

Estimating High Dimensional Monotone Index Models by Iterative Convex Optimization*

Shakeeb Khan¹, Xiaoying Lan¹, Elie Tamer², and Qingsong Yao¹

¹Dept. of Economics, Boston College

{ shakeeb.khan Xiaoyang.lan qingsong.yao}@bc.edu

²Dept. of Economics, Harvard University

{elietamer}@fas.harvard.edu

First Version: 6/2021, This Version: March 7, 2023

Abstract

We propose new approaches to estimating large dimensional semiparametric monotone index models. This class of models has been popular in the applied and theoretical econometrics literatures as it includes discrete choice, nonparametric transformation, and duration models. A main advantage of our approach is computational. For instance, rank estimation procedures such as those proposed in [Han \(1987\)](#) and [Cavanagh and Sherman \(1998\)](#) that optimize a nonsmooth, non convex objective function are difficult to optimize with more than a few regressors. Some recent progress is the work by [Ahn, Ichimura, Powell, and Ruud \(2018\)](#), but it too is only suitable for small dimensional models. Thus for such monotone index models with large, or even increasing dimension, we propose a new class of semiparametric sieve and kernel based estimators based on *batched gradient descent (BGD)*, and study their asymptotic properties. The BGD algorithm uses an iterative procedure where the key step exploits a strictly convex objective function, resulting in computational advantages.

Key Words Monotone Index models, Convex Optimization, Kernel and Sieve Estimation.

*We are grateful to conference participants at the BC/BU 2020 Econometrics Workshop, the 2019 Mid-western Econometrics Study group, the 2021 NASM of the Econometric Society, 2022 CIRAQ Econometrics conference, 2022 ISNPS conference, 2022 Advanced Methods Conference at TSE, and seminar participants from Georgetown, UC Berkeley, UC Louvain, UC Riverside, University of Bristol, UVA, University of Warwick and Yale for helpful comments.

1 Introduction

Monotone index models have received a great deal of attention in both the theoretical and applied econometrics literature, as many economic variables of interest are of a limited or qualitative nature. A leading special case in this class is the binary choice model which is usually represented by some variation of the following equation:

$$y_i = I[\mathbf{X}_{e,i}^T \boldsymbol{\beta}_e^* - u_i \geq 0] \quad (1)$$

where $I[\cdot]$ is the usual indicator function, y_i is the observed response variable, taking the values 0 or 1 and $\mathbf{X}_{e,i} = (X_{0,i}, \mathbf{X}_i^T)^T$ is an observed $p + 1$ dimensional vector of covariates which effect the behavior of y_i . Both the scalar disturbance term u_i with distribution function denoted by $G(\cdot)$, and the $(p + 1)$ - dimensional vector $\boldsymbol{\beta}_e^* = (\beta^*, \boldsymbol{\beta}^{*T})^T$ are unobserved, the latter often being the parameter estimated from a random sample $(y_i, \mathbf{X}_{e,i})$, $i = 1, 2, \dots, n$.

The disturbance term u_i is restricted in ways that ensure identification of $\boldsymbol{\beta}_e^*$. Parametric restrictions specify the distribution of u_i up to a finite dimensional parameter and assume that u_i distributed independently of the covariates \mathbf{X}_i . Under such a restriction, $\boldsymbol{\beta}_e^*$ can be estimated (up to scale) using maximum likelihood or nonlinear least squares. Estimators that are robust to these parametric distributional assumptions have been proposed and analyzed resulting in a variety of estimation procedures for $\boldsymbol{\beta}_e^*$.

An important class of semiparametric restrictions used in the literature were based on independence/index restrictions. Estimation procedures under this restriction include those proposed by Han (1987), Ichimura (1993), Klein and Spady (1993). These cover but are not limited to the above binary response model. This class of index models have a robustness advantage over parametric approaches, but estimators within this class are difficult to compute¹ due to nonconvexity and in some cases also nonsmoothness of their respective objective functions. For these objective functions, even looking for a local optimum is generally NP-Hard, let alone the global optimum (Murty and Kabadi, 1987). Furthermore the difficulty increases with the dimension of \mathbf{X}_i . Recent work which is motivated by computational concerns is Ahn, Ichimura, Powell, and Ruud (2018). However, their two step procedure involves a fully nonparametric estimator in the first stage, so is also not suitable for models with a large number of regressors.

A related drawback of all these procedures is that they are designed to estimate parameters in models of a small and *fixed* dimension. A relatively recent and thriving literature in

¹Other estimation of index models includes Stoker (1986) and Powell et al. (1989). While these are relatively easy to compute, such derivative based estimators cannot be applied unless all components of $\mathbf{X}_{e,i}$ are continuously distributed.

econometrics and machine learning is recognizing the many advantages of allowing for large dimensional models or models with a large set of controls. This class is a special case of models that consider the situation when the dimension of x_i is large, and this is now often modeled with its dimension increasing with the sample size. Due primarily to its empirical relevance, there has been a burgeoning literature on estimation and inference in certain econometric and statistics models with a large number of regressors or a large number of moment conditions. For a survey of examples in economics and finance, see [Fan et al. \(2020\)](#). Recent papers include [Newey and Windmeijer \(2009\)](#), [Chernozhukov et al. \(2017\)](#), [Belloni et al. \(2018\)](#), [Cattaneo et al. \(2018\)](#).

Related to our work is the recent literature on estimating large dimensional binary choice or monotone index models in [Sur and Candès \(2019\)](#) and [Fan et al. \(2020\)](#). [Sur and Candès \(2019\)](#) considers inference in a large dimensional logit model, where it is shown that χ^2 asymptotic approximations to the LR statistic are suspect when the dimension of x is large. [Fan, Han, Li, and Zhou \(2020\)](#) on the other hand estimate parameters by optimizing the objective function introduced in [Han \(1987\)](#), but with the number parameters increasing with the sample size. Optimizing these rank based objective functions is unfortunately hard even with recent developments in algorithms and search methods for optimizing non smooth and/or non convex objective functions. See for example important recent work based on mixed integer programming (MIP) as in, e.g. [Fan et al. \(2020\)](#) and [Shin and Todorov \(2021\)](#).

Therefore, in light of the drawbacks in the existing literature, this paper proposes a new estimation procedure that is amenable to easier computation. Specifically we aim to construct a computationally feasible estimator for a semiparametric binary choice and monotone index models with *increasing* dimension based on a convex objective function and then establish its asymptotic properties. As we will discuss in detail in the next section, our algorithm uses an iterative estimator based on a batched gradient descent (BGD) method, and we show how to use nonparametric methods to approximate the distribution in each stage of the iteration. One is the method of sieves², and the other is kernel regression. Finally, our proof in the semiparametric case requires development of approaches to handle estimators that are defined recursively while at the same time allow for an unknown link function. The paper starts out analyzing properties of the BGD estimator in parametric models with increasing dimensions. The following section considers the main case when we allow for the link function to be estimated via kernel or sieve methods. Finally a Monte Carlo Section examines the computational advantages of our approach. Also, we provide an empirical illustration that highlights the behavior of our estimator with real data.

²See [Chen \(2007\)](#) who pioneered the use of sieve methods in econometrics.

Notation: Throughout the rest of this paper, to facilitate the description and properties of estimation procedures we will be using the following notation. For any real sequences $\{a_n\}_{n=1}^\infty$ and $\{b_n\}_{n=1}^\infty$, we write $a_n = o(b_n)$ if $\limsup_{n \rightarrow \infty} |a_n/b_n| = 0$, $a_n = O(b_n)$ if $\limsup_{n \rightarrow \infty} |a_n/b_n| < \infty$, and $a_n \sim b_n$ if both $a_n = O(b_n)$ and $b_n = O(a_n)$. For any random sequences $\{a_n\}_{n=1}^\infty$ and $\{b_n\}_{n=1}^\infty$, we write $a_n = O_p(b_n)$ if for any $0 < \tau < 1$ there are N and $C > 0$ such that $P\{|a_n/b_n| > C\} < \tau$ holds for all $n \geq N$, we write $a_n = o_p(b_n)$ if for any $C > 0$, $\lim_{n \rightarrow \infty} P\{|a_n/b_n| > C\} \rightarrow 0$. For any Borel sets $A \subseteq \mathbb{R}^k$, denote its Lebesgue measure as $m(A)$. For any symmetric matrix A , we write $A \succ 0$ if A is positive definite, and $A \succeq 0$ if A is positive semi-definite. For any symmetric matrices A and B , we write $A \succ B$ if $A - B \succ 0$ and $A \succeq B$ if $A - B \succeq 0$. For any matrix A , we denote $\sigma(A)$ as its singular value, and denote $\bar{\sigma}(A)$ and $\underline{\sigma}(A)$ as its largest and smallest singular value. For any symmetric matrix A , we denote $\lambda(A)$ as its eigenvalue, and denote $\bar{\lambda}(A)$ and $\underline{\lambda}(A)$ as its largest and smallest eigenvalue. For any vector $\mathbf{x} = (x_1, \dots, x_p)^\top$, we denote its Euclidean norm as $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^p x_i^2}$. For any matrices $A = (a_{ij})_{n \times m}$, we denote $\|A\| = \sqrt{\sum_{i=1}^n \sum_{j=1}^m a_{ij}^2}$. Note that when A is positive semi-definite, there holds $\|A\mathbf{x}\| \leq \bar{\lambda}(A) \cdot \|\mathbf{x}\|$; for general square matrix A , there holds $\|A\mathbf{x}\| \leq \bar{\sigma}(A) \cdot \|\mathbf{x}\|$. Finally, for any function $f(\mathbf{x})$ with domain D , define $\|f\|_\infty = \sup_{\mathbf{x} \in D} f(\mathbf{x})$.

2 The BGD Estimator

To provide some intuition for our semiparametric estimators that will be introduced in the following sections, we first consider here a simplified version of the model where the cumulative distribution function $G(\cdot)$ is completely known. Under such setup, we explore the *batch gradient descent estimator* (BGD estimator) of β_e^* when its dimensionality p may increase, which is also important on its own right. Throughout the following analysis we assume that the data set satisfies the following assumption.

Assumption 1. An i.i.d. data set $\mathcal{D}_n = \{(\mathbf{X}_{e,i}, y_i)\}_{i=1}^n$ of sample size n is observed, where y_i is generated³ by $y_i = I(X_{0,i}\beta_0^* + \mathbf{X}_i^\top \beta^* - u_i > 0)$ with unobserved shock u_i that is independent of $\mathbf{X}_{e,i}$ and has CDF $G(\cdot)$.

Given any loss function $\ell_G(\beta_e, \mathbf{X}_e, y)$ that depends on G and is a.s. differentiable with

³Here we are decomposing the vector $\mathbf{X}_{e,i}$ into a scalar component $X_{0,i}$ and the vector \mathbf{X}_i , and decomposing the vector of parameters β_e^* into the scalar term β_0^* and the vector β^* . As we will see this is done for notational convenience when imposing scale normalizations.

respect to $\beta_e \in \mathcal{B}_e$, the BGD estimator of β_e^* is based on the following iteration,

$$\beta_{e,k+1} = \beta_{e,k} - \frac{\delta_k}{n} \sum_{i=1}^n \partial \ell_G(\beta_{e,k}, \mathbf{X}_{e,i}, y_i) / \partial \beta_e, \quad (2)$$

where $\delta_k > 0$ is the learning rate. Note that $n^{-1} \sum_{i=1}^n \partial \ell_G(\beta_e, \mathbf{X}_{e,i}, y_i) / \partial \beta_e$ constitutes a sample analogue of the derivative $\partial \mathbb{E}[\ell_G(\beta_e, \mathbf{X}_e, y)] / \partial \beta_e$. Unlike the stochastic gradient descent (SGD) algorithm, in the BGD algorithm, in each round of update we evaluate the derivative of the loss function over all data points. This increases the computational burden but provides a more accurate estimator for the derivative of the expected loss function. Given the initial guess of the parameter, $\beta_{e,1}$, we iterate based on (2) until some terminating conditions are reached.

In this paper, we consider the following loss function

$$\ell_G(\beta_e, \mathbf{X}_e, y) = \int_{-A}^{\mathbf{X}_e^T \beta_e} G(z) dz - y \mathbf{X}_e^T \beta_e, \quad (3)$$

for some sufficiently large positive constant A . The loss function (3) was also considered in Agarwal et al. (2014) and has many properties. For instance, under some mild conditions, it is easy to show that at the truth,

$$\begin{aligned} \frac{\partial \mathbb{E}(\ell_G(\beta_e^*, \mathbf{X}_e, y))}{\partial \beta_e} &= \mathbb{E} \{ (G(\mathbf{X}_e^T \beta_e^*) - y) \mathbf{X}_e \} \\ &= \mathbb{E} \{ (G(\mathbf{X}_e^T \beta_e^*) - \mathbb{E}(y | \mathbf{X}_e)) \mathbf{X}_e \} = 0, \end{aligned}$$

and

$$\frac{\partial^2 \mathbb{E}(\ell_G(\beta_e, \mathbf{X}_e, y))}{\partial \beta_e \partial \beta_e^T} = \mathbb{E} \{ G'(\mathbf{X}_e^T \beta_e) \mathbf{X}_e \mathbf{X}_e^T \} \succ 0, \forall \beta_e \in \mathcal{B}_e.$$

So β_e^* uniquely minimizes $\mathbb{E} \ell_G(\beta_e, \mathbf{X}_e, y)$ over \mathcal{B}_e . Another desirable property of the loss function (3) is that the derivative of (3) with respect to β_e , which is $(G(\mathbf{X}_e^T \beta_e) - y) \mathbf{X}_e$, depends only on $G(\cdot)$ instead of on its derivatives. So when we conduct a semiparametric iteration in the following sections, we only need to nonparametrically approximate $G(\cdot)$, which is generally more robust compared with approximating its derivatives. Based on loss function (3), the BGD estimator is obtained by using the following iteration procedure:

$$\beta_{e,k+1} = \beta_{e,k} - \frac{\delta_k}{n} \sum_{i=1}^n (G(\mathbf{X}_{e,i}^T \beta_{e,k}) - y_i) \mathbf{X}_{e,i}. \quad (4)$$

We summarize our algorithm as follows in [algorithm 1](#).

Algorithm 1: The BGD Estimator

input : Data set $\{(\mathbf{X}_{e,i}, y_i)\}_{i=1}^n$, sequence of learning rate $\{\delta_k\}_{k=1}^\infty$, initial guess $\beta_{e,1}$, CDF $G(\cdot)$, and terminating condition \mathcal{T}

output: The BGD estimator $\hat{\beta}_e$

```

1  $k \leftarrow 1$ ;
2 while The terminating condition  $\mathcal{T}$  is not satisfied do
3    $\beta_{e,k+1} \leftarrow \beta_{e,k} - \frac{\delta_k}{n} \sum_{i=1}^n (G(\mathbf{X}_{e,i}^\top \beta_{e,k}) - y_i) \mathbf{X}_{e,i}$ ;
4    $k \leftarrow k + 1$ ;
5  $\hat{\beta}_e \leftarrow \beta_{e,k}$ ;
```

Remark 1. Key to the above approach is the construction of a convex objective function that facilitates computation even with high dimensions. This transformed convex objective works for any monotone model. In particular, for any model of the form $y_i = G(\mathbf{X}_{e,i}^\top \beta_e) + \varepsilon_i$ with $E[\varepsilon_i | \mathbf{X}_{e,i}] = 0$ and monotone $G(\cdot)$, a similar *convex* criterion as in (3) can be used for inference on β_e .

We now describe the asymptotic properties of $\beta_{e,k}$. We first make the following assumption.

Assumption 2. (i) $\mathcal{X}_e = [-1, 1]^{p+1}$; (ii) \mathcal{B}_e is convex, and there exists some constant $B_0 > 0$ such that for any $\beta_e \in \mathcal{B}_e$, $|\beta_j| \leq B_0$ for any $0 \leq j \leq p$; (iii) there exists integer v_G such that G has up to v_G -th bounded derivatives; (iv) Define $M_n(\beta_e) = \frac{1}{n} \sum_{i=1}^n G'(\mathbf{X}_{e,i}^\top \beta_e) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^\top$ and $M(\beta_e) = \mathbb{E}[M_n(\beta_e)]$. For any $\beta_e \in \mathcal{B}_e$, there holds $0 < \underline{\lambda}_e \leq \underline{\lambda}(M(\beta_e)) \leq \bar{\lambda}(M(\beta_e)) \leq \bar{\lambda}_e < \infty$.

Remark 2. [Assumption 2\(i\)](#) and [Assumption 2\(ii\)](#) are convenient normalizations that facilitate the assessment of our model. Note that to ensure that $\beta_{e,k}$ falls into a compact set for each k , some form of truncation on $\beta_{e,k+1}$ in (4) is needed. While according to our results below, as long as \mathcal{B}_e is sufficiently large, it can be shown that $\beta_{e,k}$ will fall into \mathcal{B}_e for all k with probability going to 1. We then assume that $\beta_{e,k} \in \mathcal{B}_e$ for all k . [Assumption 2\(iii\)](#) imposes some smoothness conditions on G , where the requirement on v_G will be stated in the following propositions and theorems. [Assumption 2\(iv\)](#) requires that the eigenvalue of $M_n(\beta_e)$ is bounded from both below and above uniformly over \mathcal{B}_e .

For any $\beta_e \in \mathcal{B}_e$, define $\Delta\beta_e = \beta_e - \beta_e^*$. Also define $\varepsilon_i = y_i - G(\mathbf{X}_{e,i}^\top \beta_e^*)$, where $\mathbb{E}[\varepsilon_i | \mathbf{X}_{e,i}] = 0$. When [Assumption 1](#) and [Assumption 2](#) hold, we have the following result.

Theorem 1. Suppose that [Assumption 1](#) and [Assumption 2](#) hold with $v_G = 3$, that $p^5 (\log p)^2 n^{-1} \rightarrow 0$, that the learning rate is chosen such that $\delta_k = \delta \leq 2 / (3\bar{\lambda}_e)$, and that β_e is updated based on [algorithm 1](#). We have that

(i) Define

$$k_{1,n}^{BGD} = \frac{\log \|\Delta \beta_{e,1}\| + \frac{1}{2} \log (n / (p \log p))}{-\log (1 - \underline{\lambda}_e \delta / 2)},$$

we then have

$$\sup_{k \geq k_{1,n}^{BGD} + 1} \|\Delta \beta_{e,k}\| = O_p \left(\sqrt{p (\log p) / n} \right);$$

(ii) Define $k_{2,n}^{BGD}$ such that $(1 - \underline{\lambda}_e \delta)^{k_{2,n}^{BGD}} \sqrt{p \log p} \rightarrow 0$, we have

$$\sup_{k \geq k_{2,n}^{BGD} + 1} \left\| \Delta \beta_{e,k+k_{1,n}^{BGD}} - M^{-1}(\beta_e^*) \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| = o_p(1/\sqrt{n});$$

(iii) For any $k \geq k_{1,n}^{BGD} + k_{2,n}^{BGD} + 1$, define $\hat{\beta}_e = \hat{\beta}_k$. Also define

$$\Sigma_1^* = M^{-1}(\beta_e^*) \mathbb{E} [G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T] M^{-1}(\beta_e^*),$$

and

$$\hat{\Sigma}_{1,n} = M_n^{-1}(\hat{\beta}_e) \left\{ \frac{1}{n} \sum_{i=1}^n \hat{G}_i (1 - \hat{G}_i) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T \right\} M_n^{-1}(\hat{\beta}_e),$$

where $G_i^* = G(\mathbf{X}_{e,i}^T \beta_e^*)$ and $\hat{G}_i = G(\mathbf{X}_{e,i}^T \hat{\beta}_e)$. Suppose further that $\mathbb{E}(\mathbf{X}_{e,i} \mathbf{X}_{e,i}^T)$ has uniformly (with respect to p) upper bounded eigenvalues, there holds

$$\|\hat{\Sigma}_{1,n} - \Sigma_1^*\| \rightarrow_p 0.$$

(iv) For any $p+1$ vector ρ such that $\lim_{n \rightarrow \infty} \|\rho\| < \infty$, $\lim_{n \rightarrow \infty} \rho^T \Sigma_1^* \rho = \sigma^2(\rho)$, and that $\rho^T M^{-1}(\beta_e^*) \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \rightarrow_d N(0, \sigma^2(\rho))$, we have that

$$\rho^T \Delta \hat{\beta}_e / \sqrt{\hat{\sigma}^2(\rho) / n} \rightarrow_d N(0, 1),$$

where $\hat{\sigma}^2(\rho) = \rho^T \hat{\Sigma}_{1,n} \rho$.

Proof of [Theorem 1](#). See [Appendix B](#). □

When p is fixed, [Theorem 1](#)(i) implies that $\sup_{k \geq k_{1,n}^{BGD} + 1} \|\Delta \beta_{e,k}\| = O_p(1/\sqrt{n})$, and [Theorem 1](#)(ii) implies that for k sufficiently large, the BGD estimator is an asymptotically

linear estimator, so there holds $\sqrt{n}\Delta\boldsymbol{\beta}_{e,k+k_{1,n}^{BGD}} \rightarrow_d N(0, \Sigma_1^*)$ by the central limit theorem. The asymptotic variance can be estimated based on [Theorem 1\(iii\)](#). The number of iterations required to obtain root- n consistency, $k_{1,n}^{BGD}$, is determined by many factors including the sample size n , the distance between the true parameter and the initial guess $\|\Delta\boldsymbol{\beta}_{e,1}\|$, as well as the lower bound of the eigenvalues of $M_n(\boldsymbol{\beta}_e)$. In general, $k_{1,n}^{BGD}$ is of order $O(\log n)$, but in practice when we apply the above algorithm, the specific number of iteration is difficult to determine. For detailed discussion of the number of iterations, see [Remark 5](#) at the end of [Section 4](#). The inference on $\boldsymbol{\beta}_e^*$ based on the BGD estimator is given by [Theorem 1\(iv\)](#). Note that for any given vector ρ , we require that $\frac{1}{\sqrt{n}}\rho^T M^{-1}(\boldsymbol{\beta}_e^*) \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i}$ is asymptotically normally distributed. An alternative approach is to apply the high-dimensional central limit theorem to $\frac{1}{n} \sum_{i=1}^n M^{-1}(\boldsymbol{\beta}_e^*) \mathbf{X}_{e,i} \varepsilon_i$ (e.g., [Chernozhukov et al., 2017](#)).

Before we conclude this section and move to semiparametric estimation, we further comment on [Theorem 1](#). Different from the stochastic gradient descent algorithm (e.g., [Toulis and Airolidi, 2017](#)), we show in [Theorem 1](#) that the learning rate δ_k can be selected as a sufficiently small constant. Indeed, in the following results, we show that δ_k can decay to zero at any rate as long as $\sum_{k=1}^{\infty} \delta_k = \infty$ holds, and the choice of δ_k will not change the asymptotic results displayed in [Theorem 1](#). In particular, we have the following proposition.

Theorem 2. *Suppose that all the conditions in [Theorem 1](#) hold and that $\boldsymbol{\beta}_e$ is updated based on [algorithm 1](#). For any sequence of tuning parameters $\{\delta_k\}_{k=1}^{\infty}$ satisfying $\delta_k \geq 0$, $\delta_k \rightarrow 0$, $\limsup_{k \rightarrow \infty} \delta_{k-1}/\delta_k < \infty$, and $\sum_{k=1}^{\infty} \delta_k = \infty$, we have that*

(i) *Define $\tilde{k}_{1,n}^{BGD}$ such that $\sum_{k=1}^{\tilde{k}_{1,n}^{BGD}} \delta_k \geq \underline{\lambda}_e^{-1} \{\log(n/p(\log p)) + 2 \log \|\Delta\boldsymbol{\beta}_{e,1}\|\}$, and that $\sup_{k \geq \tilde{k}_{1,n}^{BGD}+1} \delta_k \leq 2/\underline{\lambda}_e$, then there holds*

$$\sup_{k \geq \tilde{k}_{1,n}^{BGD}+1} \|\Delta\boldsymbol{\beta}_{e,k}\| = O_p\left(\sqrt{p(\log p)/n}\right);$$

(ii) *Define $\tilde{k}_{2,n}^{BGD}$ such that $\sum_{k=\tilde{k}_{1,n}^{BGD}+1}^{k=\tilde{k}_{2,n}^{BGD}} \delta_k / \log p \rightarrow \infty$, then we have that*

$$\sup_{k \geq \tilde{k}_{2,n}^{BGD}+1} \left\| \Delta\boldsymbol{\beta}_{e,k+\tilde{k}_{1,n}^{BGD}} - M^{-1}(\boldsymbol{\beta}_e^*) \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| = o_p(1/\sqrt{n});$$

(iii) *For any $k \geq \tilde{k}_{1,n}^{BGD} + \tilde{k}_{2,n}^{BGD} + 1$, define $\hat{\boldsymbol{\beta}}_e = \hat{\boldsymbol{\beta}}_k$. We have that [Theorem 1\(iii\)](#) and (iv) hold.*

Proof of [Theorem 2](#). See [Appendix B](#). □

Theorem 2 shows that the choice of the learning rate basically does not affect the convergence rate as well as the asymptotic distribution of the BGD estimators. The main advantage of using a sequence of decaying learning rates is that we do not need to choose the constant δ as required in **Theorem 1**, since for k sufficiently large, $\delta_k \leq 2/(3\bar{\lambda}_e)$ will automatically hold. However, the disadvantage of using decaying learning rates is that such procedure takes much longer time to converge because the magnitude of the update in the k -th round decreases as k increases. For instance, suppose that we choose $\delta_k \sim k^{-v}$ for some $0 \leq v < 1$, we have that $\sum_{j=1}^k \delta_j \sim k^{1-v}$. Then to ensure that $\sum_{j=1}^{\tilde{k}_{1,n}^{BGD}} \delta_j \geq \lambda_e^{-1} (\log n + 2 \log \|\Delta \beta_{e,1}\|)$, we need $\tilde{k}_{1,n}^{BGD} \sim (\log n)^{\frac{1}{1-v}}$. Obviously, setting $v = 0$ leads to $k \sim \log n$, which corresponds to the requirement in **Theorem 1**(i); when $v > 0$, we can see that more rounds of iteration is needed compared with required in **Theorem 1**(i).

3 Semiparametric BGD Estimation

In the previous section, we focused on iterative estimators based on the BGD algorithm for the parametric binary choice models. We show that when the CDF of the error term is known, the iterative estimators based on the BGD algorithm are consistent and attain asymptotic normality under mild conditions. However, having prior knowledge of the form of G is generally too strong an assumption. In most applications, the source of the individual shock u in **Assumption 1** is difficult to justify, which makes it quite difficult, if not completely impossible, to know the exact expression of G . In this scenario, the algorithm proposed in the previous section is infeasible. To overcome such problem, this section generalizes the BGD estimator proposed in Section 2 to the semiparametric setting where G is unknown.

In this setup, to ensure identification we set β_0^* to be 1, so our estimation target is β^* . To simplify our notation, we denote the space of \mathbf{X} as \mathcal{X} , and the corresponding parameter space of β as \mathcal{B} . Suppose that an initial guess for β^* is given by β_1 . In the k -th round of iteration, to update β based on the BGD algorithm, we require the knowledge of G as in Section 2, which is infeasible when G is unknown. A natural idea is that we can construct an estimator for G based on the index constructed from the updated parameter in the previous round. More intuitively, suppose for a moment that in the k -th round of iteration, β_k happens to be identical to the unknown true parameter β^* , then we have that $G(z) = \mathbb{E}[y | X_0 + \mathbf{X}^T \beta^* = z] = \mathbb{E}[y | X_0 + \mathbf{X}^T \beta_k = z]$ for any $z \in R$.

This motivates semiparametric estimation by using nonparametric methods to estimate $G(\cdot)$. We consider kernel estimation and the method of sieves in each of the following subsections.

3.1 The KBGD Estimator

In this section we consider kernel estimation to estimate $G(\cdot)$. The Nadaraya-Watson kernel estimator of $G(\cdot)$ is of the form

$$\widehat{G}(z|\boldsymbol{\beta}_k) = \frac{\sum_{j=1}^n K_{h_n}(z - X_{0,j} - \mathbf{X}_j^T \boldsymbol{\beta}_k) y_j}{\sum_{j=1}^n K_{h_n}(z - X_{0,j} - \mathbf{X}_j^T \boldsymbol{\beta}_k)}, z \in R, \quad (5)$$

where $K_h(\cdot) = h^{-1}K(\cdot/h)$, $K(\cdot)$ is some kernel function, and h_n is some bandwidth parameter depending on n . Given the estimated CDF $\widehat{G}(\cdot|\boldsymbol{\beta}_k)$, we can update the parameter as if it were the true CDF $G(\cdot)$. In particular, $\boldsymbol{\beta}_k$ is updated as

$$\boldsymbol{\beta}_{k+1} = \boldsymbol{\beta}_k - \frac{\delta_k}{n} \sum_{i=1}^n \left(\widehat{G}(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k | \boldsymbol{\beta}_k) - y_i \right) \mathbf{X}_i. \quad (6)$$

Keep updating $\boldsymbol{\beta}_k$ based on (5) and (6), until some terminating conditions are reached. The resulting estimator is labeled as the *kernel-based batch gradient descent estimator* (KBGD estimator). We summarize our algorithm as follows in [algorithm 2](#).

Algorithm 2: The KBGD Estimator

input : Data set $\{(\mathbf{X}_{e,i}, y_i)\}_{i=1}^n$, sequence of learning rate $\{\delta_k\}_{k=1}^\infty$, initial guess $\boldsymbol{\beta}_1$, kernel function K , bandwidth h_n , and terminating condition \mathcal{T}

output: The KBGD estimator $\widehat{\boldsymbol{\beta}}$

```

1  $k \leftarrow 1$ ;
2 while The terminating condition  $\mathcal{T}$  is not satisfied do
3   for  $i \leftarrow 1$  to  $n$  do
4      $\widehat{G}(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k | \boldsymbol{\beta}_k) \leftarrow \frac{\sum_{j=1}^n K_{h_n}(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k - X_{0,j} - \mathbf{X}_j^T \boldsymbol{\beta}_k) y_j}{\sum_{j=1}^n K_{h_n}(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k - X_{0,j} - \mathbf{X}_j^T \boldsymbol{\beta}_k)}$ ;
5      $\boldsymbol{\beta}_{k+1} \leftarrow \boldsymbol{\beta}_k - \frac{\delta_k}{n} \sum_{i=1}^n \left( \widehat{G}(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k | \boldsymbol{\beta}_k) - y_i \right) \mathbf{X}_{e,i}$ ;
6      $k \leftarrow k + 1$ ;
7  $\widehat{\boldsymbol{\beta}} \leftarrow \boldsymbol{\beta}_k$ ;
```

Remark 3. In essence, the KBGD estimator can not be classified as a BGD estimator based on a semiparametric loss function. In the semiparametric setup, given any loss function $\ell_G(\boldsymbol{\beta}, \mathbf{X}_e, y)$ (quadratic distance in [Ichimura \(1993\)](#), log-likelihood in [Klein and Spady \(1993\)](#), or loss function given in (3)) with unknown function G , it's a common practice to replace G with its nonparametric estimator \widehat{G} and then minimize (or maximize) the resulting loss function to obtain the estimator of $\boldsymbol{\beta}$. Note that under the single-index framework, \widehat{G} usually

involves the unknown parameter β , which is of the form $\widehat{G}(\cdot) = \widehat{G}(\cdot|\beta)$. In this scenario, the BGD estimator is constructed by the following iteration

$$\beta_{k+1}^{BGD} = \beta_k^{BGD} - \frac{\delta_k}{n} \sum_{i=1}^n \frac{\partial \ell_{\widehat{G}(\cdot|\beta_k^{BGD})}(\beta_k^{BGD}, \mathbf{X}_{e,i}, y_i)}{\partial \beta},$$

where $\partial \ell_{\widehat{G}(\cdot|\beta_k^{BGD})}(\beta_k^{BGD}, \mathbf{X}_{e,i}, y_i) / \partial \beta$ involves $\partial \widehat{G}(\cdot|\beta_k) / \partial \beta$, a complicated functions of β_k . In particular, the BGD estimator under loss function (3) is given by

$$\beta_{k+1} = \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n \left(\widehat{G}(X_{0,i} + \mathbf{X}_i^T \beta_k | \beta_k) + \int_{-\infty}^{X_{0,i} + \mathbf{X}_i^T \beta_k} \frac{\partial \widehat{G}(z | \beta_k)}{\partial \beta} dz - y_i \right) \mathbf{X}_i.$$

Obviously, an additional term is introduced compared with (6). On the contrary, during the construction (6), we take G as given when taking the first order derivative of the loss function and then replace the unknown G with its non-parametric estimator in the derivative. More specifically, the KBGD estimator is updated as follows

$$\beta_{k+1} = \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n \frac{\partial \ell_G(\beta_k, \mathbf{X}_{e,i}, y_i)}{\partial \beta} \Big|_{G(\cdot) = \widehat{G}(\cdot|\beta_k)},$$

so additional terms involving $\partial \widehat{G}(\cdot|\beta_k) / \partial \beta$ are avoided. Finally, as we discussed in Section 2, the derivative of loss function (3) with respect to β depends only on G , so we also avoid approximating the derivative of G , which has poorer finite-sample performance compared with approximating G . Such update also ensures contraction map under some conditions, see Assumption 5.

For any fixed z and β , under mild conditions there holds $\widehat{G}(z|\beta) \rightarrow_p \mathbb{E}[y | X_0 + \mathbf{X}^T \beta = z]$. Denote such limit as $L(z, \beta)$. Obviously, $L(z, \beta^*) = G(z)$ holds for any $z \in \mathbb{R}$. Before we move to a formal description of the statistical properties of the KBGD estimator based on (6), we first provide some further discussion on $L(z, \beta)$. For simplicity, in the following we only focus on the case where all the covariates are continuous which permit continuous joint density function. We leave further discussion of the case where some covariates are discrete to Remark 6. We point that when there are discrete covariates, our algorithm can be directly applied without any modification, although some further assumptions will be required.

When all the covariates are continuous, denote the joint density of \mathbf{X}_e and \mathbf{X} as $f_e(\mathbf{X}_e) = f_e(X_0, \mathbf{X})$ and $f(\mathbf{X}) = \int f_e(X_0, \mathbf{X}) dX_0$, respectively. Denote $z(\mathbf{X}_e, \beta) = X_0 + \mathbf{X}^T \beta$. Also denote $f_{\mathbf{X},z}(\mathbf{X}, z|\beta)$ as the joint density of \mathbf{X} and $z(\mathbf{X}_e, \beta)$ given β . Note that for any \mathbf{x}

and z ,

$$\begin{aligned} P[\mathbf{X} \leq \mathbf{x}, z(\mathbf{X}_e, \boldsymbol{\beta}) \leq z] &= \int_{\tilde{\mathbf{X}} \leq \mathbf{x}, \tilde{X}_0 + \tilde{\mathbf{X}}^\top \boldsymbol{\beta} \leq z} f_e(\tilde{X}_0, \tilde{\mathbf{X}}) d\tilde{X}_0 d\tilde{\mathbf{X}} \\ &= \int_{\tilde{\mathbf{X}} \leq \mathbf{x}} \left[\int_{\tilde{X}_0 \leq z - \tilde{\mathbf{X}}^\top \boldsymbol{\beta}} f_e(\tilde{X}_0, \tilde{\mathbf{X}}) d\tilde{X}_0 \right] d\tilde{\mathbf{X}}. \end{aligned}$$

This implies that the joint density of \mathbf{X} and $z(\mathbf{X}_e, \boldsymbol{\beta})$ given $\boldsymbol{\beta}$ is given by

$$f_{\mathbf{X}, z}(\mathbf{X}, z | \boldsymbol{\beta}) = f_e(z - \mathbf{X}^\top \boldsymbol{\beta}, \mathbf{X}), \quad (7)$$

and the marginal density of $z(\mathbf{X}_e, \boldsymbol{\beta})$ is given by

$$f_z(z | \boldsymbol{\beta}) = \int_{\mathcal{X}} f_{\mathbf{X}, z}(\mathbf{X}, z | \boldsymbol{\beta}) d\mathbf{X} = \int_{\mathcal{X}} f_e(z - \mathbf{X}^\top \boldsymbol{\beta}, \mathbf{X}) d\mathbf{X}. \quad (8)$$

Define $f_{\mathbf{X}|z}(\mathbf{X} | z, \boldsymbol{\beta}) = f_{\mathbf{X}, z}(\mathbf{X}, z | \boldsymbol{\beta}) / f_z(z | \boldsymbol{\beta})$ as the conditional density of \mathbf{X} given z and $\boldsymbol{\beta}$, we have that

$$\begin{aligned} L(z, \boldsymbol{\beta}) &= \mathbb{E}(G(z - \mathbf{X}^\top \Delta \boldsymbol{\beta}) | z(\mathbf{X}_e, \boldsymbol{\beta}) = z) \\ &= \int_{\mathcal{X}} G(z - \mathbf{X}^\top \Delta \boldsymbol{\beta}) f_{\mathbf{X}|z}(\mathbf{X} | z, \boldsymbol{\beta}) d\mathbf{X}, \end{aligned} \quad (9)$$

where $\Delta \boldsymbol{\beta} = \boldsymbol{\beta} - \boldsymbol{\beta}^*$.

Based on the above notations, now we formally study the asymptotic properties of the KBGD estimator under increasing dimensions. We first introduce some further assumptions.

Assumption 3. *The kernel function $K(\cdot)$ satisfies: (i) K is bounded and twice continuously differentiable with bounded first and second derivatives, and the second derivative satisfies Lipschitz condition on the whole real line; (ii) $\int K(s) ds = 1$; (iii) there exists positive integer v_K such that $\int s^v K(s) du = 0$ for $1 \leq v \leq v_K - 1$ and $\int u^{v_K} K(u) du \neq 0$; (iv) $K(s) = 0$ for $|s| > 1$.*

Assumption 4. *(i) There exists some constant $\zeta > 1$ such that $\zeta^{-1} \leq f_e(\mathbf{X}_e) \leq \zeta$ holds for all $\mathbf{X}_e \in \mathcal{X}_e$; (ii) there exists positive integer v_f such that $f_e(\mathbf{X}_e)$ has bounded up to v_f -th derivatives.*

Remark 4. Assumption 4(i) together with Assumption 2(i) is a commonly-used assumption in the machine learning literature (e.g., Wager and Athey, 2018). It basically requires that the joint density of \mathbf{X}_e is uniformly bounded from both above and below over \mathcal{X}_e , so the density

does not degenerate over \mathcal{X}_e . [Assumption 4](#)(i) basically allows us to construct a subset of \mathcal{X}_e such that $f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta})$ is uniformly lower bounded from zero over such subset.

The following lemma will be useful in the proof of our theorem.

Lemma 1. Suppose that [Assumption 1](#), [Assumption 2](#)(i)-(iii), [Assumption 3](#), and [Assumption 4](#) hold with $v_G = 3$, $v_K = 2$, and $v_f = 3$. Define $\psi(n, p, h) = h^{-1} \sqrt{\log(pnh^{-1})/n} + h^2$. If $h_n \rightarrow 0$ and $p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n) \rightarrow 0$ further hold, we have that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \widehat{G}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta}) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) \mathbf{X}_i] \right\| = O_p \left(p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n) \right).$$

Proof of Lemma 1. See [Appendix A](#). □

[Lemma 1](#) implies that $\frac{1}{n} \sum_{i=1}^n \widehat{G}(Z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta}) \mathbf{X}_i$ will be closer to $\mathbb{E}[L(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) \mathbf{X}_i]$ uniformly with respect to $\boldsymbol{\beta}$ as n increases. Note that such uniform convergence results are free of trimming; we do not need to trim $\mathbf{X}_{e,i}$ even when the density of $z(\mathbf{X}_{e,i}, \boldsymbol{\beta})$ is small. So even when $\widehat{G}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta})$ is a poor estimator for $L(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta})$ for some $\mathbf{X}_{e,i}$ and $\boldsymbol{\beta}$, our results are still valid. While on the same time, the cost of not conducting any trimming is that our guaranteed convergence rate depends heavily on the dimensionality. As is required in [Lemma 1](#), the dimension p must satisfy $p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n) \rightarrow 0$. Suppose that $p/n \rightarrow 0$ and we choose $h_n = ((\log n)/n)^{1/6}$, we have that $\psi(n, p, h_n) \sim ((\log n)/n)^{1/3}$. This implies that when p is fixed, the convergence rate in [Lemma 1](#) is $((\log n)/n)^{1/3(p+1)}$. When p increases with n , the dimension p should satisfy $p \log p = O(\log n)$, implying that p is allowed to increase only mildly with n . The restriction on p basically comes from the fact that as $\mathbf{X}_{e,i}$ moves towards the boundary of \mathcal{X}_e , the density of random variable $z(\mathbf{X}_{e,i}, \boldsymbol{\beta})$ decreases faster towards zero given a larger p , which makes the convergence rate sensitive to the increase of p .

For notational simplicity, in the following we denote $z(\mathbf{X}_{e,i}, \boldsymbol{\beta}_k)$ and $z(\mathbf{X}_{e,i}, \boldsymbol{\beta}^*)$ as $z_{i,k}$ and z_i^* . Based on the results in [Lemma 1](#), we have that under all conditions as imposed in [Lemma 1](#), there holds

$$\boldsymbol{\beta}_{k+1} = \boldsymbol{\beta}_k - \delta_k \mathbb{E}[(L(z_{i,k}, \boldsymbol{\beta}_k) - G(z_i^*)) \cdot \mathbf{X}_i] + \delta_k \cdot (\text{small order terms}). \quad (10)$$

Note that $z_{i,k} = z_i^* + \mathbf{X}_i^T \Delta \boldsymbol{\beta}_k$ and $L(z_{i,k}, \boldsymbol{\beta}_k) = \int_{\mathcal{X}} G(z_{i,k} - \mathbf{X}^T \Delta \boldsymbol{\beta}_k) f_{\mathbf{X}|z}(\mathbf{X} | z_{i,k}, \boldsymbol{\beta}_k) d\mathbf{X}$, so

$(L(z_{i,k}, \beta_k) - G(z_i^*)) \cdot \mathbf{X}_i$ equals to

$$\begin{aligned} & \left\{ \int_{\mathcal{X}} [G(z_i^* + \mathbf{X}_i^T \Delta \beta_k - \mathbf{X}^T \Delta \beta_k) - G(z_i^*)] f_{\mathbf{X}|z}(\mathbf{X} | z_{i,k}, \beta_k) d\mathbf{X} \right\} \cdot \mathbf{X}_i \\ &= \int_0^1 \int_{\mathcal{X}} \left[G' \left(z_i^* + t(\mathbf{X}_i - \mathbf{X})^T \Delta \beta_k \right) f_{\mathbf{X}|z}(\mathbf{X} | z_{i,k}, \beta_k) (\mathbf{X}_i \mathbf{X}_i^T - \mathbf{X}_i \mathbf{X}^T) \right] \Delta \beta_k d\mathbf{X} dt, \quad (11) \end{aligned}$$

where the integration is understood to be element-wise. To further simplify our notation, define

$$W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \beta) = G' \left(z(\mathbf{X}_e, \beta^*) + (\mathbf{X} - \tilde{\mathbf{X}})^T \Delta \beta \right) f_{\mathbf{X}|z}(\tilde{\mathbf{X}} | z(\mathbf{X}_e, \beta), \beta),$$

$$V(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \beta) = (\mathbf{X} \mathbf{X}^T - \mathbf{X} \tilde{\mathbf{X}}^T) W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \beta),$$

and

$$\Lambda(\beta) = \mathbb{E} \left[\int_{\mathcal{X}} V(\mathbf{X}_{e,i}, \mathbf{X}_e, \beta) d\mathbf{X} \right],$$

we have that

$$\mathbb{E}[(L(z_{i,k}, \beta_k) - G(z_i^*)) \cdot \mathbf{X}_i] = \int_0^1 \Lambda(\beta^* + t \Delta \beta_k) \Delta \beta_k dt,$$

which indicates that

$$\Delta \beta_{k+1} = \left\{ \int_0^1 (I_p - \delta_k \Lambda(\beta^* + t \Delta \beta_k)) dt \right\} \Delta \beta_k + \delta_k \cdot (\text{small order terms}).$$

To ensure that with probability going to 1 the above iteration shrinks $\|\Delta \beta_k\|$, we make the following assumption.

Assumption 5. *There hold*

$$\sup_{\beta \in \mathcal{B}} \bar{\lambda}(\Lambda(\beta) + \Lambda^T(\beta)) \leq \bar{\lambda}_A < \infty,$$

and

$$\inf_{\beta \in \mathcal{B}} \underline{\lambda}(\Lambda(\beta) + \Lambda^T(\beta)) \geq \underline{\lambda}_A > 0.$$

Based on the above assumptions, we have the following result.

Theorem 3. *Suppose that [Assumption 1](#), [Assumption 2\(i\)–\(iii\)](#), [Assumption 3](#)–[Assumption 5](#) hold with $v_G = 3$, $v_K = 2$, and $v_f = 3$, $\delta_k = \delta$ such that $\delta < \min \{1/(2\underline{\lambda}_A), 1/(4p^2 \|G'\|_\infty)\}$,*

and that β is updated based on [algorithm 2](#). Define

$$k_{1,n}^{KBGD} = \frac{\log(\|\Delta\beta_1\|) - \log\left(p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n)\right)}{-\log(1 - \delta\lambda_A/4)}.$$

Then if $h_n \rightarrow 0$ and $p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n) \rightarrow 0$ hold, we have that

$$\sup_{k \geq k_{1,n}^{KBGD} + 1} \|\Delta\beta_k\| = O_p\left(p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n)\right).$$

In particular, if h_n is chosen such that $h_n = ((\log n)/n)^{1/6}$, then

$$\sup_{k \geq k_{1,n}^{KBGD} + 1} \|\Delta\beta_k\| = O_p\left(p^{\frac{5p+1}{2(p+1)}} \left(\frac{\log n}{n}\right)^{\frac{1}{3p+3}}\right).$$

Proof of [Theorem 3](#). See [Appendix B](#). □

[Theorem 3](#) implies that the iterative estimator based on (5) and (6) is consistent under increasing dimensions, no matter whether the starting point is close to the unknown true parameter or not. However, the convergence speed heavily depends on the dimensionality of the problem, p , even when p is fixed. This is not ideal under our single-index setup but is not surprising since our algorithm does not involve any trimming procedure as we have discussed in [Lemma 7](#).

We proceed to establish the asymptotic normality of the KBGD estimator. Due to technical difficulties, throughout the following analysis in this section we only consider the case where p is fixed. As we can see in [Theorem 3](#), even in the case of fixed dimensionality, the guaranteed convergence rate of the KBGD estimator based on (5) and (6) is at best $((\log n)/n)^{\frac{1}{3p+3}}$, which still depends on p . To obtain asymptotic normality, we need to slightly modify our algorithm to get rid of the dependence on dimensionality. In particular, we introduce trimming to our algorithm. When updating the parameter, we only use observations that fall into a pre-selected region as did in [Ichimura \(1993\)](#). In particular, the algorithm is modified as,

$$\beta_{k+1} = \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n I_i^\phi \cdot \left(\widehat{G}(z_{i,k} | \beta_k) - y_i \right) \mathbf{X}_i, \quad (12)$$

where $\widehat{G}(z_{i,k} | \beta_k) = \widehat{G}(z(\mathbf{X}_{e,i}, \beta_k) | \beta_k)$ is defined in (5), $I_i^\phi = I(\mathbf{X}_{e,i} \in \mathcal{X}_e^\phi)$, and \mathcal{X}_e^ϕ is a

subset of \mathcal{X}_e given by

$$\mathcal{X}_e^\phi = \{\mathbf{X}_e \in \mathcal{X}_e : |X_j| \leq 1 - \phi, 0 \leq j \leq p\} \quad (13)$$

for some $\phi > 0$ whose value will be determined later. Different from (6), the update of β_k based on (12) uses only a subset of the whole sample for which the covariate vector $\mathbf{X}_{e,i}$ falls into \mathcal{X}_e^ϕ . The reason why we choose the trimming set as in (13) is that, as we show in the Appendix A, for any $0 < \phi < 1$, there holds $\inf_{(\mathbf{X}_e, \beta) \in \mathcal{X}_e^\phi \times \mathcal{B}} f_z(z(\mathbf{X}_e, \beta) | \beta) \geq C\phi^p p^{-p}$ for some constant $C > 0$ that depends on ϕ . When p and ϕ are both fixed, $f_z(z(\mathbf{X}_e, \beta) | \beta)$ is uniformly lower bounded from zero for any combination $(\mathbf{X}_e, \beta) \in \mathcal{X}_e^\phi \times \mathcal{B}$, so the uniform estimation accuracy of $L(z(\mathbf{X}_{e,i}, \beta), \beta)$ over $\mathbf{X}_{e,i}$ and β will be improved. Note that trimming will cause some efficiency loss by dropping some observations, but such loss can be controlled to be small if we choose ϕ to be close to zero. We also point that trimming is only applied to the update of the parameter; when nonparametrically estimating G , we still use all the data points.

To simplify our following notation, given the trimming parameter ϕ , we denote $I^\phi \cdot \mathbf{X}$ as \mathbf{X}^ϕ . We also define

$$\Lambda_\phi(\beta) = \mathbb{E} \left[I_i^\phi \cdot \int_{\mathcal{X}} V(\mathbf{X}_{e,i}, \mathbf{X}_e, \beta) d\mathbf{X} \right].$$

The following theorem provides a counterpart to the results in Theorem 3.

Theorem 4. *Suppose that all the assumptions and conditions on v_G , v_K , and v_f in Theorem 3 hold. Suppose moreover that $h_n \rightarrow 0$, $\delta_k = \delta < \min\{1/(2\lambda_A), 1/(4p^2 \|G'\|_\infty)\}$, $\phi < \delta\lambda_A/(16p^2 \|G'\|_\infty \zeta)$, and that β is updated under (5) and (12) (The trimmed version of algorithm 2). Define*

$$\tilde{k}_{1,n}^{KBGD} = \frac{\log(\|\Delta\beta_1\|) - \log(\psi(n, p, h_n))}{-\log(1 - \delta\lambda_A/8)},$$

then there holds

$$\sup_{k \geq \tilde{k}_{1,n}^{KBGD} + 1} \|\Delta\beta_k\| = O_p(\psi(n, p, h_n)).$$

Proof of Theorem 4. See Appendix B. □

Note that when p is fixed, $\psi(n, p, h_n)$ no longer depends on p asymptotically. The improvement over the convergence rate basically comes from the improvement of the uniform convergence rate of the kernel estimator due to trimming. Also note that under trimming, the minimum number of iteration in Theorem 3(i), $\tilde{k}_{1,n}^{KBGD}$, is of order $\log n$ as long as $nh_n \rightarrow \infty$.

This implies that under trimming, a faster convergence rate is guaranteed with the minimum number of iterations being of the same magnitude as that of the estimator without trimming.

We now proceed to establish the asymptotic normality of β_k . Define

$$\xi_n^\phi = \frac{1}{n} \sum_{i=1}^n \left(\widehat{G}(z_i^* | \beta^*) - y_i \right) \mathbf{X}_i^\phi.$$

We note that

$$\begin{aligned} \Delta \beta_{k+1} &= \Delta \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n \left(\widehat{G}(z_{i,k} | \beta_k) - y_i \right) \mathbf{X}_i^\phi, \\ &= \Delta \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n \left(\widehat{G}(z_{i,k} | \beta_k) - \widehat{G}(z_i^* | \beta^*) \right) \mathbf{X}_i^\phi - \delta_k \xi_n^\phi \\ &= \int_0^1 \left\{ I_p - \frac{\delta_k}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\phi \frac{\partial \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta)}{\partial \beta^T} \Big|_{\beta = \beta^* + t \Delta \beta_k} \right] \right\} dt \Delta \beta_k - \delta_k \xi_n^\phi, \end{aligned} \quad (14)$$

where the integration is understood to be element-wise. To understand the properties of the above algorithm, we need the following lemmas.

Lemma 2. Suppose that all the assumptions in [Theorem 3](#) hold with $v_G = 4$, $v_K = 3$, and $v_f = 4$. For any sequence of subset $\{\mathcal{B}_n\}_{n=1}^\infty$ with $\mathcal{B}_n \subseteq \mathcal{B}$, we have that

$$\sup_{\beta \in \mathcal{B}_n} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^\phi \frac{\partial \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta)}{\partial \beta^T} - \Lambda_\phi(\beta) \right\| = O_p \left(h_n^{-2} \sqrt{(\log(nh_n^{-1}))/n} + h_n^3 + \sup_{\beta \in \mathcal{B}_n} \|\Delta \beta\| \right).$$

Proof of Lemma 2. See [Appendix A](#). □

Lemma 3. Suppose that all the assumptions in [Theorem 3](#) hold with $v_G = 4$, $v_K = 3$, and $v_f = 4$. If h_n is chosen such that $h_n^6 \rightarrow 0$, we have that $\sqrt{n} \xi_n^\phi \rightarrow_d N(0, \Sigma_\xi^\phi)$, where

$$\Sigma_\xi^\phi = \mathbb{E} \left[(1 - G(z_i^*)) G(z_i^*) \left(\mathbf{X}_i^\phi - \mathbb{E}(\mathbf{X}_i^\phi | z_i^*) \right) \left(\mathbf{X}_i^\phi - \mathbb{E}(\mathbf{X}_i^\phi | z_i^*) \right)^T \right].$$

Proof of Lemma 3. See [Appendix A](#). □

Now we are in a position to illustrate the results of the asymptotic normality of our KBGD estimator.

Theorem 5. Suppose that all the assumptions in [Theorem 3](#) hold with $v_G = 4$, $v_K = 3$, and $v_f = 4$. Suppose moreover that $\delta_k = \delta < \min\{1/(2\lambda_A), 1/(4p^2 \|G'\|_\infty)\}$, $\phi <$

$\delta\lambda_\Lambda / (16p^2 \|G'\|_\infty \zeta)$, h_n is chosen such that $nh_n^6 \rightarrow 0$ and $h_n^4 n / (\log n)^2 \rightarrow \infty$, and that β is updated under (5) and (12). Then

(i) There holds

$$\sup_{k \geq \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} \|\Delta\beta_k\| = O_p(n^{-1/2}),$$

where $k_{2,n}^{KBGD}$ is given by

$$k_{2,n}^{KBGD} = \frac{\log(n^{1/2}) + \log(\psi(n, p, h_n))}{-\log(1 - \delta\lambda_\Lambda/16)};$$

(ii) Define $\hat{\beta} = \hat{\beta}_k$ for any $k - \tilde{k}_{1,n}^{KBGD} - k_{2,n}^{KBGD} \rightarrow \infty$, we have that

$$\sqrt{n}(\hat{\beta} - \beta^*) \rightarrow N(0, \Sigma_\beta^\phi),$$

where $\Sigma_\beta^\phi = \Lambda_\phi^{-1}(\beta^*) \Sigma_\xi^\phi (\Lambda_\phi^{-1}(\beta^*))^T$.

Proof of Theorem 5. See Appendix B. □

We introduce the estimator for the variance matrix, based on which the confidence interval of β^* can be then constructed.

Theorem 6. Suppose that all the assumptions and conditions in Theorem 5 hold. Suppose also that $\hat{\beta}$ is defined as in Theorem 5. Define

$$\hat{\Sigma}_\xi^\phi = \frac{1}{n} \sum_{i=1}^n \left(\hat{G}_i (1 - \hat{G}_i) (\mathbf{X}_i^\phi - \hat{\mathbb{E}}(\mathbf{X}_i^\phi | \hat{z}_i)) (\mathbf{X}_i^\phi - \hat{\mathbb{E}}(\mathbf{X}_i^\phi | \hat{z}_i))^T \right),$$

and

$$\hat{\Lambda}_\phi(\hat{\beta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^\phi \frac{\partial \hat{G}(z(\mathbf{X}_{e,i}, \hat{\beta}) | \hat{\beta})}{\partial \beta^T},$$

where

$$\hat{G}_i = \frac{\sum_{j=1}^n K_{h_n}(\hat{z}_i - \hat{z}_j) y_j}{\sum_{j=1}^n K_{h_n}(\hat{z}_i - \hat{z}_j)}, \quad \hat{\mathbb{E}}(\mathbf{X}_i^\phi | \hat{z}_i) = \frac{\sum_{j=1}^n K_{h_n}(\hat{z}_i - \hat{z}_j) \mathbf{X}_j^\phi}{\sum_{j=1}^n K_{h_n}(\hat{z}_i - \hat{z}_j)},$$

and $\hat{z}_i = X_{0,i} + \mathbf{X}_i^T \hat{\beta}$. Then we have that

$$\left\| \hat{\Lambda}_\phi^{-1}(\hat{\beta}) \hat{\Sigma}_\xi^\phi (\hat{\Lambda}_\phi^{-1}(\hat{\beta}))^T - \Sigma_\beta^\phi \right\| \rightarrow_p 0.$$

Proof of Theorem 6. See Appendix B. □

We finally provide some remarks for the KBGD estimators.

Remark 5. We first provide some remarks on the implementation of our KBGD estimator. The KBGD estimator might be sensitive to the data magnitude. So when implementing such an estimator, we recommend first standardizing the data so that each covariate has zero mean and unit variance. Note that when constructing the KBGD estimator, we normalize the coefficient of $X_{0,i}$ to 1, indicating that the coefficients of $\mathbf{X}_{e,i}$ can not all be zeros. So we need to test whether at least one covariate affects the conditional probability of $y_i = 1$. One option is to run a Logit or Probit regression and test whether all the coefficients are equal to zero.

When applying our algorithm, it is also crucial to determine the learning rate δ , bandwidth of kernel estimator h_n , and terminating conditions of the algorithm. In [Theorem 5](#), the tuning parameter δ is required to be smaller than $1/(2\lambda_A)$ and $1/(4p^2 \|G'\|_\infty)$, neither of which is known. So we recommend setting δ to be 1 in the first place, and gradually shrink it if the iteration does not converge. For the choice of the bandwidth h_n , [Theorem 5](#) requires that h_n is chosen such that $nh_n^6 \rightarrow 0$ and $nh_n^4/(\log n)^2 \rightarrow \infty$. As a rule of thumb, we recommend choosing $h_n = C \cdot n^{-1/5}$. For the choice of the constant C , we can choose $C = C_k = \text{std}(z_{i,k})$ for the k -th round of iteration and $C = \text{std}(\hat{z}_i)$ when estimating the variance Σ_{β}^ϕ . We finally discuss the terminating conditions. As we show in [Theorem 5](#), to obtain root- n consistency and asymptotic normality, the iteration number is required to be only of order $\log(n)$. However, such rule can not be directly applied to determine the number of iterations since the initial distance $\|\Delta\beta_1\|$ as well as the lower bounded on the eigenvalues λ_A are both unknown. We recommend the terminating condition $\max_{1 \leq j \leq p} |\hat{\beta}_{j,k+1} - \hat{\beta}_{j,k}| < \varrho$ for some predetermined tolerance ϱ . During the simulation, we choose $\varrho = 10^{-5}$. Note that in many cases, $\max_{1 \leq j \leq p} |\hat{\beta}_{j,k+1} - \hat{\beta}_{j,k}|$ may not be monotonically decreasing with k ; in some extreme cases, $\max_{1 \leq j \leq p} |\hat{\beta}_{j,k+1} - \hat{\beta}_{j,k}|$ may even be oscillating and does not shrink to zero. On these condition, we recommend decreasing δ or choosing $h_n = C \cdot n^{-1/5}$ with $C = 1$ when iterating. If the maximum distance still oscillates, we recommend stop iteration when the maximum distance achieves its minimum value.

Remark 6. Our previous discussion has be confined to the case where all the covariates are continuously distributed, while our algorithm can be directly applied to the case where there are discrete covariates without any modifications. The basic reason is that, in contrast to the average derivative approach ([Stoker, 1986](#); [Powell et al., 1989](#)) that uses the differentiation with respect to covariates, the KBGD estimator performs differentiation with respect to the parameters, so it does not impose requirements on the continuity of the covariates. It should be noted that we do require at least one continuous covariate to guarantee identification

of the parameters. For simplicity, we recommend choosing a continuous covariate as the standardization covariate X_0 . Finally, we point out that stronger assumption should be imposed to make our results valid when there are discrete covariates. In particular, suppose that $\mathbf{X}_e = (\mathbf{X}_c^T, \mathbf{X}_d^T)^T$, where \mathbf{X}_c is the collection of all the continuous covariates, whereas \mathbf{X}_d is the collection of all the discrete covariates. Also denote the density function of \mathbf{X}_c conditional on \mathbf{X}_d as $f_{\mathbf{X}_c|\mathbf{X}_d}(\mathbf{X}_c|\mathbf{X}_d)$. Then we require that all the conditions imposed on the $f_e(\mathbf{X}_e)$ hold for $f_{\mathbf{X}_c|\mathbf{X}_d}(\mathbf{X}_c|\mathbf{X}_d)$ for any realizations of \mathbf{X}_d .

3.2 The SBGD Estimator

In the previous section, we introduced the KBGD algorithm, where the update of the parameter is based on a BGD-type procedure while the unknown CDF is replaced with its Nadaraya-Watson kernel estimator constructed by the initial parameter. In this section, we consider an alternative nonparametric approximation for the unknown CDF based on the method of sieves. Given a set of basis functions $\{r_j(z)\}_{j=0}^\infty$ that is complete in $C(\mathbb{R})$ space, any smooth CDF G can be represented by $G(z) = \sum_{j=0}^\infty \pi_j^* r_j(z)$ for any $z \in R$, where $\{\pi_j^*\}_{j=0}^\infty$ is the unknown coefficients of the basis functions. In practice, to make our algorithm tractable, we truncate the sequence of the basis functions and only use the first $q+1$ basis functions for approximation, where q increases with sample size n at some rate. To approximate G , it then remains to provide an estimator for the unknown coefficients of the basis functions $\{\pi_j^*\}_{j=0}^q$. Our estimation procedure for $\{\pi_j^*\}_{j=0}^q$ shares similar intuition as the one that motivates the Nadaraya-Watson kernel estimator in the previous section. In particular, suppose for a moment that in the k -th round of update, we start with $\boldsymbol{\beta}_k$, which happens to be identical to the unknown true parameter $\boldsymbol{\beta}^*$. In this case, define $\mathbf{r}_q(z) = (r_0(z), \dots, r_q(z))^T$ and $\boldsymbol{\pi}_q^* = (\pi_1^*, \dots, \pi_q^*)^T$, we have that

$$y_i = G(z_{i,k}) + \varepsilon_i \approx \mathbf{r}_q^T(z_{i,k}) \boldsymbol{\pi}_q^* + \varepsilon_i,$$

where recall that $z_{i,k} = X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k$. The above relationship motivates the following OLS estimator for the sieve coefficients

$$\hat{\boldsymbol{\pi}}_{q,n,k} = \left(\sum_{i=1}^n \mathbf{r}_q(z_{i,k}) \mathbf{r}_q^T(z_{i,k}) \right)^{-1} \left(\sum_{i=1}^n \mathbf{r}_q(z_{i,k}) y_i \right). \quad (15)$$

Given the estimator of the sieve coefficients $\hat{\pi}_{q,n,k}$, the unknown CDF G in the k -th round of update is approximated by

$$\hat{G}(z|\beta_k) = \mathbf{r}_q^T(z) \hat{\pi}_{q,n,k}, \quad -\infty < z < \infty. \quad (16)$$

Based on the estimated CDF $\hat{G}(z|\beta_k)$, the update of the parameter can be carried out based on (6). We iterate sequentially based on (15), (16) and (6) until some terminating conditions are satisfied. The resulting estimator is then labeled as the *sieve-based batch gradient descent estimator* (SBGD estimator). We summarize our algorithm as follows in [algorithm 3](#).

Algorithm 3: The SBDG Estimator

input : Data set $\{(\mathbf{X}_{e,i}, y_i)\}_{i=1}^n$, sequence of learning rate $\{\delta_k\}_{k=1}^\infty$, initial guess β_1 , the order of sieves q , sieve functions $\mathbf{r}(z) = r_0(z), \dots, r_q(z)$, and terminating condition \mathcal{T}
output: The SBDG estimator $\hat{\beta}$

```

1  $k \leftarrow 1$ ;
2 while The terminating condition  $\mathcal{T}$  is not satisfied do
3    $\hat{\pi}_{q,n,k} \leftarrow$ 
      $(\sum_{i=1}^n \mathbf{r}_q(X_{0,i} + \mathbf{X}_i^T \beta_k) \mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta_k))^{-1} (\sum_{i=1}^n \mathbf{r}_q(X_{0,i} + \mathbf{X}_i^T \beta_k) y_i);$ 
4   for  $i \leftarrow 1$  to  $n$  do
5      $\hat{G}(X_{0,i} + \mathbf{X}_i^T \beta_k | \beta_k) \leftarrow \mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta_k) \hat{\pi}_{q,n,k};$ 
6      $\beta_{k+1} \leftarrow \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n (\hat{G}(X_{0,i} + \mathbf{X}_i^T \beta_k | \beta_k) - y_i) \mathbf{X}_{e,i};$ 
7      $k \leftarrow k + 1$ ;
8  $\hat{\beta} \leftarrow \beta_k$ ;
```

Remark 7. In the above SBDG procedure, we update the sieve parameter based on the OLS-type estimation. An alternative procedure can be based on the flexible Logit regression proposed by [Hirano et al. \(2003\)](#). The advantage of using flexible Logit regression is that the estimated CDF $\hat{G}(z|\beta_k)$ always falls between 0 and 1 for all z , which makes the update more stable. While the disadvantage of such update is that the flexible Logit regression is based on MLE, which does not allow for an analytical solution. Using numerical optimization to solve for the sieve coefficients in each round of update will add to additional computational burdens.

Remark 8. Compared with the KBGD algorithm, the SBDG procedure has at least two advantages. On the one side, the sieve-based approximation for the unknown CDF is global and guarantees uniform approximation error rate. This allows us to update the parameter

without performing any form of trimming as we did for the KBGD estimator. Moreover, this allows us to develop the asymptotic distribution of the SBGD estimator for the case of increasing dimensionality. On the otherhand, the KBGD procedure relies on the kernel estimation of CDF G at n data points, whose computational complexity of each update is of order $O(n^2)$. While the most time-consuming part of the SBGD procedure is the OLS procedure (15), whose computational complexity is of order $O(nq^2 + q^3)$. When $q/\sqrt{n} \rightarrow 0$, the computational burden of SBGD estimator will be substantially lower than that of KBGD estimator.

Define $R_q(z) = G(z) - \mathbf{r}^T(z) \boldsymbol{\pi}_q^*$, $\Gamma_{q,n}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}) \mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta})$, $\Gamma_{q,n,k} = \Gamma_{q,n}(\boldsymbol{\beta}_k)$, and $\mathfrak{X}_{q,n}(z, \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}) \Gamma_{q,n}^{-1}(\boldsymbol{\beta}) \mathbf{r}_q(z) \mathbf{X}_i)$. Through tedious algebra, we can show that the SBGD procedure has the following representation,

$$\begin{aligned} \boldsymbol{\beta}_{k+1} = & \boldsymbol{\beta}_k - \frac{\delta_k}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,n}(z_{i,k}, \boldsymbol{\beta}_k)) (G(z_{i,k}) - G(z_i^*)) \\ & - \frac{\delta_k}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T(z_{i,k}) \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q(z_{j,k}) R_q(z_{j,k}) + \frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(z_{j,k}) \varepsilon_j \right) \\ & + \frac{\delta_k}{n} \sum_{i=1}^n (R_q(z_{i,k}) \mathbf{X}_i + \varepsilon_i \mathbf{X}_i), \end{aligned} \quad (17)$$

where recall that $z_i^* = X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*$. To study the properties of the above procedure, we introduce some additional assumptions.

Assumption 6. (i) There holds $\max_{0 \leq j \leq q} \|r_j\|_\infty \leq D_{q,0}$, $\max_{0 \leq j \leq q} \|r'_j\|_\infty \leq D_{q,1}$, and $\max_{0 \leq j \leq q} \|r''_j\|_\infty \leq D_{q,2}$; (ii) Define $\Gamma_q(\boldsymbol{\beta}) = \mathbb{E}(\mathbf{r}_q(X_0 + \mathbf{X}^T \boldsymbol{\beta}) \mathbf{r}_q^T(X_0 + \mathbf{X}^T \boldsymbol{\beta}))$, there hold $\inf_{\boldsymbol{\beta} \in \mathcal{B}} \underline{\lambda}(\Gamma_q(\boldsymbol{\beta})) \geq \underline{\lambda}_\Gamma > 0$ and $\sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\lambda}(\Gamma_q(\boldsymbol{\beta})) \leq \bar{\lambda}_\Gamma < \infty$ for all q ; (iii) There hold $\sup_{z \in R} |G(z) - \mathbf{r}^T(z) \boldsymbol{\pi}_q^*| \leq \mathcal{E}_{q,0}$ and $\sup_{z \in R} |G'(z) - (\mathbf{r}'(z))^T \boldsymbol{\pi}_q^*| \leq \mathcal{E}_{q,1}$, where $\mathbf{r}'(z) = (r'_0(z), \dots, r'_q(z))^T$.

For any $-\infty < z < \infty$, define the population counterpart of $\mathfrak{X}_{q,n}(z, \boldsymbol{\beta})$ as

$$\mathfrak{X}_q(z, \boldsymbol{\beta}) = \mathbb{E}(\mathbf{r}_q^T(z(\mathbf{X}_e, \boldsymbol{\beta})) \Gamma_q^{-1}(\boldsymbol{\beta}) \mathbf{r}_q(z) \mathbf{X}).$$

Then we have the following lemma.

Lemma 4. Define $\chi_{1,n} = \sqrt{pq^2 D_{q,0}^4 \log(pq D_{q,0} D_{q,1} n) / n}$, and $\chi_{2,n} = \sqrt{pq D_{q,0}^2 (\chi_{1,n} + \mathcal{E}_{q,0})}$. Suppose that [Assumption 1](#), [Assumption 2\(i\)-\(iii\)](#), and [Assumption 6](#) hold, and moreover, $v_G \geq 1$ and the combination of p , q and v_G guarantees that $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$. Then the

following holds,

$$\beta_{k+1} = \beta_k - \delta_k \mathbb{E}[(\mathbf{X} - \mathfrak{X}_q(z(\mathbf{X}_e, \beta_k), \beta_k))(G(z(\mathbf{X}_e, \beta_k)) - G(z(\mathbf{X}_e, \beta^*)))] + \delta_k \mathfrak{R}_{n,k},$$

where $\sup_{k \geq 1} \|\mathfrak{R}_{n,k}\| = O_p(\chi_{2,n})$.

Proof of Lemma 4. See Appendix A. □

Obviously, Lemma 4 provides a parallel result to (10). In particular, define

$$\Psi_q(t, \beta) = \mathbb{E}[G'(z(\mathbf{X}_e, \beta^*) + t\mathbf{X}^T \Delta \beta)(\mathbf{X}\mathbf{X}^T - \mathfrak{X}_q(z(\mathbf{X}_e, \beta), \beta)\mathbf{X}^T)],$$

under all the conditions imposed in Lemma 4, we have that

$$\Delta \beta_{k+1} = \left\{ \int_0^1 (I_p - \delta_k \Psi_q(t, \beta_k)) dt \right\} \Delta \beta_k + \delta_k \mathfrak{R}_{n,k}. \quad (18)$$

Obviously, (18) is also a parallel result to (11). As a result, to ensure that (18) actually constitutes a contraction for $\|\Delta \beta_k\|$, we impose the following assumption that is similar to Assumption 5.

Assumption 7. For any $q \geq 0$, there hold

$$\inf_{0 \leq t \leq 1, \beta \in \mathcal{B}} \underline{\lambda}(\Psi_q(t, \beta) + \Psi_q^T(t, \beta)) \geq \underline{\lambda}_\Psi > 0,$$

$$\sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \underline{\lambda}(\Psi_q(t, \beta) + \Psi_q^T(t, \beta)) \geq \bar{\lambda}_\Psi < \infty.$$

Based on the above assumptions, we have the following result.

Theorem 7. Suppose that Assumption 1, Assumption 2(i)-(iii), Assumption 6 and Assumption 7 hold, $v_G \geq 1$, and the combination of p , q and v_G guarantees that $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$. Suppose moreover that the learning rate is chosen such that $\delta_k = \delta$ with $0 < \delta < \min \left\{ 1/(2\underline{\lambda}_\Psi), \underline{\lambda}_\Psi / \left(2\|\mathbf{G}'\|_\infty^2 p^2 \{1 + \underline{\lambda}_\Gamma^{-1} q D_{q,0}^2\}^2 \right) \right\}$, and that β is updated based on algorithm 3. Define

$$k_{1,n}^{SBGD} = \frac{\log(\|\Delta \beta_1\|) - \log(\chi_{2,n})}{-\log(1 - \underline{\lambda}_\Psi \delta / 4)},$$

then we have that

$$\sup_{k \geq k_{1,n}^{SBGD} + 1} \|\Delta \beta_k\| = O_p(\chi_{2,n}).$$

Proof of Theorem 7. See Appendix B. □

According to [Theorem 7](#), when $\chi_{2,n} \rightarrow 0$ as $n \rightarrow \infty$, the SBGD estimator is consistent as long as the number of updates exceeds $k_{1,n}^{SBGD}$. Based on such consistent estimator, we are ready to establish the asymptotic normality of our SBGD estimator. Apply the mean value theorem to [\(17\)](#), we have that

$$\begin{aligned} \Delta \beta_{k+1} = & \left\{ I_p - \delta_k \int_0^1 \frac{1}{n} \sum_{i=1}^n G' (z_i^* + t \mathbf{X}_i^T \Delta \beta_k) (\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n}(z_{i,k}, \beta_k) \mathbf{X}_i^T) dt \right\} \Delta \beta_k \\ & - \frac{\delta_k}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T(z_{i,k}) \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q(z_{j,k}) R_q(z_{j,k}) + \frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(z_{j,k}) \varepsilon_j \right) \\ & + \frac{\delta_k}{n} \sum_{i=1}^n (R_q(z_{i,k}) \mathbf{X}_i + \varepsilon_i \mathbf{X}_i). \end{aligned}$$

Define $\Psi_q^* = \mathbb{E} [G' (z(\mathbf{X}_e, \beta^*)) (\mathbf{X} \mathbf{X}^T - \mathfrak{X}_q(z(\mathbf{X}_e, \beta^*), \beta^*) \mathbf{X}^T)]$ and $\mathfrak{V}_q = \mathbb{E} (\mathbf{X}_i \mathbf{r}_q^T(z_i^*) \Gamma_q^{-1}(\beta^*))$. Similar to [Lemma 2](#) and [Lemma 3](#), we provide two additional lemmas that are useful to understand the above algorithm.

Lemma 5. Suppose that [Assumption 1](#), [Assumption 2\(i\)-\(iii\)](#), and [Assumption 6](#) hold, $v_G \geq 2$ and the combination of p , q and v_G guarantees that $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$. Then for any sequence $\{\mathcal{B}_n\}_{n=1}^\infty$ with $\mathcal{B}_n \subseteq \mathcal{B}$ we have that

$$\begin{aligned} & \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}_n} \left\| \frac{1}{n} \sum_{i=1}^n G' (z_i^* + t \mathbf{X}_i^T \Delta \beta) (\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i^T) - \Psi_q^* \right\| \\ & = O_p \left(pq D_{q,0}^2 \chi_{1,n} + \sqrt{p^3 q^2 D_{q,0}^3 D_{q,1}} \sup_{\beta \in \mathcal{B}_n} \|\Delta \beta\| \right). \end{aligned}$$

Proof of Lemma 5. See [Appendix A](#). □

Lemma 6. Suppose that [Assumption 1](#), [Assumption 2\(i\)-\(iii\)](#), [Assumption 6](#), and [Assumption 7](#) hold, and the combination of p , q and v_G guarantees that $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$. Define $\mathbf{r}_{q,i,k} = \mathbf{r}_q(z_{i,k})$, and $R_{q,i,k} = R_q(z_{i,k})$. Also define

$$\chi_{3,n} = \sqrt{p^2 q D_{q,1}^2 \log(pq D_{q,2} n) / n},$$

then we have that

$$\sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_{q,j,k} R_{q,j,k} + \frac{1}{n} \sum_{j=1}^n \mathbf{r}_{q,j,k} \varepsilon_j \right) + \frac{1}{n} \sum_{i=1}^n R_q(z_{i,k}) \mathbf{X}_i - \frac{1}{n} \sum_{i=1}^n \mathfrak{X}_q(z_i^*, \boldsymbol{\beta}^*) \varepsilon_i \right\| = O_p(\chi_{4,n}),$$

where $\chi_{4,n} = \sqrt{pq} D_{q,0}^2 \mathcal{E}_{q,0} + \sqrt{pq} D_{q,0} \chi_{2,n} \chi_{3,n} + \chi_{2,n} \sqrt{p^2 q^4 D_{q,0}^6 D_{q,1}^2 (\log q) / n}$.

Proof of Lemma 6. See Appendix A. □

Based on the above two lemmas, we are now ready to study the asymptotic distribution of the SBGD estimator.

Theorem 8. Suppose that Assumption 1, Assumption 2(i)-(iii), Assumption 6 and Assumption 7 hold, $v_G \geq 2$, the combination of p , q and v_G guarantees that $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$, and that $\boldsymbol{\beta}$ is updated based on algorithm 3. We have that

(i) There holds

$$\Delta \boldsymbol{\beta}_{k+1} = (I_p - \delta \Psi_q^*) \Delta \boldsymbol{\beta}_k + \frac{\delta}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \boldsymbol{\beta}^*)) \varepsilon_i + \tilde{\mathfrak{R}}_{n,k},$$

where $\sup_{k \geq k_{1,n}^{SBGD}+1} \|\tilde{\mathfrak{R}}_{n,k}\| = O_p(\chi_{5,n})$ with

$$\chi_{5,n} = \sqrt{pq} D_{q,0}^2 (p + q D_{q,0} D_{q,1}) \chi_{2,n}^2 + \chi_{4,n};$$

(ii) Define $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_{k+k_{1,n}^{SBGD}+k_{2,n}^{SBGD}+1}$ with

$$k_{2,n}^{SBGD} = \frac{-\log \chi_{2,n} + \log \sqrt{n}}{-\log(1 - \underline{\lambda}_\Psi \delta / 4)},$$

and any $k \geq 1$. If the combination of p , q and v_G further guarantees that $\sqrt{n} \chi_{5,n} \rightarrow 0$ as $n \rightarrow \infty$, we have that

$$\sqrt{n} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) = \Psi_q^{*-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \boldsymbol{\beta}^*)) \varepsilon_i + o_p(n^{-\frac{1}{2}}).$$

Then for any $p \times 1$ vector ρ such that $\|\rho\| < \infty$ and $\frac{1}{\sqrt{n}} \sum_{i=1}^n \rho^T \Psi_q^{*-1} (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \boldsymbol{\beta}^*)) \varepsilon_i \rightarrow_d$

$N(0, \sigma_S^2(\rho))$ with

$$\sigma_S^2(\rho) = \lim_{n \rightarrow \infty} \rho^T \Psi_q^{*-1} \mathbb{E} \left\{ G(z_i^*) (1 - G(z_i^*)) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*)) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*))^T \right\} (\Psi_q^{*-1})^T \rho,$$

there holds

$$\sqrt{n} \rho^T (\hat{\beta} - \beta^*) \rightarrow_d N(0, \sigma_S^2(\rho))$$

Proof of Theorem 8. See Appendix B. □

We now provide the estimator for the variance.

Theorem 9. Suppose that all the conditions listed in Theorem 8 hold and $pq^2 D_{q,0}^4 \mathcal{E}_{q,1} \rightarrow 0$ as $n \rightarrow \infty$. Let $\hat{\beta}$ be as defined as in Theorem 8. Define $\hat{\mathbf{r}}_{q,i} = \mathbf{r}_q(z(\mathbf{X}_{e,i}, \hat{\beta}))$, $\hat{\mathbf{r}}'_{q,i} = \mathbf{r}'_q(z(\mathbf{X}_{e,i}, \hat{\beta}))$, $\hat{\pi}_q = (\sum_{i=1}^n \hat{\mathbf{r}}_{q,i} \hat{\mathbf{r}}_{q,i}^T)^{-1} (\sum_{i=1}^n \hat{\mathbf{r}}_{q,i} y_i)$, $\hat{G}_i = \hat{\mathbf{r}}_{q,i}^T \hat{\pi}_q$, $\hat{G}'_i = \hat{\mathbf{r}}'_{q,i} \hat{\pi}_q$, $\hat{\Psi}_{q,i}^* = \frac{1}{n} \sum_{i=1}^n \hat{G}'_i \cdot (\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n}(\hat{z}_i, \hat{\beta}) \mathbf{X}_i^T)$, $\hat{\mathfrak{X}}_{q,i} = \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j \hat{\mathbf{r}}_{q,j}^T \Gamma_{q,n}^{-1}(\hat{\beta}) \hat{\mathbf{r}}_{q,i}$, and

$$\hat{\sigma}_S^2(\rho) = \rho^T \hat{\Psi}_q^{*-1} \frac{1}{n} \sum_{i=1}^n \left\{ \hat{G}_i (1 - \hat{G}_i) (\mathbf{X}_i - \hat{\mathfrak{X}}_{q,i}) (\mathbf{X}_i - \hat{\mathfrak{X}}_{q,i})^T \right\} (\hat{\Psi}_q^{*-1})^T \rho,$$

Then for any $p \times 1$ vector ρ such that $\|\rho\| < \infty$, there holds

$$|\hat{\sigma}_S^2(\rho) - \sigma_S^2(\rho)| \rightarrow_p 0.$$

Proof of Theorem 9. See Appendix B. □

We finally provide some remarks on the empirical applications of the SBGD estimator.

Remark 9. For the choice of sieve functions, we can use polynomial series for the case where the error term u_i has bounded support and Hermite polynomials for the case where u_i has unbounded support. Note that when using polynomial series $\{1, z, z^2, \dots, z^q\}$, the correlation between the sieve functions increases as the approximation order q increases, which may lead to a violation of Assumption 6(ii). To improve the finite sample performance of our method, we recommend using Chebyshev or Legendre polynomials. Moreover, in the case where u_i has unbounded support, following Bierens (2014), we recommend first conducting the following transformation $G(z) = \tilde{G}(T(z))$, where $T: R \mapsto [-1, 1]$ is a differentiable function, and then using standard Chebyshev or Legendre polynomials to approximate \tilde{G} . For example, in our following simulations and empirical applications in Section 5, we use $T(z) = 2\pi^{-1} \arctan(z)$. For the uniform error bound of truncated Legendre polynomials, see Wang and Xiang (2012).

4 Monte Carlo Experiments

This section conducts Monte Carlo simulations to study the performance of our KBGD and SBGD estimators. We focus on two aspects of our estimators. First we study the finite-sample properties of the KBGD estimator, including the bias and the root mean squared error (RMSE). Let the j -th argument of the true parameter be β_j^* , and the simulation is repeated R times, where its estimator in the r -th round of simulation is $\hat{\beta}_j^r$, then the bias and RMSE are respectively given by $\text{Bias} = |\frac{1}{R} \sum_{r=1}^R (\hat{\beta}_j^r - \beta_j^*)|$ and $\text{RMSE} = \sqrt{\sum_{r=1}^R (\hat{\beta}_j^r - \beta_j^*)^2 / R}$. We also investigate whether the confidence interval based on the asymptotic distribution has good coverage rate. We consider nominal coverage rate $\alpha = 0.95$, so the confidence interval for β_j^* in the r -th round of repetition is given by $CI_j^r = [\hat{\beta}_j^r - 1.96 \cdot \widehat{\text{std}}_j^r, \hat{\beta}_j^r + 1.96 \cdot \widehat{\text{std}}_j^r]$, where $\widehat{\text{std}}_j^r$ is the estimated standard deviation of $\hat{\beta}_j^r$. The actual coverage rate is then given by $CR = \frac{1}{R} \sum_{r=1}^R I(\beta_j^* \in CI_j^r)$.

We are also interested in how sensitive our estimators are to the initial guess of the true parameter. In each repetition of our simulation, we consider three different initial guesses: the true parameter vector, the parameter vector estimated based on the Logit regression, and the parameter with all elements being zeros. If the estimation results starting from different initial guesses are close or even identical to each other, the estimation methods are insensitive to the initial guesses and thus are robust in terms of computation. Denote $\hat{\beta}_T^r$, $\hat{\beta}_L^r$, and $\hat{\beta}_Z^r$ as the estimators with starting points being true parameter, Logit estimator, and vector of zeros. We use $S_L = \sqrt{\frac{1}{R} \sum_{i=1}^n \|\hat{\beta}_L^r - \hat{\beta}_T^r\|^2}$ and $S_Z = \sqrt{\frac{1}{R} \sum_{i=1}^n \|\hat{\beta}_Z^r - \hat{\beta}_T^r\|^2}$ as the measurement of the sensitivity. To compare the performance of our method with the existing estimators, we also consider Ichimura's semiparametric least squares (SLS) estimator (Ichimura, 1993) and Klein and Spady's semiparametric maximum likelihood (SMLE) estimator (Klein and Spady, 1993).

We consider data generating process $y_i = I(X_{0,i} + \beta_1^* X_{1,i} \cdots + \beta_{10}^* X_{10,i} - u_i > 0)$, $i = 1, 2, \dots, n$, where data are i.i.d over i , and $X_{0,i}, X_{1,i}, \dots, X_{10,i}, u_i$ are also independent. We set $\beta^* = (1, 0.5, -0.5, 1, -1, 2, -2, 4, -4, 1.5, -1.5)^T$, $X_{j,i} \sim N(0, 1)$ for $0 \leq j \leq 8$, $X_{9,i} \sim \text{Bernoulli}(1/2)$, $X_{10,i} \sim \text{Poisson}(2)$, and $u_i \sim \text{Cauchy}$. We consider two sample sizes $n = 2500$ and 5000 . Finally, for finite-sample performance, we repeat the simulation 500 times; for sensitivity analysis, we repeat 100 times.

Table 1 reports the finite-sample properties of our estimators. It can be seen that our estimators work well in finite sample cases. Both estimators have small bias, whose RMSE decrease with the increase of sample size. Moreover, the confidence interval constructed based on the asymptotic variance and normal approximation has actual coverage rate that is quite close to the nominal rate 0.95.

Table 1: Finite Sample Performance of KBGD and SBGD Estimators

	Bias		RMSE		CR			Bias		RMSE		CR	
	KBGD	SBGD	KBGD	SBGD	KBGD	SBGD		KBGD	SBGD	KBGD	SBGD	KBGD	SBGD
	$n = 2500$							$n = 5000$					
β_1	0.0024	0.0031	0.1193	0.1240	0.9600	0.9680		0.0047	0.0005	0.0844	0.0867	0.9500	0.9600
β_2	0.0002	0.0055	0.1255	0.1336	0.9480	0.9500		0.0031	0.0074	0.0846	0.0878	0.9520	0.9540
β_3	0.0136	0.0260	0.1544	0.1791	0.9480	0.9460		0.0004	0.0074	0.1053	0.1112	0.9320	0.9320
β_4	0.0093	0.0213	0.1551	0.1706	0.9500	0.9440		0.0012	0.0095	0.1035	0.1117	0.9600	0.9500
β_5	0.0257	0.0482	0.2511	0.2968	0.9540	0.9400		0.0007	0.0168	0.1648	0.1889	0.9400	0.9480
β_6	0.0236	0.0477	0.2502	0.2860	0.9480	0.9580		0.0121	0.0269	0.1723	0.1931	0.9540	0.9360
β_7	0.0500	0.0964	0.4513	0.5416	0.9640	0.9420		0.0051	0.0352	0.3083	0.3525	0.9440	0.9420
β_8	0.0447	0.0920	0.4662	0.5441	0.9360	0.9520		0.0098	0.0394	0.3121	0.3477	0.9420	0.9440
β_9	0.0242	0.0454	0.2921	0.3303	0.9480	0.9500		0.0072	0.0048	0.1840	0.1909	0.9540	0.9560
β_{10}	0.0168	0.0338	0.1881	0.2223	0.9520	0.9440		0.0030	0.0147	0.1247	0.1402	0.9440	0.9380

NOTE: For KBGD estimator, we use fourth-order Epanechnikov kernel to construct the Nadaraya-Watson estimator. We choose $\delta = 1$. In each round of iteration, the bandwidth h_n is chosen as $h_n = \sigma_{\hat{z}} \cdot n^{-1/5}$, where n is sample size, $\sigma_{\hat{z}}$ is the standard deviation of $z_{i,k}$, and $z_{i,k} = X_{0,i} + \mathbf{X}_i^T \beta_k$. For SBGD estimator, we choose $q = 9$ and use Legendre polynomials with transformation discussed in [Remark 9](#). For both estimators, the stopping rule is either $\max_{1 \leq j \leq p} |\hat{\beta}_{j,k+1} - \hat{\beta}_{j,k}| < 10^{-5}$ or $k \geq 20000$. The above also applies to our empirical analysis in [Section 5](#). Trimming is ignored during all the simulations. Due to the outliers of the simulation, we trim out the lower and upper 2% simulation results and calculate the bias and RMSE.

Table 2: Sensitivity of KBGD and SBGD Estimators: Fixed Coefficients

		Sensitivity		Running Time		
	Method	S_L	S_Z	True	Logit	Zeros
$n = 2500$	KBGD	0.0242	0.0198	113.21	79.120	158.91
	SBGD	0.0175	0.0259	0.9504	0.9482	1.1587
	SLS	0.8732	251.58	35.695	37.210	35.104
	SMLE	0.9362	318.41	34.515	33.704	31.078
$n = 5000$	KBGD	0.0241	0.0175	157.48	87.954	230.07
	SBGD	0.0189	0.0282	1.4644	1.4722	1.9074
	SLS	0.6870	871.58	46.402	44.647	41.486
	SMLE	0.7343	507.69	44.563	43.256	35.904

NOTE: SLS refers to semiparametric least squares estimator, and SMLE refers to semiparametric maximum likelihood estimator. The running time is all in seconds. Due to the outliers of the simulation, we trim out the lower and upper 2% simulation results and calculate the corresponding results. The above also applies to [Table 3](#).

Table 3: Sensitivity of KBGD and SBGD Estimators: Random Coefficients

		Sensitivity		Running Time			
		Method	S_L	S_Z	True	Logit	Zeros
$n = 2500$	KBGD	0.0270	0.0214	122.00	74.433	166.94	
	SBGD	0.0123	0.0246	1.0132	0.8252	1.2044	
	SLS	0.9178	500.24	34.864	35.571	34.065	
	SMLE	0.9956	533.58	34.334	32.520	29.473	
$n = 5000$	KBGD	0.0234	0.0232	163.74	91.449	247.49	
	SBGD	0.0077	0.0234	1.5529	1.4377	1.9217	
	SLS	0.6796	10737	43.935	41.420	46.449	
	SMLE	0.6821	698.63	43.616	44.825	37.763	

Table 2 reports the sensitivity of our estimators to the starting points. We can see that for both KBGD and SBGD estimators, S_L and S_Z are close to zero, indicating that the resulting estimators starting from Logit estimator or zeros are almost identical to the ones starting from the unknown true parameter. Such a result demonstrates that our algorithms are robust to different initial guesses. On the contrary, the SLS and SMLE are both sensitive to the initial guess. As we can see, the estimators starting from parametric Logit regression differ significantly from those starting from the unknown true parameter, and such difference even explodes when we consider estimators starting from the origin point. The above results highlight the numerical robustness of our estimators.

The robustness of our algorithm might also be sensitive to the setups of coefficients. To check whether this is the case, instead of using the fixed parameters specified before, in each round of simulation we randomly draw true parameter β^* as follows $\beta_1^*, \beta_2^*, \beta_9^*, \beta_{10}^* \sim N(0, 1)$, $\beta_3^*, \beta_4^*, \beta_5^*, \beta_6^* \sim 2N(0, 1)$, and $\beta_7^*, \beta_8^* \sim 4N(0, 1)$. The simulation results are reported in Table 3. We can see that the results are similar to those under fixed parameters, indicating that our algorithm is robust to initial point under different parameter setups.

5 Empirical Illustration

As an empirical illustration of our new methods, this section applies our KBGD and SBGD estimation procedures to study how education affects the risk aversion. In the existing researches it's extensively documented that, on the individual level, risk aversion is significantly correlated with the level of education, although the directions of correlation are mixed, see Outreville (2015) for a comprehensive review. In this study, we investigate how educational background of the family affects the risk aversion of the household as well as household-level investing behaviors. We use the national survey data from 2019 China Household Financial

Survey Project (CHFS) (Gan et al., 2014), which provides household-level information over demographics, asset and debt, income and consumption, social security and insurance, and various household’s subjective preferences. The dependent variable we are interested in is the degree of risk aversion of the household. In particular, y_i is constructed to take value of 0 if the i -th household is completely against any form of risks and thus is described as being extremely risk averse; it takes value of 1 if the family is willing to bear some form of risks when making investments. We study how the probability of $y_i = 1$ is affected by a set of factors based on the binary choice model. We have a total of 11 explanatory variables in our model. The key factor that we are particularly interested in is the educational backgrounds, which is defined year of education of the head of the household. We also consider a set of other control variables including gender, ethnicity, health conditions, marital status, region of residence, economic knowledge, and total asset, whose impacts on the risk aversion are of interest on their own right. When conducting semiparametric estimation, we normalize the coefficient of total asset to 1. See Yao (2023) for detailed discussion on the construction of the data sets.

Before estimation, we normalize all the continuous variables so that the resulting variables all have zero mean and unity variance. To provide a comparison to the semiparametric estimation results, we first conduct parametric Logit regression and report the normalized coefficients in regression (I) in Table 4. We then conduct KBGD and SBGD estimation and report the estimated coefficients of education in (II) and (III). As we can see from Table 4, no matter which estimation methods we use, the coefficient of educational background is estimated to be positive with significance at 1% level. This implies that, holding other conditions fixed, on average an increase in the year of education of the head in the households leads to the increase of willingness to bear risks. Comparing the semiparametric estimation results with that of Logit regression, we can see that the KBGD and SBGD estimators are close to each other. We finally compare the computation time of each method. We can see that both KBGD and SBGD estimators take much longer to converge compared with the parametric estimation. Comparatively, the SBGD algorithm is significantly faster than the KBGD algorithm, which takes over two hours to converge. This result supports the use of SBGD algorithm when there are data of large scale.

Table 4: Estimation Results

	(I)	(II)	(III)
Estd. Coefficients	2.5543*** (0.1070)	2.4832*** (0.3638)	2.4647*** (0.3239)
Num. of Obs.	26906	26906	26906
Estimation Methods	Logit	KBGD	SBGD
Running Time	1.4276	8573.1	40.9941
Num. of Iteration	—	14996	12986

Note: For Logit regression, we report the coefficient of education divided by that of total asset. For semiparametric estimation, we normalize the coefficient of total asset to be 1. The standard deviations are reported in the brackets below the coefficients. *** indicates significance at 1% level. For both KBGD and SBGD estimators, we choose $\delta_k = 1$. For KBGD estimator, we choose $h_n = C \cdot n^{-1/5}$ with $C = C_k = \text{std}(z_{i,k})$, and use the fourth-order Epanechnikov kernel. For SBGD estimator, we choose $q = 9$ and use Legendre polynomials with transformation discussed in [Remark 9](#). The starting point of iteration for both KBGD and SBGD estimators is chosen as the origin point with all arguments being 0. The stopping rule is set as $\max_{1 \leq j \leq p} |\hat{\beta}_{j,k+1} - \hat{\beta}_{j,k}| < \varrho$ with $\varrho = 10^{-5}$. Finally, the running time is in second.

6 Conclusions

In this paper, we proposed new estimation procedures for binary choice and monotonic index models with increasing dimensions. Existing semiparametric estimation procedures for this model cannot be implemented in practice when the number of regressors is large. In contrast, our algorithmic based procedures can be used for many regressor models as it involves convex optimization at each iteration of the procedure. We show this iterative procedure also has desirable asymptotic properties when the number of regressors increases with the sample size in ways that are standard in big data literature.

A Lemmas and Proofs

This part provides some lemmas that will be used during the establishment of our results in the main context. If not otherwise stated, the dimension p of covariate \mathbf{X} is allowed to increase with sample size n .

Lemma A.1. *Consider i.i.d. random variables $\{U_i\}_{i=1}^n$ on probability space (Ω, \mathcal{A}, P) and $d_1 \times d_2$ matrix $A(U, \theta) : \Omega \times \Theta \rightarrow R^{d_1 \times d_2}$ with $\Theta \subseteq R^p$ being compact, $\sup_{U \in \Omega, \theta \in \Theta} \|A_{s,t}(U, \theta)\| \leq D_{A,0}$ and $\sup_{U \in \Omega} \|A_{s,t}(U, \theta_1) - A_{s,t}(U, \theta_2)\| \leq D_{A,1} \|\theta_1 - \theta_2\|$ uniformly for all $1 \leq s \leq d_1$ and $1 \leq t \leq d_2$. Then there holds*

$$\sup_{\theta \in \Theta} \left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta) - \mathbb{E}A(U_i, \theta) \right\| = O_p \left(\sqrt{\frac{pd_1 d_2 D_{A,0}^2 \log(d_1 d_2 D_{A,1} n)}{n}} \right).$$

Proof of Lemma A.1. Note that

$$\begin{aligned} \sup_{\theta \in \Theta} \left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta) - \mathbb{E}A(U_i, \theta) \right\| &\leq \max_{1 \leq b \leq B} \left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta_b) - \mathbb{E}A(U_i, \theta_b) \right\| \\ &\quad + \max_{1 \leq b \leq B} \sup_{\|\theta - \theta_b\| \leq \frac{C}{\sqrt[4]{B}}} \left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta) - \frac{1}{n} \sum_{i=1}^n A(U_i, \theta_b) \right\| \\ &\quad + \max_{1 \leq b \leq B} \sup_{\|\theta - \theta_b\| \leq \frac{C}{\sqrt[4]{B}}} \|\mathbb{E}A(U_i, \theta) - \mathbb{E}A(U_i, \theta_b)\|. \end{aligned}$$

For the first term, we have that

$$\begin{aligned} &P \left(\max_{1 \leq b \leq B} \left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta_b) - \mathbb{E}A(U_i, \theta_b) \right\| > \tau \right) \\ &\leq \sum_{b=1}^B P \left(\left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta_b) - \mathbb{E}A(U_i, \theta_b) \right\| > \tau \right) \\ &\leq \sum_{b=1}^B P \left(\max_{1 \leq s \leq d_1} \max_{1 \leq t \leq d_2} \left\| \frac{1}{n} \sum_{i=1}^n A_{s,t}(U_i, \theta_b) - \mathbb{E}A_{s,t}(U_i, \theta_b) \right\| > \frac{\tau}{\sqrt{d_1 d_2}} \right) \\ &\leq \sum_{b=1}^B \sum_{s=1}^{d_1} \sum_{t=1}^{d_2} P \left(\left\| \frac{1}{n} \sum_{i=1}^n A_{s,t}(U_i, \theta_b) - \mathbb{E}A_{s,t}(U_i, \theta_b) \right\| > \frac{\tau}{\sqrt{d_1 d_2}} \right) \\ &\leq \sum_{b=1}^B \sum_{s=1}^{d_1} \sum_{t=1}^{d_2} 2 \exp(-Cn\tau^2 / (d_1 d_2 D_{A,0}^2)) = 2 \exp(C \log(Bd_1 d_2) - Cn\tau^2 / (d_1 d_2 D_{A,0}^2)), \end{aligned}$$

indicating that

$$\max_{1 \leq b \leq B} \left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta_b) - \mathbb{E}A(U_i, \theta_b) \right\| = O_p \left(\sqrt{\frac{d_1 d_2 D_{A,0}^2 \log(B d_1 d_2)}{n}} \right).$$

On the other side, for the second term we have that

$$\begin{aligned} & \max_{1 \leq b \leq B} \sup_{\|\theta - \theta_b\| \leq \frac{C}{\sqrt[3]{B}}} \left\| \frac{1}{n} \sum_{i=1}^n A(U_i, \theta) - \frac{1}{n} \sum_{i=1}^n A(U_i, \theta_b) \right\| \\ & \leq \sqrt{d_1 d_2} \max_{1 \leq s \leq d_1} \max_{1 \leq t \leq d_2} \sup_{U \in \Omega} \sup_{\|\theta - \theta_b\| \leq \frac{C}{\sqrt[3]{B}}} |A_{s,t}(U, \theta) - A_{s,t}(U, \theta_b)| \leq \frac{\sqrt{d_1 d_2} D_{A,1}}{\sqrt[3]{B}}. \end{aligned}$$

The same bound holds for the third term. Then let $B = (\sqrt{n} D_{A,1})^3$, we finish the proof. \square

Lemma A.2. *If [Assumption 1](#), [Assumption 2\(i\)-\(iii\)](#), and [Assumption 4](#) hold with $\min\{v_G, v_f\} \geq 2$, then there exists a constant C that does not depend on \mathbf{X}, z, β such that the following hold*

- (i) $\sup_{\mathbf{X}, z, \beta} |\partial^s f_{\mathbf{X}, z}(\mathbf{X}, z | \beta) / \partial z^s| \leq C$ for $0 \leq s \leq v_f$;
- (ii) $\sup_{z, \beta} |\partial^s f_z(z | \beta) / \partial z^s| \leq C$ for $0 \leq s \leq v_f$;
- (iii) $\sup_{\mathbf{X}, z, \beta} \|\partial f_{\mathbf{X}, z}(\mathbf{X}, z | \beta) / \partial \beta\| \leq C \sqrt{p}$;
- (iv) $\sup_{\mathbf{X}, z, \beta} \|\partial^2 f_{\mathbf{X}, z}(\mathbf{X}, z | \beta) / \partial \beta \partial \beta^T\| \leq Cp$;
- (v) $\|\partial f_z(z | \beta) / \partial \beta\| \leq C \sqrt{p}$;
- (vi) $\|\partial^2 f_z(z | \beta) / \partial \beta \partial \beta^T\| \leq Cp$;
- (vii) $\sup_{z, \beta, f_z(z | \beta) \neq 0} |\partial^s L(z, \beta) / \partial z^s| \leq C$ for $0 \leq s \leq \min\{v_G, v_f\}$;
- (viii) $\sup_{z, \beta, f_z(z | \beta) \neq 0} \|\partial L(z, \beta) / \partial \beta\| \leq C \sqrt{p}$;
- (ix) $\sup_{z, \beta, f_z(z | \beta) \neq 0} \|\partial^2 L(z, \beta) / \partial \beta \partial \beta^T\| \leq Cp$;
- (x) $\sup_{\mathbf{X}_e, \beta, f_z(z(\mathbf{X}_e, \beta) | \beta) \neq 0} \int_{\mathcal{X}} \left\| \partial W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \beta) / \partial \beta \right\| d\tilde{\mathbf{X}} \leq C \sqrt{p}$.

Proof. To prove [Lemma A.2\(i\)](#) and [Lemma A.2\(ii\)](#), we note that for any $0 \leq s \leq v_f$,

$$\frac{\partial^s f_{\mathbf{X}, z}(\mathbf{X}, z | \beta)}{\partial z^s} = \frac{\partial^s f_e(X_0, \mathbf{X})}{\partial X_0^s} \Big|_{X_0 = z - \mathbf{X}^T \beta},$$

and

$$\frac{\partial^s f_z(z | \beta)}{\partial z^s} = \int_{\mathcal{X}} \left[\frac{\partial^s f_{\mathbf{X}, z}(\mathbf{X}, z | \beta)}{\partial X_0^s} \right] d\mathbf{X}.$$

Since $f_e(\mathbf{X}_e)$ has up to v_f -th bounded derivatives over \mathcal{X}_e according to [Assumption 4\(ii\)](#) and X_j is bounded by 1 for all $1 \leq j \leq p$ according to [Assumption 2\(i\)](#), [Lemma A.2\(i\)](#) and [Lemma A.2\(ii\)](#) hold.

Similarly, note that

$$\begin{aligned}\frac{\partial f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= - \left[\frac{\partial f_e(X_0, \mathbf{X})}{\partial X_0} \Big|_{X_0=z-\mathbf{X}^\top \boldsymbol{\beta}} \right] \mathbf{X}, \\ \frac{\partial^2 f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} &= \left[\frac{\partial^2 f_e(X_0, \mathbf{X})}{\partial X_0^2} \Big|_{X_0=z-\mathbf{X}^\top \boldsymbol{\beta}} \right] \mathbf{X} \mathbf{X}^\top, \\ \frac{\partial f_z(z|\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= - \int_{\mathcal{X}} \left[\frac{\partial f_e(X_0, \mathbf{X})}{\partial X_0} \Big|_{X_0=z-\mathbf{X}^\top \boldsymbol{\beta}} \right] \mathbf{X} d\mathbf{X}, \\ \frac{\partial^2 f_z(z|\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} &= \int_{\mathcal{X}} \left[\frac{\partial^2 f_e(X_0, \mathbf{X})}{\partial X_0^2} \Big|_{X_0=z-\mathbf{X}^\top \boldsymbol{\beta}} \right] \mathbf{X} \mathbf{X}^\top d\mathbf{X},\end{aligned}$$

we validate [Lemma A.2\(iii\)](#)-[Lemma A.2\(vi\)](#).

To prove [Lemma A.2\(vii\)](#), note that

$$\begin{aligned}\left| \frac{\partial^s L(z, \boldsymbol{\beta})}{\partial z^s} \right| &\leq C \sum_{j=0}^s \left| \int_{\mathcal{X}} G^{(j)}(z - \mathbf{X}^\top \Delta \boldsymbol{\beta}) \frac{\partial^{s-j} f_{\mathbf{X}|z}(\mathbf{X}|z, \boldsymbol{\beta})}{\partial z^{s-j}} d\mathbf{X} \right| \\ &\leq C \sum_{j=0}^s \|G^{(j)}\|_\infty \cdot \left(\int_{\mathcal{X}} \left| \frac{\partial^{s-j} f_{\mathbf{X}|z}(\mathbf{X}|z, \boldsymbol{\beta})}{\partial z^{s-j}} \right| d\mathbf{X} \right).\end{aligned}$$

According to [Assumption 2\(iii\)](#), $\|G^{(j)}\|_\infty$ is bounded for all $0 \leq j \leq v_G$. Then it remains to show that $\int_{\mathcal{X}} |\partial^{s-j} f_{\mathbf{X}|z} / \partial z^{s-j}| d\mathbf{X}$ is also upper bounded for all $0 \leq j \leq v_f$. When $j = s$, we have that $\int_{\mathcal{X}} |\partial^{s-j} f_{\mathbf{X}|z}(\mathbf{X}|z, \boldsymbol{\beta}) / \partial z^{s-j}| d\mathbf{X} = 1$. When $j = s - 1$, define $\mathbb{X}(z, \boldsymbol{\beta}) = \{\mathbf{X} : (z - \mathbf{X}^\top \boldsymbol{\beta}, \mathbf{X}) \in \mathcal{X}_e\}$. We have that

$$\begin{aligned}&\int_{\mathcal{X}} \left| \frac{\partial f_{\mathbf{X}|z}(\mathbf{X}|z, \boldsymbol{\beta})}{\partial z} \right| d\mathbf{X} \\ &= \int_{\mathcal{X}} \left| \frac{\partial f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta}) / \partial z}{\int_{\mathcal{X}} f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta}) d\mathbf{X}} - \frac{f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta}) \int_{\mathcal{X}} (\partial f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta}) / \partial z) d\mathbf{X}}{(\int_{\mathcal{X}} f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta}) d\mathbf{X})^2} \right| d\mathbf{X} \\ &\leq \frac{2 \int_{\mathcal{X}} |\partial f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta}) / \partial z| d\mathbf{X}}{\int_{\mathcal{X}} f_{\mathbf{X},z}(\mathbf{X}, z|\boldsymbol{\beta}) d\mathbf{X}} \leq \frac{2 \|\partial f_{\mathbf{X},z} / \partial z\|_\infty m(\mathbb{X}(z, \boldsymbol{\beta}))}{\zeta^{-1} m(\mathbb{X}(z, \boldsymbol{\beta}))} \leq C\end{aligned}$$

according to part (i) of this lemma. The proof of the case when $j = s - 2, \dots, 0$ are similar, so is omitted.

To prove Lemma A.2(viii), note that

$$\begin{aligned} \left\| \frac{\partial L(z, \beta)}{\partial \beta} \right\| &\leq \int_{\mathcal{X}} \|G'(z - \mathbf{X}^T \Delta \beta) f_{\mathbf{X}|z}(\mathbf{X}|z, \beta) \mathbf{X}\| d\mathbf{X} \\ &\quad + \int_{\mathcal{X}} \left\| G(Z - \mathbf{X}^T \Delta \beta) \frac{\partial f_{\mathbf{X}|z}(\mathbf{X}|z, \beta)}{\partial \beta} \right\| d\mathbf{X}. \end{aligned}$$

Obviously, the first term on the RHS is bounded by $\|G'\|_{\infty} \sqrt{p}$, and the second term is bounded by $\|G\|_{\infty} \int_{\mathcal{X}} \|\partial f_{\mathbf{X}|z}(\mathbf{X}|z, \beta) / \partial \beta\| d\mathbf{X}$. Note that

$$\begin{aligned} \int_{\mathcal{X}} \|\partial f_{\mathbf{X}|z}(\mathbf{X}|z, \beta) / \partial \beta\| d\mathbf{X} &\leq \frac{2 \int_{\mathcal{X}} \|\partial f_{\mathbf{X},z}(\mathbf{X}, z | \beta) / \partial \beta\| d\mathbf{X}}{\int_{\mathcal{X}} f_{\mathbf{X},z}(\mathbf{X}, z | \beta) d\mathbf{X}} \\ &\leq \frac{2C\sqrt{pm}(\mathbb{X}(z, \beta))}{\zeta^{-1}m(\mathbb{X}(z, \beta))} \leq C\sqrt{p}, \end{aligned}$$

according to part (iii) of this lemma. This proves Lemma A.2(viii). Lemma A.2(ix) can be similarly proved.

Finally, to show Lemma A.2(x), we note that

$$\begin{aligned} &\int_{\mathcal{X}} \left\| \frac{\partial W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \beta)}{\partial \beta} \right\| d\tilde{\mathbf{X}} \\ &\leq \int_{\mathcal{X}} \left\| G''\left(z(\mathbf{X}_e, \beta^*) + (\mathbf{X} - \tilde{\mathbf{X}})^T \Delta \beta\right) (\mathbf{X} - \tilde{\mathbf{X}}) \right\| f_{\mathbf{X}|z}(\tilde{\mathbf{X}}|z(\mathbf{X}_e, \beta), \beta) d\tilde{\mathbf{X}} \\ &\quad + \int_{\mathcal{X}} \left\| G'\left(z(\mathbf{X}_e, \beta^*) + (\mathbf{X} - \tilde{\mathbf{X}})^T \Delta \beta\right) \right\| \left\| \frac{\partial f_{\mathbf{X}|z}(\tilde{\mathbf{X}}|z(\mathbf{X}_e, \beta), \beta)}{\partial \beta} \right\| d\tilde{\mathbf{X}}. \end{aligned}$$

Obviously, the first term is bounded by $2\sqrt{p}\|G''\|_{\infty}$, and the second term is bounded by $\|G'\|_{\infty} \int_{\mathcal{X}} \left\| \partial f_{\mathbf{X}|z}(\tilde{\mathbf{X}}|z(\mathbf{X}_e, \beta), \beta) / \partial \beta \right\| d\tilde{\mathbf{X}}$. Note that

$$\int_{\mathcal{X}} \left\| \frac{\partial f_{\mathbf{X}|z}(\tilde{\mathbf{X}}|z(\mathbf{X}_e, \beta), \beta)}{\partial \beta} \right\| d\tilde{\mathbf{X}} \leq \frac{2 \int_{\mathcal{X}} \left\| \partial f_{\mathbf{X},z}(\tilde{\mathbf{X}}, z(\mathbf{X}_e, \beta) | \beta) / \partial \beta \right\| d\tilde{\mathbf{X}}}{f_z(z(\mathbf{X}_e, \beta) | \beta)}.$$

We can see that

$$\begin{aligned} \frac{\partial f_{\mathbf{X},z}(\tilde{\mathbf{X}}, z(\mathbf{X}_e, \beta) | \beta)}{\partial \beta} &= \frac{\partial f_{\mathbf{X},z}(\tilde{\mathbf{X}}, z | \beta)}{\partial z} \bigg|_{z=z(\mathbf{X}_e, \beta)} \mathbf{X} \\ &\quad + \frac{\partial f_{\mathbf{X},z}(\tilde{\mathbf{X}}, z | \beta)}{\partial \beta} \bigg|_{z=z(\mathbf{X}_e, \beta)}, \end{aligned}$$

according to (i) and (ii), we know that $\left\| \partial f_{\mathbf{X},z} \left(\tilde{\mathbf{X}}, z \mid \boldsymbol{\beta} \right) / \partial z \Big|_{z=z(\mathbf{X}_e, \boldsymbol{\beta})} \right\|$ is bounded, and $\left\| \partial f_{\mathbf{X},z} \left(\tilde{\mathbf{X}}, z \mid \boldsymbol{\beta} \right) / \partial \boldsymbol{\beta} \Big|_{z=z(\mathbf{X}_e, \boldsymbol{\beta})} \right\|$ is bounded by $C\sqrt{p}$, so $\left\| \partial f_{\mathbf{X},z} \left(\tilde{\mathbf{X}}, z(\mathbf{X}_e, \boldsymbol{\beta}) \mid \boldsymbol{\beta} \right) / \partial \boldsymbol{\beta} \right\|$ is bounded by $C\sqrt{p}$. So

$$\frac{\int_{\mathcal{X}} \left\| \partial f_{\mathbf{X},z} \left(\tilde{\mathbf{X}}, z(\mathbf{X}_e, \boldsymbol{\beta}) \mid \boldsymbol{\beta} \right) / \partial \boldsymbol{\beta} \right\| d\tilde{\mathbf{X}}}{f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) \mid \boldsymbol{\beta})} \leq \frac{C\sqrt{p} \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}))}{\zeta^{-1} \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}))} = C\sqrt{p}.$$

This finishes the proof of [Lemma A.2](#)(xii). \square

Lemma A.3. Suppose that [Assumption 1](#), [Assumption 2](#)(i)-(iii), [3](#) and [Assumption 4](#) hold with $v_G = 3$, $v_K = 2$, and $v_f = 3$. Define

$$A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) = \frac{1}{nh_n} \sum_{j=1}^n K((z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})) / h_n) \cdot (\cdot_j),$$

where \cdot is y or 1 . Also define $A(\mathbf{X}_e, \boldsymbol{\beta}) = \lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{D}_n} A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta})$, where the expectation $\mathbb{E}_{\mathcal{D}_n}$ is taken with respect to the data set \mathcal{D}_n . Then

(i) There holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} |A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - \mathbb{E}_{\mathcal{D}_n} A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta})| = O_p \left(h_n^{-1} \sqrt{p \log(nph_n^{-1}) / n} \right);$$

(ii) There holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} |\mathbb{E}_{\mathcal{D}_n} A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - A(\mathbf{X}_e, \boldsymbol{\beta})| = O_p(h_n^2);$$

(iii) Define $\psi(n, p, h_n) = h_n^{-1} \sqrt{p \log(nph_n^{-1}) / n} + h_n^2$, there holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} |A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - A(\mathbf{X}_e, \boldsymbol{\beta})| = O_p \left(h_n^{-1} \sqrt{p \log(nph_n^{-1}) / n} + h_n^2 \right).$$

Proof. [Lemma A.3](#)(i) is a direct result of [Lemma A.1](#) if we note that

$$|K((z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})) / h_n) \cdot (\cdot_j)| \leq Ch_h^{-1}$$

and

$$\|\partial(K((z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})) / h_n) \cdot (\cdot_j)) / \partial \boldsymbol{\beta}\| \leq C\sqrt{p}h_h^{-2}.$$

To prove [Lemma A.3\(ii\)](#), we only need to note that

$$\begin{aligned}
& \mathbb{E}_{\mathcal{D}_n} [A_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})] \\
&= \frac{1}{h_n} \mathbb{E}_{\mathcal{D}_n} \left[K \left(\frac{z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})}{h_n} \right) y_j \right] \\
&= \frac{1}{h_n} \mathbb{E}_{\mathcal{D}_n} \left[K \left(\frac{z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})}{h_n} \right) G(z(\mathbf{X}_{e,j}, \boldsymbol{\beta}) - \mathbf{X}_j^T \Delta \boldsymbol{\beta}) \right] \\
&= \frac{1}{h_n} \int K \left(\frac{z(\mathbf{X}_e, \boldsymbol{\beta}) - z}{h_n} \right) G(z - \mathbf{X}_j^T \Delta \boldsymbol{\beta}) f_{\mathbf{X},z}(\mathbf{X}_j, z | \boldsymbol{\beta}) d\mathbf{X}_j dz \\
&= \frac{1}{h_n} \int K \left(\frac{z(\mathbf{X}_e, \boldsymbol{\beta}) - z}{h_n} \right) f_z(z | \boldsymbol{\beta}) dz \int_{\mathcal{X}} G(z - \mathbf{X}_j^T \Delta \boldsymbol{\beta}) \frac{f_{\mathbf{X},z}(\mathbf{X}_j, z | \boldsymbol{\beta})}{f_z(z | \boldsymbol{\beta})} d\mathbf{X}_j \\
&= \frac{1}{h_n} \int K \left(\frac{z(\mathbf{X}_e, \boldsymbol{\beta}) - z}{h_n} \right) f_z(z | \boldsymbol{\beta}) L(z, \boldsymbol{\beta}) dz \\
&= \int K(z) L(z(\mathbf{X}_e, \boldsymbol{\beta}) - h_n z, \boldsymbol{\beta}) f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) - h_n z | \boldsymbol{\beta}) dz \\
&= L(z(\mathbf{X}_e, \boldsymbol{\beta})) f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta}) + \frac{h_n^2}{2} \left[\frac{\partial^2 L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial z^2} \right] \left[\int K(z) z^2 dz \right] \\
&+ \frac{h_n^3}{6} \left\{ \int K(z) z^3 \left[\frac{\partial^3 L(\tilde{z}, \boldsymbol{\beta}) f_z(\tilde{z} | \boldsymbol{\beta})}{\partial z^3} \right] dz \right\},
\end{aligned}$$

and similarly,

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}_n} [A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})] &= \frac{1}{h_n} \mathbb{E}_{\mathcal{D}_n} \left[K \left(\frac{z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})}{h_n} \right) \right] \\
&= \frac{1}{h_n} \int \left[K \left(\frac{z(\mathbf{X}_e, \boldsymbol{\beta}) - z}{h_n} \right) f_z(z | \boldsymbol{\beta}) \right] dz \\
&= \int K(z) f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) - h_n z | \boldsymbol{\beta}) dz \\
&= f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta}) + \frac{h_n^2}{2} \left[\frac{\partial^2 f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial z^2} \right] \left[\int K(z) z^2 dz \right] \\
&+ \frac{h_n^3}{6} \left\{ \int K(z) z^3 \left[\frac{\partial^3 f_z(\tilde{z} | \boldsymbol{\beta})}{\partial z^3} \right] dz \right\},
\end{aligned}$$

where \tilde{z} lies between $z(\mathbf{X}_e, \boldsymbol{\beta})$ and z . Note that according to [Lemma A.2 \(i\)](#) and [\(ii\)](#), $f_z(z | \boldsymbol{\beta})$ and $L(z, \boldsymbol{\beta}) f_z(z | \boldsymbol{\beta}) = \int_{\mathcal{X}} G(z - \mathbf{X}^T \Delta \boldsymbol{\beta}) f_{\mathbf{X},z}(\mathbf{X}, z | \boldsymbol{\beta}) d\mathbf{X}$ both have up to third bounded derivatives with respect to z , so the results hold.

Finally, [Lemma A.3 \(iii\)](#) is a combination of [Lemma A.3 \(i\)](#) and [Lemma A.3 \(ii\)](#). \square

Lemma A.4. Suppose that [Assumption 1](#), [Assumption 2\(i\)-\(iii\)](#), [Assumption 3](#), and [Assumption 4](#) hold. Given any positive sequence $\{\phi_n\}_{n=1}^{\infty}$ satisfying $p\phi_n \downarrow 0$, define

$$\mathcal{X}_{e,n} = \{\mathbf{X}_e \in \mathcal{X}_e : |X_j| \leq 1 - \phi_n, 0 \leq j \leq p\}.$$

Then

- (i) $1 - P(\mathbf{X}_e \in \mathcal{X}_{e,n}) = O(p\phi_n)$, and $\inf_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} f_Z(z(\mathbf{X}_e, \beta) | \beta) \sim \phi_n^p p^{-p}$;
- (ii) If $\psi(n, p, h_n) = o(\phi_n^p p^{-p})$, there holds

$$\sup_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} \left| \widehat{G}(z(\mathbf{X}_e, \beta) | \beta) - L(z(\mathbf{X}_e, \beta), \beta) \right| = O_p(p^p \phi_n^{-p} \psi(n, p, h_n)).$$

Proof. To prove Lemma A.4(i), note that for $p\phi_n < 1$, $m(\mathcal{X}_e - \mathcal{X}_{e,n}) = 1 - (1 - \phi_n)^p \leq p\phi_n$. So $\int_{\mathcal{X}_e - \mathcal{X}_{e,n}} f_e(\mathbf{X}_e) d\mathbf{X}_e \leq \zeta p\phi_n = O(p\phi_n)$ due to Assumption 4(i). To show the lower bound, note that given any $\beta \in \mathcal{B}$ and $\mathbf{X}_e \in \mathcal{X}_{e,n}$, there holds $|z(\mathbf{X}_e, \beta) - \tilde{\mathbf{X}}^T \beta - X_0| \leq \sum_{j=1}^p |\beta_j| |X_j - \tilde{X}_j|$. This implies that for any $\tilde{\mathbf{X}}, \tilde{\mathbf{X}} \in \mathbb{X}(z(\mathbf{X}_e, \beta), \beta)$ if

$$\tilde{\mathbf{X}} \in \left\{ \tilde{\mathbf{X}} \in [0, 1]^p : \left(\sup_{\beta \in \mathcal{B}} |\beta_j| \right) |X_j - \tilde{X}_j| \leq \phi_n/p \right\}.$$

Since the above set has Lebesgue measure of order $O(\phi_n^p/p^p)$, we have that

$$\begin{aligned} & \inf_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} f_Z(z(\mathbf{X}_e, \beta) | \beta) \\ & \geq \inf_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} \int_{\tilde{\mathbf{X}} \in \mathbb{X}(z(\mathbf{X}_e, \beta), \beta)} f_e(z(\mathbf{X}_e, \beta) - \tilde{\mathbf{X}}^T \beta, \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \sim \phi_n^p/p^p, \end{aligned}$$

due to Assumption 4(i). This proves Lemma A.4(i).

To prove Lemma A.4(ii), note that for any \mathbf{X}_e and β , we have $\widehat{G}(z(\mathbf{X}_e, \beta) | \beta) = A_{n,y}(\mathbf{X}_e, \beta) / A_{n,1}(\mathbf{X}_e, \beta)$ and $L(z(\mathbf{X}_e, \beta), \beta) = A_y(\mathbf{X}_e, \beta) / A_1(\mathbf{X}_e, \beta)$. So

$$\begin{aligned} & \sup_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} \left| \widehat{G}(z(\mathbf{X}_e, \beta) | \beta) - L(z(\mathbf{X}_e, \beta), \beta) \right| \\ & \leq \sup_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} \frac{|A_{n,y}(\mathbf{X}_e, \beta) - A_y(\mathbf{X}_e, \beta)|}{A_{n,1}(\mathbf{X}_e, \beta)} \\ & \quad + \sup_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} L(z(\mathbf{X}_e, \beta), \beta) \frac{|A_{n,1}(\mathbf{X}_e, \beta) - A_1(\mathbf{X}_e, \beta)|}{A_1(\mathbf{X}_e, \beta)}. \end{aligned}$$

Obviously, since $\psi_1(n, p, h_n) = o(\phi_n^p/p^p)$,

$$\sup_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} |A_{n,1}(\mathbf{X}_e, \beta) - A_1(\mathbf{X}_e, \beta)| = o_p(\phi_n^p/p^p),$$

so $\inf_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} A_{n,1}^{-1}(\mathbf{X}_e, \beta) = O_p(p^p \phi_n^{-p})$. Moreover, $L(z(\mathbf{X}_e, \beta), \beta)$ is upper bounded by Lemma A.2(vii). Then the results hold according to Lemma A.3. \square

Proof of Lemma 1.

Proof. Note that

$$\begin{aligned} & \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i] \right\| \\ & \leq \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \left(\widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \right) \right\| \end{aligned} \quad (1)$$

$$+ \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i] \right\|. \quad (2)$$

Obviously, (1) is bounded by

$$\begin{aligned} & \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \left(\widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \right) \right\| \\ & \leq \frac{1}{n} \sum_{i=1}^n \sup_{\beta \in \mathcal{B}} \left\| \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \right\| \cdot I_{n,i} \end{aligned} \quad (3)$$

$$+ \frac{1}{n} \sum_{i=1}^n \sup_{\beta \in \mathcal{B}} \left\| \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \right\| \cdot (1 - I_{n,i}), \quad (4)$$

where $I_{n,i} = I(\mathbf{X}_{e,i} \in \mathcal{X}_{e,n})$ and $\mathcal{X}_{e,n}$ is chosen as in [Lemma A.4](#). Note that (3) is bounded by

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \sup_{\beta \in \mathcal{B}} \left\| \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \right\| \cdot I_{n,i} \\ & \leq \sup_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} \left\| \widehat{G}(z(\mathbf{X}_e, \beta) | \beta) \mathbf{X} - L(z(\mathbf{X}_e, \beta), \beta) \mathbf{X} \right\| \\ & = O_p(p^{p+1/2} \phi_n^{-p} \psi_1(n, p, h_n)), \end{aligned}$$

according to [Lemma A.4](#). For (4), we have that

$$\begin{aligned} & \mathbb{E} \frac{1}{n} \sum_{i=1}^n \sup_{\beta \in \mathcal{B}} \left\| \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \right\| \cdot (1 - I_{n,i}) \\ & \leq C \sqrt{p} \mathbb{E} I(\mathbf{X}_{e,i} \notin \mathcal{X}_{e,n}) = O(p^{3/2} \phi_n), \end{aligned}$$

according to [Lemma A.4\(i\)](#). Then we have that (3) is of order $O_p(p^{p+1/2} \phi_n^{-p} \psi_1(n, p, h_n) + p^{3/2} \phi_n)$.

Now we go to (2). Similar to the above truncation, we have that

$$\begin{aligned} & \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i] \right\| \\ & \leq \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \cdot I_{n,i} - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \cdot I_{n,i}] \right\| \end{aligned} \quad (5)$$

$$+ \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \cdot (1 - I_{n,i}) - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i \cdot (1 - I_{n,i})] \right\|. \quad (6)$$

Obviously, (6) is $O_p(p^{3/2}\phi_n)$. For (5), note that $\|L(z(\mathbf{X}_{e,i}, \beta), \beta) X_{j,i} \cdot I_{n,i}\|$ is bounded by C and $\partial \|L(z(\mathbf{X}_{e,i}, \beta), \beta) X_{j,i} \cdot I_{n,i} / \partial \beta\|$ is bounded by $C\sqrt{p}$ by Lemma A.2(vii) and (viii), we have that (5) is of order $O_p\left(\sqrt{p^2 n \log(pn)/n}\right)$ using Lemma A.1. Then

$$\begin{aligned} & \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i] \right\| \\ & = O_p\left(\sqrt{p^2 \log(pn)/n} + p^{3/2}\phi_n\right). \end{aligned}$$

Together, we have that

$$\begin{aligned} & \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i] \right\| \\ & = O_p\left(p^{p+1/2}\phi_n^{-p}\psi_1(n, p, h_n) + \sqrt{p^2 \log(pn)/n} + p^{3/2}\phi_n\right). \end{aligned}$$

Then if we set $\phi_n = p^{\frac{p-1}{p+1}}\psi_1^{\frac{1}{p+1}}(n, p, h_n)$, we have that

$$p\phi_n = p^p\phi_n^{-p}\psi_1(n, p, h_n) = p^{\frac{2p}{p+1}}\psi_1^{\frac{1}{p+1}}(n, p, h_n) \leq p^{\frac{5p+1}{2(p+1)}}\psi_1^{\frac{1}{p+1}}(n, p, h_n) \rightarrow 0,$$

and

$$\sqrt{p^2 \log(pn)/n} = o\left(p^{\frac{5p+1}{2(p+1)}}\psi_1^{\frac{1}{p+1}}(n, p, h_n)\right),$$

so

$$\sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i] \right\| = O_p\left(p^{\frac{5p+1}{2(p+1)}}\psi_1^{\frac{1}{p+1}}(n, p, h_n)\right).$$

This finishes the whole proof. \square

Lemma A.5. Suppose that p is fixed. If all the assumptions in Lemma A.3 hold with $v_G = 4$, $v_K = 3$, and $v_f = 4$, we have that Lemma A.3(i) holds. Moreover,

(i) There holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} |\mathbb{E}_{\mathcal{D}_n} A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - A(\mathbf{X}_e, \boldsymbol{\beta})| = O_p(h_n^3);$$

(ii) There holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} |A_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - A(\mathbf{X}_e, \boldsymbol{\beta})| = O_p\left(h^{-1} \sqrt{\log(nh^{-1})/n} + h^3\right).$$

Proof. The proof is similar to the proof of [Lemma A.3](#) so is omitted. \square

Lemma A.6. Suppose that p is fixed. For any $\mathbf{X}_e \in \mathcal{X}_e$ and $\boldsymbol{\beta} \in \mathcal{B}$, define

$$A'_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) = \frac{1}{nh_n^2} \sum_{j=1}^n K'((z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta}))/h_n) (\mathbf{X} - \mathbf{X}_j) \cdot (\cdot)_j,$$

where $\cdot = 1$ or $\cdot = y$. If all the assumptions in [Lemma A.3](#) hold with $v_G = 4$, $v_K = 3$, and $v_f = 4$, then

(i) There holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} \|A'_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - \mathbb{E}_{\mathcal{D}_n} A'_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta})\| = O_p\left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n}\right);$$

(ii) Define $A'_y(\mathbf{X}_e, \boldsymbol{\beta}) = \lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{D}_n} A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})$ and $A'_1(\mathbf{X}_e, \boldsymbol{\beta}) = \lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{D}_n} A'_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})$.

We have that $A'_y(\mathbf{X}_e, \boldsymbol{\beta}) = \partial H_1(z, \mathbf{X} | \boldsymbol{\beta}) / \partial z|_{z=z(\mathbf{X}_e, \boldsymbol{\beta})}$ and $A'_1(\mathbf{X}_e, \boldsymbol{\beta}) = \partial H_2(z, \mathbf{X} | \boldsymbol{\beta}) / \partial z|_{z=z(\mathbf{X}_e, \boldsymbol{\beta})}$, where

$$H_1(z, \mathbf{X} | \boldsymbol{\beta}) = \int_{\mathcal{X}} G\left(z - \tilde{\mathbf{X}}^T \Delta \boldsymbol{\beta}\right) f_e\left(z - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}\right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}},$$

$$H_2(z, \mathbf{X} | \boldsymbol{\beta}) = \int_{\mathcal{X}} f_e\left(z - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}\right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}},$$

and the differentiation of H_1 and H_2 are element-wise. Moreover, there holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} \|\mathbb{E}_{\mathcal{D}_n} A'_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - A'(\mathbf{X}_e, \boldsymbol{\beta})\| = O_p(h_n^3),$$

(iii) There holds

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e \times \mathcal{B}} \|A'_{n,\cdot}(\mathbf{X}_e, \boldsymbol{\beta}) - A'(\mathbf{X}_e, \boldsymbol{\beta})\| = O_p\left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3\right).$$

Proof. [Lemma A.6](#)(i) is a direct result of [Lemma A.1](#) if we note that for each $1 \leq l \leq p$, $h_n^{-2} K'((z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta}))/h_n) (X_l - X_{l,j}) \cdot (\cdot)_j$ is bounded by Ch_n^{-2} and its derivatives with respect to $\boldsymbol{\beta}$ and \mathbf{X} are both upper bounded since p is fixed.

To prove [Lemma A.6\(ii\)](#), we note that

$$\begin{aligned}
& \mathbb{E}_{\mathcal{D}_n} A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta}) \\
&= \frac{1}{h_n^2} \mathbb{E}_{\mathcal{D}_n} [K'((z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})) / h_n) (\mathbf{X} - \mathbf{X}_j) \cdot G(X_{0,j} + \mathbf{X}_j^T \boldsymbol{\beta}^*)] \\
&= \frac{1}{h_n^2} \mathbb{E}_{\mathcal{D}_n} [K'((z(\mathbf{X}_e, \boldsymbol{\beta}) - z(\mathbf{X}_{e,j}, \boldsymbol{\beta})) / h_n) (\mathbf{X} - \mathbf{X}_j) \cdot G(z(\mathbf{X}_{e,j}, \boldsymbol{\beta}) - \mathbf{X}_j^T \Delta \boldsymbol{\beta})] \\
&= \frac{1}{h_n^2} \int K'((z(\mathbf{X}_e, \boldsymbol{\beta}) - z) / h_n) dz \int_{\mathcal{X}} [G(z - \tilde{\mathbf{X}}^T \Delta \boldsymbol{\beta}) f_{\mathbf{X},z}(\tilde{\mathbf{X}}, z | \boldsymbol{\beta}) (\mathbf{X} - \tilde{\mathbf{X}})] d\tilde{\mathbf{X}} \\
&= \frac{1}{h_n^2} \int [K'((z(\mathbf{X}_e, \boldsymbol{\beta}) - z) / h_n) H_1(z, \mathbf{X} | \boldsymbol{\beta})] dz \\
&= \frac{1}{h_n} \int [K'(z) H_1(z(\mathbf{X}_e, \boldsymbol{\beta}) - h_n z, \mathbf{X} | \boldsymbol{\beta})] dz
\end{aligned}$$

Note that both G and f_e have up to fourth bounded derivatives with respect to z , and the upper bounds hold uniformly with respect to z , \mathbf{X} and $\boldsymbol{\beta}$. This implies that each element of $H_1(z, \mathbf{X} | \boldsymbol{\beta})$ has up to fourth bounded derivatives with respect to z . Also note that $\int K'(v) dv = K(v)|_{-\infty}^{\infty} = 0$, $\int v K'(v) dv = K(v)|_{-\infty}^{\infty} - \int K(v) dv = -1$, $\int v^s K'(v) dv = v^s K(v)|_{-\infty}^{\infty} - s \int v^{s-1} K(v) dv = 0$ for $s = 2, 3$, and $|\int v^4 K'(v) dv| < \infty$. This implies that

$$\|\mathbb{E}_{\mathcal{D}_n} A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta}) - A'_y(\mathbf{X}_e, \boldsymbol{\beta})\| = O_p(h_n^3)$$

uniform with respect to \mathbf{X}_e and $\boldsymbol{\beta}$. The proof of the uniform distance between $\mathbb{E}_{\mathcal{D}_n} A'_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})$ and $A'_1(\mathbf{X}_e, \boldsymbol{\beta})$ is similar. So we finish the proof of [Lemma A.6\(ii\)](#).

Finally, [Lemma A.6\(iii\)](#) is a combination of [Lemma A.6\(i\)](#) and [Lemma A.6\(ii\)](#). \square

Lemma A.7. Suppose that p is fixed. If all the assumptions in [Lemma A.3](#) hold with $v_G = 4$, $v_K = 3$, and $v_f = 4$, we have that

$$\begin{aligned}
& \sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e^\phi \times \mathcal{B}} \left\| \frac{\partial \widehat{G}(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} - \frac{\partial H_1(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z}{f_z(z(\mathbf{X}_e, \boldsymbol{\beta}))} \right. \\
& \quad \left. + L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \frac{\partial H_2(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z}{f_z(z(\mathbf{X}_e, \boldsymbol{\beta}))} \right\| = O_p\left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3\right),
\end{aligned}$$

where \mathcal{X}_e^ϕ is defined in [\(13\)](#) in the main text.

Proof. Note that

$$\begin{aligned}
\frac{\partial \widehat{G}(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= \frac{\partial A_{n,y}(\mathbf{X}_e, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} - \frac{A_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} \cdot \frac{\partial A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} \\
&= \frac{A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} - \frac{A_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} \frac{A'_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})}.
\end{aligned}$$

Then

$$\begin{aligned} \left\| \frac{A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} - \frac{\partial H_1(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z}{f_z(z(\mathbf{X}_e, \boldsymbol{\beta}))} \right\| &= \left\| \frac{A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} - \frac{A'_y(\mathbf{X}_e, \boldsymbol{\beta})}{A_1(\mathbf{X}_e, \boldsymbol{\beta})} \right\| \\ &\leq \left\| \frac{A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta}) - A'_y(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} \right\| \end{aligned} \quad (7)$$

$$+ \left\| \frac{A'_y(\mathbf{X}_e, \boldsymbol{\beta})}{A_1(\mathbf{X}_e, \boldsymbol{\beta})} \frac{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta}) - A_1(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} \right\|. \quad (8)$$

Now for any $(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e^\phi \times \mathcal{B}$, $A_1(\mathbf{X}_e, \boldsymbol{\beta})$ is uniformly lower-bounded according to [Lemma A.4](#), so $A_{n,1}^{-1}(\mathbf{X}_e, \boldsymbol{\beta}) = O_p(1)$ also uniformly holds. Moreover, $\|A'_y(\mathbf{X}_e, \boldsymbol{\beta})\|$ is upper bounded, so $\|A'_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})\| = O_p(1)$ also uniformly holds. Then (7) is $O_p\left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3\right)$ and (8) is $O_p\left(h_n^{-1} \sqrt{\log(nh_n^{-1})/n} + h_n^3\right)$. Similar method can be used to show that

$$\frac{A_{n,y}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} \frac{A'_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})}{A_{n,1}(\mathbf{X}_e, \boldsymbol{\beta})} - \frac{L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \partial H_2(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z}{f_z(z(\mathbf{X}_e, \boldsymbol{\beta}))}$$

is also $O_p\left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3\right)$. This finishes the proof. \square

Lemma A.8. Suppose that p is fixed. If all the assumptions in [Lemma A.3](#) hold with $v_G = 4$, $v_K = 3$, and $v_f = 4$, then for any $\bar{\mathcal{B}} \subseteq \mathcal{B}$, we have that

$$\sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \mathcal{X}_e^\phi \times \bar{\mathcal{B}}} \left\| \frac{\partial \hat{G}(Z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} - \int W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \boldsymbol{\beta}) (\mathbf{X} - \tilde{\mathbf{X}}_e) d\tilde{\mathbf{X}}_e \right\| \leq \alpha_{1,n} + \alpha_2,$$

where $\alpha_{1,n} = O_p\left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3\right)$ and $\alpha_2 = O_p\left(\sup_{\boldsymbol{\beta} \in \bar{\mathcal{B}}} \|\Delta \boldsymbol{\beta}\|\right)$.

Proof. We only need to show that

$$\begin{aligned} \sup_{(\mathbf{X}_e, \boldsymbol{\beta}) \in \bar{\mathcal{X}}_e \times \mathcal{B}} \left\| \frac{\partial H_1(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z}{f_z(z(\mathbf{X}_e, \boldsymbol{\beta}))} - L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \frac{\partial H_2(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z}{f_z(z(\mathbf{X}_e, \boldsymbol{\beta}))} \right. \\ \left. - \int W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \boldsymbol{\beta}) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \right\| = O(\|\Delta \boldsymbol{\beta}\|). \end{aligned}$$

Note that

$$\begin{aligned}
& \partial H_1(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z - L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \partial H_2(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z \\
&= \int G'(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}} \Delta \boldsymbol{\beta}) f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \\
&+ \int G(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \Delta \boldsymbol{\beta}) \left(\partial f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) / \partial z \right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \\
&- L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \int \left(\partial f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) / \partial z \right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \\
&= \int G'(z(\mathbf{X}_e, \boldsymbol{\beta}) - \mathbf{X}^T \Delta \boldsymbol{\beta}) f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \\
&+ \int [G(z(\mathbf{X}_e, \boldsymbol{\beta}) - \mathbf{X}^T \Delta \boldsymbol{\beta}) - G(z(\mathbf{X}_e, \boldsymbol{\beta}))] \left(\partial f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) / \partial z \right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \\
&- (L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) - G(z(\mathbf{X}_e, \boldsymbol{\beta}))) \int \left(\partial f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) / \partial z \right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}}.
\end{aligned}$$

Note that

$$\begin{aligned}
& \left\| \int [G(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \Delta \boldsymbol{\beta}) - G(z(\mathbf{X}_e, \boldsymbol{\beta}))] \left(\partial f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) / \partial z \right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \right\| \\
&\leq C \cdot \sup_{\tilde{\mathbf{X}} \in \mathcal{X}} \left| G(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \Delta \boldsymbol{\beta}) - G(z(\mathbf{X}_e, \boldsymbol{\beta})) \right| \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X})) \\
&\leq C \cdot \|\Delta \boldsymbol{\beta}\| \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X})),
\end{aligned}$$

and according to our choice of \mathcal{X}_e^ϕ , we know that $m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X})) > 0$. On the other side,

$$\begin{aligned}
& \left\| (L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) - G(z(\mathbf{X}_e, \boldsymbol{\beta}))) \int \left(\partial f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^T \boldsymbol{\beta}, \tilde{\mathbf{X}}) / \partial z \right) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \right\| \\
&\leq C \cdot |L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) - G(z(\mathbf{X}_e, \boldsymbol{\beta}))| \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X})) \\
&= C \cdot |L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) - L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}^*)| \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X})) \\
&\leq C \cdot \left(\sup_{z, \boldsymbol{\beta}} \|\partial L(z, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}\| \right) \cdot \|\Delta \boldsymbol{\beta}\| \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X})) \\
&\leq C \cdot \|\Delta \boldsymbol{\beta}\| \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}))
\end{aligned}$$

due to the upper boundedness of $\|\partial L(z, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}\|$ according to [Lemma A.2\(viii\)](#). Note that

$$f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta}) > C \cdot m(\mathbb{X}(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}))$$

for some $C > 0$ due to [Assumption 4\(i\)](#) and the choice of \mathcal{X}_e^ϕ , so we have that

$$\begin{aligned}
& \left\| (\partial H_1(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z - L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \partial H_2(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) / \partial z) / f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta}) \right. \\
& \quad \left. - \int G'(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^\top \Delta \boldsymbol{\beta}) f_e(z(\mathbf{X}_e, \boldsymbol{\beta}) - \tilde{\mathbf{X}}^\top \boldsymbol{\beta}, \tilde{\mathbf{X}}) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} / f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta}) \right\| \\
& = \left\| (\partial_z H_1(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e) - L(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \partial_z H_2(z(\mathbf{X}_e, \boldsymbol{\beta}), \mathbf{X}_e)) / f_z(z(\mathbf{X}_e, \boldsymbol{\beta}) | \boldsymbol{\beta}) \right. \\
& \quad \left. - \int W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \boldsymbol{\beta}) (\mathbf{X} - \tilde{\mathbf{X}}) d\tilde{\mathbf{X}} \right\| \leq C \cdot \|\Delta \boldsymbol{\beta}\|.
\end{aligned}$$

This proves the results. \square

Now we prove [Lemma 2](#) in the main text.

Proof of Lemma 2. Note that

$$\begin{aligned}
& \sup_{\boldsymbol{\beta} \in \mathcal{B}_n} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^\phi \frac{\partial \hat{G}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} - \Lambda_\phi(\boldsymbol{\beta}) \right\| \\
& \leq \sup_{\boldsymbol{\beta} \in \mathcal{B}_n} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^\phi \left(\frac{\partial \hat{G}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} - \int W(\mathbf{X}_{e,i}, \mathbf{X}_e, \boldsymbol{\beta}) (\mathbf{X}_i - \mathbf{X}) d\mathbf{X} \right) \right\| \quad (9)
\end{aligned}$$

$$+ \sup_{\boldsymbol{\beta} \in \mathcal{B}_n} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^\phi \left(\int W(\mathbf{X}_{e,i}, \mathbf{X}_e, \boldsymbol{\beta}) (\mathbf{X}_i - \mathbf{X}) d\mathbf{X} \right) - \Lambda_\phi(\boldsymbol{\beta}) \right\|. \quad (10)$$

Obviously, (9) is of order $O_p \left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3 + \sup_{\boldsymbol{\beta} \in \mathcal{B}_n} \|\Delta \boldsymbol{\beta}\| \right)$ according to [Lemma A.8](#).

Using [Lemma A.1](#), we can show that (10) is $O_p \left(\sqrt{(\log n)/n} \right)$ by noting that each element of $\int W(\mathbf{X}_{e,i}, \mathbf{X}_e, \boldsymbol{\beta}) (\mathbf{X}_i - \mathbf{X}) d\mathbf{X}$ is bounded and that $\int_{\mathcal{X}} \left\| \partial W(\mathbf{X}_e, \tilde{\mathbf{X}}_e, \boldsymbol{\beta}) / \partial \boldsymbol{\beta} \right\| d\tilde{\mathbf{X}}$ is uniformly upper bounded according to [Lemma A.2\(x\)](#). This finishes the proof of [Lemma 2](#). \square

Now we prove [Lemma 3](#) in the main text.

Proof of Lemma 3. We first show that

$$\boldsymbol{\xi}_n^\phi = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \right) \mathbf{X}_i^\phi + o_p \left(\frac{1}{\sqrt{n}} \right),$$

Define $f_z^*(z_i^*) = f_z(z|\beta^*)$ and $f_{\mathbf{X},z}^*(\mathbf{X}, z) = f_{\mathbf{X},z}(\mathbf{X}, z|\beta^*)$. Recall that $z_i^* = z(\mathbf{X}_{e,i}, \beta^*)$, so

$$\begin{aligned} & \boldsymbol{\xi}_n^\phi - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \right) \mathbf{X}_i^\phi \\ &= \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - y_i) \right] \left[\frac{1}{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)} - \frac{1}{f_z^*(z_i^*)} \right] \mathbf{X}_i^\phi \\ &= \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \right] \left[\frac{1}{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)} - \frac{1}{f_z^*(z_i^*)} \right] \mathbf{X}_i^\phi(i) \\ &\quad - \frac{1}{n} \sum_{i=1}^n \varepsilon_i \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) \right] \left[\frac{1}{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)} - \frac{1}{f_z^*(z_i^*)} \right] \mathbf{X}_i^\phi(ii). \end{aligned}$$

For term (i), due to truncation, we have that

$$\max_{1 \leq i \leq n} \left\| \left[\frac{1}{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)} - \frac{1}{f_z^*(z_i^*)} \right] \mathbf{X}_i^\phi \right\| = O_p \left(h_n^{-1} \sqrt{\log(n)/n} + h_n^3 \right).$$

We further provide a uniform bound for $\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \mathbf{X}_i^\phi$ over i . We first note that

$$\mathbb{E}_{\mathcal{D}_n} \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \mathbf{X}_i^\phi \right] = \mathbb{E}_{\mathcal{D}_n} \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (G(z_j^*) - G(z_i^*)) \mathbf{X}_i^\phi \right],$$

where the RHS is equivalent to

$$\begin{aligned} & \mathbb{E} \left\{ \mathbb{E} \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (G(z_j^*) - G(z_i^*)) \mathbf{X}_i^\phi \middle| \mathbf{X}_{e,i} \right] \right\} \\ &= \frac{n-1}{n} \mathbb{E} \left\{ \mathbf{X}_i^\phi \int [K_{h_n}(z - z_i^*) (G(z) - G(z_i^*)) f_z^*(z)] dz \right\} \\ &= \frac{n-1}{n} \mathbb{E} \left\{ \mathbf{X}_i^\phi 0 \int [K(z) (G(z_i^* + zh_n) - G(z_i^*)) f_z^*(z_i + zh_n)] dz \right\}. \end{aligned}$$

Now note that since G and f_z^* both have up to fourth order bounded derivatives, we have that

$$\begin{aligned} & (G(z_i^* + zh_n) - G(z_i^*)) f_z^*(z_i + zh_n) \\ &= \left(G'(z_i^*) zh_n + \frac{1}{2} G''(z_i^*) z^2 h_n^2 + \frac{1}{6} G'''(z_i^*) z^3 h_n^3 + O(z^4 h_n^4) \right) (f_z^*(z_i^*) + O(zh_n)) \\ &= G'(z_i^*) f_z^*(z_i^*) zh_n + \frac{1}{2} G''(z_i^*) f_z^*(z_i^*) z^2 h_n^2 + \frac{1}{6} G'''(z_i^*) f_z^*(z_i^*) z^3 h_n^3 + O(z^4 h_n^4). \end{aligned}$$

So

$$\int [K(z) (G(z_i^* + zh_n) - G(z_i^*)) f_z^*(z_i + zh_n)] dz = O(h_n^3),$$

where the bound does not depend on i . So

$$\max_{1 \leq i \leq n} \left\| \mathbb{E} \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (G(z_j^*) - G(z_i^*)) \mathbf{X}_i^\phi \right] \right\| = O(h_n^3).$$

On the other side, we have that we have that

$$\begin{aligned} \max_{1 \leq i \leq n} \left\| \frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \mathbf{X}_i^\phi \right. \\ \left. - \mathbb{E}_{\mathcal{D}_n} \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \mathbf{X}_i^\phi \right] \right\| = O_p \left(\sqrt{(\log n)/nh_n^2} \right). \end{aligned}$$

Together we have that

$$\max_{1 \leq i \leq n} \left\| \frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \mathbf{X}_i^\phi \right\| = O_p \left(h_n^{-1} \sqrt{(\log n)/n} + h_n^3 \right).$$

So

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \right] \left[\frac{1}{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)} - \frac{1}{f_z(z_i^*)} \right] \mathbf{X}_i^\phi \right\| \\ & \leq \max_{1 \leq i \leq n} \left\| \frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) (y_j - G(z_i^*)) \mathbf{X}_i^\phi \right\| \max_{1 \leq i \leq n} \left| \frac{1}{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)} - \frac{1}{f_z(z_i^*)} \right| \\ & = O_p(h_n^{-2} (\log n)/n + h_n^6) = o_p(1/\sqrt{n}), \end{aligned}$$

according to our choice of h_n , so term (i) is $o_p(1/\sqrt{n})$.

For term (ii), without of loss of generality, we assume that $\mathbf{X}_i^\phi = X_i^\phi$ is a scalar; the general case can be proved similarly. We note that

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^n \varepsilon_i \left[\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) \right] \left[\frac{1}{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)} - \frac{1}{f_z^*(z_i^*)} \right] X_i^\phi \right] \\ & = \mathbb{E} \sum_{i=1}^n \mathbb{E} \left\{ \varepsilon_i \left[1 - \frac{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)}{f_z^*(z_i^*)} \right] X_i^\phi \middle| X_i \right\} = 0 \end{aligned}$$

due to the fact that the data is i.i.d. and that $\mathbb{E}(\varepsilon_i | \mathbf{X}_{e,i}) = 0$ for all i . Moreover,

$$\begin{aligned}
& \mathbb{V} \left[\frac{1}{n} \sum_{i=1}^n \varepsilon_i \left[1 - \frac{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)}{f_z^*(z_i^*)} \right] X_i^{\phi 2} \right] \\
&= \frac{1}{n} \mathbb{E} \left\{ G(z_i^*) (1 - G(z_i^*)) \left[1 - \frac{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)}{f_z^*(z_i^*)} \right]^2 X_i^{\phi 2} \right\} \\
&\leq \frac{C}{n} \mathbb{E} \left\{ \left(\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) - f_z^*(z_i^*) \right)^2 X_i^{\phi 2} \right\} \\
&= \frac{C}{n^3} \mathbb{E} X_i^{\phi 2} \left(\sum_{j \neq i, k \neq i, j \neq k}^n \mathbb{E} \left[(K_{h_n}(z_j^* - z_i^*) - f_z^*(z_i^*)) (K_{h_n}(z_k^* - z_i^*) - f_z^*(z_i^*)) | X_i^{\phi} \right] + O(nh_n^{-1}) \right)
\end{aligned}$$

Note that $\mathbb{E} \left[(K_{h_n}(z_j^* - z_i^*) - f_z^*(z_i^*)) (K_{h_n}(z_k^* - z_i^*) - f_z^*(z_i^*)) | X_i^{\phi} \right]$ is $O(h_n^6)$ for all $k \neq j$, $j \neq i$, and $k \neq i$. So the above term is of order $O(h_n^6/n + h_n^{-1}/n^2)$, implying that

$$\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \left[1 - \frac{\frac{1}{n} \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*)}{f_z^*(z_i^*)} \right] \mathbf{X}_i^{\phi} \right\| = O_p \left(h_n^3/\sqrt{n} + 1/(n\sqrt{h_n}) \right) = o_p(1/\sqrt{n}),$$

according to the choice of h_n . This proves the first result.

Now we obtain the asymptotic distribution of

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \right) \mathbf{X}_i^{\phi}.$$

First note that

$$\begin{aligned}
& \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \right) \mathbf{X}_i^{\phi} \\
&= \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \mathbf{X}_i^{\phi} + \frac{y_i - y_j}{f_z^*(z_j^*)} \mathbf{X}_j^{\phi} \right) \\
&= \frac{1}{n^2} \sum_{i=1}^{n-1} \sum_{j=i+1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \mathbf{X}_i^{\phi} + \frac{y_i - y_j}{f_z^*(z_j^*)} \mathbf{X}_j^{\phi} \right) \\
&= \frac{n(n-1)}{2n^2} \binom{n}{2}^{-1} \sum_{i=1}^{n-1} \sum_{j=i+1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \mathbf{X}_i^{\phi} + \frac{y_i - y_j}{f_z^*(z_j^*)} \mathbf{X}_j^{\phi} \right).
\end{aligned}$$

Let $\mathbb{E}_{j|i}$ be the expectation with respect to the j -th observation conditional on the i -th

observation. Note that

$$\begin{aligned}
& \mathbb{E}_{j|i} \left[K_{h_n} (z_j^* - z_i^*) \frac{y_j - y_i}{f_z^* (z_i^*)} \mathbf{X}_i^\phi \right] \\
&= \frac{\mathbf{X}_i^\phi}{f_z^* (z_i^*)} \mathbb{E}_{j|i} [K_{h_n} (z_j^* - z_i^*) (G (z_j^*) - y_i)] \\
&= \frac{\mathbf{X}_i^\phi}{f_z^* (z_i^*)} \int K (z) (G (z_i^* + h_n z) - y_i) f_z^* (z_i^* + h_n z) dz \\
&= \frac{\mathbf{X}_i^\phi}{f_z^* (z_i^*)} \int K (z) \left(G (z_i^*) + G' (z_i^*) z h_n + \frac{1}{2} G'' (z_i^*) z^2 h_n^2 + O (z^3 h_n^3) - y_i \right) (f_z^* (z_i^*) + O (z h_n)) dz \\
&= \frac{\mathbf{X}_i^\phi}{f_z^* (z_i^*)} \int K (z) (G (z_i^*) - y_i) f_z^* (z_i^*) dz + O (h_n^3) = \mathbf{X}_i^\phi (G (z_i^*) - y_i) + O (h_n^3),
\end{aligned}$$

and

$$\begin{aligned}
& \mathbb{E}_{j|i} \left[K_{h_n} (z_j^* - z_i^*) \left(\frac{y_i - y_j}{f_z^* (z_j^*)} \right) \mathbf{X}_j^\phi \right] \\
&= \int \frac{1}{h_n} K \left(\frac{z - z_i^*}{h_n} \right) \left(\frac{y_i - G (z)}{f_z^* (z)} \right) \mathbf{X}^\phi f_{\mathbf{X},z}^* (\mathbf{X}, z) dz d\mathbf{X} \\
&= \int K (z) \frac{y_i - G (z_i^* + h_n z)}{f_z^* (z_i^* + h_n z)} \mathbf{X}^\phi f_{\mathbf{X},z}^* (\mathbf{X}, z_i^* + h_n z) dz d\mathbf{X} \\
&= (y_i - G (z_i^*)) \int \mathbf{X}^\phi f_{\mathbf{X}|z}^* (\mathbf{X} | z_i^*) d\mathbf{X} + O (h_n^2) \\
&= (y_i - G (z_i^*)) \mathbb{E} (\mathbf{X}^\phi | z_i^*) + O (h_n^3).
\end{aligned}$$

So

$$\mathbb{E}_{j|i} \left[K_{h_n} (z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^* (z_i^*)} \mathbf{X}_i^\phi + \frac{y_i - y_j}{f_z^* (z_j^*)} \mathbf{X}_j^\phi \right) \right] = -\varepsilon_i (\mathbf{X}_i^\phi - \mathbb{E} (\mathbf{X}_i^\phi | z_i^*)) + O (h_n^3).$$

We also note that

$$\begin{aligned}
& \mathbb{E} \left\| K_{h_n} (z_j^* - z_i^*) \left(\frac{y_i - y_j}{f_z^* (z_j^*)} \right) \mathbf{X}_j^\phi \right\|^2 \leq C \mathbb{E} (K_{h_n}^2 (z_j^* - z_i^*)) = O (h_n^{-2}) = o(n), \\
& \mathbb{E}_i \mathbb{E}_{j|i} \left[K_{h_n} (z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^* (z_i^*)} \mathbf{X}_i^\phi + \frac{y_i - y_j}{f_z^* (z_j^*)} \mathbf{X}_j^\phi \right) \right] = O (h_n^3) = o \left(\frac{1}{\sqrt{n}} \right),
\end{aligned}$$

so according to [Powell et al. \(1989\)](#), we have that

$$\begin{aligned} & \sqrt{n} \binom{n}{2}^{-1} \sum_{i=1}^{n-1} \sum_{j=i+1}^n K_{h_n}(z_j^* - z_i^*) \left(\frac{y_j - y_i}{f_z^*(z_i^*)} \mathbf{X}_i^\phi + \frac{y_i - y_j}{f_z^*(z_j^*)} \mathbf{X}_j^\phi \right) \\ &= -\frac{2}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \left(\mathbf{X}_i^\phi - \mathbb{E} \left(\mathbf{X}_i^\phi \middle| z_i^* \right) \right) + o_p(1). \end{aligned}$$

This implies that

$$\sqrt{n} \boldsymbol{\xi}_n^\phi = -\frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \left(\mathbf{X}_i^\phi - \mathbb{E} \left(\mathbf{X}_i^\phi \middle| z_i^* \right) \right) + o_p(1) \rightarrow_d N \left(0, \Sigma_{\boldsymbol{\xi}}^\phi \right).$$

□

Lemma A.9. Suppose that [Assumption 1](#), [Assumption 2\(i\)](#) and (ii), and [Assumption 6](#) hold, we have that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\Gamma_{q,n}(\boldsymbol{\beta}) - \Gamma_q(\boldsymbol{\beta})\| = O_p(\chi_{1,n}).$$

Proof of Lemma A.9. This is a direct result of [Lemma A.1](#) by noting that $|r_s(z) r_s(z)| \leq D_{q,0}^2$ and $\|\partial(r_s(X_0 + \mathbf{X}^T \boldsymbol{\beta}) r_s(X_0 + \mathbf{X}^T \boldsymbol{\beta})) / \partial \boldsymbol{\beta}\| \leq C \sqrt{p} D_{q,0} D_{q,1}$. □

Lemma A.10. Suppose that [Assumption 1](#), [Assumption 2\(i\)](#) and (ii), and [Assumption 6](#) hold, and $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$. We have that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\Gamma_{q,n}^{-1}(\boldsymbol{\beta}) - \Gamma_q^{-1}(\boldsymbol{\beta})\| = O_p(\chi_{1,n}).$$

Proof of Lemma A.10. First note that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} |\underline{\lambda}(\Gamma_{q,n}(\boldsymbol{\beta})) - \underline{\lambda}(\Gamma_q(\boldsymbol{\beta}))| \leq \sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\Gamma_{q,n}(\boldsymbol{\beta}) - \Gamma_q(\boldsymbol{\beta})\| = O_p(\chi_{1,n}),$$

and

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} |\bar{\lambda}(\Gamma_{q,n}(\boldsymbol{\beta})) - \bar{\lambda}(\Gamma_q(\boldsymbol{\beta}))| \leq \sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\Gamma_{q,n}(\boldsymbol{\beta}) - \Gamma_q(\boldsymbol{\beta})\| = O_p(\chi_{1,n}).$$

Since $\chi_{1,n} \rightarrow 0$, we have that with probability going to 1, there holds

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\lambda}(\Gamma_{q,n}(\boldsymbol{\beta})) \leq \frac{3\bar{\lambda}_\Gamma}{2}, \quad \inf_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\lambda}(\Gamma_{q,n}(\boldsymbol{\beta})) \geq \frac{\bar{\lambda}_\Gamma}{2},$$

indicating that $\sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\lambda}(\Gamma_{q,n}^{-1}(\boldsymbol{\beta})) = O_p(1)$.

Note that for any positive semi-definite matrices A and B , there holds $\min\{\lambda_A \|B\|, \lambda_B \|A\|\} \leq$

$\|AB\| \leq \max \{ \bar{\lambda}_A \|B\|, \bar{\lambda}_B \|A\| \}$, so we have that

$$\begin{aligned} \sup_{\beta \in \mathcal{B}} \|\Gamma_{q,n}^{-1}(\beta) - \Gamma_q^{-1}(\beta)\| &= \sup_{\beta \in \mathcal{B}} \|\Gamma_{q,n}^{-1}(\beta) (\Gamma_{q,n}(\beta) - \Gamma_q(\beta)) \Gamma_q^{-1}(\beta)\| \\ &\leq \left(\sup_{\beta \in \mathcal{B}} \bar{\lambda}(\Gamma_{q,n}^{-1}(\beta)) \right) \left(\sup_{\beta \in \mathcal{B}} \bar{\lambda}(\Gamma_q^{-1}(\beta)) \right) \sup_{\beta \in \mathcal{B}} \|\Gamma_{q,n}(\beta) - \Gamma_q(\beta)\| = O_p(\chi_{1,n}). \end{aligned}$$

□

Lemma A.11. Suppose that [Assumption 1](#), [Assumption 2\(i\)](#) and [\(ii\)](#), and [Assumption 6](#) hold, and moreover $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$. Define

$$\mathcal{Z} = \{z : z = X_0 + \mathbf{X}^T \beta \text{ for some } \mathbf{X}_e \in \mathcal{X}_e \text{ and } \beta \in \mathcal{B}\}.$$

We have that

$$\sup_{z \in \mathcal{Z}} \sup_{\beta \in \mathcal{B}} \|\mathfrak{X}_{q,n}(z, \beta) - \mathfrak{X}_q(z, \beta)\| = O_p(\sqrt{pq} D_{q,0}^2 \chi_{1,n}).$$

Proof of [Lemma A.11](#). Note that

$$\begin{aligned} &\sup_{z \in \mathcal{Z}} \sup_{\beta \in \mathcal{B}} \|\mathfrak{X}_{q,n}(z, \beta) - \mathfrak{X}_q(z, \beta)\| \\ &\leq \sup_{z \in \mathcal{Z}} \sup_{\beta \in \mathcal{B}} \left\| \mathfrak{X}_{q,n}(z, \beta) - \frac{1}{n} \sum_{i=1}^n (\mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q(z) \mathbf{X}_i) \right\| \\ &\quad + \sup_{z \in \mathcal{Z}} \sup_{\beta \in \mathcal{B}} \left\| \mathfrak{X}_q(z, \beta) - \frac{1}{n} \sum_{i=1}^n (\mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q(z) \mathbf{X}_i) \right\|. \end{aligned}$$

For the first term, we have that

$$\begin{aligned} &\sup_{z \in \mathcal{Z}} \sup_{\beta \in \mathcal{B}} \left\| \mathfrak{X}_{q,n}(z, \beta) - \frac{1}{n} \sum_{i=1}^n (\mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q(z) \mathbf{X}_i) \right\| \\ &= \frac{1}{n} \sum_{i=1}^n \sup_{z \in \mathcal{Z}} \sup_{\beta \in \mathcal{B}} \left\| \mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta) (\Gamma_{q,n}^{-1}(\beta) - \Gamma_q^{-1}(\beta)) \mathbf{r}_q(z) \mathbf{X}_i \right\| \\ &\leq C \sqrt{pq} D_{q,0}^2 \|\Gamma_{q,n}^{-1}(\beta) - \Gamma_q^{-1}(\beta)\| = O_p(\sqrt{pq} D_{q,0}^2 \chi_{1,n}). \end{aligned}$$

For the second term, we note that

$$\begin{aligned} &\sup_{z \in \mathcal{Z}} \sup_{\beta \in \mathcal{B}} \left\| \mathfrak{X}_q(z, \beta) - \frac{1}{n} \sum_{i=1}^n (\mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q(z) \mathbf{X}_i) \right\| \\ &\leq \sup_{\beta \in \mathcal{B}} \sup_{\tilde{\beta} \in \mathcal{B}} \sup_{\mathbf{X}_e \in \mathcal{X}_e} \left\| \mathfrak{X}_q(X_0 + \mathbf{X}^T \tilde{\beta}, \beta) - \frac{1}{n} \sum_{i=1}^n (\mathbf{r}_q^T(X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q(X_0 + \mathbf{X}^T \tilde{\beta}) \mathbf{X}_i) \right\|, \end{aligned}$$

where uniformly for all $\beta, \beta_1, \beta_2, \tilde{\beta} \in \mathcal{B}$, $\mathbf{X}_e \in \mathcal{X}_e$, and $\mathbf{X}_i \in \mathcal{X}$, there hold

$$\left| \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q (X_0 + \mathbf{X}^T \tilde{\beta}) X_{i,j} \right| \leq C q D_{q,0}^2,$$

and

$$\left\| \frac{\partial \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q (X_0 + \mathbf{X}^T \tilde{\beta}) X_{i,j}}{\partial \mathbf{X}_e} \right\| \leq C \sqrt{p} q D_{q,0} D_{q,1},$$

$$\left\| \frac{\partial \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q (X_0 + \mathbf{X}^T \tilde{\beta}) X_{i,j}}{\partial \tilde{\beta}} \right\| \leq C \sqrt{p} q D_{q,0} D_{q,1},$$

$$\begin{aligned} & \left\| \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta_1) \Gamma_q^{-1}(\beta_1) \mathbf{r}_q (X_0 + \mathbf{X}^T \tilde{\beta}) - \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta_2) \Gamma_q^{-1}(\beta_2) \mathbf{r}_q (X_0 + \mathbf{X}^T \tilde{\beta}) \right\| \\ & \leq \left\| (\mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta_1) - \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta_2)) \Gamma_q^{-1}(\beta_1) \mathbf{r}_q (X_0 + \mathbf{X}^T \tilde{\beta}) \right\| \\ & + \left\| \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta_2) (\Gamma_q^{-1}(\beta_1) - \Gamma_q^{-1}(\beta_2)) \mathbf{r}_q (X_0 + \mathbf{X}^T \tilde{\beta}) \right\| \\ & \leq C \sqrt{p} q D_{q,0} D_{q,1} \|\beta_1 - \beta_2\| + C q D_{q,0}^2 \|\Gamma_q(\beta_1) - \Gamma_q(\beta_2)\| \leq C \sqrt{p} q^2 D_{q,0}^3 D_{q,1} \|\beta_1 - \beta_2\|. \end{aligned}$$

So we have that the second term is of order $O_p(\sqrt{p} \chi_{1,n})$. This finishes the proof. \square

Lemma A.12. Suppose that [Assumption 1](#), [Assumption 2\(i\)-\(iii\)](#), and [Assumption 6](#) hold with $v_G \geq 1$, and that $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$, then we have that

$$\begin{aligned} & \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_q(X_{0,i} + \mathbf{X}_i^T \beta, \beta)) (G(X_{0,i} + \mathbf{X}_i^T \beta) - G(X_{0,i} + \mathbf{X}_i^T \beta^*)) \right. \\ & \quad \left. - \mathbb{E}((\mathbf{X}_i - \mathfrak{X}_q(X_{0,i} + \mathbf{X}_i^T \beta, \beta)) (G(X_{0,i} + \mathbf{X}_i^T \beta) - G(X_{0,i} + \mathbf{X}_i^T \beta^*))) \right\| = O_p(\sqrt{p} \chi_{1,n}). \end{aligned}$$

Proof of Lemma A.12. We only need to note that uniformly for all $\mathbf{X}_{e,i}$, $1 \leq j \leq p$, and $\beta, \beta_1, \beta_2 \in \mathcal{B}$, there hold

$$\begin{aligned} & |(X_{i,j} - \mathbb{E}_{\mathbf{X}_e}(\mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta) \Gamma_q^{-1}(\beta) \mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta) X_j)) (G(X_{0,i} + \mathbf{X}_i^T \beta) - G(X_{0,i} + \mathbf{X}_i^T \beta^*))| \\ & \leq C q D_{q,0}^2, \end{aligned}$$

and

$$\begin{aligned}
& \left\| G(X_{0,i} + \mathbf{X}_i^T \beta_1) \mathbb{E}_{\mathbf{X}_e} (\mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta_1) \Gamma_q^{-1}(\beta_1) \mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta_1) X_j) \right. \\
& \quad \left. - G(X_{0,i} + \mathbf{X}_i^T \beta_2) \mathbb{E}_{\mathbf{X}_e} (\mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta_2) \Gamma_q^{-1}(\beta_2) \mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta_2) X_j) \right\| \\
& \leq \left\| (G(X_{0,i} + \mathbf{X}_i^T \beta_1) - G(X_{0,i} + \mathbf{X}_i^T \beta_2)) \mathbb{E}_{\mathbf{X}_e} (\mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta_1) \Gamma_q^{-1}(\beta_1) \mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta_1) X_j) \right\| \\
& \quad + \left\| G(X_{0,i} + \mathbf{X}_i^T \beta_2) \mathbb{E}_{\mathbf{X}_e} ((\mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta_1) - \mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta_2)) \Gamma_q^{-1}(\beta_1) \mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta_1) X_j) \right\| \\
& \quad + \left\| G(X_{0,i} + \mathbf{X}_i^T \beta_2) \mathbb{E}_{\mathbf{X}_e} (\mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta_2) (\Gamma_q^{-1}(\beta_1) - \Gamma_q^{-1}(\beta_2)) \mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta_1) X_j) \right\| \\
& \quad + \left\| G(X_{0,i} + \mathbf{X}_i^T \beta_2) \mathbb{E}_{\mathbf{X}_e} (\mathbf{r}_q^T (X_0 + \mathbf{X}^T \beta_2) \Gamma_q^{-1}(\beta_2) (\mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta_1) - \mathbf{r}_q (X_{0,i} + \mathbf{X}_i^T \beta_2)) X_j) \right\| \\
& \leq C \sqrt{pq}^2 D_{q,0}^3 D_{q,1} \|\beta_1 - \beta_2\|.
\end{aligned}$$

□

Lemma A.13. Suppose that [Assumption 1](#), [Assumption 2\(i\)-\(iii\)](#), and [Assumption 6](#) hold with $v_G \geq 1$, and that $\chi_{1,n} \rightarrow 0$ as $n \rightarrow \infty$, then we have that

$$\sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T(z(\mathbf{X}_{e,i}, \beta)) \Gamma_{q,n}^{-1}(\beta) \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q^T(z(\mathbf{X}_{e,j}, \beta)) R_q(z(\mathbf{X}_{e,j}, \beta)) \right) \right\| = O_p(\sqrt{pq} D_{q,0}^2 \mathcal{E}_{q,0}),$$

$$\sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T(z(\mathbf{X}_{e,i}, \beta)) \Gamma_{q,n}^{-1}(\beta) \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q^T(z(\mathbf{X}_{e,j}, \beta)) \varepsilon_j \right) \right\| = O_p(\sqrt{p} \chi_{1,n}),$$

and

$$\sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n (R_q(z(\mathbf{X}_{e,i}, \beta)) \mathbf{X}_i + \varepsilon_i \mathbf{X}_i) \right\| = O_p(\sqrt{p} \mathcal{E}_{q,0} + \sqrt{p(\log p)/n}).$$

Proof of Lemma A.13. For the first result, we note that

$$\begin{aligned}
& \sup_{\beta \in \mathcal{B}} \left\| \mathbf{X}_i \mathbf{r}_q^T(z(\mathbf{X}_{e,i}, \beta)) \Gamma_{q,n}^{-1}(\beta) \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q^T(z(\mathbf{X}_{e,j}, \beta)) R_q(z(\mathbf{X}_{e,j}, \beta)) \right) \right\| \\
& = O_p \left(\sqrt{p} \sup_{\beta \in \mathcal{B}, \mathbf{X}_e \in \mathcal{X}_e} \|\mathbf{r}_q(z(\mathbf{X}_e, \beta))\| \sup_{\beta \in \mathcal{B}, \mathbf{X}_e \in \mathcal{X}_e} \|\mathbf{r}_q(z(\mathbf{X}_e, \beta)) R_q(z(\mathbf{X}_e, \beta))\| \right) \\
& = O_p(\sqrt{pq} D_{q,0}^2 \mathcal{E}_{q,0}).
\end{aligned}$$

For the second result, we first have that

$$\sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{j=1}^n \mathbf{r}_q^T(z(\mathbf{X}_{e,j}, \beta)) \varepsilon_j \right\| = O_p \left(\sqrt{pq D_{q,0}^2 \log(pq D_{q,1} n) / n} \right),$$

due to the fact that $|r_l(z(\mathbf{X}_{e,j}, \beta)) \varepsilon_j| \leq C D_{q,0}$ and $\|(\partial r_l(z(\mathbf{X}_{e,j}, \beta)) / \partial \beta) \varepsilon_j\| \leq C \sqrt{p} D_{q,1}$

for all $0 \leq l \leq q$. So

$$\begin{aligned} & \sup_{\beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T(z(\mathbf{X}_{e,i}, \beta)) \Gamma_{q,n}^{-1}(\beta) \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q^T(z(\mathbf{X}_{e,j}, \beta)) \varepsilon_j \right) \right\| \\ &= O_p \left(\sqrt{pq} D_{q,0} \sqrt{pq D_{q,0}^2 \log(pq D_{q,1} n) / n} \right) = O_p(\sqrt{p} \chi_{1,n}). \end{aligned}$$

Finally for the third result, we have that $\left\| \frac{1}{n} \sum_{i=1}^n R_q(z(\mathbf{X}_{e,i}, \beta)) \mathbf{X}_i \right\| = O_p(\sqrt{p} \mathcal{E}_{q,0})$ and $\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_i \right\| = O_p(\sqrt{p(\log p) / n})$.

Combine the above results, we finish the proof. \square

Now we are ready to prove [Lemma 4](#) in the main text.

Proof of Lemma 4. We note that

$$\begin{aligned} \beta_{k+1} &= \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n \left(\widehat{G}(z_{i,k} | \beta_k) - y_i \right) \mathbf{X}_i \\ &= \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n \left(\mathbf{r}_q^T(z_{i,k}) \widehat{\pi}_{q,n,k} - \mathbf{r}_q^T(z_{i,k}) \pi_q^* \right) \mathbf{X}_i - \frac{\delta_k}{n} \sum_{i=1}^n (G(z_{i,k}) - G(z_i^*)) \mathbf{X}_i \\ &\quad + \frac{\delta_k}{n} \sum_{i=1}^n R_q(z_{i,k}) \mathbf{X}_i + \frac{\delta_k}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_i. \end{aligned}$$

Now we look at the $\widehat{\pi}_{q,n,k} - \pi_q^*$. Define $\Gamma_{q,n,k} = \Gamma_{q,n}(\beta_k)$, we have that

$$\begin{aligned} \widehat{\pi}_{q,n,k} &= \left(\frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(z_{i,k}) \mathbf{r}_q^T(z_{i,k}) \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(z_{i,k}) y_i \right) \\ &= \pi_q^* - \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(z_{i,k}) (G(z_{i,k}) - G(z_i^*)) \right) + \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(z_{i,k}) R_q(z_{i,k}) \right) \\ &\quad + \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(z_{i,k}) \varepsilon_i \right). \end{aligned}$$

Take the above expression of $\widehat{\pi}_{q,n,k} - \pi_q^*$ into the update of β_k , we have that

$$\begin{aligned} \beta_{k+1} &= \beta_k - \frac{\delta_k}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,n}(z_{i,k}, \beta_k)) (G(z_{i,k}) - G(z_i^*)) \\ &\quad - \frac{\delta_k}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T(z_{i,k}) \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q(z_{j,k}) R_q(z_{j,k}) + \frac{1}{n} \sum_{j=1}^n \mathbf{r}_q(z_{j,k}) \varepsilon_j \right) \\ &\quad + \frac{\delta_k}{n} \sum_{i=1}^n (R_q(z_{i,k}) \mathbf{X}_i + \varepsilon_i \mathbf{X}_i). \end{aligned}$$

If we define

$$\begin{aligned}
\mathfrak{R}_{n,k} &= \mathbb{E} (\mathbf{X} - \mathfrak{X}_q(z(\mathbf{X}_e, \boldsymbol{\beta}_k), \boldsymbol{\beta}_k)) (G(z(\mathbf{X}_e, \boldsymbol{\beta}_k)) - G(z(\mathbf{X}_e, \boldsymbol{\beta}^*))) \\
&\quad - \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_q(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}_k), \boldsymbol{\beta}_k)) (G(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}_k)) - G(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}^*))) \\
&\quad + \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_q(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}_k), \boldsymbol{\beta}_k)) (G(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}_k)) - G(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}^*))) \\
&\quad - \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}_k), \boldsymbol{\beta}_k)) (G(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}_k)) - G(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}^*))) \\
&\quad - \frac{\delta_k}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^\top(z_{i,k}) \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_q(z_{j,k}) R_q(z_{j,k}) + \frac{1}{n} \sum_{j=1}^n \mathbf{r}_q(z_{j,k}) \varepsilon_j \right) \\
&\quad + \frac{\delta_k}{n} \sum_{i=1}^n (R_q(z_{i,k}) \mathbf{X}_i + \varepsilon_i \mathbf{X}_i),
\end{aligned}$$

we have that

$$\boldsymbol{\beta}_{k+1} = \boldsymbol{\beta}_k - \delta_k \mathbb{E} [(\mathbf{X} - \mathfrak{X}_q(z(\mathbf{X}_e, \boldsymbol{\beta}_k), \boldsymbol{\beta}_k)) (G(z(\mathbf{X}_e, \boldsymbol{\beta}_k)) - G(z(\mathbf{X}_e, \boldsymbol{\beta}^*)))] + \delta_k \mathfrak{R}_{n,k}.$$

It remains to verify the order of $\sup_{k \geq 1} \|\mathfrak{R}_{n,k}\|$, which is done based on [Lemma A.11](#), [Lemma A.12](#), and [Lemma A.13](#). \square

Now we prove [Lemma 5](#) and [Lemma 6](#) in the main text.

Proof of Lemma 5. Recall that

$$\Psi_q(t, \boldsymbol{\beta}) = \mathbb{E} [G'(z(\mathbf{X}_e, \boldsymbol{\beta}^*) + t \mathbf{X}^\top \Delta \boldsymbol{\beta}) (\mathbf{X} \mathbf{X}^\top - \mathfrak{X}_q(z(\mathbf{X}_e, \boldsymbol{\beta}), \boldsymbol{\beta}) \mathbf{X}^\top)] .$$

We have that

$$\begin{aligned}
&\sup_{0 \leq t \leq 1, \boldsymbol{\beta} \in \mathcal{B}_n} \left\| \frac{1}{n} \sum_{i=1}^n G'(z_i^* + t \mathbf{X}_i^\top \Delta \boldsymbol{\beta}) (\mathbf{X}_i \mathbf{X}_i^\top - \mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) \mathbf{X}_i^\top) - \Psi_q^* \right\| \\
&\leq \sup_{0 \leq t \leq 1, \boldsymbol{\beta} \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n G'(z_i^* + t \mathbf{X}_i^\top \Delta \boldsymbol{\beta}) (\mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) - \mathfrak{X}_q(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta})) \mathbf{X}_i^\top \right\| \\
&\quad + \sup_{0 \leq t \leq 1, \boldsymbol{\beta} \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n G'(z_i^* + t \mathbf{X}_i^\top \Delta \boldsymbol{\beta}) (\mathbf{X}_i \mathbf{X}_i^\top - \mathfrak{X}_q(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) \mathbf{X}_i^\top) - \Psi_q(t, \boldsymbol{\beta}) \right\| \\
&\quad + \sup_{0 \leq t \leq 1, \boldsymbol{\beta} \in \mathcal{B}_n} \|\Psi_q(t, \boldsymbol{\beta}) - \Psi_q^*\|.
\end{aligned}$$

From [Lemma A.11](#), we know that

$$\sup_{z \in \mathcal{Z}} \sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\mathfrak{X}_{q,n}(z, \boldsymbol{\beta}) - \mathfrak{X}_q(z, \boldsymbol{\beta})\| = O_p(\sqrt{pq} D_{q,0}^2 \chi_{1,n}),$$

and as a result,

$$\begin{aligned} & \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n G' \left(z_i^* + t \mathbf{X}_i^T \Delta \beta \right) \left(\mathfrak{X}_{q,n} \left(z \left(\mathbf{X}_{e,i}, \beta \right), \beta \right) - \mathfrak{X}_q \left(z \left(\mathbf{X}_{e,i}, \beta \right), \beta \right) \right) \mathbf{X}_i^T \right\| \\ &= O_p \left(pq D_{q,0}^2 \chi_{1,n} \right). \end{aligned}$$

For the second term, we have that

$$\begin{aligned} & \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n G' \left(z_i^* + t \mathbf{X}_i^T \Delta \beta \right) \left(\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n} \left(z \left(\mathbf{X}_{e,i}, \beta \right), \beta \right) \mathbf{X}_i^T \right) - \Psi_q \left(t, \beta \right) \right\| \\ &= O_p \left(\sqrt{p^3 q^2 D_{q,0}^4 \log(pq D_{q,0} D_{q,1} n) / n} \right) = O_p(p \chi_{1,n}), \end{aligned}$$

due to the fact that

$$\left| G' \left(z_i^* + t \mathbf{X}_i^T \Delta \beta \right) \left(X_{i,s} X_{i,t} - \left(\mathfrak{X}_q \left(z \left(\mathbf{X}_{e,i}, \beta \right), \beta \right) \right)_s X_{i,t} \right) \right| \leq C q D_{q,0}^2,$$

and

$$\begin{aligned} & \left| G' \left(z_i^* + t \mathbf{X}_i^T \Delta \beta_1 \right) \left(X_{i,s} X_{i,t} - \left(\mathfrak{X}_q \left(z \left(\mathbf{X}_{e,i}, \beta_1 \right), \beta_1 \right) \right)_s X_{i,t} \right) \right. \\ & \quad \left. - G' \left(z_i^* + t \mathbf{X}_i^T \Delta \beta_2 \right) \left(X_{i,s} X_{i,t} - \left(\mathfrak{X}_q \left(z \left(\mathbf{X}_{e,i}, \beta_2 \right), \beta_2 \right) \right)_s X_{i,t} \right) \right| \\ & \leq C \sqrt{p} q^2 D_{q,0}^3 D_{q,1} \|\beta_1 - \beta_2\|. \end{aligned}$$

Finally,

$$\begin{aligned} & \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}_n} \left\| \Psi_q \left(t, \beta \right) - \Psi_q^* \right\| \\ & \leq \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}_n} \left\| \mathbb{E} \left[G' \left(z \left(\mathbf{X}_e, \beta^* \right) + t \mathbf{X}^T \Delta \beta \right) - G' \left(z \left(\mathbf{X}_e, \beta^* \right) \right) \left(\mathbf{X} \mathbf{X}^T - \mathfrak{X}_q \left(z \left(\mathbf{X}_e, \beta \right), \beta \right) \mathbf{X}^T \right) \right] \right\| \\ & \quad + \sup_{\beta \in \mathcal{B}_n} \left\| \mathbb{E} \left[G' \left(z \left(\mathbf{X}_e, \beta^* \right) \right) \left(\mathfrak{X}_q \left(z \left(\mathbf{X}_e, \beta \right), \beta \right) - \mathfrak{X}_q \left(z \left(\mathbf{X}_e, \beta^* \right), \beta^* \right) \mathbf{X}^T \right) \right] \right\|. \end{aligned}$$

Obviously the first term is bounded by $C \sqrt{p^3} q D_{q,0}^2 \sup_{\beta \in \mathcal{B}_n} \|\Delta \beta\|$, while the second term is

bounded by

$$\begin{aligned}
& Cp \sup_{\mathbf{X}_e, \tilde{\mathbf{X}}_e} \left\| \left(\mathbf{r}_q^T \left(z \left(\tilde{\mathbf{X}}_e, \boldsymbol{\beta} \right) \right) \Gamma_q^{-1}(\boldsymbol{\beta}) \mathbf{r}_q \left(z \left(\mathbf{X}_e, \boldsymbol{\beta} \right) \right) \right) - \left(\mathbf{r}_q^T \left(z \left(\tilde{\mathbf{X}}_e, \boldsymbol{\beta}^* \right) \right) \Gamma_q^{-1}(\boldsymbol{\beta}^*) \mathbf{r}_q \left(z \left(\mathbf{X}_e, \boldsymbol{\beta}^* \right) \right) \right) \right\| \\
& \leq Cp \sup_{\mathbf{X}_e, \tilde{\mathbf{X}}_e} \left\| \left(\mathbf{r}_q \left(z \left(\tilde{\mathbf{X}}_e, \boldsymbol{\beta} \right) \right) - \mathbf{r}_q \left(z \left(\tilde{\mathbf{X}}_e, \boldsymbol{\beta}^* \right) \right) \right)^T \Gamma_q^{-1}(\boldsymbol{\beta}) \mathbf{r}_q \left(z \left(\mathbf{X}_e, \boldsymbol{\beta} \right) \right) \right\| \\
& + Cp \sup_{\mathbf{X}_e, \tilde{\mathbf{X}}_e} \left\| \mathbf{r}_q \left(z \left(\tilde{\mathbf{X}}_e, \boldsymbol{\beta}^* \right) \right)^T \left(\Gamma_q^{-1}(\boldsymbol{\beta}) - \Gamma_q^{-1}(\boldsymbol{\beta}^*) \right) \mathbf{r}_q \left(z \left(\mathbf{X}_e, \boldsymbol{\beta} \right) \right) \right\| \\
& + Cp \sup_{\mathbf{X}_e, \tilde{\mathbf{X}}_e} \left\| \mathbf{r}_q \left(z \left(\tilde{\mathbf{X}}_e, \boldsymbol{\beta}^* \right) \right)^T \Gamma_q^{-1}(\boldsymbol{\beta}^*) \left(\mathbf{r}_q \left(z \left(\mathbf{X}_e, \boldsymbol{\beta} \right) \right) - \mathbf{r}_q \left(z \left(\mathbf{X}_e, \boldsymbol{\beta}^* \right) \right) \right) \right\| \\
& \leq C \sqrt{p^3 q^2 D_{q,0}^3 D_{q,1}} \sup_{\boldsymbol{\beta} \in \mathcal{B}_n} \|\Delta \boldsymbol{\beta}\|.
\end{aligned}$$

So

$$\sup_{0 \leq t \leq 1, \boldsymbol{\beta} \in \mathcal{B}_n} \|\Psi_q(t, \boldsymbol{\beta}) - \Psi_q^*\| = O_p \left(\sqrt{p^3 q^2 D_{q,0}^3 D_{q,1}} \sup_{\boldsymbol{\beta} \in \mathcal{B}_n} \|\Delta \boldsymbol{\beta}\| \right).$$

Combine the above results, we have that

$$\begin{aligned}
& \sup_{0 \leq t \leq 1, \boldsymbol{\beta} \in \mathcal{B}_n} \left\| \frac{1}{n} \sum_{i=1}^n G' \left(z_i^* + t \mathbf{X}_i^T \Delta \boldsymbol{\beta} \right) \left(\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n} \left(z \left(\mathbf{X}_{e,i}, \boldsymbol{\beta} \right), \boldsymbol{\beta} \right) \mathbf{X}_i^T \right) - \Psi_q^* \right\| \\
& = O_p \left(pq D_{q,0}^2 \chi_{1,n} + \sqrt{p^3 q^2 D_{q,0}^3 D_{q,1}} \sup_{\boldsymbol{\beta} \in \mathcal{B}_n} \|\Delta \boldsymbol{\beta}\| \right).
\end{aligned}$$

□

Proof of Lemma 6. According to Theorem 7, we have that $\sup_{k \geq k_{1,n}^{SBGD}+1} \|\Delta \boldsymbol{\beta}_k\| = O_p(\chi_{2,n})$.

To prove the lemma, we first show that

$$\begin{aligned}
& \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_{q,j,k} R_{q,j,k} + \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j,k} \varepsilon_j - \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j}^* \varepsilon_j \right) \right\| \\
& = O_p \left(\sqrt{pq} D_{q,0}^2 \mathcal{E}_{q,0} + \sqrt{pq} D_{q,0} \chi_{2,n} \chi_{3,n} \right),
\end{aligned}$$

where $\chi_{3,n} = \sqrt{p^2 q D_{q,1}^2 \log(pq D_{q,2} n) / n}$. Note that

$$\begin{aligned}
\sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,i,k} \varepsilon_i - \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,i}^* \varepsilon_i \right\| &= \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \left\{ \int_0^1 \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{r}'_q \left(z_i^* + t \mathbf{X}_i^T \Delta \boldsymbol{\beta}_k \right) \mathbf{X}_i^T dt \right\} \Delta \boldsymbol{\beta}_k \right\| \\
&\leq \sup_{\boldsymbol{\beta} \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{r}'_q \left(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta} \right) \mathbf{X}_i^T \right\| \sup_{k \geq k_{1,n}^{SBGD}+1} \|\Delta \boldsymbol{\beta}_k\|,
\end{aligned}$$

Obviously, we have that $\sup_{\boldsymbol{\beta} \in \mathcal{B}} \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{r}'_q \left(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta} \right) \mathbf{X}_i^T \right\| = O_p(\chi_{3,n})$ due to the fact

that $|\varepsilon_i r'_s(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}) X_t| \leq C D_{q,1}$ and $\|\partial \varepsilon_i r'_s(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}) X_t / \partial \boldsymbol{\beta}\| \leq C \sqrt{p} D_{q,2}$, so

$$\sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j,k} \varepsilon_j - \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j}^* \varepsilon_j \right\| = O_p(\chi_{2,n} \chi_{3,n}),$$

which leads to the result if we further note that

$$\begin{aligned} & \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_{q,j,k} R_{q,j,k} + \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j,k} \varepsilon_j - \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j}^* \varepsilon_j \right) \right\| \\ &= O_p \left(\sqrt{pq} D_{q,0} \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{j=1}^n \mathbf{r}_{q,j,k} R_{q,j,k} + \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j,k} \varepsilon_j - \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j}^* \varepsilon_j \right\| \right) \\ &= O_p(\sqrt{pq} D_{q,0}^2 \mathcal{E}_{q,0} + \sqrt{pq} D_{q,0} \chi_{2,n} \chi_{3,n}). \end{aligned}$$

Next we show that

$$\begin{aligned} & \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} - \mathbb{E}(\mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*) \Gamma_q^{-1}(\boldsymbol{\beta}^*)) \right\| = O_p(p \sqrt{q^3} D_{q,0}^2 D_{q,1} \chi_{2,n}). \\ & \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*) \Gamma_q^{-1}(\boldsymbol{\beta}^*) \right\| \\ & \leq \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k) \Gamma_{q,n,k}^{-1} - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*) \Gamma_{q,n,k}^{-1} \right\| \\ & + \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*) \Gamma_{q,n,k}^{-1} - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*) \Gamma_q^{-1}(\boldsymbol{\beta}^*) \right\|. \end{aligned}$$

The first term is obviously bounded in probability by

$$\begin{aligned} & C \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i (\mathbf{r}_q(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}_k) - \mathbf{r}_q(X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*))^T \right\| \\ & \leq C p \sqrt{q} D_{q,1} \|\boldsymbol{\beta}_k - \boldsymbol{\beta}^*\| = C p \sqrt{q} D_{q,1} \chi_{2,n}. \end{aligned}$$

The second term is bounded by

$$\begin{aligned} & \sup_{k \geq k_{1,n}^*} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \boldsymbol{\beta}^*) \right\| \sup_{k \geq k_{1,n}^*} \|\Gamma_{q,n,k}^{-1} - \Gamma_q^{-1}(\boldsymbol{\beta}^*)\| \\ & \leq C \sqrt{pq} D_{q,0} \sup_{k \geq k_{1,n}^{SBGD}+1} \|\Gamma_{q,n,k}^{-1} - \Gamma_q^{-1}(\boldsymbol{\beta}^*)\|. \end{aligned}$$

Now we provide an upper bound for $\sup_{k \geq k_{1,n}^{SBGD}+1} \|\Gamma_{q,n,k}^{-1} - \Gamma_q^{-1}(\beta^*)\|$. Note that

$$\begin{aligned} \sup_{k \geq k_{1,n}^{SBGD}+1} \|\Gamma_{q,n,k}^{-1} - \Gamma_q^{-1}(\beta^*)\| &= O_p \left(\sup_{k \geq k_{1,n}^{SBGD}+1} \|\Gamma_{q,n,k} - \Gamma_q(\beta^*)\| \right) \\ &= O_p \left(\sup_{k \geq k_{1,n}^{SBGD}+1} \|\Gamma_{q,n,k} - \Gamma_{q,n}(\beta^*)\| + \|\Gamma_{q,n}(\beta^*) - \Gamma_q(\beta^*)\| \right) \\ &= O_p(\sqrt{pq}D_{q,0}D_{q,1}\chi_{2,n} + \chi_{1,n}) = O_p(\sqrt{pq}D_{q,0}D_{q,1}\chi_{2,n}). \end{aligned}$$

So

$$\begin{aligned} \sup_{k \geq k_{1,n}^*} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta^*) \Gamma_{q,n,k}^{-1} - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta^*) \Gamma_q^{-1}(\beta^*) \right\| \\ = O_p \left(p\sqrt{q^3}D_{q,0}^2D_{q,1}\chi_{2,n} \right), \end{aligned}$$

and together

$$\sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta^*) \Gamma_q^{-1}(\beta^*) \right\| = O_p \left(p\sqrt{q^3}D_{q,0}^2D_{q,1}\chi_{2,n} \right).$$

Moreover, note that $\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta^*) \Gamma_q^{-1}(\beta^*) - \mathbb{E}(\mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta^*) \Gamma_q^{-1}(\beta^*)) \right\| = O_p \left(\sqrt{p^3 D_{q,0}^2 \log(pn)/n} \right)$, so we have shown the results.

Based on the above results, we have that

$$\begin{aligned} \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_{q,j,k} R_{q,j,k} + \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j,k} \varepsilon_j \right) + \frac{1}{n} \sum_{i=1}^n R_q(z_{i,k}) \mathbf{X}_i - \frac{1}{n} \sum_{i=1}^n \mathfrak{X}(z_i^*, \beta^*) \varepsilon_j \right\| \\ \leq \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{r}_{q,j,k} R_{q,j,k} + \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j,k} \varepsilon_j - \frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j}^* \varepsilon_j \right) \right\| \\ + \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{r}_{q,i,k}^T \Gamma_{q,n,k}^{-1} - \mathbb{E}(\mathbf{X}_i \mathbf{r}_q^T (X_{0,i} + \mathbf{X}_i^T \beta^*) \Gamma_q^{-1}(\beta^*)) \right) \left(\frac{1}{n} \sum_{i=1}^n \mathbf{r}_{q,j}^* \varepsilon_j \right) \right\| \\ + \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \frac{1}{n} \sum_{i=1}^n R_q(z_{i,k}) \mathbf{X}_i \right\| \\ = O_p \left(\sqrt{pq}D_{q,0}^2\mathcal{E}_{q,0} + \sqrt{pq}D_{q,0}\chi_{2,n}\chi_{3,n} + p\sqrt{q^3}D_{q,0}^2D_{q,1}\chi_{2,n}\sqrt{(qD_{q,0}^2 \log q)/n} \right) \end{aligned}$$

□

B Proofs of Theorems

Proof of Theorem 1

Proof. We first prove Theorem 1(i). Recall that $\Delta\beta_{e,k} = \beta_{e,k} - \beta_e^*$ and $\varepsilon_i = y_i - G(\mathbf{X}_{e,i}^T \beta_e^*)$. We have that

$$\Delta\beta_{e,k+1} = \Delta\beta_{e,k} - \frac{\delta}{n} \sum_{i=1}^n (G(\mathbf{X}_{e,i}^T \beta_{e,k}) - G(\mathbf{X}_{e,i}^T \beta_e^*) - \varepsilon_i) \mathbf{X}_{e,i},$$

so

$$\|\Delta\beta_{e,k+1}\| \leq \left\| \Delta\beta_{e,k} - \frac{\delta}{n} \sum_{i=1}^n (G(\mathbf{X}_{e,i}^T \beta_{e,k}) - G(\mathbf{X}_{e,i}^T \beta_e^*)) \mathbf{X}_{e,i} \right\| + \left\| \frac{\delta}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\|.$$

Note that mean value theorem leads to

$$\begin{aligned} & \Delta\beta_{e,k} - \frac{\delta}{n} \sum_{i=1}^n (G(\mathbf{X}_{e,i}^T \beta_{e,k}) - G(\mathbf{X}_{e,i}^T \beta_e^*)) \mathbf{X}_{e,i} \\ &= \Delta\beta_{e,k} - \int_0^1 \left\{ \frac{\delta}{n} \sum_{i=1}^n G'(\mathbf{X}_{e,i}^T \beta_e^* + t \mathbf{X}_{e,i}^T \Delta\beta_{e,k}) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T \Delta\beta_{e,k} \right\} dt \\ &= \int_0^1 \{ (I_{p+1} - \delta M_n(\beta_e^* + t \Delta\beta_{e,k})) \Delta\beta_{e,k} \} dt, \end{aligned}$$

where the integration is understood to be element-wise, and $\beta_e^* + t \Delta\beta_{e,k} \in \mathcal{B}_e$ due to convexity of \mathcal{B}_e .

We next provide a uniform upper bound for $\bar{\lambda}(I_{p+1} - \delta M_n(\beta_e))$ and lower bound for $\underline{\lambda}(I_{p+1} - \delta M_n(\beta_e))$ with respect to $\beta_e \in \mathcal{B}_e$ in probability. Since Assumption 2 holds, we have that $G(\mathbf{X}_{e,i}^T \beta) X_{i,t} X_{i,s}$ is bounded by $\|G\|_\infty$ and $\|\partial G(\mathbf{X}_{e,i}^T \beta) X_{i,t} X_{i,s} / \partial \beta\| \leq C \sqrt{p}$. Then according to Lemma A.1, we have that

$$\sup_{\beta_e \in \mathcal{B}} \|M_n(\beta_e) - M(\beta_e)\| = O_p \left(\sqrt{\frac{p^3 \log n}{n}} \right).$$

Since $p^5 (\log p)^2 n^{-1} \rightarrow 0$ holds, $\sqrt{p^3 (\log n) / n} \rightarrow 0$ holds, so

$$\sup_{\beta_e \in \mathcal{B}} |\bar{\lambda}(M_n(\beta_e)) - \bar{\lambda}(M(\beta_e))| = o_p(1),$$

and

$$\sup_{\beta_e \in \mathcal{B}} |\underline{\lambda}(M_n(\beta_e)) - \underline{\lambda}(M(\beta_e))| = o_p(1).$$

Due to [Assumption 2](#)(iv), with probability going to 1, there holds,

$$\underline{\lambda}_e/2 \leq \inf_{\beta_e \in \mathcal{B}} \underline{\lambda}(M_n(\beta_e)) \leq \sup_{\beta_e \in \mathcal{B}} \bar{\lambda}(M_n(\beta_e)) \leq 3\bar{\lambda}_e/2.$$

Since $\delta < 2/(3\bar{\lambda}_e)$, we have that with probability going to 1, there holds

$$0 \leq \inf_{\beta_e \in \mathcal{B}} \bar{\lambda}(I_{p+1} - \delta M_n(\beta_e)) \leq \sup_{\beta_e \in \mathcal{B}} \bar{\lambda}(I_{p+1} - \delta M_n(\beta_e)) \leq 1 - \underline{\lambda}_e \delta/2.$$

Based on the above inequality, we have that with probability going to 1, there holds

$$\begin{aligned} & \left\| \int_0^1 \{ (I_{p+1} - \delta M_n(\beta_e^* + t\Delta\beta_{e,k})) \Delta\beta_{e,k} \} dt \right\| \\ & \leq \int_0^1 \left\{ \sup_{\beta_e \in \mathcal{B}} \bar{\lambda}(I_{p+1} - \delta M_n(\beta_e)) \right\} dt \cdot \|\Delta\beta_{e,k}\| \leq (1 - \underline{\lambda}_e \delta/2) \cdot \|\Delta\beta_{e,k}\|. \end{aligned}$$

So with probability going to 1, for all k there holds

$$\begin{aligned} \|\Delta\beta_{e,k+1}\| & \leq (1 - \underline{\lambda}_e \delta/2) \|\Delta\beta_{e,k}\| + \delta \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| \\ & \leq \dots \leq (1 - \underline{\lambda}_e \delta/2)^k \|\Delta\beta_{e,1}\| + \delta \sum_{j=1}^k (1 - \underline{\lambda}_e \delta/2)^{j-1} \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| \\ & \leq (1 - \underline{\lambda}_e \delta/2)^k \|\Delta\beta_{e,1}\| + 2\underline{\lambda}_e^{-1} \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\|. \end{aligned}$$

Note that for any $\tau > 0$,

$$\begin{aligned} P \left(\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| > \tau \right) & \leq \sum_{j=0}^p P \left(\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i X_{e,j,i} \right| > \frac{\tau}{\sqrt{p+1}} \right) \\ & \leq \sum_{j=0}^p 2 \exp(Cn\tau^2/p) = 2 \exp(C_1 \log p - C_2 n\tau^2/p), \end{aligned}$$

so

$$\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| = O_p \left(\sqrt{p(\log p)/n} \right).$$

Then for k such that

$$(1 - \underline{\lambda}_e \delta/2)^k \|\Delta\beta_{e,1}\| \leq \sqrt{p(\log p)/n},$$

or equivalently,

$$k \geq k_{1,n}^{BGD} = \frac{\log \|\Delta\beta_{e,1}\| + \frac{1}{2} \log(n/(p \log p))}{-\log(1 - \underline{\lambda}_e \delta/2)},$$

we have that

$$\|\Delta\beta_{e,k+1}\| = O_p\left(\sqrt{p(\log p)/n}\right).$$

This proves [Theorem 1\(i\)](#).

Next we prove [Theorem 1\(ii\)](#). For any $k \geq k_{1,n}^{BGD} + 1$, there holds

$$\begin{aligned}\Delta\beta_{e,k+1} &= \Delta\beta_{e,k} - \frac{\delta}{n} \sum_{i=1}^n (G(\mathbf{X}_{e,i}^T \beta_{e,k}) - G(\mathbf{X}_{e,i}^T \beta_e^*) - \varepsilon_i) \mathbf{X}_{e,i}, \\ &= (I_{p+1} - \delta M_n(\bar{\beta}_{e,k})) \Delta\beta_{e,k} + \frac{\delta}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i},\end{aligned}$$

where $\bar{\beta}_{e,k}$ is element-wise and lies between $\beta_{e,k}$ and β_e^* . Since $\|\Delta\beta_{e,k}\| = O_p\left(\sqrt{p(\log p)/n}\right)$ for $k \geq k_{1,n}^{BGD} + 1$, $\|\Delta\bar{\beta}_{e,k}\| = O_p\left(\sqrt{p(\log p)/n}\right)$ also holds. Note that

$$\|M_n(\bar{\beta}_{e,k}) - M(\beta_e^*)\| \leq \|M_n(\bar{\beta}_{e,k}) - M_n(\beta_e^*)\| + \|M_n(\beta_e^*) - M(\beta_e^*)\|.$$

For the second term, $\|M_n(\beta_e^*) - M(\beta_e^*)\| = O_p\left(\sqrt{p^2(\log p)/n}\right)$ obviously holds. For the first term, since G is twice differentiable with bounded derivatives, we have that

$$\begin{aligned}\sup_{k \geq k_{1,n}^{BGD} + 1} \|M_n(\bar{\beta}_{e,k}) - M_n(\beta_e^*)\| &\leq \sup_{k \geq k_{1,n}^{BGD} + 1} \frac{1}{n} \sum_{i=1}^n \|\mathbf{X}_{e,i} \mathbf{X}_{e,i}^T\| |G''(\mathbf{X}_{e,i}^T \check{\beta}_{e,k})| |\mathbf{X}_{e,i}^T \Delta\bar{\beta}_{e,k}|, \\ &\leq C\sqrt{p^3} \sup_{k \geq k_{1,n}^{BGD} + 1} \|\bar{\beta}_{e,k} - \beta_e^*\| = O_p\left(\sqrt{p^4(\log p)/n}\right),\end{aligned}$$

where $\check{\beta}_{e,k}$ lies somewhere between $\bar{\beta}_{e,k}$ and β_e^* and is also element-wise, and the second last inequality comes from the fact that $\|\mathbf{X}_{e,i} \mathbf{X}_{e,i}^T\| \leq p$ and $|\mathbf{X}_{e,i}^T \Delta\bar{\beta}_{e,k}| \leq \|\mathbf{X}_{e,i}\| \|\Delta\bar{\beta}_{e,k}\|$. This implies that

$$\sup_{k \geq k_{1,n} + 1} \|M_n(\bar{\beta}_{e,k}) - M(\beta_e^*)\| = O_p\left(\sqrt{p^4(\log p)/n}\right).$$

Define $\omega_k = (M_n(\bar{\beta}_{e,k}) - M(\beta_e^*)) \Delta\beta_{e,k}$. Obviously, there holds

$$\begin{aligned}\sup_{k \geq k_{1,n}^{BGD} + 1} \|\omega_k\| &\leq \left(\sup_{k \geq k_{1,n}^{BGD} + 1} \|M_n(\bar{\beta}_{e,k}) - M(\beta_e^*)\| \right) \left(\sup_{k \geq k_{1,n}^{BGD} + 1} \|\Delta\beta_{e,k}\| \right) \\ &= O_p\left(\sqrt{p^5(\log p)^2/n^2}\right),\end{aligned}$$

which is $o_p(n^{-1/2})$ according to Assumption 2.

Based on the above result, we have that for any $k \geq 1$,

$$\begin{aligned}
\Delta \beta_{e,k+k_{1,n}^{BGD}+1} &= \left(I_{p+1} - \delta M_n \left(\bar{\beta}_{e,k+k_{1,n}^{BGD}} \right) \right) \Delta \beta_{e,k+k_{1,n}^{BGD}} - \frac{\delta}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \\
&= (I_{p+1} - \delta M(\beta_e^*)) \Delta \beta_{e,k+k_{1,n}^{BGD}} - \delta \omega_{k+k_{1,n}^{BGD}} - \frac{\delta}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \\
&= (I_{p+1} - \delta M(\beta_e^*))^k \Delta \beta_{e,k_{1,n}^{BGD}+1} - \delta \sum_{j=0}^{k-1} (I_{p+1} - \delta M(\beta_e^*))^j \omega_{k+k_{1,n}^{BGD}-j} \\
&\quad - \delta \left(\sum_{j=0}^{k-1} (I_{p+1} - \delta M(\beta_e^*))^j \right) \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right).
\end{aligned}$$

For the first part on the RHS of the last equality, we have that

$$\begin{aligned}
\left\| (I_{p+1} - \delta M(\beta_e^*))^k \Delta \beta_{e,k_{1,n}^{BGD}+1} \right\| &\leq (1 - \underline{\lambda}_e \delta)^k \left\| \Delta \beta_{e,k_{1,n}^{BGD}+1} \right\| \\
&= (1 - \underline{\lambda}_e \delta)^k O_p \left(\sqrt{p(\log p)/n} \right).
\end{aligned}$$

For the second part, we have that

$$\begin{aligned}
\left\| \delta \sum_{j=0}^{k-1} (I_{p+1} - \delta M(\beta_e^*))^j \omega_{k+k_{1,n}^{BGD}-j} \right\| &\leq \delta \sum_{j=0}^{\infty} (1 - \underline{\lambda}_e \delta)^j \left\| \omega_{k+k_{1,n}^{BGD}-j} \right\| \\
&\leq \underline{\lambda}_e^{-1} \sup_{k \geq 1} \left\| \omega_{k+k_{1,n}^{BGD}} \right\| = O_p \left(\sqrt{p^5 (\log p)^2 / n^2} \right) \\
&= o_p \left(n^{-1/2} \right).
\end{aligned}$$

For the third part, we have that

$$\begin{aligned}
&\left\| \left(\sum_{j=0}^{k-1} \delta (I_{p+1} - \delta M(\beta_e^*))^j \right) \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right) - M_n^{-1}(\beta_e^*) \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right) \right\| \\
&\leq \sum_{j=k}^{\infty} \delta (1 - \underline{\lambda}_e \delta)^j \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| = (1 - \underline{\lambda}_e \delta)^k \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| \\
&= (1 - \underline{\lambda}_e \delta)^k O_p \left(\sqrt{p(\log p)/n} \right).
\end{aligned}$$

This implies that when $(1 - \underline{\lambda}_e \delta)^{k_{2,n}^{BGD}} \sqrt{p \log p} \rightarrow 0$, we have that

$$\sup_{k \geq k_{2,n}^{BGD}+1} \left\| \sqrt{n} \Delta \beta_{e,k+k_{1,n}^{BGD}} - M^{-1}(\beta_e^*) \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| = o_p(1).$$

This proves [Theorem 1\(ii\)](#)

Now we prove [Theorem 1\(iii\)](#). We first note that for any square matrices A , B , and C ,

there hold $\|AB\| \leq \bar{\sigma}(A) \|B\|$ and $\|ABC\| \leq \bar{\sigma}(A) \|BC\| \leq \bar{\sigma}(A) \bar{\sigma}(B) \|C\|$. So

$$\begin{aligned} \left\| M^{-1}(\beta_e^*) - M_n^{-1}(\hat{\beta}_e) \right\| &= \left\| M^{-1}(\beta_e^*) \left(M_n(\hat{\beta}_e) - M(\beta_e^*) \right) M_n^{-1}(\hat{\beta}_e) \right\| \\ &\leq \bar{\sigma}(M^{-1}(\beta_e^*)) \cdot \bar{\sigma}(M_n^{-1}(\hat{\beta}_e)) \cdot \left\| M_n(\hat{\beta}_e) - M(\beta_e^*) \right\|, \end{aligned}$$

due to the fact that $M_n^{-1}(\hat{\beta}_e)$ and $M_n(\hat{\beta}_e) - M(\beta_e^*)$ are both symmetric. Due to [Assumption 2\(iv\)](#), we have that $\bar{\sigma}(M^{-1}(\beta_e^*)) = \bar{\lambda}(M^{-1}(\beta_e^*)) \leq \underline{\lambda}_e^{-1}$. Since $\left\| M_n(\hat{\beta}_e) - M(\beta_e^*) \right\| = o_p(1)$ holds according to the previous proof, we have that with probability going to 1, $\bar{\sigma}(M_n^{-1}(\hat{\beta}_e)) = \bar{\lambda}(M_n^{-1}(\hat{\beta}_e)) \leq 2\underline{\lambda}_e^{-1}$. Then with probability going to 1, we have that

$$\left\| M^{-1}(\beta_e^*) - M_n^{-1}(\hat{\beta}_e) \right\| \leq 2\underline{\lambda}_e^{-2} \left\| M_n(\hat{\beta}_e) - M(\beta_e^*) \right\| = O_p\left(\sqrt{p^4(\log p)/n}\right).$$

On the other side, we have that

$$\begin{aligned} &\left\| \frac{1}{n} \sum_{i=1}^n \hat{G}_i (1 - \hat{G}_i) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T - \mathbb{E}[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T] \right\| \\ &\leq \left\| \frac{1}{n} \sum_{i=1}^n \hat{G}_i (1 - \hat{G}_i) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T - \frac{1}{n} \sum_{i=1}^n G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T \right\| \\ &\quad + \left\| \frac{1}{n} \sum_{i=1}^n G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T - \mathbb{E}[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T] \right\| \\ &\leq C\sqrt{p^3} \left\| \hat{\beta}_e - \beta_e^* \right\| + O_p\left(\sqrt{p^2(\log p)/n}\right) = O_p\left(\sqrt{p^4(\log p)/n}\right). \end{aligned}$$

Together, we have that

$$\begin{aligned} \left\| \hat{\Sigma}_1 - \Sigma_1^* \right\| &\leq \left\| M^{-1}(\beta_e^*) \mathbb{E}[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T] \left(M^{-1}(\beta_e^*) - M_n^{-1}(\hat{\beta}_e) \right) \right\| \\ &\quad + \left\| M^{-1}(\beta_e^*) \left(\frac{1}{n} \sum_{i=1}^n \hat{G}_i (1 - \hat{G}_i) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T - \mathbb{E}[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T] \right) M_n^{-1}(\hat{\beta}_e) \right\| \\ &\quad + \left\| \left(M^{-1}(\beta_e^*) - M_n^{-1}(\hat{\beta}_e) \right) \left(\frac{1}{n} \sum_{i=1}^n \hat{G}_i (1 - \hat{G}_i) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T \right) M_n^{-1}(\hat{\beta}_e) \right\| \\ &\leq \bar{\lambda}(M^{-1}(\beta_e^*)) \bar{\lambda}(\mathbb{E}[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T]) \left\| M^{-1}(\beta_e^*) - M_n^{-1}(\hat{\beta}_e) \right\| \\ &\quad + \bar{\lambda}(M^{-1}(\beta_e^*)) \bar{\lambda}(M_n^{-1}(\hat{\beta}_e)) \left\| \frac{1}{n} \sum_{i=1}^n \hat{G}_i (1 - \hat{G}_i) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T - \mathbb{E}[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T] \right\| \\ &\quad + \bar{\lambda}(M_n^{-1}(\hat{\beta}_e)) \bar{\lambda}\left(\frac{1}{n} \sum_{i=1}^n \hat{G}_i (1 - \hat{G}_i) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T\right) \left\| M^{-1}(\beta_e^*) - M_n^{-1}(\hat{\beta}_e) \right\|. \end{aligned}$$

Note that $\bar{\lambda}(M^{-1}(\beta_e^*)) \leq \underline{\lambda}_e^{-1}$, $\bar{\lambda}(\mathbb{E}[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T]) \leq \frac{1}{4} \bar{\lambda}(\mathbb{E}[\mathbf{X}_{e,i} \mathbf{X}_{e,i}^T]) \leq C$, $\bar{\lambda}(M_n^{-1}(\hat{\beta}_e)) \leq$

$2\lambda_e^{-1}$ with probability going to 1, and $\bar{\lambda} \left(\frac{1}{n} \sum_{i=1}^n \widehat{G}_i \left(1 - \widehat{G}_i \right) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T \right) \leq C$ with probability going to 1, we have that

$$\begin{aligned} \left\| \widehat{\Sigma}_1 - \Sigma_1^* \right\| &\leq C \left\| M^{-1}(\boldsymbol{\beta}_e^*) - M_n^{-1}(\widehat{\boldsymbol{\beta}}_e) \right\| \\ &\quad + C \left\| \frac{1}{n} \sum_{i=1}^n \widehat{G}_i \left(1 - \widehat{G}_i \right) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T - \mathbb{E} \left[G_i^* (1 - G_i^*) \mathbf{X}_{e,i} \mathbf{X}_{e,i}^T \right] \right\| \\ &= O_p \left(\sqrt{p^4 (\log p) / n} \right) = o_p(1), \end{aligned}$$

which validates the result.

To prove (iv), we only need to show that $\widehat{\sigma}^2(\rho) - \sigma^2(\rho) = o_p(1)$. Note that

$$\left| \widehat{\sigma}_n^2(\rho) - \sigma^2(\rho) \right| = \left| \rho^T \left(\widehat{\Sigma}_1 - \Sigma_1^* \right) \rho \right| \leq \|\rho\| \left\| \left(\widehat{\Sigma}_1 - \Sigma_1^* \right) \rho \right\| \leq \|\rho\|^2 \left\| \widehat{\Sigma}_1 - \Sigma_1^* \right\| \rightarrow_p 0$$

given that $\|\rho\| < \infty$ for all n , which validates the result. \square

Proof of Theorem 2

Proof. We first show Theorem 2(i). Note that from the proof in Theorem 1, we know that with probability going to 1, we have that

$$\begin{aligned} \left\| \Delta \boldsymbol{\beta}_{e,k+1} \right\| &\leq \sup_{\boldsymbol{\beta}_e \in \mathcal{B}} \bar{\lambda} (I_{p+1} - \delta_k M_n(\boldsymbol{\beta}_e)) \left\| \Delta \boldsymbol{\beta}_{e,k} \right\| + \delta_k \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| \\ &\leq (1 - \underline{\lambda}_e \delta_k / 2) \left\| \Delta \boldsymbol{\beta}_{e,k} \right\| + \delta_k \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| \leq \dots \\ &\leq \left(\prod_{j=1}^k (1 - \underline{\lambda}_e \delta_j / 2) \right) \left\| \Delta \boldsymbol{\beta}_{e,1} \right\| + \left\{ \sum_{j=0}^{k-1} \delta_{k-j} \left(\prod_{l=0}^{j-1} (1 - \underline{\lambda}_e \delta_{k-l} / 2) \right) \right\} \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\|, \end{aligned} \tag{11}$$

where $\prod_{l=0}^{j-1} (1 - \underline{\lambda}_e \delta_{k-l} / 2) = 1$ if $j = 0$.

For the first term on the RHS of (11), since $e^x \geq 1 + x$ for all x , we have $1 - \underline{\lambda}_e \delta_j / 2 \leq \exp(-\underline{\lambda}_e \delta_j / 2)$ for all j . Define $S_0 = 0$ and $S_j = \sum_{l=1}^j \delta_l$ for $j \geq 1$, we have that

$$\left(\prod_{j=1}^k (1 - \underline{\lambda}_e \delta_j / 2) \right) \left\| \Delta \boldsymbol{\beta}_{e,1} \right\| \leq \exp \left(-\frac{\underline{\lambda}_e}{2} \sum_{j=1}^k \delta_j \right) \left\| \Delta \boldsymbol{\beta}_{e,1} \right\| = \exp \left(-\frac{\underline{\lambda}_e S_k}{2} \right) \left\| \Delta \boldsymbol{\beta}_{e,1} \right\|.$$

Next we show that $\sum_{j=0}^{k-1} \delta_{k-j} \left(\prod_{l=0}^{j-1} (1 - \underline{\lambda}_e \delta_{k-l} / 2) \right)$ is upper bounded by $\exp(\underline{\lambda}_e \delta_{k+1} / 2)$

up to some constant scale that is independent of k . Since $\limsup_k \delta_{k-1}/\delta_k < \infty$, we have that

$$\begin{aligned}
& \sum_{j=0}^{k-1} \delta_{k-j} \left(\prod_{l=0}^{j-1} (1 - \underline{\lambda}_e \delta_{k-l}/2) \right) \leq \sum_{j=0}^{k-1} \delta_{k-j} \exp \left(-\frac{\underline{\lambda}_e}{2} \sum_{l=0}^{j-1} \delta_{k-l} \right) \\
& \leq C \sum_{j=0}^{k-1} \delta_{k-j+1} \exp \left(-\frac{\underline{\lambda}_e (S_k - S_{k-j})}{2} \right) \\
& = C \exp \left(-\frac{\underline{\lambda}_e S_k}{2} \right) \sum_{j=0}^{k-1} (S_{k-j+1} - S_{k-j}) \exp \left(\frac{\underline{\lambda}_e S_{k-j}}{2} \right) \\
& \leq 2C \underline{\lambda}_e^{-1} \exp \left(-\frac{\underline{\lambda}_e S_k}{2} \right) \sum_{j=0}^{k-1} \left\{ \exp \left(\frac{\underline{\lambda}_e S_{k-j+1}}{2} \right) - \exp \left(\frac{\underline{\lambda}_e S_{k-j}}{2} \right) \right\} \leq C \exp \left(\frac{\underline{\lambda}_e \delta_{k+1}}{2} \right).
\end{aligned}$$

Then we have that

$$\|\Delta \boldsymbol{\beta}_{e,k+1}\| = O_p \left(\exp \left(-\frac{\underline{\lambda}_e S_k}{2} \right) \|\Delta \boldsymbol{\beta}_{e,1}\| \right) + O_p \left(\exp \left(\frac{\underline{\lambda}_e \delta_{k+1}}{2} \right) \sqrt{p(\log p)/n} \right).$$

When $k \geq \tilde{k}_{1,n}^{BGD} + 1$, we have that

$$\exp \left(-\frac{\underline{\lambda}_e S_k}{2} \right) \|\Delta \boldsymbol{\beta}_{e,1}\| \leq \sqrt{p(\log p)/n},$$

and

$$\exp \left(\frac{\underline{\lambda}_e \delta_{k+1}}{2} \right) \leq e,$$

so $\|\Delta \boldsymbol{\beta}_{e,k+1}\| = O_p \left(\sqrt{p(\log p)/n} \right)$. This validates [Theorem 2\(i\)](#).

For [Theorem 2\(ii\)](#), we know that for $k \geq \tilde{k}_{1,n}^{BGD} + 1$, $\|\Delta \boldsymbol{\beta}_{e,k}\| = O_p \left(\sqrt{p(\log p)/n} \right)$ holds, so we have that

$$\begin{aligned}
\Delta \boldsymbol{\beta}_{e,k+1} &= (I_{p+1} - \delta_k M_n(\bar{\boldsymbol{\beta}}_{e,k})) \Delta \boldsymbol{\beta}_{e,k} - \frac{\delta_k}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \\
&= (I_{p+1} - \delta_k M(\boldsymbol{\beta}_e^*)) \Delta \boldsymbol{\beta}_{e,k} - \delta_k (M_n(\bar{\boldsymbol{\beta}}_{e,k}) - M(\boldsymbol{\beta}_e^*)) \Delta \boldsymbol{\beta}_{e,k} - \frac{\delta_k}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i},
\end{aligned}$$

where $\bar{\boldsymbol{\beta}}_{e,k}$ lies between $\boldsymbol{\beta}_{e,k}$ and $\boldsymbol{\beta}_e^*$ and is element-wise. Following the proof of [Theorem 1](#), we can easily show that

$$\sup_{k \geq \tilde{k}_{1,n}^{BGD} + 1} \|M_n(\bar{\boldsymbol{\beta}}_{e,k}) - M(\boldsymbol{\beta}_e^*)\| = O_p \left(\sqrt{p^4(\log p)/n} \right).$$

Recall that $\omega_k = (M_n(\bar{\beta}_{e,k}) - M(\beta_e^*)) \Delta \beta_{e,k}$, so

$$\sup_{k \geq \tilde{k}_{1,n}^{BGD} + 1} \|\omega_k\| = O_p \left(\sqrt{p^5 (\log p)^2 / n^2} \right) = o_p(n^{-1/2}).$$

We have that

$$\begin{aligned} \Delta \beta_{e,k+\tilde{k}_{1,n}^{BGD}+1} &= \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k} M(\beta_e^*) \right) \Delta \beta_{e,k+\tilde{k}_{1,n}^{BGD}} - \delta_{\tilde{k}_{1,n}^{BGD}+k} \omega_{\tilde{k}_{1,n}^{BGD}+k} - \delta_{\tilde{k}_{1,n}^{BGD}+k} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \\ &= \prod_{j=0}^{k-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-j} M(\beta_e^*) \right) \Delta \beta_{e,\tilde{k}_{1,n}^{BGD}+1} \\ &\quad - \sum_{j=0}^{k-1} \left\{ \delta_{\tilde{k}_{1,n}^{BGD}+k-j} \prod_{l=0}^{j-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-l} M(\beta_e^*) \right) \right\} \omega_{\tilde{k}_{1,n}^{BGD}+k-j} \\ &\quad - \sum_{j=0}^{k-1} \left\{ \delta_{\tilde{k}_{1,n}^{BGD}+k-j} \prod_{l=0}^{j-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-l} M(\beta_e^*) \right) \right\} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i}, \end{aligned}$$

where $\prod_{l=0}^{j-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-l} M(\beta_e^*) \right) = 1$ if $j = 0$. For the first part, define $S_{\tilde{k}_{1,n}^{BGD},k} = \sum_{j=\tilde{k}_{1,n}^{BGD}+1}^{\tilde{k}_{1,n}^{BGD}+k} \delta_j$, we have that

$$\begin{aligned} \left\| \prod_{j=0}^{k-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-j} M(\beta_e^*) \right) \Delta \beta_{e,\tilde{k}_{1,n}^{BGD}+1} \right\| &\leq \prod_{j=0}^{k-1} \left(1 - \lambda_e \delta_{\tilde{k}_{1,n}^{BGD}+k-j} / 2 \right) \left\| \Delta \beta_{e,\tilde{k}_{1,n}^{BGD}+1} \right\| \\ &\leq \exp \left(-\lambda_e S_{\tilde{k}_{1,n}^{BGD},k} / 2 \right) \left\| \Delta \beta_{e,\tilde{k}_{1,n}^{BGD}+1} \right\| \\ &= O_p \left(\exp \left(-\lambda_e S_{\tilde{k}_{1,n}^{BGD},k} / 2 \right) \sqrt{p (\log p) / n} \right). \end{aligned}$$

For the second term, we have that

$$\begin{aligned} &\left\| \sum_{j=0}^{k-1} \left\{ \delta_{\tilde{k}_{1,n}^{BGD}+k-j} \prod_{l=0}^{j-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-l} M(\beta_e^*) \right) \right\} \omega_{\tilde{k}_{1,n}^{BGD}+k-j} \right\| \\ &\leq \left\{ \sum_{j=0}^{k-1} \delta_{\tilde{k}_{1,n}^{BGD}+k-j} \prod_{l=0}^{j-1} \left(1 - \lambda_e \delta_{\tilde{k}_{1,n}^{BGD}+k-l} / 2 \right) \right\} \left\{ \sup_{k \geq 1} \left\| \omega_{\tilde{k}_{1,n}^{BGD}+k} \right\| \right\} \\ &\leq \exp \left(-\lambda_e S_{\tilde{k}_{1,n}^{BGD},k} / 2 \right) \left\{ \sum_{j=0}^{k-1} \delta_{\tilde{k}_{1,n}^{BGD}+k-j} \exp \left(\lambda_e S_{\tilde{k}_{1,n}^{BGD}+k-j} / 2 \right) \right\} \left\{ \sup_{k \geq 1} \left\| \omega_{\tilde{k}_{1,n}^{BGD}+k} \right\| \right\} \\ &= O_p \left(\sqrt{p^5 (\log p)^2 / n^2} \right) \end{aligned}$$

according to the proof of [Theorem 2\(i\)](#). Now we look at the last term. Note that

$$\begin{aligned}\mathcal{M}_{k,n} &\equiv: \sum_{j=0}^{k-1} \left\{ \delta_{\tilde{k}_{1,n}^{BGD}+k-j} \prod_{l=0}^{j-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-l} M(\beta_e^*) \right) \right\} \\ &= \delta_{\tilde{k}_{1,n}^{BGD}+k} I_{p+1} + \delta_{\tilde{k}_{1,n}^{BGD}+k-1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k} M(\beta_e^*) \right) + \dots \\ &\quad + \delta_{\tilde{k}_{1,n}^{BGD}+1} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k} M(\beta_e^*) \right) \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k-1} M(\beta_e^*) \right) \dots \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+2} M(\beta_e^*) \right),\end{aligned}$$

so

$$\mathcal{M}_{k+1,n} = \delta_{\tilde{k}_{1,n}^{BGD}+k+1} I_{p+1} + \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k} M(\beta_e^*) \right) \mathcal{M}_{k,n}.$$

Note that

$$\begin{aligned}\mathcal{M}_{k+1,n} - M^{-1}(\beta_e^*) &= \mathcal{M}_{k,n} - M^{-1}(\beta_e^*) + \delta_{\tilde{k}_{1,n}^{BGD}+k+1} M(\beta_e^*) (M^{-1}(\beta_e^*) - \mathcal{M}_{k,n}) \\ &= \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k} M(\beta_e^*) \right) (\mathcal{M}_{k,n} - M^{-1}(\beta_e^*)),\end{aligned}$$

so

$$\begin{aligned}\|\mathcal{M}_{k+1,n} - M^{-1}(\beta_e^*)\| &\leq \bar{\lambda} \left(I_{p+1} - \delta_{\tilde{k}_{1,n}^{BGD}+k} M(\beta_e^*) \right) \|\mathcal{M}_{k,n} - M^{-1}(\beta_e^*)\| \\ &\leq \left(1 - \delta_{\tilde{k}_{1,n}^{BGD}+k} \underline{\lambda}_e \right) \|\mathcal{M}_{k,n} - M^{-1}(\beta_e^*)\| \\ &\leq \exp \left(-\underline{\lambda}_e S_{\tilde{k}_{1,n}^{BGD},k} \right) \|\mathcal{M}_{1,n} - M^{-1}(\beta_e^*)\|.\end{aligned}$$

Then

$$\begin{aligned}&\sum_{j=0}^{k-1} \left\{ \delta_{\tilde{k}_{1,n}^{BGD}+k-1-j} \prod_{l=0}^{j-1} \left(I - \delta_{\tilde{k}_{1,n}^{BGD}+k-1-l} M(\beta_e^*) \right) \right\} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \\ &= M^{-1}(\beta_e^*) \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} + O_p \left(\exp \left(-\underline{\lambda}_e S_{\tilde{k}_{1,n}^{BGD},k} \right) \sqrt{p(\log p)/n} \right).\end{aligned}$$

So we have

$$\begin{aligned}\left\| \sqrt{n} \Delta \beta_{e,k+\tilde{k}_{1,n}^{BGD}} - M^{-1}(\beta_e^*) \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{e,i} \right\| &= O_p \left(\exp \left(-\underline{\lambda}_e S_{\tilde{k}_{1,n}^{BGD},k}/2 \right) \sqrt{p(\log p)/n} \right) \\ &\quad + O_p \left(\sqrt{p^5 (\log p)^2 / n^2} \right) \\ &\quad + O_p \left(\exp \left(-S_{\tilde{k}_{1,n}^{BGD},k} \right) \sqrt{p(\log p)/n} \right).\end{aligned}$$

According to the definition of $\tilde{k}_{2,n}^{BGD}$, we have that for $k \geq \tilde{k}_{2,n}^{BGD}$, there holds $S_{\tilde{k}_{1,n}^{BGD},k}/\log p \rightarrow \infty$, this proves [Theorem 2\(ii\)](#).

The proof of [Theorem 2\(iii\)](#) and [Theorem 2\(iv\)](#) is the same as that in the proof of [Theorem 1](#), so is left out. \square

Proof of [Theorem 3](#)

Proof. Define

$$\eta_{1,n}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \widehat{G}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta}) \mathbf{X}_i - \mathbb{E}[L(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) \mathbf{X}_i],$$

$$\eta_{2,n} = \left(\frac{1}{n} \sum_{i=1}^n G(z_i^*) \mathbf{X}_i - \mathbb{E}[G(z_i^*) \mathbf{X}_i] \right) + \frac{1}{n} \sum_{i=1}^n \varepsilon_i \cdot \mathbf{X}_i.$$

Note that when $\boldsymbol{\beta}^* \in \mathcal{B}$ and $\boldsymbol{\beta}_k \in \mathcal{B}$, we have that $\boldsymbol{\beta}^* + t\Delta\boldsymbol{\beta}_k \in \mathcal{B}$ for all $0 \leq t \leq 1$, so

$$\begin{aligned} \|\Delta\boldsymbol{\beta}_{k+1}\| &\leq \left\| \int_0^1 (I_p - \delta\Lambda(\boldsymbol{\beta}^* + t\Delta\boldsymbol{\beta}_k)) dt \Delta\boldsymbol{\beta}_k \right\| + \delta \|\eta_{1,n}(\boldsymbol{\beta}_k)\| + \delta \|\eta_{2,n}\| \\ &\leq \left\{ \sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\sigma}(I_p - \delta\Lambda(\boldsymbol{\beta})) \right\} \|\Delta\boldsymbol{\beta}_k\| + \delta \|\eta_{1,n}(\boldsymbol{\beta}_k)\| + \delta \|\eta_{2,n}\|. \end{aligned}$$

Note that for any $1 \leq s, t \leq p$,

$$\begin{aligned} |(\Lambda(\boldsymbol{\beta}))_{s,t}| &= \left| \mathbb{E} \left[\int_{\mathcal{X}} (X_{s,i} X_{t,i} - X_{s,i} X_t) W(\mathbf{X}_{e,i}, \mathbf{X}_e, \boldsymbol{\beta}) d\mathbf{X} \right] \right| \\ &\leq 2 \|G'\|_{\infty} \mathbb{E} \left[\int_{\mathcal{X}} f_{\mathbf{X}|z}(\mathbf{X} | z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) d\mathbf{X} \right] = 2 \|G'\|_{\infty}, \end{aligned}$$

so each element of $\Lambda^T(\boldsymbol{\beta}) \Lambda(\boldsymbol{\beta})$ is bounded by $2p\|G'\|_{\infty}$, and we have that

$$\begin{aligned} \sup_{\boldsymbol{\beta} \in \mathcal{B}} \left| \bar{\sigma}^2(I_p - \delta\Lambda(\boldsymbol{\beta})) - \bar{\lambda}(I_p - \delta(\Lambda(\boldsymbol{\beta}) + \Lambda^T(\boldsymbol{\beta}))) \right| \\ \leq \sup_{\boldsymbol{\beta} \in \mathcal{B}} \delta^2 \|\Lambda^T(\boldsymbol{\beta}) \Lambda(\boldsymbol{\beta})\| \leq 2 \|G'\|_{\infty} p^2 \delta^2. \end{aligned}$$

Then according to [Assumption 5](#), we have that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\sigma}^2(I_p - \delta\Lambda(\boldsymbol{\beta})) \leq 1 - \delta\lambda_A + 2 \|G'\|_{\infty} p^2 \delta^2.$$

When $\delta < \min\{1/(2\lambda_A), 1/(4\|G'\|_{\infty} p^2)\}$, we have that

$$0 \leq 1 - \delta\lambda_A + 2 \|G'\|_{\infty} p^2 \delta^2 \leq 1 - \delta\lambda_A/2 < 1.$$

So

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\sigma}(I_p - \delta\Lambda(\boldsymbol{\beta})) \leq \sqrt{1 - \delta\lambda_A/2} \leq 1 - \delta\lambda_A/4,$$

and

$$\begin{aligned}
\|\Delta\beta_{k+1}\| &\leq (1 - \delta\lambda_A/4) \|\Delta\beta_k\| + \delta \|\eta_{1,n}(\beta_k)\| + \delta \|\eta_{2,n}\| \\
&\leq \cdots \leq (1 - \delta\lambda_A/4)^k \|\Delta\beta_1\| + \delta \cdot \sum_{j=0}^{k-1} (1 - \delta\lambda_A/4)^j (\|\eta_{1,n}(\beta_j)\| + \|\eta_{2,n}\|) \\
&\leq (1 - \delta\lambda_A/4)^k \|\Delta\beta_1\| + \delta \cdot \sum_{j=0}^{\infty} (1 - \delta\lambda_A/4)^j \left(\sup_{\beta \in \mathcal{B}} \|\eta_{1,n}(\beta)\| + \|\eta_{2,n}\| \right) \\
&= (1 - \delta\lambda_A/4)^k \|\Delta\beta_1\| + 4\lambda_A^{-1} \left(\sup_{\beta \in \mathcal{B}} \|\eta_{1,n}(\beta)\| + \|\eta_{2,n}\| \right).
\end{aligned}$$

Note that

$$\sup_{\beta \in \mathcal{B}} \|\eta_{1,n}(\beta)\| = p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n)$$

according to [Lemma 1](#), and

$$\|\eta_{2,n}\| = O_p\left(\sqrt{p(\log p)/n}\right) = o_p\left(p^{\frac{5p+1}{2(p+1)}} (\psi(n, p, h_n))^{\frac{1}{3p+3}}\right)$$

under any choices of $h_n \rightarrow 0$. This implies that when

$$(1 - \delta\lambda_A/4)^k \|\Delta\beta_1\| \leq p^{\frac{5p+1}{2(p+1)}} (\psi(n, p, h_n))^{\frac{1}{p+1}},$$

or equivalently,

$$k \geq k_{1,n}^{KBGD} = \frac{\log(\|\Delta\beta_1\|) - \frac{5p+1}{2(p+1)} \log p - \frac{1}{p+1} \log \psi(n, p, h_n)}{-\log(1 - \delta\lambda_A/4)},$$

we have that $\sup_{k \geq k_{1,n}^{KBGD}+1} \|\Delta\beta_k\| = O_p\left(p^{\frac{5p+1}{2(p+1)}} \psi^{\frac{1}{p+1}}(n, p, h_n)\right)$. □

Proof of [Theorem 4](#)

Proof. We first note that

$$\left\| \int_{\mathcal{X}} V(\mathbf{X}_{e,i}, \mathbf{X}_e, \beta) d\mathbf{X} \right\| \leq 2p \|G'\|_{\infty} \int_{\mathcal{X}} f_{\mathbf{X}|z}(\mathbf{X}|z(\mathbf{X}_{e,i}, \beta), \beta) d\mathbf{X} = 2p \|G'\|_{\infty},$$

for all $\mathbf{X}_{e,i}$, so

$$\sup_{\beta \in \mathcal{B}} \|\Lambda_{\phi}(\beta) - \Lambda(\beta)\| \leq 2p \|G'\|_{\infty} \mathbb{E}\left(1 - I_i^{\phi}\right) \leq 2\zeta p^2 \|G'\|_{\infty} \phi,$$

where the last inequality comes from the fact that $m(\mathcal{X}_e^{\phi}) = 1 - (1 - \phi)^p \leq p\phi$. So

$$\sup_{\beta \in \mathcal{B}} \|\Lambda_{\phi}(\beta) - \Lambda(\beta)\| \leq \delta\lambda_A/8 \tag{12}$$

holds under the choice of ϕ .

Based on (12), the following proof is similar to the proof of Theorem 3. Define

$$\begin{aligned}\eta_{1,n}^\phi(\boldsymbol{\beta}) &= \frac{1}{n} \sum_{i=1}^n \widehat{G}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta}) \mathbf{X}_i^\phi - \mathbb{E} \left[L(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}), \boldsymbol{\beta}) \mathbf{X}_i^\phi \right], \\ \eta_{2,n}^\phi &= \frac{1}{n} \sum_{i=1}^n G(z_i^*) \mathbf{X}_i^\phi - \mathbb{E} \left[G(z_i^*) \mathbf{X}_i^\phi \right] + \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{X}_i^\phi.\end{aligned}$$

We have that

$$\begin{aligned}\Delta \boldsymbol{\beta}_{k+1} &= \Delta \boldsymbol{\beta}_k - \frac{\delta}{n} \sum_{i=1}^n \left(\widehat{G}(z_{i,k} | \boldsymbol{\beta}_k) - Y_i \right) \mathbf{X}_i^\phi \\ &= \Delta \boldsymbol{\beta}_k - \delta \mathbb{E} \left[(L(z_{i,k}, \boldsymbol{\beta}_k) - G(Z_i^*)) \mathbf{X}_i^\phi \right] + \delta \left(\eta_{1,n}^\phi(\boldsymbol{\beta}_k) + \eta_{2,n}^\phi \right) \\ &= \int_0^1 \{I_p - \delta \Lambda_\phi(\boldsymbol{\beta}^* + t \Delta \boldsymbol{\beta}_k)\} \Delta \boldsymbol{\beta}_k dt + \delta \left(\eta_{1,n}^\phi(\boldsymbol{\beta}_k) + \eta_{2,n}^\phi \right),\end{aligned}$$

so

$$\|\boldsymbol{\beta}_{k+1}\| \leq \sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\sigma}(I_p - \delta \Lambda_\phi(\boldsymbol{\beta})) \|\boldsymbol{\beta}_k\| + \delta \left(\sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\eta_{1,n}^\phi(\boldsymbol{\beta})\| + \|\eta_{2,n}^\phi\| \right).$$

Obviously, since p is fixed, we have that $\|\eta_{2,n}^\phi\| = O_p(n^{-1/2})$. Due to trimming, we also have that $\sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\eta_{1,n}^\phi(\boldsymbol{\beta})\| = O_p(\psi(n, p, h_n))$. Note that (12) holds, so we have that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \|\{I_p - \delta \Lambda_\phi(\boldsymbol{\beta})\} - \{I_p - \delta \Lambda(\boldsymbol{\beta})\}\| \leq \delta \underline{\lambda}_\Lambda / 8.$$

According to the proof of Theorem 3, there holds $\sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\sigma}(I_p - \delta \Lambda(\boldsymbol{\beta})) \leq 1 - \delta \underline{\lambda}_\Lambda / 4$ under the choice of δ , so we have that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \bar{\sigma}(I_p - \delta \Lambda_\phi(\boldsymbol{\beta})) \leq 1 - \delta \underline{\lambda}_\Lambda / 8.$$

Then based on the proof of Theorem 3, it remains to note that

$$\sup_{\boldsymbol{\beta} \in \mathcal{B}} \left(\|\eta_{1,n}^\phi(\boldsymbol{\beta})\| + \|\eta_{2,n}^\phi\| \right) = O_p(\psi(n, p, h_n))$$

holds under any fixed trimming parameter ϕ . □

Proof of Theorem 5

Proof. Note that under the choice of δ and ϕ , $\sup_{k \geq \tilde{k}_{1,n}^{KBGD}+1} \|\beta_k - \beta^*\| = O_p(\psi(n, p, h_n))$ according to Theorem 4. According to (14), we have that

$$\begin{aligned} & \left\| \Delta \beta_{k+\tilde{k}_{1,n}^{KBGD}+1} \right\| \\ & \leq \sup_{k \geq \tilde{k}_{1,n}^{KBGD}+1, t \in [0,1]} \bar{\sigma} \left\{ I_p - \frac{\delta}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\phi \frac{\partial \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta)}{\partial \beta^T} \right]_{\beta=\beta^*+t\Delta\beta_k} \right\} \|\Delta\beta_k\| + \delta \|\xi_n^\phi\|. \end{aligned}$$

According to Lemma 2, we have that

$$\begin{aligned} & \sup_{k \geq \tilde{k}_{1,n}^{KBGD}, t \in [0,1]} \left\| \left\{ I_p - \frac{\delta}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\phi \frac{\partial \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta)}{\partial \beta^T} \right]_{\beta=\beta^*+t\Delta\beta_k} \right\} \right. \\ & \quad \left. - \{I_p - \delta \Lambda_\phi(\beta^* + t\Delta\beta_k)\} \right\| = \delta O_p \left(h_n^{-2} \sqrt{(\log(nh_n^{-1}))/n} + h_n^3 \right), \end{aligned} \quad (13)$$

due to the fact that

$$\sup_{k \geq \tilde{k}_{1,n}^{KBGD}+1} \|\Delta\beta_k\| = O_p(\psi_1(n, p, h_n)) = o_p \left(h_n^{-2} \sqrt{(\log(nh_n^{-1}))/n} + h_n^3 \right),$$

when p is fixed and $h_n \rightarrow 0$.

When $nh_n^6 \rightarrow 0$ and $h_n^4 n / (\log n)^2 \rightarrow \infty$, we have that $h_n^{-2} \sqrt{(\log(nh_n^{-1}))/n} + h_n^3 \rightarrow 0$. So we have that (13) is smaller than $\delta \underline{\lambda}_\Lambda / 16$ with probability going to 1. According to the choice of ϕ and δ , we have that $\sup_{\beta \in \mathcal{B}} \bar{\sigma}(I_p - \delta \Lambda_\phi(\beta)) \leq 1 - \delta \underline{\lambda}_\Lambda / 8$ according to the proof of Theorem 4. So as n increases, with probability going to 1, there holds

$$\sup_{k \geq \tilde{k}_{1,n}^{KBGD}+1, t \in [0,1]} \bar{\sigma} \left(I_p - \frac{\delta}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\phi \frac{\partial \widehat{G}(z(\mathbf{X}_{e,i}, \beta) | \beta)}{\partial \beta^T} \right]_{\beta=\beta^*+t\Delta\beta_k} \right) \leq 1 - \delta \underline{\lambda}_\Lambda / 16,$$

Then as n increases, with probability going to 1 there holds

$$\begin{aligned} \left\| \Delta \beta_{k+\tilde{k}_{1,n}^{KBGD}+1} \right\| & \leq (1 - \delta \underline{\lambda}_\Lambda / 16) \left\| \Delta \beta_{k+\tilde{k}_{1,n}^{KBGD}} \right\| + \delta \|\xi_n^\phi\| \\ & \leq \dots \leq (1 - \delta \underline{\lambda}_\Lambda / 16)^k \left\| \Delta \beta_{\tilde{k}_{1,n}^{KBGD}+1} \right\| + 16 \underline{\lambda}_\Lambda^{-1} \|\xi_n^\phi\|. \end{aligned}$$

According to Lemma 3, $\|\xi_n^\phi\| = O_p(n^{-1/2})$. Also note that $\left\| \Delta \beta_{\tilde{k}_{1,n}^{KBGD}+1} \right\| = O_p(\psi(n, p, h_n))$, then if we choose $k_{2,n}^{KBGD}$ such that $(1 - \delta \underline{\lambda}_\Lambda / 16)^{k_{2,n}^{KBGD}-1} \leq n^{-1/2} \psi^{-1}(n, p, h_n)$, or equivalently,

$$k_{2,n}^{KBGD} \geq -\frac{\log(n^{1/2}) + \log(\psi(n, p, h_n))}{\log(1 - \delta \underline{\lambda}_\Lambda / 16)} + 1,$$

we have that $\sup_{k \geq k_{2,n}^{KBGD}+1} \left\| \Delta \beta_{k+\tilde{k}_{1,n}^{KBGD}} \right\| = O_p(n^{-1/2})$. This proves (i).

To prove (ii), we consider the following decomposition,

$$\Delta \boldsymbol{\beta}_{k+1} = (I_p - \delta \Lambda_\phi(\boldsymbol{\beta}^*)) \Delta \boldsymbol{\beta}_k + \delta \bar{\omega}_1(\boldsymbol{\beta}_k) + \delta \bar{\omega}_2(\boldsymbol{\beta}_k) - \delta \boldsymbol{\xi}_n^\phi,$$

where

$$\bar{\omega}_1(\boldsymbol{\beta}_k) = \int_0^1 \left\{ \Lambda_\phi(\boldsymbol{\beta}^* + t \Delta \boldsymbol{\beta}_k) - \frac{1}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\phi \frac{\partial \widehat{G}(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \right]_{\boldsymbol{\beta}=\boldsymbol{\beta}^*+t\Delta\boldsymbol{\beta}_k} \right\} dt \Delta \boldsymbol{\beta}_k,$$

and

$$\bar{\omega}_2(\boldsymbol{\beta}_k) = \int_0^1 \{ \Lambda_\phi(\boldsymbol{\beta}^*) - \Lambda_\phi(\boldsymbol{\beta}^* + t \Delta \boldsymbol{\beta}_k) \} dt \Delta \boldsymbol{\beta}_k.$$

Obviously, according to [Lemma 2](#),

$$\sup_{k \geq \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} \|\bar{\omega}_1(\boldsymbol{\beta}_k)\| = O_p \left(h_n^{-2} \sqrt{(\log(nh_n^{-1}))/n} + h_n^3 \right) O_p \left(n^{-\frac{1}{2}} \right) = o_p \left(n^{-\frac{1}{2}} \right).$$

We also note that each element of matrix $I_i^\phi \cdot \int_{\mathcal{X}} V(\mathbf{X}_{e,i}, \mathbf{X}_e, \boldsymbol{\beta}) d\mathbf{X}$ has bounded derivative with respect to $\boldsymbol{\beta}$ for any $\mathbf{X}_{e,i}$. This is because, if $\mathbf{X}_{e,i} \notin \mathcal{X}_e^\phi$, $I_i^\phi = 0$ so each element will be zero and the results hold; if $\mathbf{X}_{e,i} \in \mathcal{X}_e^\phi$, then $f_z(z(\mathbf{X}_{e,i}, \boldsymbol{\beta}) | \boldsymbol{\beta}) > 0$, so $\int_{\mathcal{X}} \|\partial W(\mathbf{X}_{e,i}, \mathbf{X}_e, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}\| d\mathbf{X}$ is bounded according to [Lemma A.2\(x\)](#). This implies that

$$\sup_{k \geq \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} \|\bar{\omega}_2(\boldsymbol{\beta}_k)\| \leq C \|\Delta \boldsymbol{\beta}_k\|^2 = o_p \left(n^{-\frac{1}{2}} \right).$$

Then

$$\begin{aligned} & \Delta \boldsymbol{\beta}_{k + \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} \\ &= (I_p - \delta \Lambda_\phi(\boldsymbol{\beta}^*))^k \Delta \boldsymbol{\beta}_{\tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} + \delta \sum_{j=1}^k (I_p - \delta \Lambda_\phi(\boldsymbol{\beta}^*))^{k-j} \bar{\omega}_1 \left(\boldsymbol{\beta}_{\tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + j} \right) \\ &+ \sum_{j=1}^k (I_p - \delta \Lambda_\phi(\boldsymbol{\beta}^*))^{k-j} \bar{\omega}_2 \left(\boldsymbol{\beta}_{\tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + j} \right) - \delta \sum_{j=1}^k (I_p - \delta \Lambda_\phi(\boldsymbol{\beta}^*))^{k-j} \boldsymbol{\xi}_n^\phi. \end{aligned}$$

Note that $\sup_{\beta \in \mathcal{B}} \bar{\sigma}(I_p - \delta \Lambda_\phi(\beta)) \leq 1 - \delta \underline{\lambda}_\Lambda / 8$, so

$$\begin{aligned}
& \left\| (I_p - \delta \Lambda_\phi(\beta^*))^k \Delta \beta_{\tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} \right\| \leq (1 - \delta \underline{\lambda}_\Lambda / 8)^k \left\| \Delta \beta_{\tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} \right\|, \\
& \delta \left\| \sum_{j=1}^k (I_p - \delta \Lambda_\phi(\beta^*))^{k-j} \omega_1 \left(\beta_{\tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + j} \right) \right\| \leq \delta \sum_{j=0}^{\infty} (1 - \delta \underline{\lambda}_\Lambda / 8)^j \sup_{k \geq \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + 1} \|\bar{\omega}_1(\beta_k)\| \\
& \quad = o_p(n^{-1/2}), \\
& \delta \left\| \sum_{j=1}^k (I_p - \delta \Lambda_\phi(\beta^*))^{k-j} \omega_2 \left(\beta_{\tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD} + j} \right) \right\| \leq \delta \sum_{j=0}^{\infty} (1 - \delta \underline{\lambda}_\Lambda / 8)^j \sup_{k \geq \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD}} \|\bar{\omega}_2(\beta_k)\| \\
& \quad = o_p(n^{-1/2}), \\
& \left\| \Lambda_\phi^{-1}(\beta^*) \boldsymbol{\xi}_n^\phi - \delta \sum_{j=1}^k (I_p - \delta \Lambda_\phi(\beta^*))^{k-j} \boldsymbol{\xi}_n^\phi \right\| \leq 8 \lambda_\Lambda^{-1} (1 - \delta \underline{\lambda}_\Lambda / 8)^{k+1} \|\boldsymbol{\xi}_n^\phi\|.
\end{aligned}$$

As $k \rightarrow \infty$, we have that $\lambda_\Lambda^{-1} (1 - \delta \underline{\lambda}_\Lambda / 8)^{k+1} \|\boldsymbol{\xi}_n^\phi\| = o_p(n^{-1/2})$, so

$$\Delta \beta_{k + \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD}} = \Lambda_\phi^{-1}(\beta^*) \boldsymbol{\xi}_n^\phi + o_p(n^{-1/2}).$$

According to [Lemma 3](#), we have that $\sqrt{n} \boldsymbol{\xi}_n^\phi \rightarrow N(0, \Sigma_\xi^\phi)$, so we have that

$$\sqrt{n} \Delta \beta_{k + \tilde{k}_{1,n}^{KBGD} + k_{2,n}^{KBGD}} = \Lambda_\phi^{-1}(\beta^*) \sqrt{n} \boldsymbol{\xi}_n^\phi + o_p(1) \rightarrow_d N(0, \Lambda_\phi^{-1}(\beta^*) \Sigma_\xi^\phi (\Lambda_\phi^{-1}(\beta^*))^\top).$$

□

Proof of [Theorem 6](#)

Proof. We only need to show that $\left\| \hat{\Lambda}_{\phi,n}^{-1}(\hat{\beta}) - \Lambda_\phi^{-1}(\beta^*) \right\| \rightarrow_p 0$ and $\left\| \hat{\Sigma}_\xi^\phi - \Sigma_\xi^\phi \right\| \rightarrow_p 0$ both hold. Note that [Lemma 2](#) indicates that $\left\| \hat{\Lambda}_{\phi,n}(\hat{\beta}) - \Lambda_\phi(\beta^*) \right\| \rightarrow_p 0$, which implies that $\left\| \hat{\Lambda}_{\phi,n}^{-1}(\hat{\beta}) - \Lambda_\phi^{-1}(\beta^*) \right\| \rightarrow_p 0$ also holds.

Now we show that $\left\| \hat{\Sigma}_\xi^\phi - \Sigma_\xi^\phi \right\| \rightarrow_p 0$ holds. Our basic proof method is similar to that of [Lemma 1](#). In particular, let $\phi_n \downarrow 0$ and $\mathcal{X}_{e,n}$ be as defined as in the proof of [Lemma 1](#). Then

we have that $f_z^*(z_i^*) \geq C\phi_n^p$ as long as $\mathbf{X}_{e,i} \in \mathcal{X}_{e,n}$. Denote $G_i^* = G(z_i^*)$, we have

$$\begin{aligned} \left\| \widehat{\Sigma}_{\xi}^{\phi} - \Sigma_{\xi}^{\phi} \right\| \leq & \left\| \frac{1}{n} \sum_{i=1}^n \left(I_{n,i} \cdot \widehat{G}_i (1 - \widehat{G}_i) \left(\mathbf{X}_i^{\phi} - \widehat{\mathbb{E}} \left(\mathbf{X}_i^{\phi} \middle| \widehat{z}_i \right) \right) \left(\mathbf{X}_i^{\phi} - \widehat{\mathbb{E}} \left(\mathbf{X}_i^{\phi} \middle| \widehat{z}_i \right) \right)^{\top} \right) \right. \\ & \left. - \mathbb{E} \left(I_{n,i} \cdot G_i^* (1 - G_i^*) \left(\mathbf{X}_i^{\phi} - \mathbb{E} \left(\mathbf{X}_i^{\phi} \middle| z_i^* \right) \right) \left(\mathbf{X}_i^{\phi} - \mathbb{E} \left(\mathbf{X}_i^{\phi} \middle| z_i^* \right) \right)^{\top} \right) \right\| \end{aligned} \quad (14)$$

$$\begin{aligned} & + \left\| \frac{1}{n} \sum_{i=1}^n \left((1 - I_{n,i}) \cdot \widehat{G}_i (1 - \widehat{G}_i) \left(\mathbf{X}_i^{\phi} - \widehat{\mathbb{E}} \left(\mathbf{X}_i^{\phi} \middle| \widehat{z}_i \right) \right) \left(\mathbf{X}_i^{\phi} - \widehat{\mathbb{E}} \left(\mathbf{X}_i^{\phi} \middle| \widehat{z}_i \right) \right)^{\top} \right) \right. \\ & \left. - \mathbb{E} \left((1 - I_{n,i}) \cdot G_i^* (1 - G_i^*) \left(\mathbf{X}_i^{\phi} - \mathbb{E} \left(\mathbf{X}_i^{\phi} \middle| z_i^* \right) \right) \left(\mathbf{X}_i^{\phi} - \mathbb{E} \left(\mathbf{X}_i^{\phi} \middle| z_i^* \right) \right)^{\top} \right) \right\|. \end{aligned} \quad (15)$$

Note that \widehat{G}_i , G_i^* , \mathbf{X}_i^{ϕ} , $\widehat{\mathbb{E}} \left(\mathbf{X}_i^{\phi} \middle| \widehat{z}_i \right)$, and $\mathbb{E} \left(\mathbf{X}_i^{\phi} \middle| z_i^* \right)$ are all upper bounded, so (15) is $O_p(\phi_n)$. Now we look at (14). Note that

$$\widehat{G}_i - \frac{\sum_{j=1}^n K_{h_n}(z_i^* - z_j^*) y_j}{\sum_{j=1}^n K_{h_n}(z_i^* - z_j^*)} = \frac{\partial \widehat{G} \left(z \left(\mathbf{X}_{e,i}, \widetilde{\beta} \right) \middle| \widetilde{\beta} \right)}{\partial \beta^{\top}} \Delta \widetilde{\beta},$$

where $\widetilde{\beta}$ lies somewhere between $\widehat{\beta}$ and β^* . According to the proof of Lemma A.7, we have that

$$\sup_{(\mathbf{X}_e, \beta) \in \mathcal{X}_{e,n} \times \mathcal{B}} \left\| \frac{\partial \widehat{G} \left(z \left(\mathbf{X}_e, \beta \right) \middle| \beta \right)}{\partial \beta^{\top}} \right\| = O_p(1)$$

if $\phi_n^{-p} \left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right) \rightarrow 0$, since $\|f_z^{-1}(z(\mathbf{X}_e, \beta)) \partial H_1(z(\mathbf{X}_e, \beta), \mathbf{X}_e) / \partial z\|$ and $\|L(z(\mathbf{X}_e, \beta), \beta) f_z^{-1}(z(\mathbf{X}_e, \beta)) \partial H_2(z(\mathbf{X}_e, \beta), \mathbf{X}_e) / \partial z\|$ are both bounded for all $\beta \in \mathcal{B}$ and $\mathbf{X}_e \in \mathcal{X}_{e,n}$. So

$$\max_{1 \leq i \leq n} \left| \left(\widehat{G}_i - \frac{\sum_{j=1}^n K_{h_n}(z_i^* - z_j^*) y_j}{\sum_{j=1}^n K_{h_n}(z_i^* - z_j^*)} \right) \cdot I_{n,i} \right| = O_p(n^{-1/2}).$$

Also note that when $\phi_n^{-p} \left(h_n^{-2} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right) \rightarrow 0$,

$$\max_{1 \leq i \leq n} \left| \left(\frac{\sum_{j=1}^n K_{h_n}(z_i^* - z_j^*) y_j}{\sum_{j=1}^n K_{h_n}(z_i^* - z_j^*)} - G(z_i^*) \right) \cdot I_{n,i} \right| = O_p \left(\phi_n^{-p} \left(h_n^{-1} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right) \right),$$

this indicates that

$$\max_{1 \leq i \leq n} I_{n,i} \cdot \left| \widehat{G}_i - G(z_i^*) \right| = O_p \left(\phi_n^{-p} \left(h_n^{-1} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right) \right),$$

due to $n^{1/2} \left(h_n^{-1} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right) \rightarrow \infty$ under the choice of h_n . Using similar argument,

we can also show that

$$\max_{1 \leq i \leq n} \left\| \left(\widehat{\mathbb{E}} \left(\mathbf{X}_i^\phi \middle| \widehat{z}_i \right) - \mathbb{E} \left(\mathbf{X}_i^\phi \middle| z_i^* \right) \right) \cdot I_{n,i} \right\| = O_p \left(\phi_n^{-p} \left(h_n^{-1} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right) \right).$$

So we have that (14) is of order $O_p \left(\phi_n^{-p} \left(h_n^{-1} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right) + n^{-1/2} \right)$. It remains to choose

$$\phi_n = O \left(\left(h_n^{-1} \sqrt{\log(nh_n^{-1})/n} + h_n^3 \right)^{\frac{1}{p+1}} \right)$$

to conclude the proof. \square

Proof of Theorem 7

Proof. The proof is similar to that of Theorem 3. Note that

$$\begin{aligned} & \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \left| \bar{\sigma}^2(I_p - \delta \Psi_q(t, \beta)) - \bar{\lambda} (I_p - \delta (\Psi_q(t, \beta) + \Psi_q^T(t, \beta))) \right| \\ & \leq \delta^2 \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \|\Psi_q(t, \beta)\|^2 \leq \delta^2 \|G'\|_\infty^2 p^2 \{1 + \underline{\lambda}_\Gamma^{-1} q D_{q,0}^2\}^2. \end{aligned}$$

So if $\delta^2 \|G'\|_\infty^2 p^2 \{1 + \underline{\lambda}_\Gamma^{-1} q D_{q,0}^2\}^2 \leq \frac{1}{2} \underline{\lambda}_\Psi \delta$, or equivalently, $\delta \leq \underline{\lambda}_\Psi / \left(2 \|G'\|_\infty^2 p^2 \{1 + \underline{\lambda}_\Gamma^{-1} q D_{q,0}^2\}^2 \right)$, we have that

$$\sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \left| \bar{\sigma}^2(I_p - \delta \Psi_q(t, \beta)) - \bar{\lambda} (I_p - \delta (\Psi_q(t, \beta) + \Psi_q^T(t, \beta))) \right| \leq \underline{\lambda}_\Psi \delta / 2,$$

so

$$\sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \bar{\sigma}^2(I_p - \delta \Psi_q(t, \beta)) \leq 1 - \underline{\lambda}_\Psi \delta / 2 < 1,$$

and

$$\sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \bar{\sigma}(I_p - \delta \Psi_q(t, \beta)) \leq 1 - \underline{\lambda}_\Psi \delta / 4.$$

Then we have that

$$\begin{aligned} \|\Delta \beta_{k+1}\| & \leq \left\| \int_0^1 (I_p - \delta \Psi_q(t, \beta_k)) \Delta \beta dt + \delta \mathfrak{R}_{n,k} \right\| \\ & \leq \sup_{0 \leq t \leq 1, \beta \in \mathcal{B}} \bar{\sigma}(I_p - \delta \Psi_q(t, \beta)) \|\Delta \beta_k\| + \delta_k \|\mathfrak{R}_{n,k}\| \leq (1 - \underline{\lambda}_\Psi \delta / 4) \|\Delta \beta_k\| + \delta \|\mathfrak{R}_{n,k}\| \leq \dots \\ & \leq (1 - \underline{\lambda}_\Psi \delta / 4)^k \|\Delta \beta_1\| + \delta \sum_{j=1}^k (1 - \underline{\lambda}_\Psi \delta / 4)^{k-j} \|\mathfrak{R}_{n,j}\| \\ & \leq (1 - \underline{\lambda}_\Psi \delta / 4)^k \|\Delta \beta_1\| + 4 / \underline{\lambda}_\Psi O_p \left(\sup_{k \geq 1} \|\mathfrak{R}_{n,k}\| \right). \end{aligned}$$

When $(1 - \underline{\lambda}_\Psi \delta / 4)^k \|\Delta \beta_1\| \leq \chi_{2,n}$, or equivalently, $k \geq \frac{\log(\|\Delta \beta_1\|) - \log(\chi_{2,n})}{-\log(1 - \underline{\lambda}_\Psi \delta / 4)} = k_{1,n}^{SBGD}$, there holds $\|\Delta \beta_{k+1}\| = O_p(\chi_{2,n})$. \square

Proof of Theorem 8

Proof. We first prove Theorem 8 (i). Note that

$$\begin{aligned} \Delta \beta_{k+1} &= \left\{ \int_0^1 (I_p - \delta \Psi_q^*) dt \right\} \Delta \beta_k + \delta \mathfrak{R}_{n,k} \\ &= (I_p - \delta \Psi_q^*) \Delta \beta_k + \frac{\delta}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i + \delta \left\{ \mathfrak{R}_{n,k} - \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i \right. \\ &\quad \left. + \int_0^1 \left(\Psi_q^* - \frac{1}{n} \sum_{i=1}^n G'(z_i^* + t \mathbf{X}_i^T \Delta \beta) (\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i^T) \right) dt \Delta \beta_k \right\}. \end{aligned}$$

Define

$$\begin{aligned} \tilde{\mathfrak{R}}_{n,k} &= \mathfrak{R}_{n,k} - \frac{\delta}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i + \\ &\quad \int_0^1 \left(\Psi_q^* - \frac{1}{n} \sum_{i=1}^n G'(z_i^* + t \mathbf{X}_i^T \Delta \beta) (\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i^T) \right) dt \Delta \beta_k. \end{aligned}$$

According to Lemma 5, we have that

$$\begin{aligned} &\sup_{k \geq k_{1,n}^{SBGD} + 1} \left\| \int_0^1 \left(\Psi_q^* - \frac{1}{n} \sum_{i=1}^n G'(z_i^* + t \mathbf{X}_i^T \Delta \beta) (\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i^T) \right) dt \Delta \beta_k \right\| \\ &\leq \sup_{k \geq k_{1,n}^{SBGD} + 1, 0 \leq t \leq 1} \left\| \Psi_q^* - \frac{1}{n} \sum_{i=1}^n G'(z_i^* + t \mathbf{X}_i^T \Delta \beta) (\mathbf{X}_i \mathbf{X}_i^T - \mathfrak{X}_{q,n}(z(\mathbf{X}_{e,i}, \beta), \beta) \mathbf{X}_i^T) \right\| \sup_{k \geq k_{1,n}^{SBGD} + 1} \|\Delta \beta_k\| \\ &= O_p \left(\sqrt{pq} D_{q,0}^2 (p + q D_{q,0} D_{q,1}) \sup_{k \geq k_{1,n}^{SBGD} + 1} \|\Delta \beta\|^2 \right) \\ &= O_p \left(\sqrt{pq} D_{q,0}^2 (p + q D_{q,0} D_{q,1}) \chi_{2,n}^2 \right). \end{aligned}$$

According to Lemma 6, we have that

$$\sup_{k \geq k_{1,n}^{SBGD} + 1} \left\| \mathfrak{R}_{n,k} - \frac{\delta}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i \right\| = O_p(\chi_{4,n}).$$

This shows the result.

To prove [Theorem 8\(ii\)](#), we note that

$$\begin{aligned}
\Delta\beta_{k+k_{1,n}^{SBGD}+1} &= (I_p - \delta\Psi_q^*) \Delta\beta_{k+k_{1,n}^{SBGD}} + \frac{\delta}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i + \tilde{\mathfrak{R}}_{n,k+k_{1,n}^{SBGD}}, \\
&= (I_p - \delta\Psi_q^*)^k \Delta\beta_{k_{1,n}^{SBGD}+1} + \sum_{j=1}^k (I_p - \delta\Psi_q^*)^{j-1} \left(\frac{\delta}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i \right) \\
&\quad + \sum_{j=1}^k (I_p - \delta\Psi_q^*)^{j-1} \tilde{\mathfrak{R}}_{n,k+k_{1,n}^{SBGD}+1-j} \\
&= \Psi_q^{\star-1} \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i + (I_p - \delta\Psi_q^*)^k \Delta\beta_{k_{1,n}^{SBGD}+1} + \sum_{j=1}^k (I_p - \delta\Psi_q^*)^{j-1} \tilde{\mathfrak{R}}_{n,k+k_{1,n}^{SBGD}+1-j} \\
&\quad + \sum_{j=k+1}^{\infty} (I_p - \delta\Psi_q^*)^{j-1} \left(\frac{\delta}{n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i \right).
\end{aligned}$$

Then since

$$\left\| (I_p - \delta\Psi_q^*)^k \Delta\beta_{k_{1,n}^{SBGD}+1} \right\| = O_p \left((1 - \underline{\lambda}_\Psi \delta/4)^k \chi_{2,n} \right),$$

$$\left\| \sum_{j=1}^k (I_p - \delta\Psi_q^*)^{j-1} \tilde{\mathfrak{R}}_{n,k+k_{1,n}^{SBGD}+1-j} \right\| \leq \sum_{j=1}^{\infty} (1 - \underline{\lambda}_\Psi \delta/4)^{j-1} \sup_{k \geq k_{1,n}^{SBGD}+1} \left\| \tilde{\mathfrak{R}}_{n,k} \right\| = O_p(\chi_{5,n}),$$

and

$$\begin{aligned}
\left\| \sum_{j=k+1}^{\infty} (I_p - \delta\Psi_q^*)^{j-1} \left(\frac{\delta}{n} \sum_{i=1}^n (\mathfrak{V}_q \mathbf{r}_q(z_i^*) + \mathbf{X}_i) \varepsilon_i \right) \right\| &\leq (1 - \underline{\lambda}_\Psi \delta/4)^k \left\| \frac{4}{\underline{\lambda}_\Psi n} \sum_{i=1}^n (\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i \right\| \\
&= O_p \left((1 - \underline{\lambda}_\Psi \delta/4)^k \sqrt{\frac{pq D_{q,0}^2 (\log p)}{n}} \right) \\
&= O_p \left((1 - \underline{\lambda}_\Psi \delta/4)^k \chi_{2,n} \right).
\end{aligned}$$

So as long as $(1 - \underline{\lambda}_\Psi \delta/4)^k \chi_{2,n} \leq n^{-1/2}$, or equivalently, $k \geq k_{2,n}^{SBGD} = \frac{-\log \chi_{2,n} + \log \sqrt{n}}{-\log(1 - \underline{\lambda}_\Psi \delta/4)}$, we have that

$$\sup_{k \geq k_{2,n}^{SBGD}+1} \left\| \Delta\beta_{k+k_{1,n}^{SBGD}+1} - \Psi_q^{\star-1} \frac{1}{n} \sum_{i=1}^n (\mathfrak{V}_q \mathbf{r}_q(z_i^*) + \mathbf{X}_i) \varepsilon_i \right\| = o_p \left(n^{-\frac{1}{2}} \right).$$

The following results hold trivially. □

Proof of Theorem 9

Proof. Note that under all the conditions imposed in Theorem 8, we have that

$$\left\| \widehat{\beta} - \beta^* \right\| = O_p \left(\sqrt{pq^2 D_{q,0}^4 (\log p) / n} \right),$$

due to the fact that each element of $(\mathbf{X}_i - \mathfrak{X}_{q,i}) \varepsilon_i$ is bounded by $CqD_{q,0}^2$ and Assumption 7 holds.

To prove the theorem, we first show that

$$\sup_{1 \leq i \leq n} \left| \widehat{G}_i - G_i(z_i^*) \right| = O_p \left(\sqrt{p^2 q^4 D_{q,0}^8 (\log p) / n} + qD_{q,0}^2 \mathcal{E}_{q,0} \right).$$

Define $\widehat{z}_i = z \left(\mathbf{X}_{e,i}, \widehat{\beta} \right)$. To show the above result, note that

$$\begin{aligned} \sup_{1 \leq i \leq n} \left| \widehat{G}_i - G(z_i^*) \right| &\leq \sup_{1 \leq i \leq n} \left| \widehat{\mathbf{r}}_{q,i}^T (\widehat{\pi}_q - \pi_q^*) \right| \\ &+ \sup_{1 \leq i \leq n} \left| \widehat{\mathbf{r}}_{q,i}^T \pi_q^* - G(\widehat{z}_i) \right| + \sup_{1 \leq i \leq n} |G(\widehat{z}_i) - G(z_i^*)|. \end{aligned}$$

Obviously, the second and third terms on RHS are of order $O_p(\mathcal{E}_{q,0})$ and $O_p \left(\sqrt{p^2 q^2 D_{q,0}^4 (\log p) / n} \right)$, while the first term is bounded by $\sqrt{q} D_{q,0} \left\| \widehat{\pi}_q - \pi_q^* \right\|$. Note that

$$\begin{aligned} \widehat{\pi}_q - \pi_q^* &= \Gamma_{q,n}^{-1}(\widehat{\beta}) \left(\frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{r}}_{q,i} (G(\widehat{z}_i) - G(z_i^*)) \right) + \Gamma_{q,n}^{-1}(\widehat{\beta}) \left(\frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{r}}_{q,i} R_q(\widehat{z}_i) \right) \\ &+ \Gamma_{q,n}^{-1}(\widehat{\beta}) \left(\frac{1}{n} \sum_{i=1}^n \mathbf{r}_q(\widehat{z}_i) \varepsilon_i \right). \end{aligned}$$

So we have that $\left\| \widehat{\pi}_q - \pi_q^* \right\| = O_p \left(\sqrt{p^2 q^3 D_{q,0}^6 (\log p) / n} + \sqrt{q} D_{q,0} \mathcal{E}_{q,0} \right)$ and the third term is of order $O_p \left(\sqrt{p^2 q^4 D_{q,0}^8 (\log p) / n} + qD_{q,0}^2 \mathcal{E}_{q,0} \right)$. This proves the first result.

We also note that according to the proof of Lemma A.11, we have that

$$\sup_{1 \leq i \leq n} \left\| \mathfrak{X}_{q,n}(\widehat{z}_i, \widehat{\beta}) - \mathfrak{X}_q(z_i^*, \beta^*) \right\| = O_p \left(\sqrt{p^3 q^6 D_{q,0}^{10} D_{q,1}^2 \log(pn) / n} \right).$$

Then we show that

$$\begin{aligned} \max_{1 \leq i \leq n} &\left\| \widehat{G}_i (1 - \widehat{G}_i) \left(\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta}) \right) \left(\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta}) \right)^T - G_i (1 - G_i) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*)) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*))^T \right\| \\ &= O_p \left(\sqrt{p^4 q^8 D_{q,0}^{14} (\log pn) / n} (D_{q,0} + D_{q,1}) + pq^3 D_{q,0}^6 \mathcal{E}_{q,0} \right). \end{aligned}$$

Note that the above is bounded by

$$\begin{aligned} & \max_{1 \leq i \leq n} \left\| \left(\widehat{G}_i (1 - \widehat{G}_i) - G_i (1 - G_i) \right) \left(\mathbf{X}_i - \mathfrak{X}_{q,n} (\widehat{z}, \widehat{\beta}) \right) \left(\mathbf{X}_i - \mathfrak{X}_{q,n} (\widehat{z}, \widehat{\beta}) \right)^\top \right\| \\ & + \max_{1 \leq i \leq n} \left\| G_i (1 - G_i) \left(\left(\mathbf{X}_i - \mathfrak{X}_{q,n} (\widehat{z}, \widehat{\beta}) \right) \left(\mathbf{X}_i - \mathfrak{X}_{q,n} (\widehat{z}, \widehat{\beta}) \right)^\top - (\mathbf{X}_i - \mathfrak{X}_q (z_i^*, \beta^*)) (\mathbf{X}_i - \mathfrak{X}_q (z_i^*, \beta^*))^\top \right) \right\| \end{aligned}$$

where the first term is of order $O_p \left(\sqrt{p^4 q^8 D_{q,0}^{16} (\log p) / n} + p q^3 D_{q,0}^6 \mathcal{E}_{q,0} \right)$, while the second term is of order $O_p \left(\sqrt{p^4 q^8 D_{q,0}^{14} D_{q,1}^2 (\log pn) / n} \right)$. Together we show the result.

Next we show that

$$\left\| \widehat{\Psi}_q^* - \Psi_q^* \right\| = O_p \left(\sqrt{p^4 q^4 D_{q,0}^4 \log (pq D_{q,0} D_{q,1} n) / n} \right).$$

Since $v_G \geq 2$, we have that

$$\begin{aligned} & \sup_{1 \leq i \leq n} \left| \widehat{G}'_i - G' (z(\mathbf{X}_{e,i}, \beta^*)) \right| \leq \sup_{1 \leq i \leq n} \left| \widehat{\mathbf{r}}_{q,i}^\top (\widehat{\pi}_q - \pi_q^*) \right| \\ & + \sup_{1 \leq i \leq n} \left| \widehat{\mathbf{r}}_{q,i}^\top \pi_q^* - G' (\widehat{z}_i) \right| + \sup_{1 \leq i \leq n} |G' (\widehat{z}_i) - G' (z_i^*)| \\ & = O_p \left(\sqrt{p^2 q^4 D_{q,0}^8 D_{q,1}^2 (\log p) / n} + q D_{q,0} D_{q,1} \mathcal{E}_{q,0} + \mathcal{E}_{q,1} \right). \end{aligned}$$

So

$$\begin{aligned} \left\| \widehat{\Psi}_q^* - \Psi_q^* \right\| & \leq \left\| \frac{1}{n} \sum_{i=1}^n \left(\widehat{G}'_i - G' (z_i^*) \right) \cdot \left(\mathbf{X}_i \mathbf{X}_i^\top - \mathfrak{X}_{q,n} (\widehat{z}_i, \widehat{\beta}) \mathbf{X}_i^\top \right) \right\|, \\ & + \left\| \frac{1}{n} \sum_{i=1}^n G' (z_i^*) \cdot \left(\left(\mathfrak{X}_{q,n} (\widehat{z}_i, \widehat{\beta}) - \mathfrak{X}_q (z_i^*, \beta^*) \right) \mathbf{X}_i^\top \right) \right\| \\ & + \left\| \frac{1}{n} \sum_{i=1}^n G' (z_i^*) \cdot \mathfrak{X}_q (z_i^*, \beta^*) \mathbf{X}_i^\top - \Psi_q^* \right\| \\ & = O_p \left(\sqrt{p^4 q^6 D_{q,0}^{12} D_{q,1}^2 \log (pn) / n} + p q^2 D_{q,0}^3 D_{q,1} \mathcal{E}_{q,0} + p q D_{q,0}^2 \mathcal{E}_{q,1} \right), \end{aligned}$$

which also implies that $\bar{\sigma} \left(\widehat{\Psi}_q^{*-1} \right) = O_p (1)$, and

$$\left\| \widehat{\Psi}_q^{*-1} - \Psi_q^{*-1} \right\| = O_p \left(\sqrt{p^4 q^6 D_{q,0}^{12} D_{q,1}^2 (\log pn) / n} + p q^2 D_{q,0}^3 D_{q,1} \mathcal{E}_{q,0} + q D_{q,0}^2 \mathcal{E}_{q,1} \right)$$

Now we are ready to demonstrate the consistency of the variance estimator. Note that

$$\begin{aligned}
& |\widehat{\sigma}_S^2(\rho) - \sigma_S^2(\rho)| \\
& \leq \|\rho\|^2 \left\| \widehat{\Psi}_q^{\star-1} \frac{1}{n} \sum_{i=1}^n \left\{ \widehat{G}_i(1 - \widehat{G}_i) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta})) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta}))^\top \right\} (\widehat{\Psi}_q^{\star-1})^\top \right. \\
& \quad \left. - \Psi_q^{\star-1} \mathbb{E} \left\{ G(z_i^*) (1 - G(z_i^*)) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*)) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*))^\top \right\} (\Psi_q^{\star-1})^\top \right\| \\
& \leq \|\rho\|^2 \left\| \widehat{\Psi}_q^{\star-1} - \Psi_q^{\star-1} \right\| \left\| \frac{1}{n} \sum_{i=1}^n \left\{ \widehat{G}_i(1 - \widehat{G}_i) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta})) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta}))^\top \right\} (\Psi_q^{\star-1})^\top \right\| \\
& \quad + \|\rho\|^2 \left\| \Psi_q^{\star-1} \left(\frac{1}{n} \sum_{i=1}^n \left\{ \widehat{G}_i(1 - \widehat{G}_i) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta})) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta}))^\top \right\} \right. \right. \\
& \quad \left. \left. - \mathbb{E} \left\{ G(z_i^*) (1 - G(z_i^*)) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta})) (\mathbf{X}_i - \mathfrak{X}_{q,n}(\widehat{z}, \widehat{\beta}))^\top \right\} \right\} (\widehat{\Psi}_q^{\star-1})^\top \right\| \\
& \quad + \|\rho\|^2 \left\| \Psi_q^{\star-1} \mathbb{E} \left\{ G(z_i^*) (1 - G(z_i^*)) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*)) (\mathbf{X}_i - \mathfrak{X}_q(z_i^*, \beta^*))^\top \right\} (\widehat{\Psi}_q^{\star-1} - \Psi_q^{\star-1})^\top \right\|.
\end{aligned}$$

The first and the third terms are of order $O_p \left(\sqrt{p^6 q^8 D_{q,0}^{16} D_{q,1}^2 (\log pn) / n} + p^2 q^3 D_{q,0}^4 D_{q,1} \mathcal{E}_{q,0} + p q^2 D_{q,0}^4 \mathcal{E}_{q,1} \right)$, and the second term is of order $O_p \left(\sqrt{p^4 q^8 D_{q,0}^{14} (\log pn) / n} (D_{q,0} + D_{q,1}) + p q^3 D_{q,0}^6 \mathcal{E}_{q,0} \right)$. Together, we have that

$$|\widehat{\sigma}_S^2(\rho) - \sigma_S^2(\rho)| = O_p \left(\sqrt{p^6 q^8 D_{q,0}^{16} D_{q,1}^2 (\log pn) / n} + p q^3 D_{q,0}^4 (p D_{q,1} + D_{q,0}^2) \mathcal{E}_{q,0} + p q^2 D_{q,0}^4 \mathcal{E}_{q,1} \right),$$

which implies that $|\widehat{\sigma}_S^2(\rho) - \sigma_S^2(\rho)| \rightarrow_p 0$ under all the conditions. \square

References

- Alekh Agarwal, Sham Kakade, Nikos Karampatziakis, Le Song, and Gregory Valiant. Least squares revisited: Scalable approaches for multi-class prediction. In *International Conference on Machine Learning*, pages 541–549. PMLR, 2014.
- Hyungtaik Ahn, Hidehiko Ichimura, James L Powell, and Paul A Ruud. Simple estimators for invertible index models. *Journal of Business & Economic Statistics*, 36(1):1–10, 2018.
- A. Belloni, V. Chernozhukov, D. Chetverikov, and Y. Wei. Uniformly valid post-regularization confidence regions for many functional parameters in z-estimation framework. *Annals of Statistics*, 46:3643–3675, 2018.
- Herman J Bierens. Consistency and asymptotic normality of sieve ml estimators under low-level conditions. *Econometric Theory*, 30(5):1021–1076, 2014.
- M.D. Cattaneo, M. Jansson, and W.K. Newey. Alternative asymptotics and the partially linear model with many regressors. *Econometric Theory*, 34:277–301, 2018.
- Christopher Cavanagh and Robert P Sherman. Rank estimators for monotonic index models. *Journal of Econometrics*, 84(2):351–382, 1998.
- Xiaohong Chen. Large sample sieve estimation of semi-nonparametric models. *Handbook of Econometrics*, 6:5549–5632, 2007.
- Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. Central limit theorems and bootstrap in high dimensions. *The Annals of Probability*, 45(4):2309–2352, 2017.
- Yanqin Fan, Fang Han, Wei Li, and Xiao-Hua Zhou. On rank estimators in increasing dimensions. *Journal of Econometrics*, 214(2):379–412, 2020.
- Li Gan, Zhichao Yin, Nan Jia, Shu Xu, Shuang Ma, Lu Zheng, et al. Data you need to know about china. *Springer Berlin Heidelberg*. [https://doi, 10:978-3](https://doi.org/10.978-3), 2014.
- Aaron K Han. Non-parametric analysis of a generalized regression model: the maximum rank correlation estimator. *Journal of Econometrics*, 35(2-3):303–316, 1987.
- Keisuke Hirano, Guido W Imbens, and Geert Ridder. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica*, 71(4):1161–1189, 2003.

- Hidehiko Ichimura. Semiparametric least squares (sls) and weighted sls estimation of single-index models. *Journal of econometrics*, 58(1-2):71–120, 1993.
- Roger W Klein and Richard H Spady. An efficient semiparametric estimator for binary response models. *Econometrica: Journal of the Econometric Society*, pages 387–421, 1993.
- Katta G Murty and Santosh N Kabadi. Some np-complete problems in quadratic and non-linear programming. *Mathematical Programming*, 39:117–129, 1987.
- W.K. Newey and F. Windmeijer. Generalized method of moments with many weak moment conditions. *Econometrica*, 77(3):687–719, 2009.
- J François Outreville. The relationship between relative risk aversion and the level of education: A survey and implications for the demand for life insurance. *Journal of economic surveys*, 29(1):97–111, 2015.
- James L Powell, James H Stock, and Thomas M Stoker. Semiparametric estimation of index coefficients. *Econometrica: Journal of the Econometric Society*, pages 1403–1430, 1989.
- Y. Shin and Z. Todorov. Exact computation of the maximum rank correlation estimator. *Forthcoming, Econometrics Journal*, 2021.
- Thomas M Stoker. Consistent estimation of scaled coefficients. *Econometrica: Journal of the Econometric Society*, pages 1461–1481, 1986.
- Pragya Sur and Emmanuel J Candès. A modern maximum-likelihood theory for high-dimensional logistic regression. *Proceedings of the National Academy of Sciences*, 116(29):14516–14525, 2019.
- Panos Toulis and Edoardo M Airolidi. Asymptotic and finite-sample properties of estimators based on stochastic gradients. *The Annals of Statistics*, 45(4):1694–1727, 2017.
- Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.
- Haiyong Wang and Shuhuang Xiang. On the convergence rates of legendre approximation. *Mathematics of computation*, 81(278):861–877, 2012.
- Qingsong Yao. *Kernel-Based Learning of Monotone Index Models with Large Dimensionality*. PhD thesis, Boston College, 2023.