{A}_{0i}^{1/2} \mathbf{U}_i \\ \vdots \\ \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{U}_i \end{pmatrix} \right),$$

which is a principal submatrix of

$$\mathbf{D}' := \left( \frac{1}{n} \sum_i \begin{pmatrix} \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_1 \mathbf{A}_{0i}^{1/2} (\mathbf{G}(\mathbf{X}_i \beta_0), \mathbf{U}_i) \\ \vdots \\ \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} (\mathbf{G}(\mathbf{X}_i \beta_0), \mathbf{U}_i) \end{pmatrix} \right)^\top \mathbf{W}_{0n}^{-1} \left( \frac{1}{n} \sum_i \begin{pmatrix} \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_1 \mathbf{A}_{0i}^{1/2} (\mathbf{G}(\mathbf{X}_i \beta_0), \mathbf{U}_i) \\ \vdots \\ \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} (\mathbf{G}(\mathbf{X}_i \beta_0), \mathbf{U}_i) \end{pmatrix} \right).$$

Noting  $\mathbf{U}_i = \mathbf{N}_{0i} - \mathbf{G}(\mathbf{X}_i \beta_0) \mathbf{P}$ , we have

$$(\mathbf{G}(\mathbf{X}_i \beta_0), \mathbf{U}_i) = \mathbf{V}_{0i} \begin{pmatrix} \mathbf{I} & -\mathbf{P} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}.$$

Since  $\|\mathbf{P}\|_{op}$  is bounded, both

$$\begin{pmatrix} \mathbf{I} & -\mathbf{P} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}$$

and its inverse

$$\begin{pmatrix} \mathbf{I} & \mathbf{P} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}$$

have bounded eigenvalues. Furthermore, since  $\mathbf{W}_{0n}^{-1}$  has eigenvalues bounded away from zero, and the sample average in the definition of  $\mathbf{D}'$  can be approximated by its population counterpart, we only need to show that

$$\sum_{\ell} E \left[ \mathbf{V}_{0i}^\top \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \mathbf{V}_{0i} \right]^{\otimes 2}$$

has eigenvalues bounded away from zero, which is true by assumption.  $\square$

### 9.3. Proof of asymptotic properties of PQIF estimators in ultra-high dimensions

For convergence rate, we only need to show

$$\inf_{\|\alpha - \alpha_0\| = Lr_n} Q_n(\alpha) + \sum_{j=1}^p q_{\lambda_p}(|\beta_j|) + \sum_{k=1}^q q_{\lambda_q}(|\gamma_k|) > Q_n(\alpha_0) + \sum_{j=1}^p q_{\lambda_p}(|\beta_{0j}|) + \sum_{k=1}^q q_{\lambda_q}(|\gamma_{0k}|). \quad (15)$$

We already see in (7) that  $\inf_{\|\alpha - \alpha_0\| = Lr_n} Q_n(\alpha) > Q_n(\alpha_0)$ . We will show when  $\|\alpha - \alpha_0\| \leq Lr_n$ ,  $q_{\lambda_p}(|\beta_j|) \geq q_{\lambda_p}(|\beta_{0j}|)$ ,  $j \in \{1, \dots, p\}$ . Similarly we can show  $q_{\lambda_q}(|\gamma_k|) \geq q_{\lambda_q}(|\gamma_{0k}|)$ ,  $k \in \{1, \dots, q\}$ , which immediately implies (15). Indeed, when  $j > p_s$ ,  $q_{\lambda_p}(|\beta_j|) \geq 0 = q_{\lambda_p}(|\beta_{0j}|)$ . On the other hand, when  $j \leq p_s$ , since  $|\beta_{0j}| \geq C\lambda_p$  and  $|\hat{\beta}_j - \beta_{0j}| \leq \|\alpha - \alpha_0\| = o(\lambda_p)$ , both  $|\beta_{0j}|$  and  $|\hat{\beta}_j|$  are large enough to be in the region of the domain of  $q_{\lambda_p}$  that is nonzero by the specific expression of the SCAD penalty, and thus  $q_{\lambda_p}(|\beta_j|) = q_{\lambda_p}(|\beta_{0j}|)$ .



Next we consider variable selection consistency. Suppose, by way of contradiction, that  $\widehat{\beta}_{j^*} \neq 0$  for some  $j^* \in \{p_s + 1, \dots, p\}$ , and components of  $\widehat{\gamma}$  can be similarly dealt with. Define  $\check{\beta}$  such that its  $j^*$ -component is zero while other components are equal to those of  $\widehat{\beta}$ . We will show that

$$Q_n(\check{\alpha}) + \sum_{j=1}^p q_{\lambda_p}(|\check{\beta}_j|) + \sum_{k=1}^q q_{\lambda_q}(|\check{\gamma}_k|) < Q_n(\widehat{\alpha}) + \sum_{j=1}^p q_{\lambda_p}(|\widehat{\beta}_j|) + \sum_{k=1}^q q_{\lambda_q}(|\widehat{\gamma}_k|), \quad (16)$$

which leads to a contradiction. In fact, in Lemma 4, we show  $Q_n(\check{\alpha}) - Q_n(\alpha_0) = O_p(\lambda_p) \|\check{\alpha} - \widehat{\alpha}\|$ . Furthermore, by the definition of  $\check{\beta}$  which only differs from  $\widehat{\beta}$  in the  $j^*$ -th component, we have

$$\sum_{j=1}^p q_{\lambda_p}(|\check{\beta}_{0j}|) - \sum_{j=1}^p q_{\lambda_p}(|\widehat{\beta}_j|) = -q_{\lambda_p}(|\widehat{\beta}_{j^*}|) = -\lambda_p |\widehat{\beta}_{j^*}|,$$

where the last step is due to  $|\widehat{\beta}_{j^*}| \leq \|\widehat{\alpha} - \alpha_0\| = o_p(\lambda_p)$  implying  $|\widehat{\beta}_{j^*}|$  is in the region of the domain of  $q_{\lambda_p}(\cdot)$  that is a linear function by the specific expression of the SCAD penalty. This finishes the proof of (16).

**Lemma 4.** Uniformly for  $j^* \in \{p_s + 1, \dots, p\}$ , where  $\check{\alpha}$  below implicitly depends on  $j^*$ ,

$$Q_n(\widehat{\alpha}) - Q_n(\check{\alpha}) = O_p \left( \sqrt{\frac{(H^3 + H^2 p_s + q_s) \ln^2 p}{n}} + \sqrt{H + p_s + q_s} \right) r_n |\widehat{\beta}_{j^*}|.$$

**Proof of Lemma 4.** The proof is based on Taylor's expansion largely the same as in Lemma 1. We only briefly present some of the calculations for illustration. We decompose

$$\begin{aligned} Q_n(\widehat{\alpha}) - Q_n(\check{\alpha}) &= (\bar{\mathbf{g}}_n(\widehat{\alpha}) - \bar{\mathbf{g}}_n(\check{\alpha}))^\top \mathbf{W}_{0n}^{-1} \bar{\mathbf{g}}_{0n}(\alpha_0) + \bar{\mathbf{g}}_{0n}(\alpha_0)^\top (\mathbf{W}_n^{-1}(\alpha_0) - \mathbf{W}_n^{-1}(\check{\alpha})) \bar{\mathbf{g}}_{0n}(\alpha_0) \\ &\quad + \bar{\mathbf{g}}_{0n}(\alpha_0)^\top \mathbf{W}_{0n}^{-1} (\bar{\mathbf{g}}_n(\widehat{\alpha}) - \bar{\mathbf{g}}_n(\check{\alpha})) + \dots, \end{aligned} \quad (17)$$

where we omitted the higher order terms. Consider the first term as an example.  $\|\bar{\mathbf{g}}_{0n}(\alpha_0)\| = O_p(r_n)$  using (e) of Lemma 1. For  $\bar{\mathbf{g}}_n(\widehat{\alpha}) - \bar{\mathbf{g}}_n(\check{\alpha})$ , similar to the calculations in Lemma 1 we can get that, for example, the main terms of

$$\sum_i \mathbf{G}(\mathbf{X}_i \widehat{\beta}) \mathbf{A}_i^{1/2}(\widehat{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\widehat{\alpha}) (\mathbf{Y}_i - \mu_i(\widehat{\alpha})) - \sum_i \mathbf{G}(\mathbf{X}_i \check{\beta}) \mathbf{A}_i^{1/2}(\check{\alpha}) \mathbf{M}_\ell \mathbf{A}_i^{-1/2}(\check{\alpha}) (\mathbf{Y}_i - \mu_i(\check{\alpha}))$$

are

$$\begin{aligned} &\sum_i \left\{ \dot{\mathbf{G}}(\mathbf{X}_i \beta_0) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{-1/2} + \frac{1}{2} \mathbf{G}(\mathbf{X}_i \beta_0) \mathbf{A}_{0i}^{1/2} \text{diag}(\ddot{\mu}_{0i} \odot \dot{\eta}(\mathbf{X}_i \beta_0)) - \frac{1}{2} \mathbf{G}(\mathbf{X}_i \beta_0) \mathbf{A}_{0i}^{-3/2} \text{diag}(\ddot{\mu}_{0i} \odot \dot{\eta}(\mathbf{X}_i \beta_0)) \right\} \\ &(\mathbf{X}_{i(j^*)} \odot \epsilon_i) \widehat{\beta}_{j^*} + \sum_i \mathbf{G}(\mathbf{X}_i \beta_0) \mathbf{A}_{0i}^{1/2} \mathbf{M}_\ell \mathbf{A}_{0i}^{1/2} \text{diag}(\dot{\eta}(\mathbf{X}_i \beta_0)) \mathbf{X}_{i(j^*)} \widehat{\beta}_{j^*}, \end{aligned} \quad (18)$$

where the  $T$ -dimensional vector  $\mathbf{X}_{i(j^*)}$  is the  $j^*$ -th column of  $\mathbf{X}_i$ . The first term of (18) has mean zero, and is of order  $O_p(\sqrt{nH^3 \ln^2 p} |\widehat{\beta}_{j^*}|)$ , where the logarithmic term comes from applying Bernstein's inequality to get uniform bound over  $j^*$ . The second term in (18) is more easily derived to be of order  $n\sqrt{H} |\widehat{\beta}_{j^*}|$ . This and similar bounds would give

$$\|\bar{\mathbf{g}}_n(\widehat{\alpha}) - \bar{\mathbf{g}}_n(\check{\alpha})\| = O_p \left( \sqrt{\frac{(H^3 + H^2 p_s + q_s) \ln^2 p}{n}} + \sqrt{H + p_s + q_s} \right) |\widehat{\beta}_{j^*}|.$$

Then the first term in (17) would be of order  $O_p \left( \sqrt{\frac{(H^3 + H^2 p_s + q_s) \ln^2 p}{n}} + \sqrt{H + p_s + q_s} \right) r_n |\widehat{\beta}_{j^*}|$ . The third term in (17) is easily seen to be of the same order, while the second term in (17) can be shown to be of smaller order as is also the case in (10).

## CRediT authorship contribution statement

**Brittany Green:** Data curation, Methodology, Formal analysis Writing – original draft. **Heng Lian:** Methodology, Theory, Writing – review & editing. **Yan Yu:** Conceptualization, Methodology, Funding acquisition, Supervision, Writing – review & editing. **Tianhai Zu:** Data curation, Methodology, Formal analysis, Writing – review & editing.

## Acknowledgments

The Framingham Heart Study is conducted and supported by the National Heart, Lung, and Blood Institute (NHLBI), United States in collaboration with Boston University (Contract No. N01-HC-25195, HHSN268201500001I and 75N92019D-00031). This manuscript was not prepared in collaboration with investigators of the Framingham Heart Study and does not necessarily reflect the opinions or views of the Framingham Heart Study, Boston University, or NHLBI.

Funding for SHARe Affymetrix genotyping was provided by NHLBI, United States Contract N02-HL64278. SHARe Illumina genotyping was provided under an agreement between Illumina and Boston University. Funding for Affymetrix genotyping of the FHS Omni cohorts was provided by Intramural NHLBI funds from Andrew D. Johnson and Christopher J. O'Donnell.

The research of Heng Lian is partially supported by RGC GRF 11300519, 11300721 and 11311822.

## Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jmva.2023.105175>.

## References

- [1] P. An, M. Feitosa, S. Ketkar, A. Adelman, S. Lin, I. Borecki, M. Province, Epistatic interactions of CDKN2B-TCF7L2 for risk of type 2 diabetes and of CDKN2B-JAZF1 for triglyceride/high-density lipoprotein ratio longitudinal change: Evidence from the Framingham Heart Study, *BMC Proc.* 3 (Suppl. 7) (2009) S71.
- [2] Y. Bai, W.K. Fung, Z.Y. Zhu, Penalized quadratic inference functions for single-index models with longitudinal data, *J. Multivariate Anal.* 100 (1) (2009) 152–161.
- [3] L. Cai, H. Wu, D. Li, K. Zhou, F. Zou, Type 2 diabetes biomarkers of human gut microbiota selected via iterative sure independent screening method, *PLoS One* 10 (10) (2015) 1–15.
- [4] R.J. Carroll, J. Fan, I. Gijbels, M.P. Wand, Generalized partially linear single-index models, *J. Amer. Statist. Assoc.* 92 (438) (1997) 477–489.
- [5] H. Cho, A. Qu, Model selection for correlated data with diverging number of parameters, *Statist. Sinica* 23 (2) (2013) 901–927.
- [6] T.R. Dawber, G.F. Meadors, F.E. Moore Jr., Epidemiological approaches to heart disease: The Framingham Study, *Am. J. Publ. Health Nations Health* 41 (3) (1951) 279–286.
- [7] J. Fan, R. Li, Variable selection via nonconcave penalized likelihood and its oracle properties, *J. Amer. Statist. Assoc.* 96 (456) (2001) 1348–1360.
- [8] J. Fan, J. Lv, Sure independence screening for ultrahigh dimensional feature space, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 70 (5) (2008) 849–911.
- [9] Y. Fang, Y. Qin, N. Zhang, J. Wang, H. Wang, X. Zheng, DISIS: Prediction of drug response through an iterative sure independence screening, *PLoS One* 10 (3) (2015) 1–13.
- [10] P.W. Franks, Genex environment interactions in type 2 diabetes, *Curr. Diabetes Rep.* 11 (6) (2011) 552.
- [11] K.J. Gaulton, T. Ferreira, Y. Lee, A. Raimondo, R. Mägi, M.E. Reschen, A. Mahajan, et al., Genetic fine mapping and genomic annotation defines causal mechanisms at type 2 diabetes susceptibility loci, *Nature Genet.* 47 (12) (2015) 1415–1425.
- [12] B. Green, H. Lian, Y. Yu, T. Zu, Ultra high-dimensional semiparametric longitudinal data analysis, *Biometrics* (2020).
- [13] L.P. Hansen, Large sample properties of generalized method of moments estimators, *Econometrica* (1982) 1029–1054.
- [14] T. Hastie, R. Tibshirani, M. Wainwright, *Statistical Learning with Sparsity: The LASSO and Generalizations*, CRC Press, 2015.
- [15] X. He, Z.-Y. Zhu, W.-K. Fung, Estimation in a semiparametric model for longitudinal data with unspecified dependence structure, *Biometrika* 89 (3) (2002) 579–590.
- [16] J. Huang, J.L. Horowitz, F. Wei, Variable selection in nonparametric additive models, *Ann. Statist.* 38 (4) (2010) 2282–2313.
- [17] R. Karlsson Linnér, P. Biroli, E. Kong, et al., Genome-wide association analyses of risk tolerance and risky behaviors in over 1 million individuals identify hundreds of loci and shared genetic influences, *Nature Genet.* 51 (2) (2019) 245–257, Number: 2 Publisher: Nature Publishing Group.
- [18] P. Lai, G. Li, H. Lian, Quadratic inference functions for partially linear single-index models with longitudinal data, *J. Multivariate Anal.* 118 (2013) 115–127.
- [19] G. Li, P. Lai, H. Lian, Variable selection and estimation for partially linear single-index models with longitudinal data, *Stat. Comput.* 25 (3) (2015) 579–593.
- [20] J. Lindström, J. Tuomilehto, The diabetes risk score: A practical tool to predict type 2 diabetes risk, *Diabetes Care* 26 (3) (2003) 725–731, Publisher: American Diabetes Association Section: Epidemiology/Health Services/Psychosocial Research.
- [21] S. Ma, H. Liang, C.-L. Tsai, Partially linear single index models for repeated measurements, *J. Multivariate Anal.* 130 (2014) 354–375.
- [22] J.H. Macke, P. Berens, A.S. Ecker, A.S. Tolias, M. Bethge, Generating spike trains with specified correlation coefficients, *Neural Comput.* 21 (2) (2009) 397–423.
- [23] J.B. Meigs, A.K. Manning, C.S. Fox, J.C. Florez, C. Liu, L.A. Cupples, J. Dupuis, Genome-wide association with diabetes-related traits in the Framingham Heart Study, *BMC. Med. Genet.* 8 (1) (2007) 1–10.
- [24] J.R.B. Perry, M.I. McCarthy, A.T. Hattersley, E. Zeggini, Wellcome Trust Case Control Consortium, M.N. Weedon, T.M. Frayling, Interrogating type 2 diabetes genome-wide association data using a biological pathway-based approach, *Diabetes* 58 (6) (2009) 1463–1467.
- [25] R.B. Prasad, L. Groop, Genetics of Type 2 Diabetes—Pitfalls and Possibilities, *Genes* 6 (1) (2015) 87–123, <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4377835/>.
- [26] A. Qu, R. Li, Quadratic inference functions for varying-coefficient models with longitudinal data, *Biometrics* 62 (2) (2006) 379–391.
- [27] A. Qu, B.G. Lindsay, B. Li, Improving generalised estimating equations using quadratic inference functions, *Biometrika* 87 (4) (2000) 823–836.
- [28] D. Ruppert, R.J. Carroll, Theory & methods: Spatially-adaptive penalties for spline fitting, *Aust. N. Z. J. Stat.* 42 (2) (2000) 205–223.
- [29] L. Schumaker, *Spline Functions: Basic Theory*, Cambridge University Press, 2007.
- [30] M.P. Stern, K. Williams, S.M. Haffner, Identification of persons at high risk for type 2 diabetes mellitus: Do we need the oral glucose tolerance test? *Ann. Intern. Med.* 136 (8) (2002) 575–581.
- [31] J.Y. Taylor, Y.V. Sun, S.C. Hunt, S.L.R. Kardia, Gene-environment interaction for hypertension among African American women across generations, *Biol. Res. Nurs.* 12 (2) (2010) 149–155.
- [32] R. Tibshirani, Regression shrinkage and selection via the Lasso, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 58 (1) (1996) 267–288, <https://www.jstor.org/stable/2346178>.

- [33] T. Vos, C. Allen, M. Arora, R.M. Barber, Z.A. Bhutta, A. Brown, A. Carter, D.C. Casey, F.J. Charlson, A.Z. Chen, et al., Global, regional, and national incidence, prevalence, and years lived with disability for 310 diseases and injuries, 1990–2015: A systematic analysis for the global burden of disease study 2015, *Lancet* 388 (10053) (2016) 1545–1602.
- [34] H. Wang, B. Li, C. Leng, Shrinkage tuning parameter selection with a diverging number of parameters, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 71 (3) (2009) 671–683.
- [35] L. Wang, X. Liu, H. Liang, R.J. Carroll, Estimation and variable selection for generalized additive partial linear models, *Ann. Statist.* 39 (4) (2011) 1827.
- [36] P. Wang, G.-f. Tsai, A. Qu, Conditional inference functions for mixed-effects models with unspecified random-effects distribution, *J. Amer. Statist. Assoc.* 107 (498) (2012) 725–736.
- [37] L. Wang, L. Xue, A. Qu, H. Liang, Estimation and model selection in generalized additive partial linear models for correlated data with diverging number of covariates, *Ann. Statist.* 42 (2) (2014) 592–624.
- [38] L. Wang, J. Zhou, A. Qu, Penalized generalized estimating equations for high-dimensional longitudinal data analysis, *Biometrics* 68 (2) (2012) 353–360.
- [39] A.H. Young, I.N. Ferrier, S.G. Ball, M.K. Mohiuddin, C.E. Todhunter, J.C. Mansfield, et al., Genome-wide association study of 14,000 cases of seven common diseases and 3000 shared controls, *Nature* (2007).
- [40] Y. Yu, D. Ruppert, Penalized spline estimation for partially linear single-index models, *J. Amer. Statist. Assoc.* 97 (460) (2002) 1042–1054.
- [41] Y. Yu, C. Wu, Y. Zhang, Penalised spline estimation for generalised partially linear single-index models, *Stat. Comput.* 27 (2) (2017) 571–582.
- [42] S.L. Zeger, K.-Y. Liang, Longitudinal data analysis for discrete and continuous outcomes, *Biometrics* 42 (1) (1986) 121–130.
- [43] Y. Zhang, H. Lian, Y. Yu, Estimation and variable selection for quantile partially linear single-index models, *J. Multivariate Anal.* 162 (2017) 215–234.
- [44] J. Zhou, A. Qu, Informative estimation and selection of correlation structure for longitudinal data, *J. Amer. Statist. Assoc.* 107 (498) (2012) 701–710.