

CAMBRIDGE WORKING PAPERS IN ECONOMICS  
CAMBRIDGE-INET WORKING PAPERSA Unified Framework for Efficient Estimation of  
General Treatment Models

Chunrong

Ai

University of  
Florida

Oliver

Linton

University of  
Cambridge

Kaiji

Motegi

Bank of  
England

Zheng

Zhang

University of  
China

## Abstract

This paper presents a weighted optimization framework that unifies the binary, multivalued, continuous, as well as mixture of discrete and continuous treatment, under the unconfounded treatment assignment. With a general loss function, the framework includes the average, quantile and asymmetric least squares causal effect of treatment as special cases. For this general framework, we first derive the semiparametric efficiency bound for the causal effect of treatment, extending the existing bound results to a wider class of models. We then propose a generalized optimization estimation for the causal effect with weights estimated by solving an expanding set of equations. Under some sufficient conditions, we establish consistency and asymptotic normality of the proposed estimator of the causal effect and show that the estimator attains our semiparametric efficiency bound, thereby extending the existing literature on efficient estimation of causal effect to a wider class of applications. Finally, we discuss estimation of some causal effect functionals such as the treatment effect curve and the average outcome. To evaluate the finite sample performance of the proposed procedure, we conduct a small scale simulation study and find that the proposed estimation has practical value. To illustrate the applicability of the procedure, we revisit the literature on campaign advertise and campaign contributions. Unlike the existing procedures which produce mixed results, we find no evidence of campaign advertise on campaign contribution.

## Reference Details

1934 Cambridge Working Papers in Economics  
2019/28 Cambridge-INET Working Paper Series

Published 3 March 2019

Key Words Treatment effect, Semiparametric efficiency, Stabilized Weights

Websites [www.econ.cam.ac.uk/cwpe](http://www.econ.cam.ac.uk/cwpe)  
[www.inet.econ.cam.ac.uk/working-papers](http://www.inet.econ.cam.ac.uk/working-papers)

# A Unified Framework for Efficient Estimation of General Treatment Models

Chunrong Ai <sup>\*\*</sup>, Oliver Linton <sup>†\*</sup>, Kaiji Motegi <sup>‡†</sup>, and Zheng Zhang <sup>§‡</sup>

*<sup>\*</sup>Department of Economics, University of Florida*

*<sup>\*</sup>Faculty of Economics, University of Cambridge*

*<sup>†</sup>Graduate School of Economics, Kobe University*

*<sup>‡</sup>Institute of Statistics & Big Data, Renmin University of China*

March 18, 2019

## Abstract

This paper presents a weighted optimization framework that unifies the binary, multi-valued, continuous, as well as mixture of discrete and continuous treatment, under the unconfounded treatment assignment. With a general loss function, the framework includes the average, quantile and asymmetric least squares causal effect of treatment as special cases. For this general framework, we first derive the semiparametric efficiency bound for the causal effect of treatment, extending the existing bound results to a wider class of models. We then propose a generalized optimization estimation for the causal effect with weights estimated by solving an expanding set of equations. Under some sufficient conditions, we establish consistency and asymptotic normality of the proposed estimator of the causal effect and show that the estimator attains our semiparametric efficiency bound, thereby extending the existing literature on efficient estimation of causal effect to a wider class of applications. Finally, we discuss estimation of some causal effect functionals such as the treatment effect curve and the average outcome. To evaluate the finite sample performance of the proposed procedure, we conduct a small scale simulation study and find that the proposed estimation has practical value. To illustrate the applicability of the procedure, we revisit the literature on campaign advertise and campaign contributions. Unlike the existing procedures which produce mixed results, we find no evidence of campaign advertise on campaign contribution.

---

<sup>\*</sup>E-mail: chunrong.ai@warrington.ufl.edu

<sup>†</sup>E-mail: obl20@cam.ac.uk

<sup>‡</sup>E-mail: motegi@econ.kobe-u.ac.jp

<sup>§</sup>E-mail: zhengzhang@ruc.edu.cn

*Keywords:* Treatment effect; Semiparametric efficiency; Stabilized Weights.

# 1 Introduction

Modeling and estimating the causal effect of treatment have received considerable attention from both econometrics and statistics literature (see [Hirano, Imbens, and Ridder \(2003\)](#), [Imbens \(2004\)](#), [Abadie \(2005\)](#), [Heckman and Vytlacil \(2005\)](#), [Angrist and Pischke \(2008\)](#), [Imbens and Wooldridge \(2009\)](#), [Fan and Park \(2010\)](#), [Chernozhukov, Fernández-Val, and Melly \(2013\)](#), [Rothe \(2017\)](#), and [Śłoczyński and Wooldridge \(2018\)](#), [Athey, Imbens, and Wager \(2018\)](#), [Wager and Athey \(2018\)](#) for examples). Most existing studies focus on the binary treatment where an individual either receives the treatment or does not, ignoring the treatment intensity. In many applications, however, the treatment intensity is a part of the treatment, and its causal effect is also of great interest to decision makers. For example, in evaluating how financial incentives affect health care providers, the causal effect may depend on not only the introduction of incentive but also the level of incentive. For example, in evaluating the effect of corporate bond purchase schemes on market quality, the causal effect may depend not just on whether the bond is selected into the scheme but how much of it is purchased (see [Boneva, Elliot, Kaminska, Linton, Morely, and McLaren, \(2018\)](#)). Similarly, in studying how taxes affect addictive substance usages, the causal effect may depend on the imposition of tax as well as the tax rate. In recognition of the importance of the treatment intensity, the binary treatment literature has been extended to the multi-valued treatment (e.g., [Imbens \(2000\)](#) and [Cattaneo \(2010\)](#)) and continuous treatment (e.g., [Hirano and Imbens \(2004\)](#), [Imai and van Dyk \(2004\)](#), [Florens, Heckman, Meghir, and Vytlacil \(2008\)](#), [Fong, Hazlett, and Imai \(2018\)](#) and [Yiu and Su \(2018\)](#)).

The parameter of primary interest in this literature is the average causal effect of treatment, defined as the difference in response to two levels of treatment by the same individual, averaged over a set of individuals. The identification and estimation difficulty is that each individual only receives one level of treatment. To overcome the difficulty, researchers impose the *unconfounded treatment assignment* condition, which allows them to find statistical matches for each observed individual from all other treatment levels. Under this condition, the average causal effect is estimated in the binary treatment by the difference of the weighted average responses with the propensity scores as weights (e.g., [Rosenbaum and Rubin \(1983\)](#), [Hirano, Imbens, and Ridder \(2003\)](#), [Busso, DiNardo, and McCrary \(2014\)](#)). Other popular methods include regression adjustment ([Rubin, 1977](#), [Angrist and Pischke, 2008](#)), matching ([Imbens, 2004](#), [Abadie and Imbens, 2006, 2011, 2012, 2016](#)), imputation ([Heckman, Ichimura, and Todd, 1998](#), [Cattaneo and Farrell, 2011](#)), and hybrid method

(Farrell, 2015, Chernozhukov, Escanciano, Ichimura, and Newey, 2016, Słoczyński and Wooldridge, 2018). The efficiency bound of the average causal effect in this model is derived by Robins, Rotnitzky, and Zhao (1994) and Hahn (1998), and efficient estimation is proposed by Robins, Rotnitzky, and Zhao (1994), Hahn (1998), Hirano, Imbens, and Ridder (2003), Bang and Robins (2005), Qin and Zhang (2007), Cao, Tsiatis, and Davidian (2009), Tan (2010), Vansteelandt, Bekaert, and Claeskens (2010), Graham, Pinto, and Egel (2012), and Chan, Yam, and Zhang (2016). Of particular interest in this literature is the study by Hirano, Imbens, and Ridder (2003) which shows that the weighted average difference estimator attains the semiparametric efficiency bound if the weights are estimated by the empirical likelihood estimation. In the multi-valued treatment, Imbens (2000) generalizes the propensity score, and Cattaneo (2010) derives the efficiency bound and proposes an estimator that attains the efficiency bound. In the continuous treatment, Hirano and Imbens (2004) and Imai and van Dyk (2004) parameterize the generalized propensity score function and propose a consistent estimation of the average causal effect. Their estimators are not efficient and could be biased if the generalized propensity score function is misspecified. Florens, Heckman, Meghir, and Vytlacil (2008) use a control function approach to identify the average causal effect in the continuous treatment and propose a consistent estimation. It is unclear if their estimation is efficient. Galvao and Wang (2015) estimate the continuous treatment effects through stabilized weighting. They do not study how to construct the stabilized weights such that their estimation is efficient. Kennedy, Ma, McHugh, and Small (2017) propose a nonparametric kernel estimator for the treatment effects curve, again the efficient estimation is still unclear. Fong, Hazlett, and Imai (2018) propose an estimation of the average causal effect of continuous treatment but do not establish consistency of their estimation. In fact, their simulation results indicate their estimation could be seriously biased. Yiu and Su (2018) study the average causal effect of both discrete and continuous treatment by parameterizing propensity scores. Their estimation could be biased if their parameterization is incorrect.

Without unconfoundedness, several methods have also been developed to conduct causal analysis, for example, sensitivity analyses (Rosenbaum and Rubin, 1983), bounds analyses (Manski, 1990, 2003, 2009, Imbens and Manski, 2004), instrumental variables (Imbens and Angrist, 1994, Angrist, Imbens, and Rubin, 1996, Imbens and Rubin, 1997, Hong and Nekipelov, 2010, Kitagawa, 2015), regression discontinuity (Hahn, Todd, and Van der Klaauw, 2001, Imbens and Lemieux, 2008, Lee and Lemieux, 2010, Cattaneo, Idrobo, and Titiunik, 2017), difference-in-differences (Ashenfelter and Card, 1985, Bertrand, Duflo, and Mullainathan, 2004, Abadie, 2003, 2005, Donald and Lang, 2007, Athey and Imbens, 2006), and synthetic controls (Abadie and Gardeazabal, 2003, Abadie, Diamond, and

Hainmueller, 2010, 2015).

In addition to the average causal effect of treatment (ATE), it is also important to investigate the distributional impact of treatment. For instance, a decision maker may be interested in the causal effect of a treatment on the outcome dispersion or on the lower tail of the outcome distribution. Doksum (1974) and Lehmann (1975) introduce the quantile causal effect of treatment (QTE). Firpo (2007) computes the efficiency bound and proposes an efficient estimation of QTE for the binary treatment. For additional studies on QTE, we refer to Abadie, Angrist, and Imbens (1998), Chernozhukov and Hansen (2005), Angrist and Pischke (2008), Frölich and Melly (2013), and Donald and Hsu (2014).

To the best of our knowledge, we are unaware of existence of any studies that compute the efficiency bound and propose efficient estimation of the causal effect in the continuous or mixture of discrete and continuous treatment under general loss function that permits ATE and QTE. The main objective of this paper is to present a weighted optimization framework that unifies the binary, multi-valued, continuous as well as mixture of discrete and continuous treatment and allows for general loss function, where the weights are called the *stabilized weights* by Robins, Hernán, and Brumback (2000) and are defined as the ratio of the marginal probability distribution of the treatment status over the conditional probability distribution of the treatment status given covariates. For this general framework, we first apply the approach of Bickel, Klaassen, Ritov, and Wellner (1993) to compute the efficiency bound of the causal effect of treatment, extending the semiparametric efficiency bound results of Hahn (1998), Firpo (2007), and Cattaneo (2010) from the binary treatment to a variety of treatments and to the general loss function. Our bound reveals that the weighted optimization with known stabilized weights does not produce efficient estimation since it fails to account for the information restricting the stabilized weights. Similar observation is also noted by Hirano, Imbens, and Ridder (2003) in the binary treatment. Here we show that their observation holds true for much wider class of treatment models. We exploit the information that the stabilized weights satisfy some (finite but expanding number of) equations by estimating the stabilized weights from those equations and then estimate the causal effect by the generalized optimization with the true stabilized weights replaced by the estimated weights. Under some sufficient conditions, we show that our proposed estimator is consistent and asymptotically normally distributed and, more importantly, it attains our semiparametric efficiency bound, thereby extending the efficient estimation work of Hirano, Imbens, and Ridder (2003), Firpo (2007), and Cattaneo (2010) to a much wider class of treatment models.

The paper is organized as follows. Section 2 sets up the basic framework, Section 3 computes the semiparametric efficiency bound of the causal effect of treatment, Section 4

presents a generalized optimization estimator, Section 5 establishes large sample properties of the proposed estimator, Section 6 presents a consistent covariance matrix, Section 7 proposes two data-driven approaches for selecting tuning parameters, Section 8 discusses some extensions, Section 9 reports on a simulation study, Section 10 presents an application, followed by some concluding remarks in Section 11. All technical proofs and extra simulation results are relegated to the supplemental material [Ai, Linton, Motegi, and Zhang \(2019\)](#).

## 2 Basic framework and notation

Let  $T$  denote the observed treatment status variable with support  $\mathcal{T} \subset \mathbb{R}$ , where  $\mathcal{T}$  is either a discrete or a continuous or a mixture of discrete and continuous subset, and has a marginal probability distribution function  $F_T(t)$ . Let  $Y^*(t)$  denote the potential response when treatment  $T = t$  is assigned. Let  $L(\cdot)$  denote a known convex loss function whose derivative, denoted by  $L'(\cdot)$ , exists almost everywhere. For the leading part of the paper, we shall maintain that there exists a parametric causal effect function  $g(t; \beta)$  with the unknown value  $\beta_0 \in \mathbb{R}^p$  (with  $p \in \mathbb{N}$ ) uniquely solving:

$$\beta_0 = \arg \min_{\beta} \int_{\mathcal{T}} \mathbb{E} [L(Y^*(t) - g(t; \beta))] dF_T(t). \quad (2.1)$$

The parameterization of the causal effect is restrictive. Some extensions to the unspecified causal effect function shall be discussed later in the paper (see Section 8).

The generality of model (2.1) permits many important models and much more. For example, it includes the average causal effect of binary treatment studied in [Hahn \(1998\)](#) and [Hirano, Imbens, and Ridder \(2003\)](#) (i.e.,  $\mathcal{T} = \{0, 1\}$ ,  $L(v) = v^2$  and  $g(t; \beta_0) = \beta_0 + \beta_1 t$ ), the quantile causal effect of binary treatment studied in [Firpo \(2007\)](#) (i.e.,  $\mathcal{T} = \{0, 1\}$ ,  $L(v) = v(\tau - I(v \leq 0))$  is an almost everywhere differentiable function with  $\tau \in (0, 1)$  and  $g(t; \beta_0) = t\beta_1 + (1 - t)\beta_0$ ), the average causal effect of multi-valued treatment studied in [Cattaneo \(2010\)](#) (i.e.,  $\mathcal{T} = \{0, 1, \dots, J\}$  for some  $J \in \mathbb{N}$ ,  $L(v) = v^2$  and  $g(t; \beta_0) = \sum_{j=0}^J \beta_j I(t = j)$ ), and the average causal effect of continuous treatment studied in [Hirano and Imbens \(2004\)](#) (i.e.,  $L(v) = v^2$  and  $\mathbb{E}[Y^*(t)] = g(t; \beta_0)$ ). It also includes the quantile causal effect of multi-valued (i.e.,  $L(v) = v(\tau - I(v \leq 0))$  with  $\tau \in (0, 1)$  and  $g(t; \beta_0) = \sum_{j=0}^J \beta_j I(t = j)$ ) and continuous treatment (i.e.,  $L(v) = v(\tau - I(v \leq 0))$  and  $\inf \{q : \mathbb{P}(Y^*(t) \geq q) \leq \tau\} = g(t; \beta_0)$ ). The latter has never been studied in the existing literature. Moreover, with  $L(v) = v^2 |\tau - I(v \leq 0)|$ , it covers asymmetric least squares estimation of the causal effect of (binary, multi-valued, continuous, mixture of discrete and

continuous) treatment. The asymmetric least squares regression received attention from some noted econometricians (see [Newey and Powell \(1987\)](#)) but zero attention in the causal effect literature.

The problem with (2.1) is that the potential outcome  $Y^*(t)$  is not observed for all  $t$ . Let  $Y := Y^*(T)$  denote the observed response. One may attempt to solve:

$$\min_{\beta} \mathbb{E}[L(Y - g(T; \beta))].$$

However, if there exists a selection into treatment, the true value  $\beta_0$  does not solve the above minimization problem. Indeed, in this case, the observed response and treatment assignment data alone cannot identify  $\beta_0$ . To address this identification issue, most studies in the literature impose a selection on observable condition (e.g., [Hirano, Imbens, and Ridder \(2003\)](#), [Imai and van Dyk \(2004\)](#) and [Fong, Hazlett, and Imai \(2018\)](#)). Specifically, let  $\mathbf{X}$  denote a vector of covariates. The following condition shall be maintained throughout the paper.

**Assumption 1** (*Unconfounded Treatment Assignment*). For all  $t \in \mathcal{T}$ , given  $\mathbf{X}$ ,  $T$  is independent of  $Y^*(t)$ , i.e.,  $Y^*(t) \perp T | \mathbf{X}$ , for all  $t \in \mathcal{T}$ .

Let  $F_{T|\mathbf{X}}$  denote the conditional probability distribution of  $T$  given the observed covariates  $\mathbf{X}$  and let  $dF_{T|\mathbf{X}}$  denote the probability measure. In the literature,  $dF_{T|\mathbf{X}}$  is called the *generalized propensity score* ([Hirano and Imbens, 2004](#), [Imai and van Dyk, 2004](#)). Suppose that  $dF_{T|\mathbf{X}}(T | \mathbf{X})$  is positive everywhere and denote

$$\pi_0(T, \mathbf{X}) := \frac{dF_T(T)}{dF_{T|\mathbf{X}}(T | \mathbf{X})}.$$

The function  $\pi_0(T, \mathbf{X})$  is called the *stabilized weight* in [Robins, Hernán, and Brumback \(2000\)](#). Under Assumption 1, we obtain

$$\mathbb{E}[\pi_0(T, \mathbf{X})L(Y - g(T; \beta))] = \int \mathbb{E}[L(Y^*(t) - g(t; \beta))] dF_T(t) \quad (2.2)$$

(see Appendix A), and hence the true value  $\beta_0$  solves the weighted optimization problem:

$$\beta_0 = \arg \min_{\beta} \mathbb{E}[\pi_0(T, \mathbf{X})L(Y - g(T; \beta))]. \quad (2.3)$$

This result is very insightful. It tells us that the selection bias in the *unconfounded treatment assignment* can be corrected through covariate-balancing. More importantly, it says that the true value  $\beta_0$  can be identified from the observed data. The weighted optimization (2.3) provides a unified framework for estimating the causal effect of a variety of treatments, including binary, multi-level, continuous and mixture of discrete and continuous treatment, and under general loss function. The goal of this paper is to compute the semiparametric efficiency bound and present an efficient estimation of  $\beta_0$  under this general framework.



### 3 Efficiency bound

We begin by applying the approach of [Bickel, Klaassen, Ritov, and Wellner \(1993\)](#) to compute the semiparametric efficiency bound of the parameter  $\beta_0$  defined by (2.1) under Assumption 1. This gives the least possible variance achievable by a regular estimator in the semiparametric model. The result is presented in the following theorem.

**Theorem 1.** *Suppose that  $g(T; \beta)$  is twice differentiable with respect to  $\beta$  in the parameter space  $\Theta \subset \mathbb{R}^p$ , with  $m(T; \beta_0) := \nabla_{\beta} g(T; \beta_0)$ , and  $\mathbb{E}[L'(Y - g(T; \beta))|Y, \mathbf{X}]$  is differentiable with respect to  $\beta \in \Theta$ . Denote  $\varepsilon(T, \mathbf{X}; \beta_0) := \mathbb{E}[L'(Y - g(T; \beta_0))|T, \mathbf{X}]$ ,  $H_0 := -\nabla_{\beta} \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta))m(T; \beta)]|_{\beta=\beta_0}$ , and*

$$\begin{aligned} \psi(Y, T, \mathbf{X}; \beta_0) &:= \pi_0(T, \mathbf{X})m(T; \beta_0)L'(Y - g(T; \beta_0)) - \pi_0(T, \mathbf{X})m(T; \beta_0)\varepsilon(T, \mathbf{X}; \beta_0) \\ &\quad + \mathbb{E}[\varepsilon(T, \mathbf{X}; \beta_0)\pi_0(T, \mathbf{X})m(T; \beta_0)|T] + \mathbb{E}[\varepsilon(T, \mathbf{X}; \beta_0)\pi_0(T, \mathbf{X})m(T; \beta_0)|\mathbf{X}]. \end{aligned}$$

*Suppose that  $H_0$  is nonsingular and  $\mathbb{E}[\psi(Y, T, \mathbf{X}; \beta_0)\psi(Y, T, \mathbf{X}; \beta_0)^{\top}]$  exists and is finite. Under Assumption 1, namely  $Y^*(t) \perp T|\mathbf{X}$  for all  $t \in \mathcal{T}$ , and model (2.1), the efficient influence function of  $\beta_0$  is given by*

$$S_{eff}(Y, T, \mathbf{X}; \beta_0) = H_0^{-1}\psi(Y, T, \mathbf{X}; \beta_0).$$

*Consequently, the efficient variance bound of  $\beta_0$  is*

$$V_{eff} = \mathbb{E}[S_{eff}(Y, T, \mathbf{X}; \beta_0)S_{eff}(Y, T, \mathbf{X}; \beta_0)^{\top}].$$

The proof of Theorem 1 is given in the supplemental material [Ai, Linton, Motegi, and Zhang \(2019, Section 2.1\)](#). We can rewrite the influence function  $\psi(Y, T, \mathbf{X}; \beta_0)$  defined in Theorem 1 in the more intuitive form: denote  $\varrho(Y, T, \mathbf{X}; \beta) := \pi_0(T, \mathbf{X})m(T; \beta)L'(Y - g(T; \beta))$  and

$$\psi(Y, T, \mathbf{X}; \beta_0) = \varrho(T, \mathbf{X}, Y; \beta_0) - res_{add}\varrho(T, \mathbf{X}, Y; \beta_0),$$

where the operator  $res_{add}(\cdot)$  is defined by

$$\begin{aligned} res_{add}f(Y, T, \mathbf{X}) &:= \mathbb{E}[f(T, \mathbf{X}, Y)|T, \mathbf{X}] - \mathbb{E}_{add}[f(T, \mathbf{X}, Y)|T, \mathbf{X}], \\ \mathbb{E}_{add}[f(T, \mathbf{X}, Y)|T, \mathbf{X}] &:= \mathbb{E}[f(T, \mathbf{X}, Y)|T] + \mathbb{E}[f(T, \mathbf{X}, Y)|\mathbf{X}]. \end{aligned}$$

where the operator  $\mathbb{E}_{add}[\cdot]$  projects a random variable on to the space of additive functions

$$\{g(T, \mathbf{X}) : g(T, \mathbf{X}) = h_T(T) + h_X(\mathbf{X})\}$$



inside the space generated by  $T, \mathbf{X}$ , except that the projection is with respect to product measure  $dF_T(t) \times dF_X(\mathbf{x})$  (Nielsen and Linton, 1998).

In the continuous case  $\pi_0(T, \mathbf{X})$  can be written as

$$\pi_0(T, \mathbf{X}) = \frac{f_T(T)f_X(\mathbf{X})}{f_{T,X}(T, \mathbf{X})}$$

and we know that  $-\mathbb{E}[\log \pi(T, \mathbf{X})]$  is the Kullback-Leibler divergence of the joint density from the product of the marginals. The property of  $\pi_0(T, \mathbf{X})$  in Theorem 1 can also be stated as that for any function  $g(T, \mathbf{X})$ :

$$\mathbb{E}[\pi_0(T, \mathbf{X})g(T, \mathbf{X})] = \int \int g(t, \mathbf{x}) f_T(t) dt f_X(\mathbf{x}) d\mathbf{x},$$

which is the expectation of  $g(T, \mathbf{X})$  taken with respect to the product measure  $f_T(t)f_X(\mathbf{x})dtd\mathbf{x}$ . In the case where  $g(T, \mathbf{X})$  is separable the resulting moment factorizes, that is,

$$\mathbb{E}[\pi_0(T, \mathbf{X})u(T)v(\mathbf{X})] = \mathbb{E}[u(T)] \mathbb{E}[v(\mathbf{X})].$$

Kernel estimators will not satisfy the sample version of this property but they will satisfy the smoothed empirical version, that is,

$$\int \hat{\pi}_0(t, \mathbf{x}) u(t) v(\mathbf{x}) dF_N(t, \mathbf{x}) \neq \int u(t) dF_N(t) \int v(\mathbf{x}) dF_N(\mathbf{x}),$$

where  $F_N(t, \mathbf{x})$  is the joint empirical measure and  $F_N(t)$  and  $F_N(\mathbf{x})$  are the marginals. However, it will satisfy

$$\int \hat{\pi}_0(t, \mathbf{x}) u(t) v(\mathbf{x}) dF_N^*(t, \mathbf{x}) = \int u(t) dF_N^*(t) \int v(\mathbf{x}) dF_N^*(\mathbf{x}),$$

where  $F_N^*(t, \mathbf{x})$  is the smoothed empirical (i.e.,  $dF_N^*(t, \mathbf{x})$  is the kernel density estimator used in constructing  $\hat{\pi}_0(t, \mathbf{x})$ ).

It is worth noting that our bound  $V_{eff}$  is equal to the bound of Hahn (1998) for the case of binary average treatment, the bound of Cattaneo (2010) for the case of multi-valued average treatment, and the bound of Firpo (2007) for the case of binary quantile treatment (see Ai, Linton, Motegi, and Zhang, 2019, Sections 2.2-2.4). Moreover, our bound applies to a much wider class of models, including quantile causal effect of multi-valued, continuous and mixture of discrete and continuous treatment as well as the asymmetric least squares estimation of the causal effect of all kinds of treatments.

Based on the expression of the efficient influence function, many existing literature construct the efficient estimator by solving the estimated efficient score equation (Athey,

Imbens, Pham, and Wager, 2017, Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins, 2018). Such kind of estimators typically have double or multiple robustness property. However, our the efficient influence function  $S_{eff}(T, \mathbf{X}, Y; \beta)$  involves five unknown functionals  $f_T(T)$ ,  $f_{T|X}(T|\mathbf{X})$ ,  $\varepsilon(T, \mathbf{X}; \beta)$ ,  $\mathbb{E}[\pi_0(T, \mathbf{X})\varepsilon(T, \mathbf{X}; \beta)m(T, \beta)|T]$ , and  $\mathbb{E}[\pi_0(T, \mathbf{X})\varepsilon(T, \mathbf{X}; \beta)m(T, \beta)|\mathbf{X}]$ . Estimation of these functionals are quite difficult in practice, and the performance of estimated  $\beta_0$  would be poor. Instead of explicitly estimating the efficient influence function  $S_{eff}$ , we propose a simple weighted optimization estimator based on (2.3) by estimating the stabilized weights  $\pi_0(T, \mathbf{X})$ .

It is also worth noting that, if the stabilized weights are known and  $g(t; \beta_0)$  is correctly specified, one can estimate  $\beta_0$  by solving the sample analogue of the weighted optimization (2.3). The asymptotic variance of the estimator is

$$V_{ineff} = \mathbb{E} [S_{ineff}(Y, T, \mathbf{X}; \beta_0) S_{ineff}(Y, T, \mathbf{X}; \beta_0)^\top],$$

with

$$S_{ineff}(Y, T, \mathbf{X}; \beta_0) = H_0^{-1} \cdot \pi_0(T, \mathbf{X}) m(T; \beta_0) L' \{Y - g(T; \beta_0)\}.$$

It is easy to show that  $V_{ineff} > V_{eff}$  (see Proposition B.1 of Appendix B), implying that the weighted optimization estimator is not efficient. This follows because the weighted optimization does not account for the restriction on the stabilized weight  $\pi_0(t, \mathbf{x})$ :

$$\mathbb{E} [\pi_0(T, \mathbf{X}) u(T) v(\mathbf{X})] = \mathbb{E}[u(T)] \cdot \mathbb{E}[v(\mathbf{X})] \quad (3.1)$$

holds for any suitable functions  $u(t)$  and  $v(\mathbf{x})$ . Incorporating restriction (3.1) into the estimation of the causal effect can improve efficiency. Similar observation is also noted by Hirano, Imbens, and Ridder (2003) in the binary treatment. Exactly how to incorporate restriction (3.1) into the estimation is the subject of the next section.

## 4 Efficient estimation

One way to incorporate (3.1) into the estimation is to estimate the stabilized weights from (3.1) and then implement (2.3) with the estimated weights. But before doing so, we must verify that (3.1) uniquely identifies  $\pi_0(T, \mathbf{X})$ . After some manipulations, it is straightforward to show that

$$\mathbb{E} [\pi(T, \mathbf{X}) u(T) v(\mathbf{X})] = \mathbb{E}[u(T)] \cdot \mathbb{E}[v(\mathbf{X})]$$

holds for any suitable functions  $u(t)$  and  $v(\mathbf{x})$  if and only if  $\pi(T, \mathbf{X}) = \pi_0(T, \mathbf{X})$ . Therefore, condition (3.1) identifies the stabilized weights. The challenge now is that (3.1) implies infinite number of equations. With a finite sample of observations, it is impossible

to solve infinite number of equations. To overcome this difficulty, we approximate the (infinite dimensional) function space with the (finite dimensional) sieve space. Specifically, let  $u_{K_1}(T) = (u_{K_1,1}(T), \dots, u_{K_1,K_1}(T))^\top$  and  $v_{K_2}(\mathbf{X}) = (v_{K_2,1}(\mathbf{X}), \dots, v_{K_2,K_2}(\mathbf{X}))^\top$  denote the known basis functions with dimensions  $K_1 \in \mathbb{N}$  and  $K_2 \in \mathbb{N}$  respectively, and let  $K := K_1 \cdot K_2$ . The functions  $u_{K_1}(t)$  and  $v_{K_2}(\mathbf{x})$  are called the *approximation sieves* that can approximate any suitable functions  $u(t)$  and  $v(\mathbf{x})$  arbitrarily well (see [Chen \(2007\)](#) for more discussion on sieve approximation). Since the sieve approximating space is also a subspace of the functional space,  $\pi_0(T, \mathbf{X})$  satisfies

$$\mathbb{E} [\pi_0(T, \mathbf{X}) u_{K_1}(T) v_{K_2}(\mathbf{X})^\top] = \mathbb{E}[u_{K_1}(T)] \cdot \mathbb{E}[v_{K_2}(\mathbf{X})]^\top. \quad (4.1)$$

Unfortunately, it is not the only solution. Indeed, for any monotonic increasing and globally concave function  $\rho(v)$ , with

$$\Lambda_{K_1 \times K_2}^* = \arg \max_{\Lambda \in \mathbb{R}^{K_1 \times K_2}} \mathbb{E} [\rho(u_{K_1}(T)^\top \Lambda v_{K_2}(\mathbf{X}))] - \mathbb{E}[u_{K_1}(T)]^\top \Lambda \mathbb{E}[v_{K_2}(\mathbf{X})], \quad (4.2)$$

$\pi_K^*(T, \mathbf{X}) = \rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}))$  also solves (4.1), where  $\rho'(v)$  denotes the first derivative. Let  $\pi_K(T, \mathbf{X}) = \rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X}))$  denote the best approximation of  $\pi_0(T, \mathbf{X})$  under the sup norm  $L_\infty$  and suppose that  $\|\pi_K(\cdot, \cdot) - \pi_0(\cdot, \cdot)\|_\infty = O(K^{-\alpha})$  for some  $\alpha > 0$ . Then,

$$\begin{aligned} & \|\pi_K^*(T, \mathbf{X}) - \pi_0(T, \mathbf{X})\|_{L^2} \\ &= O(\max \{ \|\pi_K^*(T, \mathbf{X}) - \pi_K(T, \mathbf{X})\|_{L^2}, \|\pi_K(\cdot, \cdot) - \pi_0(\cdot, \cdot)\|_\infty \}) \\ &= \max \{ O(K^{-\alpha}), O(K^{-\alpha}) \} = O(K^{-\alpha}). \end{aligned}$$

(see [Ai, Linton, Motegi, and Zhang, 2019](#), Lemma 3.1).

Let  $\{T_i, \mathbf{X}_i, Y_i\}_{i=1}^N$  denote an independently and identically distributed sample of observations drawn from the joint distribution of  $(T, \mathbf{X}, Y)$ . We propose to estimate the stabilized weights  $\pi_i = \pi_0(T_i, \mathbf{X}_i)$  by solving the entropy maximization problem:

$$\begin{cases} \max \left\{ -\sum_{i=1}^N \pi_i \log \pi_i \right\} \\ \text{subject to } \frac{1}{N} \sum_{i=1}^N \pi_i u_{K_1}(T_i) v_{K_2}(\mathbf{X}_i)^\top = \left( \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \right) \left( \frac{1}{N} \sum_{j=1}^N v_{K_2}(\mathbf{X}_j)^\top \right). \end{cases} \quad (4.3)$$

The primal problem (4.3) is difficult to compute. We instead consider its dual problem that can be solved by numerically efficient and stable algorithms. Specifically, let  $\rho(v) := -e^{-v-1}$  for any  $v \in \mathbb{R}$ , [Tseng and Bertsekas \(1991\)](#) showed that the dual solution is given by:

$$\hat{\pi}_K(T_i, \mathbf{X}_i) := \rho' \left( u_{K_1}(T_i)^\top \hat{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right),$$

where  $\hat{\Lambda}_{K_1 \times K_2}$  is the maximizer of the strictly concave function  $\hat{G}_{K_1 \times K_2}$  defined by

$$\hat{\Lambda}_{K_1 \times K_2} = \arg \max_{\Lambda} \hat{G}_{K_1 \times K_2}(\Lambda) := \frac{1}{N} \sum_{i=1}^N \rho(u_{K_1}(T_i)^\top \Lambda v_{K_2}(\mathbf{X}_i)) - \left( \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \right)^\top \Lambda \left( \frac{1}{N} \sum_{j=1}^N v_{K_2}(\mathbf{X}_j) \right). \quad (4.4)$$

The duality between (4.3) and (4.4) is shown in Appendix C. By [Ai, Linton, Motegi, and Zhang \(2019, Corollary 3.3\)](#), we have

$$\int_{\mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})|^2 dF_{T, \mathbf{X}}(t, \mathbf{x}) = O_p \left( \sqrt{\frac{K}{N}} \right).$$

Having estimated the weights, we now estimate  $\beta_0$  by applying the generalized optimization:

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L(Y_i - g(T_i; \beta)). \quad (4.5)$$

**Remarks:**

1. Alternatively, one can estimate the stabilized weights by estimating the generalized propensity score function as well as the marginal distribution of the treatment variable nonparametrically (e.g., kernel estimation). But these alternatively estimated weights do not satisfy empirical moment in (4.3) and may not result in efficient estimation of the causal effect.
2. The primal problem (4.3) is different from the empirical likelihood ([Smith, 1997, Imbens, 2002](#)). Notice that  $\rho(v) = -e^{-v-1}$  satisfies the invariance property (i.e.,  $-\rho''(v) = \rho'(v)$ ). It turns out that this invariance property is critical for establishing consistency of the generalized optimization estimator. Any other choice of  $\rho(\cdot)$  that does not have the invariance property may result in biased causal effect estimate.
3. The proposed estimation (4.5) is a semiparametric estimation problem that contains both finite dimensional and infinite unknown parameters. The general semiparametric estimation problems have been studied by [Ai and Chen \(2003\)](#) and [Chen, Linton, and Van Keilegom \(2003\)](#). [Ai and Chen \(2003\)](#) study the large sample properties under smooth objective functions, and [Chen, Linton, and Van Keilegom \(2003\)](#) extend those to nonsmooth criterion functions. (4.5) is a special case of the general setting of [Chen, Linton, and Van Keilegom \(2003\)](#). Indeed, we will apply their result (e.g. Theorem 2 of [Chen, Linton, and Van Keilegom \(2003\)](#), page 1594) to derive the asymptotic properties of  $\hat{\beta}$ . However, there is a major difference. Our focus is on the

efficiency bound derivation and efficient estimation, whereas their focus is on deriving the asymptotic properties of the sequential estimator under high level conditions (e.g. Condition 2.6, page 1594). These high level conditions are nontrivial to verify. Most of our derivations are indeed verifying those high level conditions, see Section 4.2 of the supplemental material [Ai, Linton, Motegi, and Zhang \(2019\)](#).

## 5 Large sample properties

To establish the large sample properties of the generalized optimization estimator, we first show that the estimated weight function  $\hat{\pi}_K(t, \mathbf{x})$  is consistent and compute its convergence rates under both  $L_\infty$  norm and the  $L_2$  norm. The following conditions shall be imposed.

**Assumption 2.** (i) The support  $\mathcal{X}$  of  $\mathbf{X}$  is a compact subset of  $\mathbb{R}^r$ . The support  $\mathcal{T}$  of the treatment variable  $T$  is a compact subset of  $\mathbb{R}$ . (ii) There exist two positive constants  $\eta_1$  and  $\eta_2$  such that

$$0 < \eta_1 \leq \pi_0(t, \mathbf{x}) \leq \eta_2 < \infty, \forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}.$$

**Assumption 3.** There exist  $\Lambda_{K_1 \times K_2} \in \mathbb{R}^{K_1 \times K_2}$  and a positive constant  $\alpha > 0$  such that

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |(\rho'^{-1}(\pi_0(t, \mathbf{x})) - u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}))| = O(K^{-\alpha}).$$

**Assumption 4.** (i) For every  $K_1$  and  $K_2$ , the smallest eigenvalues of  $\mathbb{E}[u_{K_1}(T)u_{K_1}(T)^\top]$  and  $\mathbb{E}[v_{K_2}(\mathbf{X})v_{K_2}(\mathbf{X})^\top]$  are bounded away from zero uniformly in  $K_1$  and  $K_2$ . (ii) There are two sequences of constants  $\zeta_1(K_1)$  and  $\zeta_2(K_2)$  satisfying  $\sup_{t \in \mathcal{T}} \|u_{K_1}(t)\| \leq \zeta_1(K_1)$  and  $\sup_{\mathbf{x} \in \mathcal{X}} \|v_{K_2}(\mathbf{x})\| \leq \zeta_2(K_2)$ ,  $K = K_1(N)K_2(N)$  and  $\zeta(K) := \zeta_1(K_1)\zeta_2(K_2)$ , such that  $\zeta(K)K^{-\alpha} \rightarrow 0$  and  $\zeta(K)\sqrt{K/N} \rightarrow 0$  as  $N \rightarrow \infty$ .

Assumption 2 (i) restricts both the covariates and treatment level to be bounded. This condition is restrictive but convenient for computing the convergence rate under  $L_\infty$  norm. It is commonly imposed in the nonparametric regression literature. This condition can be relaxed, however, if we restrict the tail distribution of  $(\mathbf{X}, T)$ . Assumption 2 (ii) restricts the weight function to be bounded and bounded away from zero. Given Assumption 2 (i), this condition is equivalent to  $dF_{T|\mathbf{X}}(T|\mathbf{X})$  being bounded away from zero, meaning that each type of individuals (denoted by  $\mathbf{X}$ ) always have a sufficient portion participating in each level of treatment. This restriction is important for our analysis since each individual participates only in one level of treatment and this condition allows us to construct her statistical counterparts from all other treatments. Although Assumption 2 (ii) is useful

in causal analysis and establishing the convergence rates, it is not essential and could be relaxed by allowing  $\eta_1$  (resp.  $\eta_2$ ) to depend on  $N$  and go to zero (resp. infinity) slowly, as  $N \rightarrow \infty$ . Notice that  $u_{K_1}(t)^\top \Lambda v_{K_2}(\mathbf{x})$  is a linear sieve approximation to any suitable function of  $(\mathbf{X}, T)$ . Assumption 3 requires the sieve approximation error of  $\rho'^{-1}(\pi_0(t, \mathbf{x}))$  to shrink at a polynomial rate. This condition is satisfied for a variety of sieve basis functions. For example, if both  $\mathbf{X}$  and  $T$  are discrete, then the approximation error is zero for sufficient large  $K$  and in this case Assumption 3 is satisfied with  $\alpha = +\infty$ . If some components of  $(\mathbf{X}, T)$  are continuous, the polynomial rate depends positively on the smoothness of  $\rho'^{-1}(\pi_0(t, \mathbf{x}))$  in continuous components and negatively on the number of the continuous components. We will show that the convergence rate of the estimated weight function (and consequently the rate of the generalized optimization estimator) is bounded by this polynomial rate. Assumption 4 (i) essentially ensures the sieve approximation estimator is non-degenerate. Similar condition is common in the sieve regression literature (see Andrews (1991) and Newey (1997)). If the approximation error is nonzero, Assumption 4 (ii) requires it to shrink to zero at an appropriate rate as sample size increases.

Under these conditions, we are able to establish the following theorem:

**Theorem 2.** *Suppose that Assumptions 2-4 hold. Then, we obtain the following:*

$$\int_{\mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_0(t, \mathbf{x})|^2 dF_{T, \mathbf{X}}(t, \mathbf{x}) = O_p \left( \max \left\{ K^{-2\alpha}, \frac{K}{N} \right\} \right),$$

$$\frac{1}{N} \sum_{i=1}^N |\hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_0(T_i, \mathbf{X}_i)|^2 = O_p \left( \max \left\{ K^{-2\alpha}, \frac{K}{N} \right\} \right).$$

The proof of Theorem 2 immediately follows from the supplemental material Ai, Linton, Motegi, and Zhang (2019, Lemma 3.1 & Corollary 3.3).

The following additional condition is needed to establish the consistency of the proposed estimator  $\hat{\beta}$ .

**Assumption 5.** (i) *The parameter space  $\Theta \subset \mathbb{R}^p$  is a compact set and the true parameter  $\beta_0$  is in the interior of  $\Theta$ , where  $p \in \mathbb{N}$ .* (ii)  $\mathbb{E} [\sup_{\beta \in \Theta} |L(Y - g(T; \beta))|^2] < \infty$ .

Assumption 5 (i) is commonly imposed in the nonlinear regression, but can be relaxed if  $g(t; \beta)$  is linear in  $\beta$ . Assumption 5 (ii) is an envelope condition that is sufficient for the applicability of the uniform law of large numbers.

Under these and other conditions, we establish the consistency of the generalized optimization estimator.

**Theorem 3.** *Suppose that Assumptions 1-5 hold. Then,  $\|\hat{\beta} - \beta_0\| \xrightarrow{p} 0$ .*

To establish the asymptotic distribution of the proposed estimator, we need some smoothness condition on the regression function and some under-smoothing condition on the sieve approximation (i.e., larger  $K$  than needed for consistency). We also have to address the possibility of a nonsmooth loss function. These conditions are presented below.

**Assumption 6.**

- (i) *The loss function  $L(v)$  is differentiable almost everywhere,  $g(t; \beta)$  is twice continuously differentiable in  $\beta \in \Theta$  and we denote its first derivative by  $m(t; \beta) := \nabla_{\beta} g(t; \beta)$ ;*
- (ii)  *$\mathbb{E} [\pi_0(T, \mathbf{X}) L'(Y - g(T; \beta)) m(T; \beta)]$  is differentiable with respect to  $\beta$  and  $H_0 := -\nabla_{\beta} \mathbb{E} [\pi_0(T, \mathbf{X}) L'(Y - g(T; \beta)) m(T; \beta)] \Big|_{\beta=\beta_0}$  is nonsingular;*
- (iii)  *$\varepsilon(t, \mathbf{x}; \beta_0) := \mathbb{E}[L'(Y - g(T; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}]$  is continuously differentiable in  $(t, \mathbf{x})$ ;*
- (iv) *Suppose that  $N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \hat{\beta})) m(T_i; \hat{\beta}) = o_p(N^{-1/2})$  holds with probability approaching one.*

**Assumption 7.** (i)  $\mathbb{E} [\sup_{\beta \in \Theta} |L'(Y - g(T; \beta))|^{2+\delta}] < \infty$  for some  $\delta > 0$ ; (ii) *The function class  $\{L'(y - g(t; \beta)) : \beta \in \Theta\}$  satisfies:*

$$\mathbb{E} \left[ \sup_{\beta_1: \|\beta_1 - \beta\| < \delta} |L'(Y - g(T; \beta_1)) - L'(Y - g(T; \beta))|^2 \right]^{1/2} \leq a \cdot \delta^b$$

for any  $\beta \in \Theta$  and any small  $\delta > 0$  and for some finite positive constants  $a$  and  $b$ .

Assumption 6 (i) imposes sufficient regularity conditions on both regression function and loss function. These conditions permit nonsmooth loss functions and are satisfied by the example loss functions mentioned in previous sections. Assumption 6 (ii) ensures that the efficient variance to be finite. Assumption 6 (iv) is essentially the first order condition, similar to the one imposed in  $Z$ -estimation. Again, this first order condition is satisfied by popular nonsmooth loss functions, see [Pakes and Pollard \(1989\)](#). Assumption 7 is a stochastic equicontinuity condition which is needed for establishing weak convergence, see [Andrews \(1994\)](#). Again, it is satisfied by widely used nonsmooth loss functions.

Under above sufficient conditions, we have the following theorem:



**Theorem 4.** Suppose that Assumptions 1-7 hold, and strengthen Assumption 4 (ii) to

$$\textbf{Assumption 4 (ii)'} \quad \zeta(K)\sqrt{K^2/N} \rightarrow 0 \text{ and } \sqrt{N}K^{-\alpha} \rightarrow 0.$$

Then,  $\sqrt{N}(\hat{\beta} - \beta_0) \xrightarrow{d} \mathcal{N}(0, V_{eff})$ , where  $V_{eff} = \mathbb{E}[S_{eff}(T, \mathbf{X}, Y; \beta_0)S_{eff}(T, \mathbf{X}, Y; \beta_0)^\top]$ . Therefore,  $\hat{\beta}$  attains the semi-parametric efficiency bound of Theorem 1.

Assumption 4 (ii)' imposes further restriction on the smoothing parameter ( $K$ ) so that the sieve approximation is under-smoothed. This condition is stronger than Assumption 4 (ii) but it is commonly imposed in the semiparametric regression literature. The proof of Theorem 4 is given in the supplemental material [Ai, Linton, Motegi, and Zhang \(2019, Section 4\)](#).

## 6 Variance estimation

In order to conduct statistical inference, a consistent covariance matrix is needed. Theorem 1 suggests that such consistent covariance can be obtained by replacing  $H_0$  and  $\psi(Y, T, \mathbf{X}; \beta_0)$  with some consistent estimates. Since the nonsmooth loss function may invalidate the exchangeability between the expectation and derivative operator, some care in the estimation of  $H_0$  is warranted. Using the tower property of conditional expectation, we rewrite  $H_0$  as:

$$\begin{aligned} H_0 &= -\nabla_\beta \mathbb{E}[\pi_0(T, \mathbf{X}) \mathbb{E}[L'(Y - g(T; \beta)) | T, \mathbf{X}] m(T; \beta)] \Big|_{\beta=\beta_0} \\ &= -\mathbb{E} \left[ \pi_0(T, \mathbf{X}) \nabla_\beta \mathbb{E}[L'(Y - g(T; \beta)) | T, \mathbf{X}] \Big|_{\beta=\beta_0} m(T; \beta_0)^\top \right] \\ &\quad - \mathbb{E}[\pi_0(T, \mathbf{X}) \mathbb{E}[L'(Y - g(T; \beta_0)) | T, \mathbf{X}] \nabla_\beta m(T; \beta_0)]. \end{aligned}$$

Applying integration by parts (see Appendix D), we obtain

$$\begin{aligned} &\nabla_\beta \mathbb{E}[L'(Y - g(T; \beta)) | T = t, \mathbf{X} = \mathbf{x}] \Big|_{\beta=\beta_0} \\ &= \mathbb{E} \left[ L'(Y - g(T; \beta_0)) \frac{\frac{\partial}{\partial y} f_{Y,T,\mathbf{X}}(Y, T, \mathbf{X})}{f_{Y,T,\mathbf{X}}(Y, T, \mathbf{X})} \Big| T = t, \mathbf{X} = \mathbf{x} \right] m(t; \beta_0) \end{aligned} \quad (6.1)$$

and consequently

$$H_0 = -\mathbb{E} \left[ \pi_0(T, \mathbf{X}) L'(Y - g(T; \beta_0)) \left\{ \frac{\frac{\partial}{\partial y} f_{Y,T,\mathbf{X}}(Y, T, \mathbf{X})}{f_{Y,T,\mathbf{X}}(Y, T, \mathbf{X})} m(T; \beta_0) m(T; \beta)^\top + \nabla_\beta m(T; \beta_0)^\top \right\} \right].$$

The density  $f_{Y,T,\mathbf{X}}(y, t, \mathbf{x})$  can be estimated via the widely used sieve extremum estimator (Chen, 2007, Example 2.6, page 5565):

$$\hat{f}_{Y,T,\mathbf{X}}(y, t, \mathbf{x}) := \frac{\exp(\hat{a}_{K_0}^\top r_{K_0}(y, t, \mathbf{x}))}{\int_{\mathcal{Y} \times \mathcal{T} \times \mathcal{X}} \exp(\hat{a}_{K_0}^\top r_{K_0}(y, t, \mathbf{x})) dy dt d\mathbf{x}},$$

where  $\hat{a}_{K_0} \in \mathbb{R}^{K_0}$  ( $K_0 \in \mathbb{N}$ ) maximizes the following concave objective function:

$$\hat{a}_{K_0} := \arg \max_{a \in \mathbb{R}^{K_0}} \frac{1}{N} \sum_{i=1}^N \left[ a^\top r_{K_0}(Y_i, T_i, \mathbf{X}_i) - \log \int_{\mathcal{Y} \times \mathcal{T} \times \mathcal{X}} \exp(a^\top r_{K_0}(y, t, \mathbf{x})) dy dt d\mathbf{x} \right],$$

and  $r_{K_0}(t, y, \mathbf{x})$  is a  $K_0$ -dimensional sieve basis. Then  $H_0$  can be estimated by

$$\hat{H} := -\frac{1}{N} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \hat{\beta})) \left\{ \frac{\frac{\partial}{\partial y} \hat{f}_{Y,T,\mathbf{X}}(Y_i, T_i, \mathbf{X}_i)}{\hat{f}_{Y,T,\mathbf{X}}(Y_i, T_i, \mathbf{X}_i)} m(T_i; \hat{\beta}) m(T_i; \hat{\beta})^\top + \nabla_{\beta} m(T_i; \hat{\beta}) \right\}.$$

Also,  $\psi(Y, T, \mathbf{X}; \beta_0)$  can be directly estimated by the sieve estimator:

$$\begin{aligned} \hat{\psi}(Y, T, \mathbf{X}; \hat{\beta}) &= \hat{\pi}_K(T, \mathbf{X}) L'(Y - g(T; \hat{\beta})) m(T; \hat{\beta}) - \hat{\pi}_K(t, \mathbf{x}) \hat{\mathbb{E}} \left[ L'(Y - g(T; \hat{\beta})) | T, \mathbf{X} \right] m(T; \hat{\beta}) \\ &\quad + \hat{\mathbb{E}} \left[ \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y - g(T; \hat{\beta})) | T \right] m(T; \hat{\beta}) + \hat{\mathbb{E}} \left[ \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y - g(T; \hat{\beta})) | \mathbf{X} \right] m(T; \hat{\beta}), \end{aligned}$$

where

$$\begin{aligned} \hat{\mathbb{E}} \left[ L'(Y - g(T; \hat{\beta})) | T, \mathbf{X} \right] &:= \left[ \sum_{i=1}^N L'(Y_i - g(T_i; \hat{\beta})) w_{K_0}(T_i, \mathbf{X}_i)^\top \right] \left[ \sum_{i=1}^N w_{K_0}(T_i, \mathbf{X}_i) w_{K_0}(T_i, \mathbf{X}_i)^\top \right]^{-1} w_{K_0}(T, \mathbf{X}), \\ \hat{\mathbb{E}} \left[ \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y - g(T; \hat{\beta})) | T \right] &:= \left[ \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \hat{\beta})) u_{K_0}(T_i)^\top \right] \left[ \sum_{i=1}^N u_{K_0}(T_i) u_{K_0}(T_i)^\top \right]^{-1} u_{K_0}(T), \\ \hat{\mathbb{E}} \left[ \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y - g(T; \hat{\beta})) | \mathbf{X} \right] &:= \left[ \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \hat{\beta})) v_{K_0}(\mathbf{X}_i)^\top \right] \left[ \sum_{i=1}^N v_{K_0}(\mathbf{X}_i) v_{K_0}(\mathbf{X}_i)^\top \right]^{-1} v_{K_0}(\mathbf{X}). \end{aligned}$$

The asymptotic covariance is estimated by

$$\hat{V} := \hat{H}^{-1} \left\{ \frac{1}{N} \sum_{i=1}^N \hat{\psi}(Y_i, T_i, \mathbf{X}_i; \hat{\beta}) \hat{\psi}(Y_i, T_i, \mathbf{X}_i; \hat{\beta})^\top \right\} (\hat{H}^\top)^{-1}.$$

It is well known that the sieve extreme estimation is uniformly strong consistent (in almost surely sense), see Chen (2007, Theorem 3.1). Also from Theorems 2 and 3, we have  $\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_0(t, \mathbf{x})| = o_p(1)$  and  $\|\hat{\beta} - \beta_0\| \rightarrow 0$ . With these results, we obtain the consistency of  $\hat{V}$ .

**Theorem 5.** Suppose that Assumptions 1-5 hold. Then,  $\hat{V}$  converges to  $V_{eff}$  in probability.

## 7 Selection of tuning parameters

The large sample properties of the proposed estimator permit a wide range of values of  $K_1$  and  $K_2$ . This presents a dilemma for applied researchers who have only one finite sample and would like to have some guidance on the selection of smoothing parameters. In this section, we present two data-driven approaches to select  $K_1$  and  $K_2$ . The first one is simply minimizing a (penalized) loss function. Define

$$\bar{L}(K_1, K_2) = \frac{1}{N} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L(Y_i - g(T_i; \hat{\beta})).$$

There are several ways to penalize using large  $K_1$  or  $K_2$ :

**No penalty.**  $\mathcal{L}(K_1, K_2) = \bar{L}(K_1, K_2)$ .

**Additive penalty.**  $\mathcal{L}(K_1, K_2) = (1 + 2(K_1 + K_2)/N) \times \bar{L}(K_1, K_2)$ .

**Multiplicative penalty.**  $\mathcal{L}(K_1, K_2) = (1 + 2K_1K_2/N) \times \bar{L}(K_1, K_2)$ .

Choose  $(K_1^*, K_2^*)$  that minimizes  $\mathcal{L}(K_1, K_2)$  in some choice sets  $(K_1, K_2) \in \mathbb{K}_1 \times \mathbb{K}_2$ .

The second approach is the  $J$ -folder cross validation which proceeds as follows.

1. We divide  $N$  samples into  $J$  groups, (say  $J = 5$  or  $10$ ), and let  $n = N/J$ . The data in the  $j^{th}$  group is denoted by  $S_j = \{\mathbf{X}_i^{(j)}, T_i^{(j)}, Y_i^{(j)} : i = 1, \dots, n\}$  for  $j \in \{1, \dots, J\}$ .
2. For each  $j \in \{1, \dots, J\}$ , we denote the dataset  $S_{(-j)} = \{\mathbf{X}_i, T_i, Y_i\}_{i=1}^N / S_j$ . We compute the following quantities based on  $S_{(-j)}$ :

$$\begin{aligned} \hat{\Lambda}_{K_1 \times K_2}^{(-j)} &= \arg \max_{\Lambda} \hat{G}_K^{(-j)}(\Lambda) \\ &= \frac{1}{N-n} \sum_{i \in S_{(-j)}} \rho(u_{K_1}^\top(T_i) \Lambda v_{K_2}(\mathbf{X}_i)) - \left[ \frac{1}{N-n} \sum_{i \in S_{(-j)}} u_{K_1}^\top(T_i) \right] \Lambda \left[ \frac{1}{N-n} \sum_{i \in S_{(-j)}} v_{K_2}(\mathbf{X}_i) \right] \end{aligned}$$

and

$$\hat{\pi}_K^{(-j)}(T, \mathbf{X}) = \rho' \left( u_{K_1}^\top(T) \hat{\Lambda}_{K_1 \times K_2}^{(-j)} v_{K_2}(\mathbf{X}) \right)$$

and

$$\hat{\beta}_K^{(-j)} = \arg \min \sum_{i \in S_{(-j)}} \hat{\pi}_K^{(-j)}(T_i, \mathbf{X}_i) \{Y_i - g(T_i; \beta)\}^2.$$

3. We choose optimal  $K_1$  and  $K_2$  so that the following cross validation criterion is minimized:

$$CV(K_1, K_2) = \sum_{j=1}^J \left[ \sum_{k \in S_j} \hat{\pi}_K^{(-j)}(T_k, \mathbf{X}_k) \left\{ Y_k - g\left(T_k; \hat{\beta}_K^{(-j)}\right) \right\}^2 \right].$$

## 8 Some extensions

The condition (2.1) that the causal effect is parameterized may be restrictive for some applications. To relax this condition, we can consider the nonparametric specification:

$$\min_{g(\cdot)} \int_{\mathcal{T}} \mathbb{E} [L(f(Y^*(t)) - g(t))] dF_T(t),$$

where  $f(\cdot)$  is a known function while  $g(\cdot)$  is unknown. Under Assumption 1, above optimization is equivalent to

$$\min_{g(\cdot)} \mathbb{E} [\pi_0(T, \mathbf{X}) L(f(Y) - g(T))].$$

We can estimate  $g(\cdot)$  through the weighted nonparametric sieve regression:

$$\min_{g(\cdot) \in \mathcal{H}_{K_1}} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L(f(Y_i) - g(T_i)),$$

where  $\mathcal{H}_{K_1} := \{g(\cdot) : \mathcal{T} \rightarrow \mathbb{R}, g(t) = \lambda^\top u_{K_1}(t) : \lambda \in \mathbb{R}^{K_1}\}$  is a specified sieve space. Extension to the general loss function requires considerable derivation and shall be dealt with in a separate paper. In this section, we only consider three particular cases: (i) the treatment effect curve  $\theta(t) := \mathbb{E}[Y^*(t)]$ , which corresponds to  $L(v) = v^2$  and  $f(Y) = Y$ ; (ii) the average effect  $\psi := \int_{\mathcal{T}} \mathbb{E}[\theta(t)] dF_T(t)$  proposed by Kennedy, Ma, McHugh, and Small (2017).

### 8.1 Estimation of effect curve

We begin with estimation of  $\theta(t)$ . Note that, for all  $t \in \mathcal{T}$  and under Assumption 1, we can rewrite  $\theta(t)$  as

$$\theta(t) = \mathbb{E}[Y^*(t)] = \mathbb{E}[\pi_0(T, \mathbf{X}) Y | T = t].$$

With  $\pi_0(T, \mathbf{X})$  replaced by  $\hat{\pi}_K(T, \mathbf{X})$ , we estimate  $\theta(t)$  by regressing  $\hat{\pi}_K(T, \mathbf{X})Y$  on  $u_{K_1}(T)$ :

$$\hat{\theta}_K(t) = \left[ \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) Y_i u_{K_1}(T_i)^\top \right] \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} u_{K_1}(t).$$

To aid presentation of the asymptotic properties of  $\hat{\theta}_K(t)$ , we denote

$$\varepsilon_i = \pi_0(T_i, \mathbf{X}_i) Y_i - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i, \mathbf{X}_i] + \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | \mathbf{X}_i] - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i],$$

$$\sigma^2(T_i) = \text{Var}(\varepsilon_i | T_i), \quad \Sigma_{K_1 \times K_1} = \mathbb{E}[u_{K_1}(T) u_{K_1}(T)^\top \sigma^2(T)], \quad \Phi_{K_1 \times K_1} = \mathbb{E}[u_{K_1}(T) u_{K_1}(T)^\top],$$

and

$$V_K(t) = u_{K_1}^\top(t) \Phi_{K_1 \times K_1}^{-1} \Sigma_{K_1 \times K_1} \Phi_{K_1 \times K_1}^{-1} u_{K_1}(t).$$

**Theorem 6.** Suppose that for some  $\tilde{\alpha} > 0$ , there exists some  $\gamma^* \in \mathbb{R}^{K_1}$  such that  $\sup_{t \in \mathcal{T}} |\theta(t) - (\gamma^*)^\top u_{K_1}(t)| = O(K_1^{-\tilde{\alpha}})$  and  $\sigma^2(t)$  is bounded, and that Assumptions 1-6 hold. Then:

1. (Consistency)

$$\int_{\mathcal{T}} |\hat{\theta}_K(t) - \theta(t)|^2 dF_T(t) = O_p \left( \frac{\zeta(K)^2 K}{N} + \zeta(K)^2 K^{-2\alpha} + K_1^{-2\tilde{\alpha}} \right)$$

and

$$\sup_{t \in \mathcal{T}} |\hat{\theta}_K(t) - \theta(t)| = O_p \left[ \zeta_1(K_1) \left( \zeta(K) \sqrt{K/N} + \zeta(K) K^{-2\alpha} + K_1^{-\tilde{\alpha}} \right) \right].$$

2. (Asymptotic Normality) suppose  $\sqrt{N} K^{-\tilde{\alpha}} \rightarrow 0$  for any fixed  $t \in \mathcal{T}$ ,

$$\sqrt{N} V_K(t)^{-1/2} \left[ \hat{\theta}_K(t) - \theta(t) \right] \xrightarrow{d} N(0, 1).$$

See [Ai, Linton, Motegi, and Zhang \(2019, Section 5.1\)](#) for a proof of Theorem 6.

## 8.2 Efficient estimation of average effect

To estimate the average effect  $\psi$ , we notice that

$$\psi = \mathbb{E}[\pi_0(T, \mathbf{X}) Y].$$

Hence, we estimate  $\psi$  by

$$\hat{\psi}_K = \frac{1}{N} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) Y_i.$$

The following theorem establishes the asymptotic distribution of  $\hat{\psi}_K$  and shows that  $\hat{\psi}_K$  attains the efficiency bound of  $\psi$  derived by [Kennedy, Ma, McHugh, and Small \(2017\)](#).

**Theorem 7.** *Under all assumptions stated in Theorem 4. Then:*

1. (Consistency)  $\hat{\psi}_K \xrightarrow{p} \psi$ ;
2. (Asymptotic Efficiency)  $\sqrt{N}(\hat{\psi}_K - \psi) \xrightarrow{d} N(0, V_{eff}^\psi)$ , where  $V_{eff}^\psi = \mathbb{E} \left[ \left( S_{eff}^\psi \right)^2 \right]$   
and  

$$S_{eff}^\psi(T, \mathbf{X}, Y) = \pi_0(T, \mathbf{X}) \{Y - \mathbb{E}[Y|\mathbf{X}, T]\} + \{\mathbb{E}[\pi_0(T, \mathbf{X})Y|\mathbf{X}] - \mathbb{E}[\pi_0(T, \mathbf{X})Y]\} \\ + \{\mathbb{E}[\pi_0(T, \mathbf{X})Y|T] - \mathbb{E}[\pi_0(T, \mathbf{X})Y]\}.$$

See [Ai, Linton, Motegi, and Zhang \(2019, Section 5.2\)](#) for a proof of Theorem 7.

## 9 Monte Carlo simulations

The large sample properties established in previous sections do not indicate how the generalized optimization estimator behaves in finite samples. To evaluate its finite sample performance, we conduct a simulation study on a continuous treatment. We present a simulation design in Section 9.1 and results in Section 9.2.

### 9.1 Simulation design

Let  $\mathbf{X}_i = (X_{1i}, X_{2i})^\top$  be covariates, and assume that  $\mathbf{X}_i \stackrel{i.i.d.}{\sim} N(0, I_2)$ . Error terms are drawn mutually independently as  $\xi_i \stackrel{i.i.d.}{\sim} N(0, 1)$  and  $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, 1)$ . Consider four data generating processes (DGPs):

**DGP-L1**  $T = 1 + 0.2X_1 + \xi$  and  $Y = 1 + X_1 + T + \epsilon$ . ( $X_2$  does not play any role, and  $X_1$  affects  $T$  and  $Y$  linearly.)

**DGP-NL1**  $T = 0.1X_1^2 + \xi$  and  $Y = X_1^2 + T + \epsilon$ . ( $X_2$  does not play any role, and  $X_1$  affects  $T$  and  $Y$  non-linearly.)

**DGP-L2**  $T = 1 + 0.2 \sum_{j=1}^2 X_j + \xi$  and  $Y = 1 + (1/2) \sum_{j=1}^2 X_j + T + \epsilon$ . ( $X_1$  and  $X_2$  affect  $T$  and  $Y$  linearly.)

**DGP-NL2**  $T = 0.1(\sum_{j=1}^2 X_j)^2 + \xi$  and  $Y = 1/2 + [(1/2) \sum_{j=1}^2 X_j]^2 + T + \epsilon$ . ( $X_1$  and  $X_2$  affect  $T$  and  $Y$  non-linearly.)

For each DGP, the true link function is  $\mathbb{E}[Y(t)] = 1 + t$ , a simple linear function with  $\beta_1^* = \beta_2^* = 1$ . Below we use a linear link function  $g(T_i; \beta) = \beta_1 + \beta_2 T_i$ , compute the generalized optimization estimator  $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2)^\top$ , and examine its performance.

To compute the generalized optimization estimator, two approximating basis functions  $u_{K_1}(T)$  and  $v_{K_2}(\mathbf{X})$  need to be specified. For  $u_{K_1}(T)$ ,  $K_1 \in \{2, 3, 4\} \equiv \mathbb{K}_1$  is considered:

$$u_2(T) = (1, T)^\top, \quad u_3(T) = (1, T, T^2)^\top, \quad u_4(T) = (1, T, T^2, T^3)^\top.$$

For  $v_{K_2}(\mathbf{X})$ , the choice set  $\mathbb{K}_2$  depends on the number of covariates. For DGP-L1 and DGP-NL1,  $K_2 \in \{2, 3, 4\} \equiv \mathbb{K}_2^1$  is considered:

$$v_2(X_1) = (1, X_1)^\top, \quad v_3(X_1) = (1, X_1, X_1^2)^\top, \quad v_4(X_1) = (1, X_1, X_1^2, X_1^3)^\top. \quad (9.1)$$

For DGP-L2 and DGP-NL2,  $K_2 \in \{3, 6, 10\} \equiv \mathbb{K}_2^2$  is considered:

$$\begin{aligned} v_3(\mathbf{X}) &= (1, X_1, X_2)^\top, \\ v_6(\mathbf{X}) &= (1, X_1, X_2, X_1^2, X_2^2, X_1 X_2)^\top, \\ v_{10}(\mathbf{X}) &= (1, X_1, X_2, X_1^2, X_2^2, X_1 X_2, X_1^3, X_2^3, X_1^2 X_2, X_1 X_2^2)^\top. \end{aligned} \quad (9.2)$$

Besides fixed pairs of  $(K_1, K_2) \in \mathbb{K}_1 \times \mathbb{K}_2$ , the data-driven selections described in Section 7 are employed. First, the (penalized) loss function approaches are implemented with  $L(Y_i - g(T_i; \hat{\beta})) = (Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . Second, the  $J$ -folder cross validation is implemented with  $J \in \{5, 10\}$ .

We also compute the covariate balancing generalized propensity score estimator proposed by Fong, Hazlett, and Imai (2018) under the quadratic loss function. Fong, Hazlett, and Imai (2018) use a linear model specification. Hence, their specification is correct under DGP-L1 and DGP-L2 while it is incorrect under DGP-NL1 and DGP-NL2.

Our proposed estimator and Fong, Hazlett, and Imai's (2018) estimator are computed in a simulated sample with size  $N \in \{100, 500, 1000\}$ , after which another sample is generated and both estimators are computed again. This exercise is repeated  $M = 1000$  times.

To evaluate the performance of point estimation, the bias, standard deviation, and root mean squared error (RMSE) of  $\hat{\beta}_1$  and  $\hat{\beta}_2$  are calculated from (a subset of)  $M = 1000$  simulations. In a small portion of the  $M = 1000$  samples,  $\bar{\pi}_N \equiv (1/N) \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)$ , which should be equal to 1 in theory, takes a value far from 1 due to numerical instability in the computation of  $\Lambda_{K_1 \times K_2}^*$ . The maximization with respect to  $\Lambda$  should lead to a global maximizer  $\Lambda_{K_1 \times K_2}^*$  in theory, but optimizing the  $K_1 K_2$  elements of  $\Lambda$  all at once is often hard in practical computation. Hence, we calculate the bias, standard deviation, and RMSE



from Monte Carlo samples such that  $\bar{\pi}_N \in [0.5, 2]$ . Other few samples having  $\bar{\pi}_N \notin [0.5, 2]$  are simply discarded.

We also evaluate the performance of variance estimation. The true covariance matrix of  $\hat{\beta}$  is written as

$$V_{eff} = \begin{bmatrix} V_{11} & V_{12} \\ V_{12} & V_{22} \end{bmatrix}.$$

Different DGPs have different true values of  $(V_{11}, V_{12}, V_{22})$ , and they are computed in Section 6.2 of the supplemental material [Ai, Linton, Motegi, and Zhang \(2019\)](#). For each DGP, we compute  $\hat{\beta}$  based on  $(K_1, K_2)$  that leads to sharp point estimation. Then we use other sieve bases of dimension  $(K'_1, K'_2) \in \mathbb{K}_1 \times \mathbb{K}_2$  to re-estimate the propensity score  $\pi_{K'}(T, \mathbf{X})$ . We allow  $(K_1, K_2)$  and  $(K'_1, K'_2)$  to be different from each other since  $\hat{\pi}_K(T, \mathbf{X})$  that leads to sharp point estimation might be different from  $\hat{\pi}_{K'}(T, \mathbf{X})$  that leads to sharp variance estimation.

Using  $\hat{\pi}_{K'}(T, \mathbf{X})$  and variance-specific sieve bases  $v_{M_0}(\mathbf{X})$  and  $w_{K_0}(T, \mathbf{X})$ , the variance estimator  $\hat{V}_{eff}$  is computed. For  $v_{M_0}(\mathbf{X})$ ,  $M_0 \in \{2, 3\}$  is used for DGP-L1 and DGP-NL1 and  $M_0 \in \{3, 6\}$  is used for DGP-L2 and DGP-NL2 (see (9.1) and (9.2)). For  $w_{K_0}(T, \mathbf{X})$ ,  $K_0 \in \{3, 5\}$  is used for DGP-L1 and DGP-NL1:

$$w_3(T, X_1) = (1, T, X_1)^\top, \quad w_5(T, X_1) = (1, T, X_1, T^2, X_1^2)^\top.$$

$K_0 \in \{4, 8\}$  is used for DGP-L2 and DGP-NL2:

$$w_4(T, \mathbf{X}) = (1, T, X_1, X_2)^\top, \quad w_8(T, \mathbf{X}) = (1, T, X_1, X_2, T^2, X_1^2, X_2^2, TX_1)^\top.$$

Data-driven selection of  $(K'_1, K'_2, M_0, K_0)$  is beyond the scope of the present paper.

## 9.2 Simulation results

We discuss point estimation first, and then discuss variance estimation. See Tables 1-2 for point estimation results on DGP-L1; Tables 3-4 for DGP-NL1; Tables 5-6 for DGP-L2; Tables 7-8 for DGP-NL2. We also draw bar charts that depict the share of  $(K_1, K_2)$  selected by each data-driven method. See Figures 1-2 for the bar charts for DGP-L1; Figures 3-4 for DGP-NL1; Figures 5-6 for DGP-L2; Figures 7-8 for DGP-NL2.

Unless otherwise noted, we discuss results on the causal parameter  $\beta_2$  only, since it is economically more important than the constant  $\beta_1$ . Under DGP-L1, the generalized optimization estimator (labeled as GOE) has small enough RMSE for any fixed  $(K_1, K_2)$  (Table 2). It is not a surprising result since DGP-L1 has a simple linear structure. The

data-driven methods often choose  $(K_1^*, K_2^*) = (2, 2)$ , the simplest possible approximation basis (Figures 1-2). The RMSE of the covariate balancing generalized propensity score estimator (labeled as CBGPS) is even smaller than the RMSE of GOE. It is not surprising since CBGPS has a correct parametric specification under DGP-L1.

Under DGP-NL1, GOE dominates CBGPS. GOE leads to small enough RMSE for both  $\beta_1$  and  $\beta_2$  as long as  $K_2 \geq 3$ . The relatively large RMSE under  $K_2 = 2$  suggests that  $X_1^2$  needs to be included in  $v_{K_2}(X_1)$  (see (9.1)). That is a reasonable result since DGP-NL1 has a quadratic structure. As desired, any data-driven method considered often selects pairs with  $K_2 \geq 3$  (Figures 3-4). CBGPS, in contrast, fails with the bias for  $\beta_2$  being around 0.2. The bias arises because the linear specification of CBGPS is incorrect under DGP-NL1. This result highlights that GOE performs well for both linear and nonlinear scenarios while CBGPS performs well for linear scenarios only.

The two-covariate scenarios yield similar implications to the single-covariate scenarios. Under DGP-L2, GOE with any fixed  $(K_1, K_2)$  has small RMSE (Table 6). The data-driven methods often choose  $(K_1^*, K_2^*) = (2, 3)$ , the simplest possible approximation basis (Figures 5-6). The RMSE of CBGPS is even smaller than the RMSE of GOE due to the linear structures of DGP-L2.

Under DGP-NL2, GOE with  $K_2 \geq 6$  leads to small RMSE, and any data-driven method considered often selects pairs with  $K_2 \geq 6$  as desired (Tables 7-8 and Figures 7-8). CBGPS, in contrast, fails with substantial bias of around 0.17 for  $\beta_2$ . This result again highlights the remarkable advantage of GOE relative to CBGPS.

Summarizing point estimation, the generalized optimization estimator performs well in finite samples, and the performance is still good even when the true DGP is nonlinear. In contrast, the existing alternative estimator of Fong, Hazlett, and Imai (2018) is sensitive to model misspecification.

We now discuss variance estimation results. The values of true covariance matrix,  $V_{eff}$ , are also provided in Tables 9-12. See Ai, Linton, Motegi, and Zhang (2019, Section 6) for how to compute the true values. For each DGP, we compute  $\hat{\beta}$  via  $(K_1, K_2) = (2, 2)$  for DGP-L1,  $(2, 3)$  for DGP-NL1,  $(2, 3)$  for DGP-L2, and  $(2, 6)$  for DGP-NL2. Recall from Tables 1-8 that those values are optimal values that lead to smallest MSEs in point estimation. Then we present in Tables 9-12 the bias, standard deviation, and RMSE of  $\hat{V}_{eff}$  with respect to  $V_{eff}$ , where  $(K'_1, K'_2, M_0, K_0) = (3, 3, 3, 5)$  for DGP-L1,  $(3, 3, 3, 5)$  for DGP-NL1,  $(3, 3, 6, 8)$  for DGP-L2, and  $(2, 10, 3, 4)$  for DGP-NL2. Under those values, we observe desired results that  $\hat{V}_{eff}$  converges to  $V_{eff}$  as sample size  $N$  increases. When  $N = 1000$ , the bias and standard deviation are small enough. Under DGP-NL1 and DGP-NL2, CBGPS suffers from large bias in variance estimation (Tables 10 and 12). That is

reasonable since the point estimation is already biased (Tables 3-4 and 7-8).

## 10 Empirical application

To illustrate the applicability of the generalized optimization procedure, we revisit U.S. presidential campaign data analyzed by Urban and Niebler (2014) and Fong, Hazlett, and Imai (2018). The motivation of the original study, Urban and Niebler (2014), is well summarized in Fong, Hazlett, and Imai (2018, Section 2):

Urban and Niebler (2014) explored the potential causal link between advertising and campaign contributions. Presidential campaigns ordinarily focus their advertising efforts on competitive states, but if political advertising drives more donations, then it may be worthwhile for candidates to also advertise in non-competitive states. The authors exploit the fact that media markets sometimes cross state boundaries. This means that candidates may inadvertently advertise in noncompetitive states when they purchase advertisements for media markets that mainly serve competitive states. By restricting their analysis to noncompetitive states, the authors attempt to isolate the effect of advertising from that of other campaigning, which do not incur these media market spillovers.

The treatment of interest, the number of political advertisements aired in each zip code, can be regarded as a continuous variable since it takes a range of values from 0 to 22379 across  $N = 16265$  zip codes. Urban and Niebler (2014) restricted themselves to a binary treatment framework, and they dichotomized the treatment variable by examining whether a zip code received more than 1000 advertisements or not. Their empirical results suggest that advertising in non-competitive states had a significant impact on the level of campaign contributions.

Dichotomizing a continuous treatment variable requires an ad-hoc choice of a cut-off value, and it makes an empirical result hard to interpret. Fong, Hazlett, and Imai (2018) analyzed the continuous version of the treatment variable, taking advantage of their proposed CBGPS method. Their empirical results suggest, contrary to Urban and Niebler (2014), that advertising in non-competitive states did *not* have a significant impact on the level of campaign contributions (cf. Fong, Hazlett, and Imai, 2018, Table 2).

As shown in Section 9, our generalized optimization estimator has a better performance than Fong, Hazlett, and Imai’s (2018) CBGPS estimator. Our estimator exhibits a solid performance even if a DGP of treatment  $T_i$  or outcome  $Y_i$  is nonlinear in covariate  $\mathbf{X}_i$ . It is

thus of interest to apply our approach to the continuous version of the treatment variable in order to see how results change.

## 10.1 Fong, Hazlett, and Imai’s (2018) CBGPS approach

We begin with Fong, Hazlett, and Imai’s (2018) CBGPS estimator as a benchmark. It requires a choice of pre-treatment covariates  $\mathbf{X}_i$  in a generalized propensity score model. There are eight covariates

$$\mathbf{X}_1 = \begin{bmatrix} \log(\text{Population}) \\ \% \text{Over 65} \\ \log(\text{Income} + 1) \\ \% \text{Hispanic} \\ \% \text{Black} \\ \text{Population Density} \\ \% \text{College Graduates} \\ \text{Can Commute} \end{bmatrix}. \quad (10.1)$$

Subscript  $i$  is omitted for brevity, but (10.1) is defined for each zip code  $i \in \{1, \dots, N\}$ . The definition of each covariate is almost self-explanatory (see Fong, Hazlett, and Imai, 2018, Sec. 5 for more details). Following Fong, Hazlett, and Imai (2018, Table 1), we add squared terms to construct a  $15 \times 1$  vector of pre-treatment covariates:

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \{\log(\text{Population})\}^2 \\ \{\% \text{Over 65}\}^2 \\ \{\log(\text{Income} + 1)\}^2 \\ \{\% \text{Hispanic}\}^2 \\ \{\% \text{Black}\}^2 \\ \{\text{Population Density}\}^2 \\ \{\% \text{College Graduates}\}^2 \end{bmatrix}. \quad (10.2)$$

The square of “Can Commute” is not added since it is a binary indicator of whether it is possible to commute to zip code  $i$  from a competitive state so that  $\text{Can Commute} = \{\text{Can Commute}\}^2$ .

Let  $T_i$  be the treatment of interest (i.e. the number of political advertisements aired in each zip code). The CBGPS approach assumes that the standardized treatment variable

$$T_i^* = s_T^{-1/2}(T_i - \bar{T}) \quad (10.3)$$

follows the standard normal distribution, where  $\bar{T} = (1/N) \sum_{i=1}^N T_i$  and  $s_T = (1/(N - 1)) \sum_{i=1}^N (T_i - \bar{T})^2$ . Given the data of political advertisements, the normality assumption is far from satisfied (see Panel 1 of Figure 9). Fong, Hazlett, and Imai (2018) therefore run a Box-Cox transformation  $T'_i = \{(T_i + 1)^\lambda - 1\}/\lambda$  with  $\lambda = -0.16$  and then standardize  $T'_i$  according to (10.3). They choose  $\lambda = -0.16$  since it yields the greatest correlation between the sample quantiles of the standardized treatment and the corresponding theoretical quantiles of the standard normal distribution. As Fong, Hazlett, and Imai (2018, p.15) admit, the Gaussian approximation is very poor even after running the Box-Cox transformation (see Panels 2-3 of Figure 9). This result suggests that the normality of a standardized treatment is often a too strong assumption to make in practice.

For an outcome model, we consider four cases for covariates  $Z_i$ :

**Case #1.**  $Z_i = [T_i, T_i^2, 1]^\top$ .

**Case #2.**  $Z_i = [T_i, T_i^2, \mathbf{SD}_i^\top]^\top$ .

**Case #3.**  $Z_i = [T_i, T_i^2, 1, \mathbf{X}_{1i}^\top]^\top$ .

**Case #4.**  $Z_i = [T_i, T_i^2, \mathbf{SD}_i^\top, \mathbf{X}_{1i}^\top]^\top$ .

Note that  $\mathbf{SD}_i = [SD_{1i}, SD_{2i}, \dots, SD_{24i}]^\top$ , where  $SD_{ji}$  is a binary indicator that equals 1 if zip code  $i$  belongs to state  $j$  and equals 0 otherwise. Any zip code contained in the dataset belongs to one and only one of 24 states (e.g. Alabama, Arkansas, ..., Wyoming).

For each of Cases #1–#4, we compute the CBGPS estimator and its asymptotic 95% confidence bands (see Fong, Hazlett, and Imai, 2018, Sec. 3.2 for procedures). Our main interest lies in the parameters of  $(T_i, T_i^2)$  and their statistical significance. See Table 13 for results. It is evident that the empirical results depend critically on a specification of  $Z_i$ . In Case #2,  $T_i$  has a significantly positive impact on  $Y_i$  and  $T_i^2$  has a significantly negative impact on  $Y_i$ . In the other three cases, both  $T_i$  and  $T_i^2$  have *insignificant* impacts on  $Y_i$ .

## 10.2 Generalized optimization approach

A practical advantage of our proposed approach over the CBGPS approach is that we do not require the normality assumption for the treatment variable  $T$ . As indicated in Figure 9,

the normality assumption is too strong for the number of political advertisements aired in each zip code whether or not the Box-Cox transformation is implemented. The generalized optimization approach allows us to work with the original treatment variable (Panel 1 of Figure 9).

We assume that the link function is quadratic with  $p = 3$ :

$$g(T, \boldsymbol{\beta}) = \beta_1 + \beta_2 T + \beta_3 T^2.$$

Our covariates  $\mathbf{X}$  are chosen to be identical to Eq. (10.2). Given that the dimension of  $\mathbf{X}$  is as large as 15, we use simple polynomials with  $K_1 = 3$  and  $K_2 = 16$  to compute  $\hat{\pi}_K(T, \mathbf{X})$  and  $\boldsymbol{\beta}$ :

$$u_{K_1}(T) = [1, T, T^2]^\top, \quad v_{K_2}(\mathbf{X}) = [1, \mathbf{X}^\top]^\top.$$

To compute variance estimator  $\hat{V}_{eff}$ , we use the same propensity score  $\hat{\pi}_K(T, \mathbf{X})$  and variance-specific polynomials with  $M_0 = 3$  and  $K_0 = 17$ :

$$v_{M_0}(\mathbf{X}) = [1, \mathbf{X}^\top]^\top, \quad w_{K_0}(T, \mathbf{X}) = [1, T, \mathbf{X}^\top]^\top.$$

See Table 14 for results. Neither  $\hat{\beta}_2$  nor  $\hat{\beta}_3$  is different from 0 at the 5% level. Hence there do not exist statistically significant impacts of the political advertisements on the level of campaign contributions  $Y$ .

## 11 Conclusions

In this paper we present a weighted optimization framework that unifies the binary, multi-valued, continuous and a mixture of discrete and continuous treatment, under the condition of unconfounded treatment assignment. Under this general framework, we first apply the result of [Bickel, Klaassen, Ritov, and Wellner \(1993\)](#) to compute the semiparametric efficiency bound for the causal effect of treatment under a general loss function. We then propose a generalized optimization estimation with the weights estimated by solving an expanding set of equations. These equations impose restriction on the weights and extract valuable information about the causal effect. Under some sufficient conditions, we establish consistency and asymptotic normality of the generalized optimization estimator and show that it attains the semiparametric efficiency bound. Since none of the existing studies has investigated efficient estimation of the continuous treatment model, our efficiency bound result and efficient estimation extend the existing literature on the binary and multi-valued treatment to the continuous treatment model. To relax the parametric restriction, we also

discuss the nonparametric specification of the causal link function, and propose estimators for both average treatment effects curve and average effect.

There are several extensions worth pursuing in future projects. First, estimation of the nonparametric causal effect function under general loss function has not been completely dealt with in this paper. But this is an important extension since it removes the burden of parameterizing the causal effect. Second, extension of the current setting to that with high dimensional covariates is also an important project. Third, panel data are common in empirical literature. Our approach is readily applicable to those data, though efficiency issue is more difficult. All these extensions shall be taken up in future studies.

## References

- ABADIE, A. (2003): “Semiparametric instrumental variable estimation of treatment response models,” *Journal of Econometrics*, 113(2), 231–263.
- (2005): “Semiparametric difference in differences estimators,” *The Review of Economic Studies*, 72(1), 1–19.
- ABADIE, A., J. D. ANGRIST, AND G. W. IMBENS (1998): “Instrumental Variables Estimation of Quantile Treatment Effects,” *NBER Technical Working Paper No. 229*.
- ABADIE, A., A. DIAMOND, AND J. HAINMUELLER (2010): “Synthetic control methods for comparative case studies: Estimating the effect of California’s tobacco control program,” *Journal of the American statistical Association*, 105(490), 493–505.
- (2015): “Comparative politics and the synthetic control method,” *American Journal of Political Science*, 59(2), 495–510.
- ABADIE, A., AND J. GARDEAZABAL (2003): “The economic costs of conflict: A case study of the Basque Country,” *American Economic Review*, 93(1), 113–132.
- ABADIE, A., AND G. W. IMBENS (2006): “Large sample properties of matching estimators for average treatment effects,” *Econometrica*, 74(1), 235–267.
- (2011): “Bias-corrected matching estimators for average treatment effects,” *Journal of Business & Economic Statistics*, 29(1), 1–11.
- (2012): “A martingale representation for matching estimators,” *Journal of the American Statistical Association*, 107(498), 833–843.
- (2016): “Matching on the estimated propensity score,” *Econometrica*, 84(2), 781–807.



- AI, C., AND X. CHEN (2003): “Efficient estimation of models with conditional moment restrictions containing unknown functions,” *Econometrica*, 71(6), 1795–1843.
- AI, C., O. LINTON, K. MOTEGI, AND Z. ZHANG (2019): “Supplemental material for ‘A Unified Framework for Efficient Estimation of General Treatment Models’,” Discussion paper, University of Florida.
- ANDREWS, D. W. K. (1991): “Asymptotic normality of series estimators for nonparametric and semiparametric regression models,” *Econometrica: Journal of the Econometric Society*, 59, 307–345.
- (1994): “Empirical process methods in econometrics,” *Handbook of econometrics*, 4, 2247–2294.
- ANGRIST, J. D., G. W. IMBENS, AND D. B. RUBIN (1996): “Identification of causal effects using instrumental variables,” *Journal of the American statistical Association*, 91(434), 444–455.
- ANGRIST, J. D., AND J.-S. PISCHKE (2008): *Mostly harmless econometrics: An empiricist’s companion*. Princeton University Press.
- ASHENFELTER, O. C., AND D. CARD (1985): “Using the longitudinal structure of earnings to estimate the effect of training programs,” *Review of Economics and Statistics*, 67(4), 648–660.
- ATHEY, S., G. IMBENS, T. PHAM, AND S. WAGER (2017): “Estimating average treatment effects: Supplementary analyses and remaining challenges,” *American Economic Review*, 107(5), 278–81.
- ATHEY, S., AND G. W. IMBENS (2006): “Identification and inference in nonlinear difference-in-differences models,” *Econometrica*, 74(2), 431–497.
- ATHEY, S., G. W. IMBENS, AND S. WAGER (2018): “Approximate residual balancing: debiased inference of average treatment effects in high dimensions,” *Journal of the Royal Statistical Society*.
- BANG, H., AND J. M. ROBINS (2005): “Doubly robust estimation in missing data and causal inference models,” *Biometrics*, 61(4), 962–973.
- BERTRAND, M., E. DUFLO, AND S. MULLAINATHAN (2004): “How much should we trust differences-in-differences estimates?,” *The Quarterly journal of economics*, 119(1), 249–275.
- BICKEL, P. J., C. A. J. KLAASSEN, Y. RITOV, AND J. A. WELLNER (1993): *Efficient and Adaptive Estimation for Semiparametric Models*. The Johns Hopkins University Press.

- BONEVA, L., D. ELLIOT, I. KAMINSKA, O. LINTON, B. MORELY, AND N. MCLAREN (2018): “The Impact of QE on liquidity: Evidence from the UK Corporate Bond Purchase Scheme,” *Bank of England working paper*.
- BUSO, M., J. DINARDO, AND J. MCCRARY (2014): “New evidence on the finite sample properties of propensity score reweighting and matching estimators,” *Review of Economics and Statistics*, 96(5), 885–897.
- CAO, W., A. A. TSIATIS, AND M. DAVIDIAN (2009): “Improving efficiency and robustness of the doubly robust estimator for a population mean with incomplete data,” *Biometrika*, 96(3), 723–734.
- CATTANEO, M. D. (2010): “Efficient semiparametric estimation of multi-valued treatment effects under ignorability,” *Journal of Econometrics*, 155(2), 138–154.
- CATTANEO, M. D., AND M. H. FARRELL (2011): “Efficient estimation of the dose-response function under ignorability using subclassification on the covariates,” in *Missing Data Methods: Cross-Sectional Methods and Applications*, pp. 93–127. Emerald Group Publishing Limited.
- CATTANEO, M. D., N. IDROBO, AND R. TITIUNIK (2017): “A Practical Introduction to Regression Discontinuity Designs,” *Cambridge Elements: Quantitative and Computational Methods for Social Science-Cambridge University Press I*.
- CHAN, K. C. G., S. C. P. YAM, AND Z. ZHANG (2016): “Globally efficient non-parametric inference of average treatment effects by empirical balancing calibration weighting,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(3), 673–700.
- CHEN, X. (2007): “Large sample sieve estimation of semi-nonparametric models,” *Handbook of Econometrics*, 6(B), 5549–5632.
- CHEN, X., O. LINTON, AND I. VAN KEILEGOM (2003): “Estimation of Semiparametric Models When the Criterion Function Is Not Smooth,” *Econometrica*, 71(5), 1591–1608.
- CHERNOZHUKOV, V., D. CHETVERIKOV, M. DEMIRER, E. DUFLO, C. HANSEN, W. NEWAY, AND J. ROBINS (2018): “Double/debiased machine learning for treatment and structural parameters,” *The Econometrics Journal*, 21(1), C1–C68.
- CHERNOZHUKOV, V., J. C. ESCANCIANO, H. ICHIMURA, AND W. K. NEWAY (2016): “Locally robust semiparametric estimation,” *arXiv preprint arXiv:1608.00033*.
- CHERNOZHUKOV, V., I. FERNÁNDEZ-VAL, AND B. MELLY (2013): “Inference on counterfactual distributions,” *Econometrica*, 81(6), 2205–2268.

- CHERNOZHUKOV, V., AND C. HANSEN (2005): “An IV Model of Quantile Treatment Effects,” *Econometrica*, 73(1), 245–261.
- DOKSUM, K. (1974): “Empirical Probability Plots and Statistical Inference for Nonlinear Models in the Two-Sample Case,” *Annals of Statistics*, 2(2), 267–277.
- DONALD, S. G., AND Y.-C. HSU (2014): “Estimation and inference for distribution functions and quantile functions in treatment effect models,” *Journal of Econometrics*, 178(3), 383–397.
- DONALD, S. G., AND K. LANG (2007): “Inference with difference-in-differences and other panel data,” *The Review of Economics and Statistics*, 89(2), 221–233.
- FAN, Y., AND S. S. PARK (2010): “Sharp bounds on the distribution of treatment effects and their statistical inference,” *Econometric Theory*, 26(3), 931–951.
- FARRELL, M. H. (2015): “Robust inference on average treatment effects with possibly more covariates than observations,” *Journal of Econometrics*, 189(1), 1–23.
- FIRPO, S. (2007): “Efficient Semiparametric Estimation of Quantile Treatment Effects,” *Econometrica*, 75(1), 259–276.
- FLORENS, J. P., J. J. HECKMAN, C. MEGHIR, AND E. VYTLACIL (2008): “Identification of treatment effects using control functions in models with continuous, endogenous treatment and heterogeneous effects,” *Econometrica*, 76(5), 1191–1206.
- FONG, C., C. HAZLETT, AND K. IMAI (2018): “Covariate Balancing Propensity Score for a Continuous Treatment: Application to the Efficacy of Political Advertisements,” *Annals of Applied Statistics*, 12(1), 156–177.
- FRÖLICH, M., AND B. MELLY (2013): “Unconditional Quantile Treatment Effects Under Endogeneity,” *Journal of Business and Economic Statistics*, 31(3), 346–357.
- GALVAO, A. F., AND L. WANG (2015): “Uniformly semiparametric efficient estimation of treatment effects with a continuous treatment,” *Journal of the American Statistical Association*, 110(512), 1528–1542.
- GRAHAM, B. S., C. C. D. X. PINTO, AND D. EGEL (2012): “Inverse probability tilting for moment condition models with missing data,” *The Review of Economic Studies*, 79(3), 1053–1079.
- HAHN, J. (1998): “On the role of the propensity score in efficient semiparametric estimation of average treatment effects,” *Econometrica*, 66(2), 315–331.
- HAHN, J., P. TODD, AND W. VAN DER KLAUW (2001): “Identification and estimation of treatment effects with a regression-discontinuity design,” *Econometrica*, 69(1), 201–209.

- HECKMAN, J. J., H. ICHIMURA, AND P. TODD (1998): "Matching as an econometric evaluation estimator," *The review of economic studies*, 65(2), 261–294.
- HECKMAN, J. J., AND E. VYTLACIL (2005): "Structural equations, treatment effects, and econometric policy evaluation," *Econometrica*, 73(3), 669–738.
- HIRANO, K., AND G. W. IMBENS (2004): "The propensity score with continuous treatments," in *Applied Bayesian Modeling and Causal Inference from Incomplete-Data Perspectives*, ed. by A. Gelman, and X.-L. Meng, chap. 7, pp. 73–84. John Wiley & Sons Ltd.
- HIRANO, K., G. W. IMBENS, AND G. RIDDER (2003): "Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score," *Econometrica*, 71(4), 1161–1189.
- HONG, H., AND D. NEKIPELOV (2010): "Semiparametric efficiency in nonlinear LATE models," *Quantitative Economics*, 1(2), 279–304.
- IMAI, K., AND D. A. VAN DYK (2004): "Causal inference with general treatment regimes: Generalizing the propensity score," *Journal of the American Statistical Association*, 99(467), 854–866.
- IMBENS, G., AND J. ANGRIST (1994): "Identification and estimation of local average treatment effects," *Econometrica: Journal of the Econometric Society*, 62(2), 467–475.
- IMBENS, G. W. (2000): "The role of the propensity score in estimating dose-response functions," *Biometrika*, 87(3), 706–710.
- (2002): "Generalized method of moments and empirical likelihood," *Journal of Business and Economic Statistics*, 20(4), 493–506.
- (2004): "Nonparametric estimation of average treatment effects under exogeneity: A review," *The Review of Economics and Statistics*, 86(1), 4–29.
- IMBENS, G. W., AND T. LEMIEUX (2008): "Regression discontinuity designs: A guide to practice," *Journal of econometrics*, 142(2), 615–635.
- IMBENS, G. W., AND C. F. MANSKI (2004): "Confidence intervals for partially identified parameters," *Econometrica*, 72(6), 1845–1857.
- IMBENS, G. W., AND D. B. RUBIN (1997): "Estimating outcome distributions for compliers in instrumental variables models," *The Review of Economic Studies*, 64(4), 555–574.
- IMBENS, G. W., AND J. M. WOOLDRIDGE (2009): "Recent developments in the econometrics of program evaluation," *Journal of Economic Literature*, 47(1), 5–86.

- KENNEDY, E. H., Z. MA, M. D. MCHUGH, AND D. S. SMALL (2017): “Non-parametric methods for doubly robust estimation of continuous treatment effects,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(4), 1229–1245.
- KITAGAWA, T. (2015): “A test for instrument validity,” *Econometrica*, 83(5), 2043–2063.
- LEE, D. S., AND T. LEMIEUX (2010): “Regression discontinuity designs in economics,” *Journal of economic literature*, 48(2), 281–355.
- LEHMANN, E. L. (1975): *Nonparametrics: Statistical Methods Based on Ranks*. San Francisco: Holden-Day.
- MANSKI, C. F. (1990): “Nonparametric bounds on treatment effects,” *The American Economic Review*, 80(2), 319–323.
- (2003): *Partial identification of probability distributions*. Springer Science & Business Media.
- (2009): *Identification for prediction and decision*. Harvard University Press.
- NEWKEY, W. K. (1997): “Convergence Rates and Asymptotic Normality for Series Estimators,” *Journal of Econometrics*, 79(1), 147–168.
- NEWKEY, W. K., AND J. L. POWELL (1987): “Asymmetric least squares estimation and testing,” *Econometrica: Journal of the Econometric Society*, 55, 819–847.
- NIELSEN, J. P., AND O. LINTON (1998): “An optimization interpretation of integration and back-fitting estimators for separable nonparametric models,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(1), 217–222.
- PAKES, A., AND D. POLLARD (1989): “Simulation and the asymptotics of optimization estimators,” *Econometrica: Journal of the Econometric Society*, pp. 1027–1057.
- QIN, J., AND B. ZHANG (2007): “Empirical-likelihood-based inference in missing response problems and its application in observational studies,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(1), 101–122.
- ROBINS, J. M., M. A. HERNÁN, AND B. BRUMBACK (2000): “Marginal structural models and causal inference in epidemiology,” *Epidemiology*, 11(5), 550–560.
- ROBINS, J. M., A. ROTNITZKY, AND L. P. ZHAO (1994): “Estimation of regression coefficients when some regressors are not always observed,” *Journal of the American Statistical Association*, 89(427), 846–866.
- ROSENBAUM, P. R., AND D. B. RUBIN (1983): “The central role of the propensity score in observational studies for causal effects,” *Biometrika*, 70(1), 41–55.

- ROTHER, C. (2017): “Robust confidence intervals for average treatment effects under limited overlap,” *Econometrica*, 85(2), 645–660.
- RUBIN, D. B. (1977): “Assignment to treatment group on the basis of a covariate,” *Journal of educational Statistics*, 2(1), 1–26.
- SŁOCZYŃSKI, T., AND J. M. WOOLDRIDGE (2018): “A general double robustness result for estimating average treatment effects,” *Econometric Theory*, 34, 112–133.
- SMITH, R. J. (1997): “Alternative semi-parametric likelihood approaches to generalised method of moments estimation,” *The Economic Journal*, 107(441), 503–519.
- TAN, Z. (2010): “Bounded, efficient and doubly robust estimation with inverse weighting,” *Biometrika*, 97(3), 661–682.
- TSENG, P., AND D. P. BERTSEKAS (1991): “Relaxation methods for problems with strictly convex costs and linear constraints,” *Mathematics of Operations Research*, 16(3), 462–481.
- URBAN, C., AND S. NIEBLER (2014): “Dollars on the Sidewalk: Should U.S. Presidential Candidates Advertise in Uncontested States?,” *American Journal of Political Science*, 58(2), 322–336.
- VANSTEELENDT, S., M. BEKAERT, AND G. CLAESKENS (2010): “On model selection and model misspecification in causal inference,” *Statistical Methods in Medical Research*, 21(1), 7–30.
- WAGER, S., AND S. ATHEY (2018): “Estimation and inference of heterogeneous treatment effects using random forests,” *Journal of the American Statistical Association*.
- YIU, S., AND L. SU (2018): “Covariate association eliminating weights: A unified weighting framework for causal effect estimation,” *Biometrika*, 105(3), 709–722.

## Appendix

### A Proof of (2.2)

Using the law of iterated expectation and Assumption 1, we can deduce that

$$\begin{aligned}
& \mathbb{E} [\pi_0(T, \mathbf{X}) L(Y - g(T; \beta))] \\
&= \mathbb{E} [\mathbb{E} [\pi(T, \mathbf{X}) L(Y^*(T) - g(T; \beta)) | T, \mathbf{X}]] \\
&= \int \pi_0(t, \mathbf{x}) \cdot \mathbb{E} [L(Y^*(T) - g(T; \beta)) | T = t, \mathbf{X} = \mathbf{x}] dF_{T|X}(t|\mathbf{x}) dF_X(\mathbf{x})
\end{aligned}$$

$$\begin{aligned}
&= \int \mathbb{E} [L(Y^*(t) - g(t; \beta)) | T = t, \mathbf{X} = \mathbf{x}] dF_T(t) dF_X(\mathbf{x}) \\
&= \int \mathbb{E} [L(Y^*(t) - g(t; \beta)) | \mathbf{X} = \mathbf{x}] dF_T(t) dF_X(\mathbf{x}) \quad (\text{using Assumption 1}) \\
&= \int \mathbb{E} [L(Y^*(t) - g(t; \beta))] dF_T(t).
\end{aligned}$$

## B Asymptotic result when $\pi_0(T, \mathbf{X})$ is known

Suppose the stabilized weights  $\pi_0(T, \mathbf{X})$  is known, the weighted optimization estimator of  $\beta_0$ , denoted by  $\beta^*$ , is

$$\beta^* = \min_{\beta} \sum_{i=1}^N \pi_0(T_i, \mathbf{X}_i) L(Y_i - g(T_i; \beta)).$$

We also assume the asymptotic first order condition

$$\frac{1}{N} \sum_{i=1}^N \pi_0(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \beta^*)) m(T_i; \beta^*) = o_P(N^{-1/2}) \quad (\text{B.1})$$

holds with probability approaching to one.

**Proposition B.1** Suppose Assumptions 5, 6 (i-ii), and 7 hold, and (B.1) holds, then we have

1.  $\beta^* \xrightarrow{p} \beta_0$ ;
2.  $\sqrt{N}(\beta^* - \beta_0) \xrightarrow{d} \mathcal{N}(0, V_{ineff})$ , where

$$V_{ineff} := H_0^{-1} \cdot \mathbb{E} [\pi_0(T, \mathbf{X})^2 L'(Y - g(T; \beta_0))^2 m(T; \beta_0) m(T; \beta_0)^\top] \cdot H_0^{-1};$$

3. furthermore, if  $\mathbb{E} [L'(Y(t) - g(t; \beta_0))] = 0$  holds for all  $t \in \mathcal{T}$ , then  $V_{ineff} \geq V_{eff}$  in the sense of that  $c^\top \cdot V_{ineff} \cdot c \geq c^\top \cdot V_{eff} \cdot c$  for any vector  $c \in \mathbb{R}^p$ .

*Proof.* By Assumption 5 and the uniform law of large number, we obtain

$$\frac{1}{N} \sum_{i=1}^N \pi_0(T_i, \mathbf{X}_i) L\{Y_i - g(T_i; \beta)\} \rightarrow \mathbb{E} [\pi_0(T, \mathbf{X}) L\{Y - g(T; \beta)\}] \text{ in probability uniformly over } \beta,$$

which implies the consistency result  $\|\beta^* - \beta_0\| \xrightarrow{p} 0$ .

The first order condition (B.1) holds with probability approaching to one. Note that  $L'(\cdot)$  may not be a differentiable function, e.g.  $L'(v) = \tau - I(v < 0)$  in quantile regression, we cannot simply apply Mean Value Theorem on (B.1) to obtain the expression for



$\sqrt{N}(\beta^* - \beta_0)$ . To solve this problem, we resort to the empirical process theory in [Andrews \(1994\)](#). Define

$$f(\beta) := \mathbb{E} [\pi_0(T, \mathbf{X}) L'(Y - g(T; \beta)) m(T; \beta)],$$

which is a differentiable function in  $\beta$  and by (2.3)  $f(\beta_0) = 0$ . Using Mean Value Theorem, we can obtain

$$0 = \sqrt{N}f(\beta_0) = \sqrt{N}f(\beta^*) - \nabla_{\beta}f(\bar{\beta}) \cdot \sqrt{N}(\beta^* - \beta_0),$$

where  $\bar{\beta}$  lies on the line joining  $\beta^*$  and  $\beta_0$ . Because  $\nabla_{\beta}f(\beta)$  is continuous in  $\beta$  at  $\beta_0$ , and  $\|\beta^* - \beta_0\| \xrightarrow{p} 0$ , then we have

$$\sqrt{N}(\beta^* - \beta_0) = [\nabla_{\beta}f(\beta_0)]^{-1} \cdot \sqrt{N}f(\beta^*).$$

Define the empirical process

$$\nu_N(\beta) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \{ \pi_0(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \beta)) m(T_i; \beta) - \mathbb{E} [\pi_0(T, \mathbf{X}) L'(Y - g(T; \beta)) m(T; \beta)] \}.$$

By (B.1) and the definition of  $\nu_N(\beta)$ , we have

$$\begin{aligned} \sqrt{N}(\beta^* - \beta_0) &= \nabla_{\beta}f(\beta_0)^{-1} \cdot \left\{ \sqrt{N}f(\beta^*) - \frac{1}{\sqrt{N}} \sum_{i=1}^N \pi_0(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \beta^*)) m(T_i; \beta^*) \right. \\ &\quad \left. + \frac{1}{\sqrt{N}} \sum_{i=1}^N \pi_0(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \beta^*)) m(T_i; \beta^*) \right\} \\ &= -\nabla_{\beta}f(\beta_0)^{-1} \cdot \nu_N(\beta^*) + o_p(1) \\ &= H_0^{-1} \cdot \left\{ (\nu_N(\beta^*) - \nu_N(\beta_0)) + \nu_N(\beta_0) \right\} + o_p(1). \end{aligned}$$

By Assumptions 6, 7, Theorems 4 and 5 of [Andrews \(1994\)](#), we have that  $\nu_N(\cdot)$  is stochastically equicontinuous, which implies  $\nu_N(\beta^*) - \nu_N(\beta_0) \xrightarrow{p} 0$ . Therefore,

$$\sqrt{N}(\beta^* - \beta_0) = H_0^{-1} \frac{1}{\sqrt{N}} \sum_{i=1}^N \pi_0(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \beta_0)) m(T_i; \beta_0) + o_p(1),$$

then we can conclude that the asymptotic variance of  $\sqrt{N}(\beta^* - \beta_0)$  is  $V_{ineff}$ .

We next show  $V_{ineff} \geq V_{eff}$ . From Theorem 1, we have

$$V_{eff} = H_0^{-1} \cdot \left\{ \mathbb{E} \left[ \pi_0(T, \mathbf{X})^2 L'(Y - g(T; \beta_0))^2 m(T; \beta_0) m(T; \beta_0)^{\top} \right] \right\}$$

$$\begin{aligned}
& + \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}]^\top \right] \\
& + \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}]^\top \right] \\
& + \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T]^\top \right] \\
& - 2 \cdot \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)^\top \right] \\
& - 2 \cdot \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}]^\top \right] \\
& - 2 \cdot \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T]^\top \right] \\
& + 2 \cdot \mathbb{E} \left[ \pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0) \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}]^\top \right] \\
& + 2 \cdot \mathbb{E} \left[ \pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0) \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T]^\top \right] \\
& + 2 \cdot \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}]^\top \right] \Big\} H_0^{-1} \\
& = H_0^{-1} \left\{ \mathbb{E} \left[ \pi_0(T, \mathbf{X})^2 L'(Y - g(T; \beta_0))^2 m(T; \beta_0) m(T; \beta_0)^\top \right] \right. \\
& \quad - \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}]^\top \right] \\
& \quad \left. + \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}]^\top \right] \right\} H_0^{-1},
\end{aligned}$$

where the last equality holds by noting

$$\mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T = t] = \mathbb{E}[L'(Y^*(t) - g(t; \beta_0))] \cdot m(t; \beta_0) = 0,$$

since the model is correctly specified, i.e.  $\mathbb{E}[L'(Y^*(t) - g(t; \beta_0))] = 0$  for  $t \in \mathcal{T}$ . Therefore,

$$\begin{aligned}
& V_{ineff} - V_{eff} \\
& = H_0^{-1} \left\{ \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}]^\top \right] \right. \\
& \quad \left. - \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})L'(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}]^\top \right] \right\} H_0^{-1} \geq 0,
\end{aligned}$$

where the last inequality holds by using Jensen's inequality:

$$\begin{aligned}
& \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})(Y - g(T; \beta_0))m(T; \beta_0)|\mathbf{X}]^\top \right] \\
& = \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}]^\top \right] \\
& < \mathbb{E} \left[ \mathbb{E}[\pi_0(T, \mathbf{X})(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}] \cdot \mathbb{E}[\pi_0(T, \mathbf{X})(Y - g(T; \beta_0))m(T; \beta_0)|T, \mathbf{X}]^\top \right].
\end{aligned}$$

□

## C Duality of primal problem (4.3)

Let

$$\boldsymbol{\pi} := (\pi_1, \dots, \pi_N)^\top, \quad f(\boldsymbol{\pi}) := \sum_{i=1}^N \pi_i \log \pi_i,$$

the primal problem (4.3) can be written as:

$$\left\{ \begin{array}{l} \min f(\boldsymbol{\pi}) \\ \text{subject to } N^{-1} \sum_{i=1}^N \pi_i u_{K_1,k}(T_i) v_{K_2,k'}(\mathbf{X}_i) = \bar{u}_{K_1,k} \cdot \bar{v}_{K_2,k'}, \\ k \in \{1, \dots, K_1\}, k' \in \{1, \dots, K_2\}, \end{array} \right. \quad (\text{C.1})$$

where

$$\bar{u}_{K_1,k} := \frac{1}{N} \sum_{j=1}^N u_{K_1,k}(T_j), \quad \bar{v}_{K_2,k'} := \frac{1}{N} \sum_{j=1}^N v_{K_2,k'}(\mathbf{X}_j),$$

and  $u_{K_1,k}(T)$  (resp.  $v_{K_2,k'}$ ) is the  $k^{\text{th}}$  (resp.  $k'^{\text{th}}$ ) component of  $u_{K_1}(T)$  (resp.  $v_{K_2}(\mathbf{X})$ ). Let  $m_{K_1 K_2}(T, \mathbf{X})$  be a  $K = K_1 \cdot K_2$  dimensional vector function whose elements are  $\{u_{K_1,k}(T) v_{K_2,k'}(\mathbf{X}); k = 1, \dots, K_1, k' = 1, \dots, K_2\}$ , and we define

$$E_{K_1 K_2 \times N} := (m_{K_1 K_2}(T_1, \mathbf{X}_1), \dots, m_{K_1 K_2}(T_N, \mathbf{X}_N)).$$

Let  $b_{K_1 K_2}$  be a  $K = K_1 \cdot K_2$  dimensional vector function whose elements are  $\{\bar{u}_{K_1,k} \bar{v}_{K_2,k'}; k = 1, \dots, K_1, k' = 1, \dots, K_2\}$ . The optimization problem (C.1) becomes

$$\left\{ \begin{array}{l} \min_{\boldsymbol{\pi}} f(\boldsymbol{\pi}) \\ \text{subject to } E_{K_1 K_2 \times N} \cdot \boldsymbol{\pi} = N \cdot b_{K_1 K_2} \end{array} \right. . \quad (\text{C.2})$$

By Tseng and Bertsekas (1991), the conjugate convex function of  $f(\cdot)$  to be

$$f^*(\mathbf{z}) = \sup_{\boldsymbol{\pi}} \sum_{i=1}^N \{z_i \pi_i - \pi_i \log \pi_i\} = \sum_{i=1}^N \{z_i \pi_i^* - \pi_i^* \log \pi_i^*\},$$

where  $\pi_j^*$  satisfies the first order condition:

$$z_j = \log \pi_j^* + 1 \Rightarrow \pi_j^* = e^{z_j - 1};$$

then we have

$$f^*(\mathbf{z}) = \sum_{i=1}^N \{z_i e^{z_i-1} - e^{z_i-1}(z_i - 1)\} = \sum_{i=1}^N e^{z_i-1} = \sum_{i=1}^N -\rho(-z_i),$$

where  $\rho(z_j) := -e^{-z_j-1}$ . By [Tseng and Bertsekas \(1991\)](#), the dual problem of (C.2) is

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}^K} \{ \lambda^\top (N \cdot b_{K_1 K_2}) - f^*(\lambda^\top E_{K_1 K_2, N}) \} \\ &= \max_{\Lambda \in \mathbb{R}^{K_1 \times K_2}} \sum_{j=1}^N \{ \bar{u}_{K_1}^\top \Lambda \bar{v}_{K_2} + \rho(-u_{K_1}(T_j)^\top \Lambda v_{K_2}(\mathbf{X}_j)) \} \\ &= \max_{\Lambda \in \mathbb{R}^{K_1 \times K_2}} \sum_{j=1}^N \{ \rho(u_{K_1}(T_j)^\top \Lambda v_{K_2}(\mathbf{X}_j)) - \bar{u}_{K_1}^\top \Lambda \bar{v}_{K_2} \} \\ &= \max_{\Lambda \in \mathbb{R}^{K_1 \times K_2}} \hat{G}_{K_1 \times K_2}(\Lambda), \end{aligned} \tag{C.3}$$

where  $E_{K_1 K_2, N}^{(j)}$  is the  $j^{\text{th}}$  column of  $E_{K_1 K_2 \times N}$  and

$$G_{K_1 \times K_2}(\Lambda) = \frac{1}{N} \sum_{j=1}^N \rho(u_{K_1}(T_j)^\top \Lambda v_{K_2}(\mathbf{X}_j)) - \bar{u}_{K_1}^\top \Lambda \bar{v}_{K_2}.$$

Therefore, the dual solution of (4.3) is given by

$$\hat{\pi}_K(T_i, \mathbf{X}_i) = \rho' \left( u_{K_1}(T_i)^\top \hat{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right), \tag{C.4}$$

where  $\hat{\Lambda}_{K_1 \times K_2}$  is the maximizer of the strictly concave objective function  $\hat{G}_{K_1 \times K_2}$ .

## D Proof of (6.1)

$$\begin{aligned} & \nabla_\beta \mathbb{E} [L'(Y - g(T; \beta)) | T = t, \mathbf{X} = \mathbf{x}] \Big|_{\beta=\beta_0} \\ &= \nabla_\beta \left[ \int_{\mathbb{R}} L'(y - g(t; \beta)) f_{Y|T, X}(y|t, \mathbf{x}) dy \right] \Big|_{\beta=\beta_0} \\ &= \nabla_\beta \left[ \int_{\mathbb{R}} L'(z) f_{Y|T, X}(z + g(t; \beta) | t, \mathbf{x}) dz \right] \Big|_{\beta=\beta_0} \quad (\text{use } z = y - g(t; \beta)) \\ &= \int_{\mathbb{R}} L'(z) \cdot \frac{\partial}{\partial y} f_{Y|T, X}(z + g(t; \beta_0) | t, \mathbf{x}) dz \cdot m(t; \beta_0) \\ &= \int_{\mathbb{R}} L'(y - g(t; \beta_0)) \cdot \frac{\partial}{\partial y} f_{Y|T, X}(y|t, \mathbf{x}) dy \cdot m(t; \beta_0) \end{aligned}$$

$$\begin{aligned}
&= \int_{\mathbb{R}} L'(y - g(t; \boldsymbol{\beta})) \cdot \frac{\frac{\partial}{\partial y} f_{Y|T,X}(y|t, \mathbf{x})}{f_{Y|T,X}(y|t, \mathbf{x})} f_{Y|T,X}(y|t, \mathbf{x}) dy \cdot m(t; \boldsymbol{\beta}_0) \\
&= \int_{\mathbb{R}} L'(y - g(t; \boldsymbol{\beta})) \cdot \frac{\frac{\partial}{\partial y} f_{Y,T,X}(y, t, \mathbf{x})}{f_{Y,T,X}(y, t, \mathbf{x})} f_{Y|T,X}(y|t, \mathbf{x}) dy \cdot m(t; \boldsymbol{\beta}_0) \\
&= \mathbb{E} \left[ L'(Y - g(T; \boldsymbol{\beta}_0)) \frac{\frac{\partial}{\partial y} f_{Y,T,\mathbf{X}}(Y, T, \mathbf{X})}{f_{Y,T,\mathbf{X}}(Y, T, \mathbf{X})} \middle| T = t, \mathbf{X} = \mathbf{x} \right] m(t; \boldsymbol{\beta}_0).
\end{aligned}$$

Table 1: Simulation results on point estimation of intercept  $\beta_1$  under DGP-L1 ( $\beta_1^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 2)	0.002	0.108	0.108	0.001	0.047	0.047	0.002	0.034	0.034
GOE	(2, 3)	-0.002	0.106	0.106	0.001	0.046	0.046	-0.000	0.034	0.034
GOE	(2, 4)	0.009	0.117	0.117	0.008	0.048	0.049	0.009	0.035	0.036
GOE	(3, 2)	0.005	0.108	0.109	0.001	0.049	0.049	0.002	0.034	0.034
GOE	(3, 3)	0.002	0.119	0.119	0.004	0.049	0.049	0.007	0.035	0.036
GOE	(3, 4)	0.010	0.123	0.123	0.012	0.052	0.053	0.017	0.036	0.040
GOE	(4, 2)	0.013	0.114	0.115	0.004	0.048	0.048	0.005	0.034	0.034
GOE	(4, 3)	0.011	0.127	0.128	0.014	0.051	0.052	0.013	0.038	0.040
GOE	(4, 4)	0.020	0.133	0.134	0.017	0.052	0.055	0.019	0.037	0.042
GOE	MSE (none)	0.006	0.110	0.110	0.009	0.048	0.049	0.008	0.034	0.035
GOE	MSE (add)	0.005	0.105	0.105	0.004	0.047	0.048	0.005	0.034	0.035
GOE	MSE (multi)	-0.004	0.105	0.105	0.003	0.046	0.046	0.002	0.034	0.034
GOE	CV ( $J = 5$ )	0.000	0.107	0.107	0.001	0.046	0.046	0.001	0.033	0.033
GOE	CV ( $J = 10$ )	-0.000	0.104	0.104	0.001	0.046	0.046	0.001	0.032	0.032
CBGPS	-	0.003	0.144	0.144	0.001	0.065	0.065	0.003	0.047	0.047

DGP-L1:  $T = 1 + 0.2X_1 + \xi$  and  $Y = 1 + X_1 + T + \epsilon$ , where  $X_1 \sim N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, X_{1i})(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .

Table 2: Simulation results on point estimation of slope  $\beta_2$  under DGP-L1 ( $\beta_2^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 2)	-0.002	0.185	0.185	-0.000	0.081	0.081	-0.001	0.056	0.056
GOE	(2, 3)	0.010	0.178	0.178	-0.001	0.080	0.080	0.004	0.057	0.057
GOE	(2, 4)	-0.000	0.196	0.196	-0.005	0.083	0.083	-0.001	0.057	0.057
GOE	(3, 2)	-0.002	0.185	0.185	-0.001	0.081	0.081	0.002	0.057	0.057
GOE	(3, 3)	-0.001	0.190	0.190	0.000	0.080	0.080	-0.003	0.057	0.057
GOE	(3, 4)	-0.007	0.201	0.201	-0.011	0.085	0.086	-0.011	0.060	0.061
GOE	(4, 2)	-0.005	0.184	0.185	-0.002	0.080	0.080	-0.000	0.055	0.055
GOE	(4, 3)	-0.007	0.205	0.205	-0.006	0.083	0.084	-0.011	0.060	0.061
GOE	(4, 4)	-0.020	0.207	0.208	-0.012	0.084	0.084	-0.013	0.062	0.064
GOE	MSE (none)	0.002	0.171	0.171	-0.008	0.079	0.080	-0.006	0.057	0.058
GOE	MSE (add)	-0.013	0.169	0.170	-0.005	0.076	0.076	-0.002	0.057	0.057
GOE	MSE (multi)	0.003	0.165	0.165	-0.001	0.079	0.079	-0.003	0.056	0.056
GOE	CV ( $J = 5$ )	0.006	0.191	0.191	0.004	0.080	0.080	0.001	0.058	0.058
GOE	CV ( $J = 10$ )	0.005	0.182	0.182	0.001	0.079	0.079	0.001	0.057	0.057
CBGPS	-	-0.002	0.102	0.103	-0.002	0.045	0.046	-0.002	0.032	0.032

DGP-L1:  $T = 1 + 0.2X_1 + \xi$  and  $Y = 1 + X_1 + T + \epsilon$ , where  $X_1 \sim N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, X_{1i})(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .

Table 3: Simulation results on point estimation of intercept  $\beta_1$  under DGP-NL1 ( $\beta_1^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 2)	0.180	0.179	0.254	0.195	0.084	0.213	0.192	0.058	0.201
GOE	(2, 3)	-0.002	0.106	0.106	-0.001	0.045	0.045	0.002	0.032	0.032
GOE	(2, 4)	0.012	0.118	0.118	0.005	0.047	0.048	0.006	0.034	0.034
GOE	(3, 2)	0.170	0.183	0.250	0.191	0.082	0.208	0.192	0.057	0.201
GOE	(3, 3)	0.010	0.131	0.131	0.007	0.048	0.048	0.006	0.034	0.035
GOE	(3, 4)	0.013	0.132	0.133	0.010	0.051	0.052	0.011	0.037	0.039
GOE	(4, 2)	0.176	0.187	0.257	0.188	0.086	0.206	0.192	0.058	0.201
GOE	(4, 3)	0.019	0.136	0.138	0.015	0.053	0.055	0.011	0.036	0.037
GOE	(4, 4)	0.024	0.139	0.141	0.016	0.055	0.058	0.017	0.037	0.040
GOE	MSE (none)	0.006	0.132	0.132	0.012	0.055	0.056	0.011	0.041	0.043
GOE	MSE (add)	0.001	0.123	0.123	0.015	0.059	0.061	0.010	0.043	0.044
GOE	MSE (multi)	0.021	0.135	0.137	0.015	0.060	0.062	0.012	0.045	0.047
GOE	CV ( $J = 5$ )	0.089	0.166	0.188	0.053	0.091	0.105	0.039	0.076	0.086
GOE	CV ( $J = 10$ )	0.075	0.152	0.170	0.052	0.089	0.103	0.038	0.075	0.084
CBGPS	-	-0.035	0.234	0.237	-0.024	0.076	0.080	-0.022	0.052	0.057

DGP-NL1:  $T = 0.1X_1^2 + \xi$  and  $Y = X_1^2 + T + \epsilon$ , where  $X_1 \sim N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, X_{1i})(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .



Table 4: Simulation results on point estimation of slope  $\beta_2$  under DGP-NL1 ( $\beta_2^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 2)	-0.029	0.171	0.174	-0.020	0.075	0.078	-0.021	0.055	0.059
GOE	(2, 3)	0.001	0.175	0.175	-0.003	0.074	0.074	-0.000	0.054	0.054
GOE	(2, 4)	-0.004	0.177	0.177	-0.001	0.080	0.080	-0.000	0.057	0.057
GOE	(3, 2)	-0.042	0.168	0.173	-0.021	0.075	0.077	-0.021	0.054	0.058
GOE	(3, 3)	-0.017	0.186	0.187	-0.001	0.081	0.081	0.002	0.056	0.056
GOE	(3, 4)	-0.025	0.180	0.182	-0.004	0.080	0.080	-0.000	0.058	0.058
GOE	(4, 2)	-0.063	0.170	0.181	-0.028	0.076	0.081	-0.023	0.053	0.058
GOE	(4, 3)	-0.015	0.197	0.197	-0.002	0.087	0.087	0.001	0.058	0.058
GOE	(4, 4)	-0.044	0.187	0.192	-0.011	0.081	0.082	-0.003	0.058	0.058
GOE	MSE (none)	-0.061	0.175	0.185	-0.021	0.078	0.081	-0.013	0.057	0.058
GOE	MSE (add)	-0.075	0.164	0.180	-0.021	0.079	0.081	-0.015	0.054	0.056
GOE	MSE (multi)	-0.057	0.171	0.181	-0.017	0.077	0.079	-0.013	0.055	0.056
GOE	CV ( $J = 5$ )	-0.035	0.174	0.177	-0.010	0.079	0.079	-0.006	0.055	0.056
GOE	CV ( $J = 10$ )	-0.026	0.171	0.173	-0.013	0.077	0.078	-0.006	0.055	0.055
CBGPS	-	0.189	0.186	0.266	0.190	0.080	0.206	0.195	0.055	0.203

DGP-NL1:  $T = 0.1X_1^2 + \xi$  and  $Y = X_1^2 + T + \epsilon$ , where  $X_1 \sim N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, X_{1i})(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .

Table 5: Simulation results on point estimation of intercept  $\beta_1$  under DGP-L2 ( $\beta_1^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 3)	0.003	0.108	0.108	0.002	0.050	0.050	0.001	0.034	0.034
GOE	(2, 6)	0.030	0.125	0.129	0.026	0.051	0.058	0.025	0.038	0.045
GOE	(2, 10)	0.045	0.131	0.138	0.042	0.052	0.067	0.040	0.037	0.055
GOE	(3, 3)	0.015	0.116	0.117	0.020	0.052	0.056	0.015	0.036	0.040
GOE	(3, 6)	0.029	0.129	0.132	0.031	0.056	0.064	0.030	0.041	0.050
GOE	(3, 10)	0.041	0.142	0.148	0.036	0.056	0.066	0.035	0.039	0.053
GOE	(4, 3)	0.024	0.126	0.128	0.024	0.051	0.057	0.022	0.038	0.044
GOE	(4, 6)	0.034	0.137	0.141	0.037	0.057	0.067	0.037	0.038	0.053
GOE	(4, 10)	0.039	0.153	0.158	0.028	0.055	0.062	0.031	0.039	0.049
GOE	MSE (none)	0.030	0.110	0.114	0.024	0.050	0.056	0.021	0.037	0.042
GOE	MSE (add)	0.015	0.105	0.106	0.017	0.049	0.052	0.017	0.036	0.039
GOE	MSE (multi)	0.005	0.107	0.107	0.005	0.050	0.050	0.010	0.035	0.037
GOE	CV ( $J = 5$ )	0.010	0.106	0.106	0.006	0.050	0.050	0.005	0.036	0.036
GOE	CV ( $J = 10$ )	0.011	0.107	0.107	0.003	0.047	0.047	0.005	0.034	0.034
CBGPS	-	-0.003	0.156	0.156	0.003	0.070	0.070	-0.001	0.049	0.049

DGP-L2:  $T = 1 + 0.2 \sum_{j=1}^2 X_j + \xi$  and  $Y = 1 + (1/2) \sum_{j=1}^2 X_j + T + \epsilon$ , where  $X_1, X_2 \stackrel{i.i.d.}{\sim} N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .

Table 6: Simulation results on point estimation of slope  $\beta_2$  under DGP-L2 ( $\beta_2^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 3)	-0.004	0.163	0.163	-0.003	0.075	0.075	0.003	0.053	0.053
GOE	(2, 6)	-0.019	0.178	0.179	-0.013	0.078	0.079	-0.010	0.056	0.057
GOE	(2, 10)	-0.036	0.196	0.199	-0.031	0.076	0.082	-0.030	0.056	0.064
GOE	(3, 3)	-0.005	0.178	0.178	-0.004	0.078	0.079	-0.001	0.054	0.054
GOE	(3, 6)	-0.038	0.190	0.194	-0.027	0.084	0.088	-0.025	0.060	0.065
GOE	(3, 10)	-0.036	0.207	0.210	-0.033	0.082	0.088	-0.028	0.058	0.065
GOE	(4, 3)	-0.014	0.188	0.188	-0.006	0.081	0.081	-0.007	0.058	0.058
GOE	(4, 6)	-0.037	0.202	0.205	-0.034	0.082	0.089	-0.028	0.058	0.065
GOE	(4, 10)	-0.026	0.213	0.215	-0.025	0.083	0.086	-0.027	0.058	0.065
GOE	MSE (none)	-0.028	0.162	0.165	-0.019	0.072	0.075	-0.014	0.052	0.054
GOE	MSE (add)	-0.009	0.160	0.161	-0.014	0.072	0.073	-0.010	0.052	0.053
GOE	MSE (multi)	-0.006	0.163	0.163	-0.002	0.073	0.073	-0.006	0.052	0.052
GOE	CV ( $J = 5$ )	0.003	0.161	0.161	0.001	0.075	0.075	0.001	0.052	0.052
GOE	CV ( $J = 10$ )	0.003	0.164	0.164	-0.001	0.071	0.071	-0.002	0.053	0.053
CBGPS	-	-0.003	0.114	0.114	-0.001	0.050	0.050	-0.001	0.036	0.036

DGP-L2:  $T = 1 + 0.2 \sum_{j=1}^2 X_j + \xi$  and  $Y = 1 + (1/2) \sum_{j=1}^2 X_j + T + \epsilon$ , where  $X_1, X_2 \stackrel{i.i.d.}{\sim} N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .

Table 7: Simulation results on point estimation of intercept  $\beta_1$  under DGP-NL2 ( $\beta_1^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 3)	0.167	0.127	0.210	0.185	0.059	0.194	0.183	0.042	0.187
GOE	(2, 6)	0.027	0.123	0.126	0.023	0.053	0.058	0.025	0.039	0.047
GOE	(2, 10)	0.037	0.127	0.132	0.041	0.052	0.066	0.038	0.038	0.053
GOE	(3, 3)	0.160	0.132	0.208	0.181	0.061	0.190	0.181	0.041	0.186
GOE	(3, 6)	0.027	0.126	0.129	0.031	0.051	0.060	0.030	0.039	0.049
GOE	(3, 10)	0.033	0.140	0.144	0.034	0.053	0.063	0.035	0.038	0.051
GOE	(4, 3)	0.162	0.136	0.212	0.178	0.061	0.188	0.181	0.042	0.185
GOE	(4, 6)	0.031	0.132	0.136	0.034	0.054	0.064	0.036	0.039	0.053
GOE	(4, 10)	0.039	0.139	0.144	0.028	0.054	0.060	0.032	0.039	0.050
GOE	MSE (none)	0.019	0.121	0.123	0.027	0.052	0.059	0.028	0.036	0.046
GOE	MSE (add)	0.034	0.121	0.126	0.029	0.051	0.058	0.028	0.037	0.046
GOE	MSE (multi)	0.078	0.117	0.141	0.038	0.057	0.069	0.028	0.040	0.049
GOE	CV ( $J = 5$ )	0.109	0.127	0.168	0.064	0.073	0.096	0.046	0.053	0.070
GOE	CV ( $J = 10$ )	0.105	0.127	0.165	0.061	0.071	0.094	0.042	0.051	0.066
CBGPS	-	-0.042	0.129	0.136	-0.040	0.054	0.068	-0.036	0.038	0.053

DGP-NL2:  $T = 0.1(\sum_{j=1}^2 X_j)^2 + \xi$  and  $Y = 1/2 + [(1/2) \sum_{j=1}^2 X_j]^2 + T + \epsilon$ , where  $X_1, X_2 \stackrel{i.i.d.}{\sim} N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .

Table 8: Simulation results on point estimation of slope  $\beta_2$  under DGP-NL2 ( $\beta_2^* = 1$ )

		$N = 100$			$N = 500$			$N = 1000$		
	$(K_1, K_2)$	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	(2, 3)	-0.037	0.125	0.130	-0.036	0.054	0.065	-0.036	0.037	0.052
GOE	(2, 6)	-0.008	0.141	0.141	0.005	0.062	0.062	0.007	0.045	0.045
GOE	(2, 10)	-0.022	0.136	0.138	-0.007	0.059	0.060	-0.007	0.043	0.044
GOE	(3, 3)	-0.045	0.123	0.131	-0.036	0.052	0.063	-0.037	0.037	0.052
GOE	(3, 6)	-0.031	0.132	0.135	-0.012	0.060	0.061	-0.005	0.045	0.045
GOE	(3, 10)	-0.031	0.147	0.151	-0.016	0.061	0.063	-0.014	0.043	0.045
GOE	(4, 3)	-0.049	0.128	0.137	-0.039	0.055	0.068	-0.037	0.038	0.053
GOE	(4, 6)	-0.032	0.148	0.151	-0.014	0.060	0.061	-0.009	0.044	0.045
GOE	(4, 10)	-0.046	0.155	0.162	-0.016	0.059	0.061	-0.016	0.044	0.047
GOE	MSE (none)	-0.056	0.134	0.146	-0.023	0.057	0.061	-0.018	0.042	0.045
GOE	MSE (add)	-0.044	0.128	0.136	-0.022	0.056	0.060	-0.021	0.041	0.047
GOE	MSE (multi)	-0.048	0.121	0.130	-0.022	0.054	0.058	-0.017	0.040	0.043
GOE	CV ( $J = 5$ )	-0.027	0.123	0.125	-0.013	0.056	0.058	-0.007	0.043	0.044
GOE	CV ( $J = 10$ )	-0.030	0.125	0.129	-0.013	0.058	0.059	-0.009	0.044	0.044
CBGPS	-	0.168	0.139	0.218	0.177	0.058	0.186	0.183	0.041	0.188

DGP-NL2:  $T = 0.1(\sum_{j=1}^2 X_j)^2 + \xi$  and  $Y = 1/2 + [(1/2) \sum_{j=1}^2 X_j]^2 + T + \epsilon$ , where  $X_1, X_2 \stackrel{i.i.d.}{\sim} N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. “MSE (none)” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “MSE (add)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “MSE (multi)” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . “CV” signifies that we pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . The choice set of  $(K_1, K_2)$  is the nine pairs listed in the table. “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. The number of Monte Carlo iterations is  $M = 1000$ .

Table 9: Simulation results on variance estimation under DGP-L1

$N = 100$

	$V_{11}$ (truth: 3.142)			$V_{12}$ (truth: -1.097)			$V_{22}$ (truth: 1.097)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	0.172	1.273	1.285	-0.117	0.725	0.735	0.109	0.573	0.584
CBGPS	-1.109	1.298	1.707	0.113	0.421	0.436	-0.124	0.365	0.385

$N = 500$

	$V_{11}$ (truth: 3.142)			$V_{12}$ (truth: -1.097)			$V_{22}$ (truth: 1.097)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	0.025	0.458	0.458	-0.021	0.283	0.284	0.037	0.253	0.256
CBGPS	-1.043	0.318	1.091	0.038	0.212	0.215	-0.037	0.187	0.190

$N = 1000$

	$V_{11}$ (truth: 3.142)			$V_{12}$ (truth: -1.097)			$V_{22}$ (truth: 1.097)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	0.026	0.333	0.334	-0.007	0.201	0.201	0.018	0.164	0.165
CBGPS	-1.013	0.244	1.042	0.013	0.173	0.174	-0.012	0.162	0.162

DGP-L1:  $T = 1 + 0.2X_1 + \xi$  and  $Y = 1 + X_1 + T + \epsilon$ , where  $X_1 \sim N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $(K_1, K_2) = (2, 2)$  is used to compute  $\hat{\pi}_K(T, X_1)$ , which is used to estimate  $\beta$ .  $(K'_1, K'_2, M_0, K_0) = (3, 3, 3, 5)$  is used to compute  $\hat{\pi}_{K'}(T, X_1)$ , which is used to estimate  $V_{eff}$ . “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. We report the bias, standard deviation, and RMSE of each element of the variance estimator  $\hat{V}_{eff}$  across  $M = 1000$  Monte Carlo samples.

Table 10: Simulation results on variance estimation under DGP-NL1

$N = 100$

	$V_{11}$ (truth: 3.043)			$V_{12}$ (truth: $-0.118$ )			$V_{22}$ (truth: 1.074)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	$-0.427$	0.812	0.917	0.054	0.582	0.584	0.214	0.615	0.651
CBGPS	$-0.273$	1.149	1.181	0.381	0.713	0.809	1.644	1.252	2.067

$N = 500$

	$V_{11}$ (truth: 3.043)			$V_{12}$ (truth: $-0.118$ )			$V_{22}$ (truth: 1.074)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	$-0.249$	0.439	0.505	0.078	0.260	0.271	0.027	0.188	0.190
CBGPS	$-0.205$	0.338	0.395	0.501	0.387	0.633	2.147	0.885	2.323

$N = 1000$

	$V_{11}$ (truth: 3.043)			$V_{12}$ (truth: $-0.118$ )			$V_{22}$ (truth: 1.074)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	$-0.235$	0.309	0.389	0.071	0.191	0.204	0.005	0.136	0.136
CBGPS	$-0.205$	0.229	0.307	0.507	0.268	0.574	2.172	0.660	2.270

DGP-NL1:  $T = 0.1X_1^2 + \xi$  and  $Y = X_1^2 + T + \epsilon$ , where  $X_1 \sim N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $(K_1, K_2) = (2, 3)$  is used to compute  $\hat{\pi}_K(T, X_1)$ , which is used to estimate  $\beta$ .  $(K'_1, K'_2, M_0, K_0) = (3, 3, 3, 5)$  is used to compute  $\hat{\pi}_{K'}(T, X_1)$ , which is used to estimate  $V_{eff}$ . “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. We report the bias, standard deviation, and RMSE of each element of the variance estimator  $\hat{V}_{eff}$  across  $M = 1000$  Monte Carlo samples.

Table 11: Simulation results on variance estimation under DGP-L2

$N = 100$

	$V_{11}$ (truth: 2.840)			$V_{12}$ (truth: -1.236)			$V_{22}$ (truth: 1.236)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	-0.098	1.016	1.021	0.037	0.619	0.620	0.002	0.544	0.544
CBGPS	0.091	16.173	16.173	-0.159	7.307	7.308	-0.039	3.464	3.464

$N = 500$

	$V_{11}$ (truth: 2.840)			$V_{12}$ (truth: -1.236)			$V_{22}$ (truth: 1.236)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	-0.096	0.533	0.541	0.035	0.346	0.348	-0.021	0.300	0.301
CBGPS	-0.652	0.429	0.780	0.124	0.296	0.320	-0.122	0.259	0.287

$N = 1000$

	$V_{11}$ (truth: 2.840)			$V_{12}$ (truth: -1.236)			$V_{22}$ (truth: 1.236)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	-0.098	0.429	0.440	0.029	0.285	0.287	-0.012	0.240	0.241
CBGPS	-0.582	0.423	0.720	0.072	0.283	0.292	-0.071	0.271	0.280

DGP-L2:  $T = 1 + 0.2 \sum_{j=1}^2 X_j + \xi$  and  $Y = 1 + (1/2) \sum_{j=1}^2 X_j + T + \epsilon$ , where  $X_1, X_2 \stackrel{i.i.d.}{\sim} N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $(K_1, K_2) = (2, 3)$  is used to compute  $\hat{\pi}_K(T, \mathbf{X})$ , which is used to estimate  $\beta$ .  $(K'_1, K'_2, M_0, K_0) = (3, 3, 6, 8)$  is used to compute  $\hat{\pi}_{K'}(T, \mathbf{X})$ , which is used to estimate  $V_{eff}$ . “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. We report the bias, standard deviation, and RMSE of each element of the variance estimator  $\hat{V}_{eff}$  across  $M = 1000$  Monte Carlo samples.



Table 12: Simulation results on variance estimation under DGP-NL2

$N = 100$

	$V_{11}$ (truth: 1.867)			$V_{12}$ (truth: -0.476)			$V_{22}$ (truth: 1.458)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	-0.104	0.602	0.610	0.309	0.474	0.566	-0.056	0.730	0.732
CBGPS	-0.499	0.284	0.574	0.410	0.272	0.492	-0.104	0.496	0.507

$N = 500$

	$V_{11}$ (truth: 1.867)			$V_{12}$ (truth: -0.476)			$V_{22}$ (truth: 1.458)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	-0.058	0.326	0.331	0.132	0.304	0.331	0.048	0.466	0.468
CBGPS	-0.426	1.294	1.363	0.434	0.175	0.468	0.155	0.643	0.662

$N = 1000$

	$V_{11}$ (truth: 1.867)			$V_{12}$ (truth: -0.476)			$V_{22}$ (truth: 1.458)		
	Bias	Stdev	RMSE	Bias	Stdev	RMSE	Bias	Stdev	RMSE
GOE	-0.015	0.301	0.301	0.085	0.288	0.300	0.094	0.376	0.388
CBGPS	-0.467	0.134	0.486	0.438	0.123	0.455	0.185	0.307	0.358

DGP-NL2:  $T = 0.1(\sum_{j=1}^2 X_j)^2 + \xi$  and  $Y = 1/2 + [(1/2)\sum_{j=1}^2 X_j]^2 + T + \epsilon$ , where  $X_1, X_2 \stackrel{i.i.d.}{\sim} N(0, 1)$ . “GOE” is the proposed generalized optimization estimator.  $(K_1, K_2) = (2, 6)$  is used to compute  $\hat{\pi}_K(T, \mathbf{X})$ , which is used to estimate  $\beta$ .  $(K'_1, K'_2, M_0, K_0) = (2, 10, 3, 4)$  is used to compute  $\hat{\pi}_{K'}(T, \mathbf{X})$ , which is used to estimate  $V_{eff}$ . “CBGPS” is Fong, Hazlett, and Imai’s (2018) parametric covariate balancing generalized propensity score estimator. We report the bias, standard deviation, and RMSE of each element of the variance estimator  $\hat{V}_{eff}$  across  $M = 1000$  Monte Carlo samples.

Table 13: Empirical results of Fong, Hazlett, and Imai's (2018) CBGPS approach

	Covariates	Parameter of $T_i$	Parameter of $T_i^2$
Case #1	$\mathbf{Z}_i = [T_i, T_i^2, 1]^\top$	0.088 (0.456) [-0.804, 0.981]	$-8.5 \times 10^{-6}$ ( $2.3 \times 10^{-5}$ ) [ $-5.4 \times 10^{-5}, 3.7 \times 10^{-5}$ ]
Case #2	$\mathbf{Z}_i = [T_i, T_i^2, \mathbf{SD}_i^\top]^\top$	1.333 (0.444) [0.462, 2.204]	$-8.6 \times 10^{-5}$ ( $2.0 \times 10^{-5}$ ) [ $-1.3 \times 10^{-4}, -4.6 \times 10^{-5}$ ]
Case #3	$\mathbf{Z}_i = [T_i, T_i^2, 1, \mathbf{X}_{1i}^\top]^\top$	-0.545 (0.423) [-1.373, 0.284]	$-2.2 \times 10^{-5}$ ( $2.2 \times 10^{-5}$ ) [ $-6.6 \times 10^{-5}, 2.1 \times 10^{-5}$ ]
Case #4	$\mathbf{Z}_i = [T_i, T_i^2, \mathbf{SD}_i^\top, \mathbf{X}_{1i}^\top]^\top$	-0.216 (0.422) [-1.044, 0.611]	$2.7 \times 10^{-5}$ ( $2.1 \times 10^{-5}$ ) [ $-1.4 \times 10^{-5}, 6.8 \times 10^{-5}$ ]

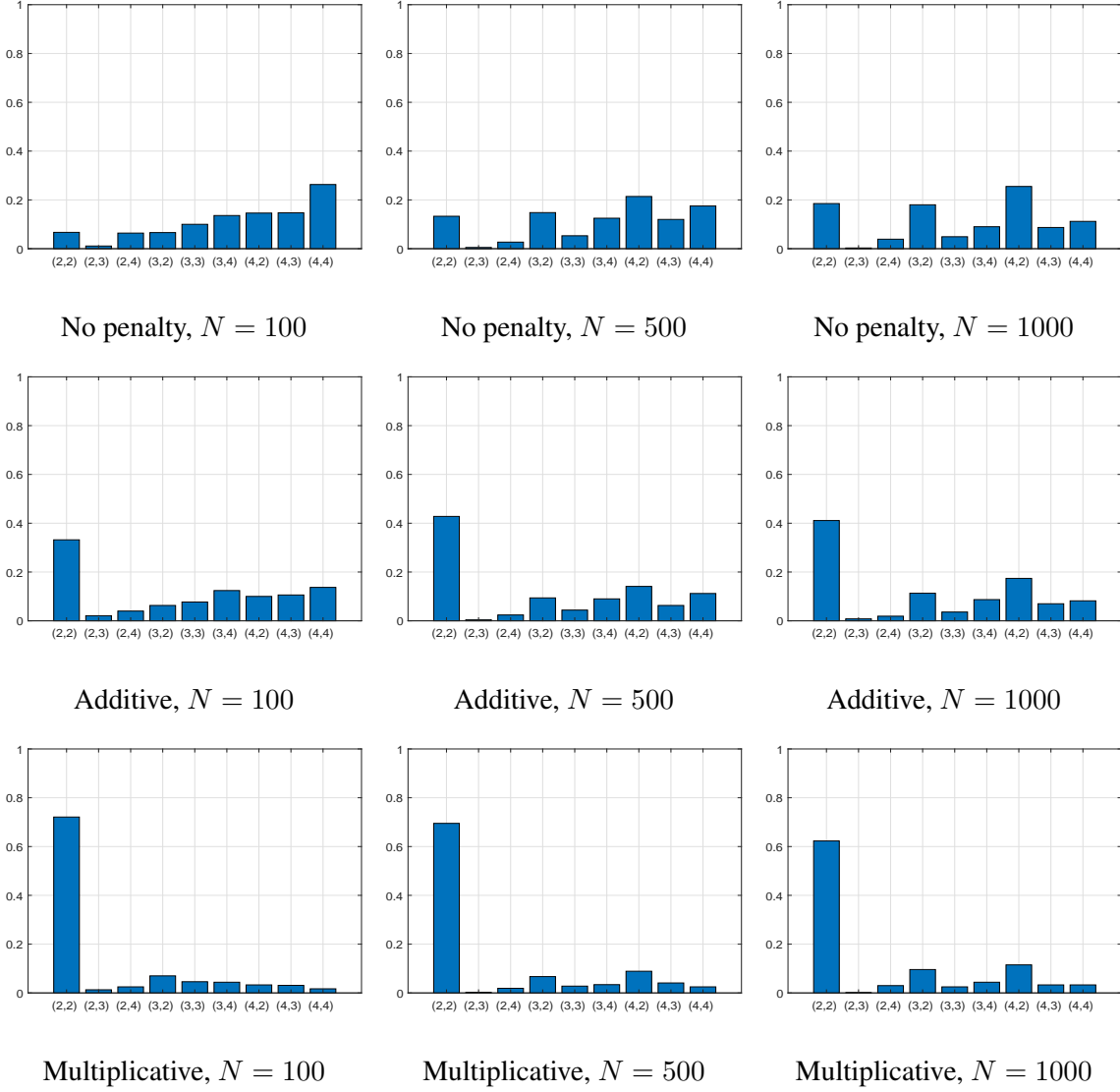
$\mathbf{X}_{1i}$  is a vector of eight covariates used in the generalized propensity score model (cf. Eq. (10.1)).  $\mathbf{SD}_i = [SD_{1i}, SD_{2i}, \dots, SD_{24i}]^\top$ , where  $SD_{ji}$  is a binary indicator that equals 1 if zip code  $i$  belongs to state  $j$  and equals 0 otherwise. Any zip code contained in the dataset belongs to one and only one of 24 states. In this table we report the CBGPS estimates for the parameters of  $T_i$  and  $T_i^2$  as well as their standard errors in round brackets and 95% confidence bands in square brackets.

Table 14: Empirical results of the generalized optimization approach

	$\beta_1$	$\beta_2$	$\beta_3$
Point estimate	22.09	$-4.7 \times 10^{-4}$	$1.5 \times 10^{-8}$
Standard error	1.214	0.001	$4.3 \times 10^{-8}$
95% confidence band	[19.71, 24.47]	[-0.002, 0.001]	$[-7.0 \times 10^{-8}, 1.0 \times 10^{-7}]$

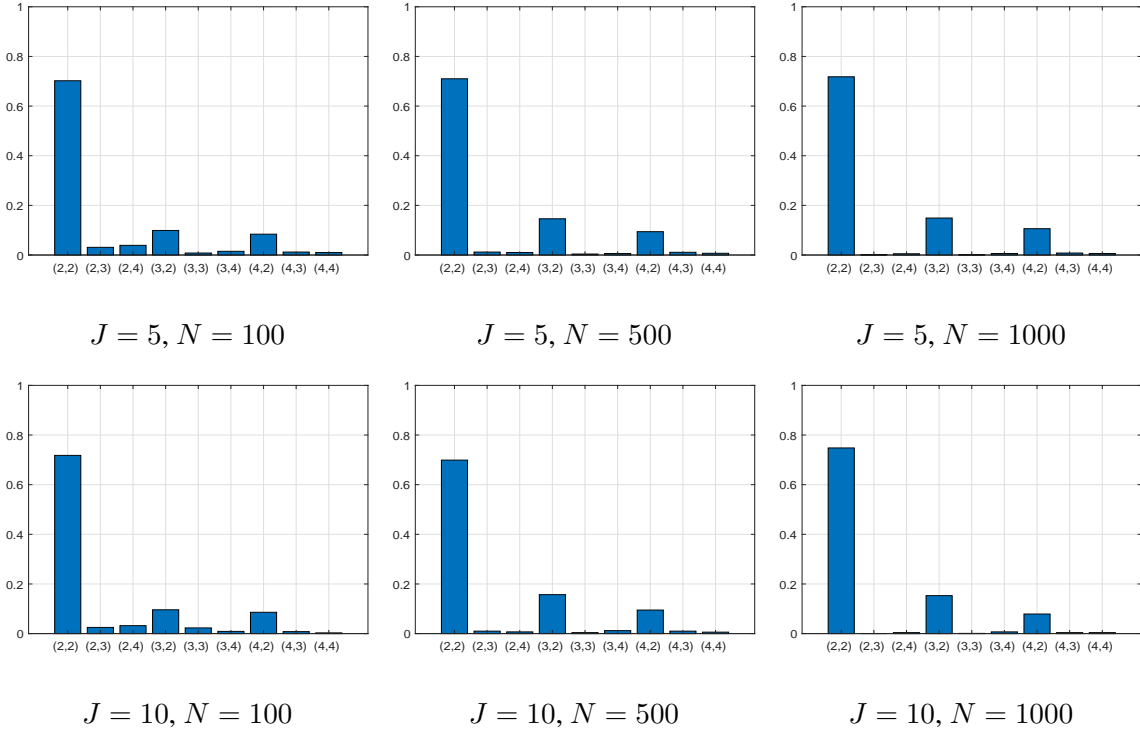
The link function is  $g(T, \beta) = \beta_1 + \beta_2 T + \beta_3 T^2$ . Covariates  $\mathbf{X}$  are defined in Eq. (10.2).

Figure 1: Share of  $(K_1, K_2)$  selected under DGP-L1 (MSE criteria)



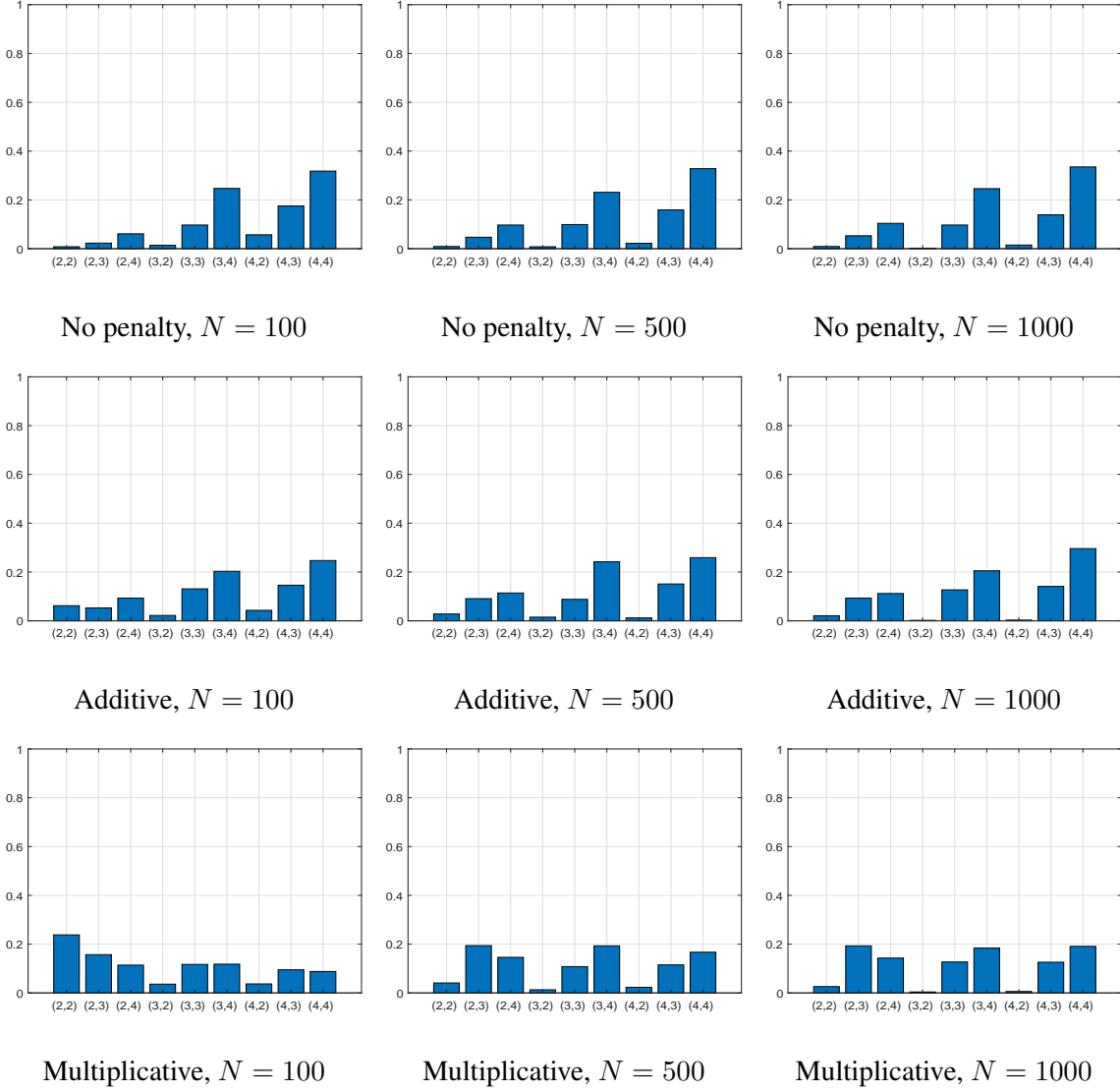
DGP-L1:  $T = 1 + X_1 + \xi$  and  $Y = 1 + X_1 + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. “No penalty” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, X_{1i})(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “Additive” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “Multiplicative” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

Figure 2: Share of  $(K_1, K_2)$  selected under DGP-L1 ( $J$ -folder cross validation)



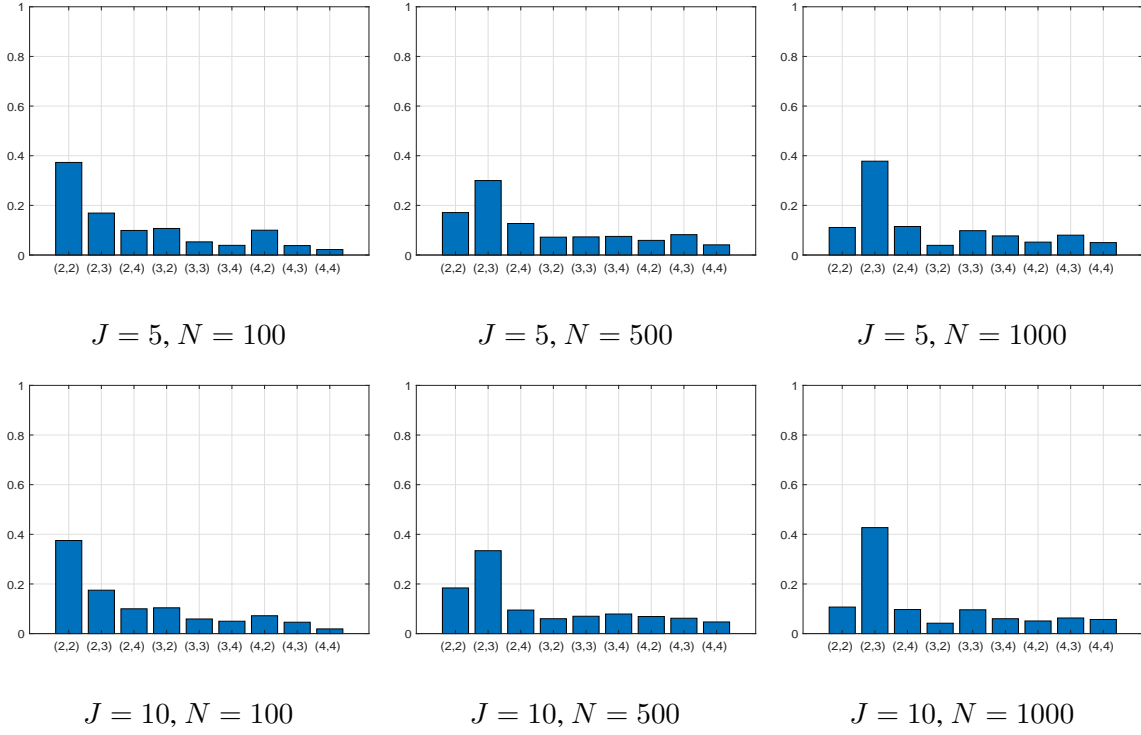
DGP-L1:  $T = 1 + X_1 + \xi$  and  $Y = 1 + X_1 + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. We pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

Figure 3: Share of  $(K_1, K_2)$  selected under DGP-NL1 (MSE criteria)



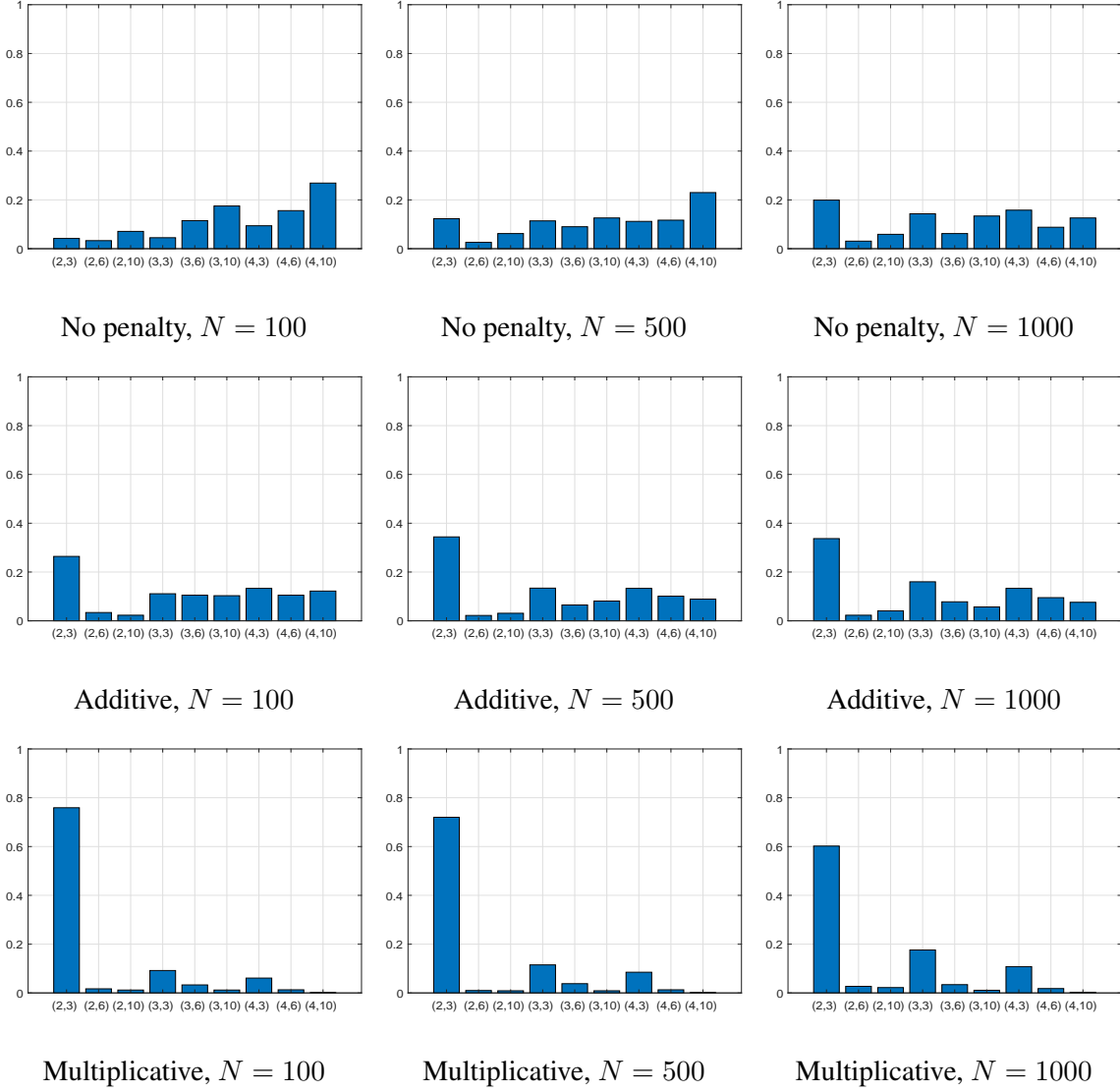
DGP-NL1:  $T = X_1^2 + \xi$  and  $Y = X_1^2 + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. "No penalty" signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, X_{1i})(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . "Additive" signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . "Multiplicative" signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

Figure 4: Share of  $(K_1, K_2)$  selected under DGP-NL1 ( $J$ -folder cross validation)



DGP-NL1:  $T = X_1^2 + \xi$  and  $Y = X_1^2 + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $X_1$ , respectively. We pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

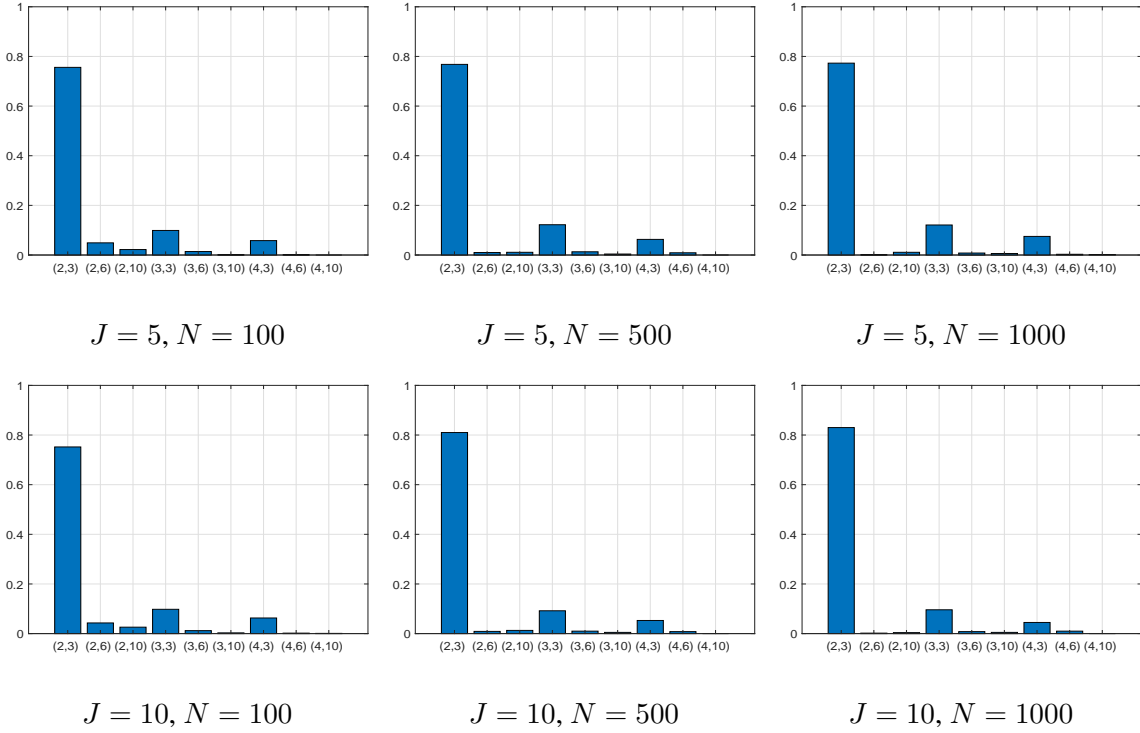
Figure 5: Share of  $(K_1, K_2)$  selected under DGP-L2 (MSE criteria)



DGP-L2:  $T = 1 + (1/2) \sum_{j=1}^2 X_j + \xi$  and  $Y = 1 + (1/2) \sum_{j=1}^2 X_j + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. “No penalty” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “Additive” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “Multiplicative” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

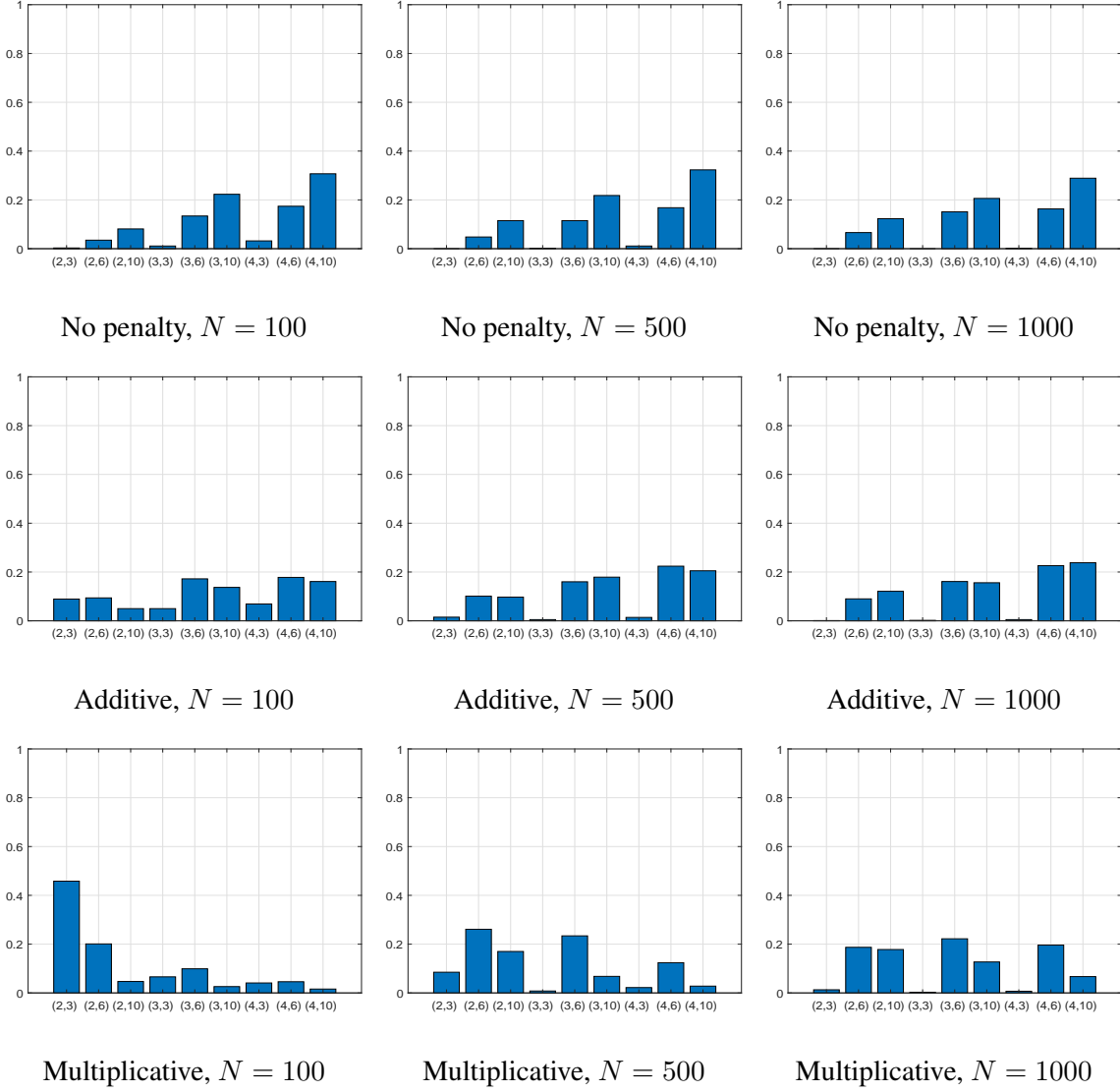


Figure 6: Share of  $(K_1, K_2)$  selected under DGP-L2 ( $J$ -folder cross validation)



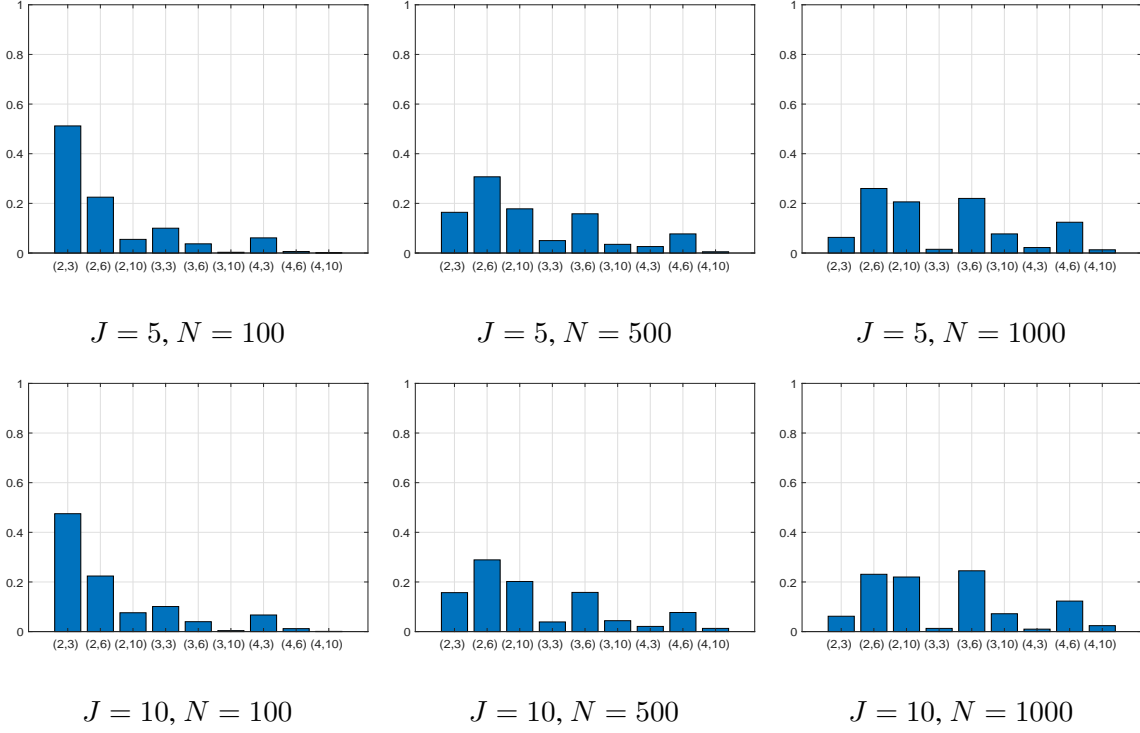
DGP-L2:  $T = 1 + (1/2) \sum_{j=1}^2 X_j + \xi$  and  $Y = 1 + (1/2) \sum_{j=1}^2 X_j + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. We pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

Figure 7: Share of  $(K_1, K_2)$  selected under DGP-NL2 (MSE criteria)



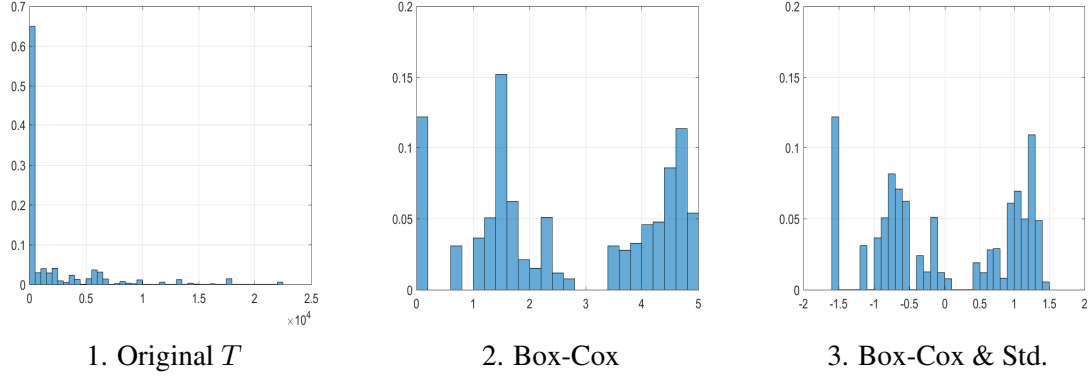
DGP-NL2:  $T = 1/2 + ((1/2) \sum_{j=1}^2 X_j)^2 + \xi$  and  $Y = 1/2 + ((1/2) \sum_{j=1}^2 X_j)^2 + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. “No penalty” signifies that we pick  $(K_1, K_2)$  that minimizes  $MSE(K_1, K_2) = N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)(Y_i - \hat{\beta}_1 - \hat{\beta}_2 T_i)^2$ . “Additive” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2(K_1 + K_2)/N) \times MSE(K_1, K_2)$ . “Multiplicative” signifies that we pick  $(K_1, K_2)$  that minimizes  $(1 + 2K_1 K_2/N) \times MSE(K_1, K_2)$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

Figure 8: Share of  $(K_1, K_2)$  selected under DGP-NL2 ( $J$ -folder cross validation)



DGP-NL2:  $T = 1/2 + ((1/2) \sum_{j=1}^2 X_j)^2 + \xi$  and  $Y = 1/2 + ((1/2) \sum_{j=1}^2 X_j)^2 + T + \epsilon$ .  $K_1$  and  $K_2$  are the dimensions of the polynomials of  $T$  and  $\mathbf{X} = (X_1, X_2)^\top$ , respectively. We pick  $(K_1, K_2)$  that minimizes the loss function of the  $J$ -folder cross validation with  $J \in \{5, 10\}$ . In this figure we plot the empirical probability of selecting each pair across  $M = 1000$  Monte Carlo samples. The choice set of  $(K_1, K_2)$  is the nine pairs put on the horizontal axis.

Figure 9: Empirical distributions of political advertisements



In this figure we draw histograms of the treatment variable studied in [Fong, Hazlett, and Imai \(2018\)](#) (i.e. the number of political advertisements aired in each zip code). Panel 1 plots the original treatment  $T$ ; Panel 2 plots  $T'$ , namely the treatment after running the Box-Cox transformation with  $\lambda = -0.16$ ; Panel 3 plots  $T^*$ , namely the standardized version of  $T'$ . For each histogram, the sum of the heights of bars is normalized to 1 so that the vertical axis is associated with empirical probability.

Supplemental Material for  
“*A Unified Framework for Efficient Estimation of  
General Treatment Models*”

Chunrong Ai\*– University of Florida  
Oliver Linton<sup>†</sup>– University of Cambridge  
Kaiji Motegi<sup>‡</sup>– Kobe University  
Zheng Zhang<sup>§</sup>– Renmin University of China

This draft: March 18, 2019

---

\*Department of Economics, University of Florida. E-mail: [chunrong.ai@warrington.ufl.edu](mailto:chunrong.ai@warrington.ufl.edu)

<sup>†</sup>Faculty of Economics, University of Cambridge. E-mail: [obl20@cam.ac.uk](mailto:obl20@cam.ac.uk)

<sup>‡</sup>Graduate School of Economics, Kobe University. E-mail: [motegi@econ.kobe-u.ac.jp](mailto:motegi@econ.kobe-u.ac.jp)

<sup>§</sup>Institute of Statistics and Big Data, Renmin University of China. E-mail: [zhengzhang@ruc.edu.cn](mailto:zhengzhang@ruc.edu.cn)

# Contents

<b>1</b>	<b>Assumptions</b>	<b>3</b>
<b>2</b>	<b>Efficiency Bound</b>	<b>4</b>
2.1	Proof of Theorem 1 . . . . .	4
2.2	Particular Case I: Binary Average Treatment Effects . . . . .	8
2.3	Particular Case II: Multiple Average Treatment Effects . . . . .	11
2.4	Particular Case III: Binary Quantile Treatment Effects . . . . .	14
<b>3</b>	<b>Convergence Rate of Estimated Stabilized Weights</b>	<b>16</b>
3.1	Lemma 3.1 . . . . .	16
3.2	Lemma 3.2 . . . . .	24
3.3	Corollary 3.3 . . . . .	31
<b>4</b>	<b>Efficient Estimation</b>	<b>33</b>
4.1	Proof of Theorem 4 . . . . .	33
4.2	Proof of (36) . . . . .	35
<b>5</b>	<b>Some Extensions</b>	<b>45</b>
5.1	Proof of Theorem 6 . . . . .	45
5.2	Proof of Theorem 7 . . . . .	51
<b>6</b>	<b>Variance Estimation in Monte Carlo Simulations</b>	<b>56</b>
6.1	Proposed Variance Estimator . . . . .	56
6.2	True Values of $V_{eff}$ in Monte Carlo Simulations . . . . .	57
6.2.1	DGP-L1 . . . . .	57
6.2.2	DGP-NL1 . . . . .	59

# 1 Assumptions

**Assumption 1.1 (Unconfounded Treatment Assignment)** For all  $t \in \mathcal{T}$ , given  $\mathbf{X}$ ,  $T$  is independent of  $Y^*(t)$ , i.e.,  $Y^*(t) \perp T | \mathbf{X}$ , for all  $t \in \mathcal{T}$ .

**Assumption 1.2** The support  $\mathcal{X}$  of  $\mathbf{X}$  is a compact subset of  $\mathbb{R}^r$ . The support  $\mathcal{T}$  of the treatment variable  $T$  is a compact subset of  $\mathbb{R}$ .

**Assumption 1.3** There exist two positive constants  $\eta_1$  and  $\eta_2$  such that

$$0 < \eta_1 \leq \pi_0(t, \mathbf{x}) \leq \eta_2 < \infty, \quad \forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}.$$

**Assumption 1.4** There exist  $\Lambda_{K_1 \times K_2} \in \mathbb{R}^{K_1 \times K_2}$  and a positive constant  $\alpha > 0$  such that

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |(\rho'^{-1}(\pi_0(t, \mathbf{x})) - u_{K_1}(t))^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})| = O(K^{-\alpha}).$$

**Assumption 1.5** For every  $K_1$  and  $K_2$ , the smallest eigenvalues of  $\mathbb{E}[u_{K_1}(T)u_{K_1}(T)^\top]$  and  $\mathbb{E}[v_{K_2}(\mathbf{X})v_{K_2}(\mathbf{X})^\top]$  are bounded away from zero uniformly in  $K_1$  and  $K_2$ .

**Assumption 1.6** There are two sequences of constants  $\zeta_1(K_1)$  and  $\zeta_2(K_2)$  satisfying

$\sup_{t \in \mathcal{T}} \|u_{K_1}(t)\| \leq \zeta_1(K_1)$  and  $\sup_{\mathbf{x} \in \mathcal{X}} \|v_{K_2}(\mathbf{x})\| \leq \zeta_2(K_2)$ ,  $K = K_1(N)K_2(N)$  and  $\zeta(K) = \zeta_1(K_1)\zeta_2(K_2)$ , such that  $\zeta(K)K^{-\alpha} \rightarrow 0$  and  $\zeta(K)\sqrt{K/N} \rightarrow 0$  as  $N \rightarrow \infty$ .

**Assumption 1.7** The parameter space  $\Theta \subset \mathbb{R}^p$  is a compact set and the true parameter  $\beta_0$  is in the interior of  $\Theta$ , where  $p \in \mathbb{N}$ .

**Assumption 1.8** There exists a unique solution  $\beta_0$  for the optimization problem

$$\min_{\beta \in \Theta} \int_{\mathcal{T}} \mathbb{E}[L(Y^*(t) - g(t; \beta))] dF_T(t).$$

**Assumption 1.9**  $\mathbb{E}[\sup_{\beta \in \Theta} |L(Y - g(T; \beta))|^2] < \infty$ .

**Assumption 1.10** The following conditions hold true:

1.  $g(t; \beta)$  is twice continuously differentiable in  $\beta \in \Theta$ ;
2.  $L(Y - g(T; \beta))$  is differentiable in  $\beta$  with probability one, i.e., for any directional vector  $\boldsymbol{\eta} \in \mathbb{R}^p$ , there exists an integrable random variable  $L'(Y - g(T; \beta))$  such that

$$\mathbb{P}\left(\lim_{\epsilon \rightarrow 0} \frac{L(Y - g(T; \beta + \epsilon \boldsymbol{\eta})) - L(Y - g(T; \beta))}{\epsilon} = L'(Y - g(T; \beta)) \cdot \langle m(T; \beta), \boldsymbol{\eta} \rangle_{\mathbb{R}^p}\right) = 1,$$

where  $\langle \cdot, \cdot \rangle_{\mathbb{R}^p}$  is the inner product in Euclidean space  $\mathbb{R}^p$ ;

$$3. \mathbb{E} [L'(Y - g(T; \beta_0))^2] < \infty.$$

**Assumption 1.11** Suppose that

$$\frac{1}{N} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) L' \left( Y_i - g(T_i; \hat{\beta}) \right) m(T_i; \hat{\beta}) = o_P(N^{-1/2})$$

holds with probability approaching one.

**Assumption 1.12**  $\mathbb{E} [\pi_0(T, \mathbf{X}) L'(Y - g(T; \beta)) m(T; \beta)]$  is differentiable with respect to  $\beta$  and  $H_0 := -\nabla_{\beta} \mathbb{E} [\pi_0(T, \mathbf{X}) L'(Y - g(T; \beta)) m(T; \beta)] \Big|_{\beta=\beta_0}$  is nonsingular.

**Assumption 1.13**  $\varepsilon(t, \mathbf{x}; \beta_0) := \mathbb{E}[L'(Y - g(T; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}]$  is continuously differentiable in  $(t, \mathbf{x})$ .

**Assumption 1.14**

1.  $\mathbb{E} [\sup_{\beta \in \Theta} |L'(Y - g(T; \beta))|^{2+\delta}] < \infty$  for some  $\delta > 0$ ;
2. The function class  $\{L'(y - g(t; \beta)) : \beta \in \Theta\}$  satisfies:

$$\mathbb{E} \left[ \sup_{\beta_1: \|\beta_1 - \beta\| < \delta} |L'(Y - g(T; \beta_1)) - L'(Y - g(T; \beta))|^2 \right]^{1/2} \leq a \cdot \delta^b$$

for any  $\forall \beta \in \Theta$  and any small  $\delta > 0$  and for some finite positive constants  $a$  and  $b$ .

**Assumption 1.15**  $\zeta(K) \sqrt{K^4/N} \rightarrow 0$  and  $\sqrt{N} K^{-\alpha} \rightarrow 0$  as  $N \rightarrow \infty$ .

## 2 Efficiency Bound

### 2.1 Proof of Theorem 1

Without loss of generality, we only consider the distribution of  $(T, \mathbf{X}, Y)$  to be absolutely continuous with respect to Lebesgue measure, i.e., there exists a density function  $f_{T,X,Y}(t, \mathbf{x}, y)$  such that  $dF_{T,X,Y}(t, \mathbf{x}, y) = f_{T,X,Y}(t, \mathbf{x}, y) dt d\mathbf{x} dy$ . For discrete cases, the proof can be established by using a similar argument.



We follow the approach of [Bickel, Klaassen, Ritov, and Wellner \(1993, Section 3.3\)](#) to derive the variance bound of  $\beta_0$ , see also [Tchetgen Tchetgen and Shpitser \(2012\)](#). Let  $\{f_{Y,T,X}^\alpha(y, t, \mathbf{x})\}_{\alpha \in \mathbb{R}}$  denote a one dimensional regular parametric submodel with  $f_{Y,T,X}^{\alpha=0}(y, t, \mathbf{x}) = f_{Y,T,X}(y, t, \mathbf{x})$ . By definition,  $\beta_0$  solves following equation:

$$\int_{\mathcal{T}} \mathbb{E}[m(t; \beta_0) L'(Y^*(t) - g(t; \beta_0))] f_T(t) dt = 0. \quad (1)$$

By Assumption 1.1, (1) is equivalent to

$$\int_{\mathcal{T}} \int_{\mathcal{X}} \mathbb{E}[m(T; \beta_0) L'(Y - g(T; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] f_X(\mathbf{x}) f_T(t) d\mathbf{x} dt = 0.$$

Therefore, the parameter  $\beta(\alpha)$  induced by the submodel  $f_{Y,T,X}^\alpha(y, t, \mathbf{x})$  satisfies:

$$\int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta(\alpha)) \cdot \mathbb{E}^\alpha[L'(Y - g(t; \beta(\alpha))) | T = t, \mathbf{X} = \mathbf{x}] f_T^\alpha(t) f_X^\alpha(\mathbf{x}) d\mathbf{x} dt = 0, \quad (2)$$

where  $\mathbb{E}^\alpha[\cdot | T = t, \mathbf{X} = \mathbf{x}]$  denotes taking expectation with respect to the submodel  $f_{Y,T,X}^\alpha(\cdot | t, \mathbf{x})$ .

Differentiating both sides of (2) with respect to  $\alpha$ , evaluating at  $\alpha = 0$  and using the condition  $Y^*(t) \perp T | \mathbf{X}$ , we can deduce that

$$\begin{aligned} 0 &= \int_{\mathcal{T}} \int_{\mathcal{X}} \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \{m(t; \beta(\alpha)) \mathbb{E}^\alpha[L'(Y - g(t; \beta(\alpha))) | T = t, \mathbf{X} = \mathbf{x}] f_T^\alpha(t) f_X^\alpha(\mathbf{x})\} d\mathbf{x} dt \\ &= \int_{\mathcal{T}} \int_{\mathcal{X}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] f_T(t) f_X(\mathbf{x}) \nabla_\beta m(t; \beta_0) d\mathbf{x} dt \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) \\ &\quad + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_X^\alpha(\mathbf{x}) \Big|_{\alpha=0} f_T(t) d\mathbf{x} dt \\ &\quad + \int_{\mathcal{Y} \times \mathcal{X} \times \mathcal{T}} m(t; \beta_0) L'(y - g(t; \beta_0)) \cdot \frac{\partial}{\partial \alpha} f_{Y|T,X}^\alpha(y | t, \mathbf{x}) \Big|_{\alpha=0} f_X(\mathbf{x}) f_T(t) dy d\mathbf{x} dt \\ &\quad + \int_{\mathcal{X} \times \mathcal{T}} m(t; \beta_0) \cdot \nabla_\beta \mathbb{E}[L'(Y^*(t) - g(t; \beta)) | T = t, \mathbf{X} = \mathbf{x}] \Big|_{\beta=\beta_0} \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) \cdot f_T(t) f_X(\mathbf{x}) d\mathbf{x} dt \\ &\quad + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_T^\alpha(t) \Big|_{\alpha=0} f_X(\mathbf{x}) d\mathbf{x} dt \\ &= \int_{\mathcal{T}} \int_{\mathcal{X}} \mathbb{E}[L'(Y^*(t) - g(t; \beta_0)) | \mathbf{X} = \mathbf{x}] f_T(t) f_X(\mathbf{x}) \nabla_\beta m(t; \beta_0) d\mathbf{x} dt \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) \\ &\quad + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_X^\alpha(\mathbf{x}) \Big|_{\alpha=0} f_T(t) d\mathbf{x} dt \end{aligned}$$

$$\begin{aligned}
& + \int_{\mathcal{Y} \times \mathcal{X} \times \mathcal{T}} m(t; \beta_0) L'(y - g(t; \beta_0)) \cdot \frac{\partial}{\partial \alpha} f_{Y|T,X}^\alpha(y|t, \mathbf{x}) \Big|_{\alpha=0} f_X(\mathbf{x}) f_T(t) dy d\mathbf{x} dt \\
& + \int_{\mathcal{X} \times \mathcal{T}} m(t; \beta_0) \cdot \nabla_\beta \mathbb{E}[L'(Y^*(t) - g(t; \beta)) | \mathbf{X} = \mathbf{x}] \Big|_{\beta=\beta_0} \cdot f_T(t) f_X(\mathbf{x}) d\mathbf{x} dt \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) \\
& + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_T^\alpha(t) \Big|_{\alpha=0} f_X(\mathbf{x}) d\mathbf{x} dt \\
& = \int_{\mathcal{T}} \mathbb{E}[L'(Y^*(t) - g(t; \beta_0))] \cdot f_T(t) \nabla_\beta m(t; \beta_0) dt \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) \\
& + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_X^\alpha(\mathbf{x}) \Big|_{\alpha=0} f_T(t) d\mathbf{x} dt \\
& + \int_{\mathcal{Y} \times \mathcal{X} \times \mathcal{T}} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) \cdot \frac{\partial}{\partial \alpha} f_{Y|T,X}^\alpha(y|t, \mathbf{x}) \Big|_{\alpha=0} f_X(\mathbf{x}) f_T(t) dy d\mathbf{x} dt \\
& + \int_{\mathcal{T}} \nabla_\beta \mathbb{E}[L'(Y^*(t) - g(t; \beta))] \Big|_{\beta=\beta_0} m(t; \beta_0) \cdot f_T(t) dt \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) \\
& + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_T^\alpha(t) \Big|_{\alpha=0} f_X(\mathbf{x}) d\mathbf{x} dt \\
& = \nabla_\beta \left\{ \int_{\mathcal{T}} \mathbb{E}[L'(Y^*(t) - g(t; \beta))] \cdot m(t; \beta) f_T(t) dt \right\} \Big|_{\beta=\beta_0} \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) \\
& + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_X^\alpha(\mathbf{x}) \Big|_{\alpha=0} f_T(t) d\mathbf{x} dt \\
& + \int_{\mathcal{Y} \times \mathcal{X} \times \mathcal{T}} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) \cdot \frac{\partial}{\partial \alpha} f_{Y|T,X}^\alpha(y|t, \mathbf{x}) \Big|_{\alpha=0} f_X(\mathbf{x}) f_T(t) dy d\mathbf{x} dt \\
& + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_T^\alpha(t) \Big|_{\alpha=0} f_X(\mathbf{x}) d\mathbf{x} dt.
\end{aligned}$$

Since  $H_0 = -\nabla_\beta \left\{ \int_{\mathcal{T}} \mathbb{E}[L'(Y^*(t) - g(t; \beta))] \cdot m(t; \beta) f_T(t) dt \right\} \Big|_{\beta=\beta_0}$  is invertible by Assumption 1.9, we get

$$\begin{aligned}
\frac{\partial}{\partial \alpha} \Big|_{\alpha=0} \beta(\alpha) & = H_0^{-1} \cdot \left\{ \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_X^\alpha(\mathbf{x}) \Big|_{\alpha=0} f_T(t) d\mathbf{x} dt \right. \\
& + \int_{\mathcal{Y} \times \mathcal{X} \times \mathcal{T}} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) \cdot \frac{\partial}{\partial \alpha} f_{Y|T,X}^\alpha(y|t, \mathbf{x}) \Big|_{\alpha=0} f_X(\mathbf{x}) f_T(t) dy d\mathbf{x} dt \\
& \left. + \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E}[L'(Y - g(t; \beta_0)) | T = t, \mathbf{X} = \mathbf{x}] m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} f_T^\alpha(t) \Big|_{\alpha=0} f_X(\mathbf{x}) d\mathbf{x} dt \right\}.
\end{aligned}$$

The efficient influence function of  $\beta_0$ , denoted by  $S_{eff}(Y, T, \mathbf{X}; \beta_0)$ , is a unique function satisfying

the following equation:

$$\left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} \beta(\alpha) = \mathbb{E} \left[ S_{eff}(Y, T, \mathbf{X}; \beta_0) \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} \log f_{Y,X,T}^\alpha(Y, \mathbf{X}, T) \right]. \quad (3)$$

Therefore, to justify our theorem, it suffices to substitute  $S_{eff}(Y, T, \mathbf{X}; \beta_0) = H_0^{-1} \psi(Y, T, \mathbf{X}; \beta_0)$  into (3) and check the validity. Note that

$$\begin{aligned} & \mathbb{E} \left[ S_{eff}(Y, T, \mathbf{X}; \beta_0) \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} \log f_{Y,X,T}^\alpha(Y, \mathbf{X}, T) \right] \\ &= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} \psi(y, t, \mathbf{x}; \beta_0) \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f_{Y|X,T}^\alpha(y|\mathbf{x}, t) f_{T,X}(t, \mathbf{x}) dy d\mathbf{x} dt \end{aligned} \quad (4)$$

$$+ H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} \psi(y, t, \mathbf{x}; \beta_0) f_{Y|X,T}(y|\mathbf{x}, t) \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f_{T|X}^\alpha(t|\mathbf{x}) f_X(\mathbf{x}) dy d\mathbf{x} dt \quad (5)$$

$$+ H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} \psi(y, t, \mathbf{x}; \beta_0) f_{Y|X,T}(y|\mathbf{x}, t) f_{T|X}(t|\mathbf{x}) \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f_X^\alpha(\mathbf{x}) dy d\mathbf{x} dt. \quad (6)$$

For the term (4), we have

$$\begin{aligned} (4) &= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} \left\{ \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) - \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot \varepsilon(t, \mathbf{x}; \beta_0) \right. \\ &\quad \left. + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X} = \mathbf{x}] + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | T = t] \right\} \\ &\quad \times \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f_{Y|X,T}^\alpha(y|\mathbf{x}, t) f_{T,X}(t, \mathbf{x}) dy d\mathbf{x} dt \\ &= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) \cdot \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f_{Y|X,T}^\alpha(y|\mathbf{x}, t) f_{T,X}(t, \mathbf{x}) dy d\mathbf{x} dt \\ &= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) \cdot \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f_{Y|X,T}^\alpha(y|\mathbf{x}, t) f_T(t) f_X(\mathbf{x}) dy d\mathbf{x} dt. \end{aligned}$$

For the term (5), we have

$$\begin{aligned} (5) &= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} \left\{ \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) - \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot \varepsilon(t, \mathbf{x}; \beta_0) \right. \\ &\quad \left. + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X} = \mathbf{x}] + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | T = t] \right\} \\ &\quad \times f_{Y|X,T}(y|\mathbf{x}, t) \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f_{T|X}^\alpha(t|\mathbf{x}) f_X(\mathbf{x}) dy d\mathbf{x} dt \\ &= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \left\{ \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X} = \mathbf{x}] + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | T = t] \right\} \end{aligned}$$

$$\begin{aligned}
& \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_{T|X}^\alpha(t|\mathbf{x}) f_X(\mathbf{x}) d\mathbf{x} dt \\
&= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | T = t] \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_{T|X}^\alpha(t|\mathbf{x}) f_X(\mathbf{x}) d\mathbf{x} dt \\
&= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | T = t] \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_T^\alpha(t) dt \\
&= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \varepsilon(t, \mathbf{x}; \beta_0) \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_T^\alpha(t) \cdot f_{X|T}(\mathbf{x}|t) d\mathbf{x} dt \\
&= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \varepsilon(t, \mathbf{x}; \beta_0) m(t; \beta_0) \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_T^\alpha(t) \cdot f_X(\mathbf{x}) d\mathbf{x} dt,
\end{aligned}$$

where the first equality holds in accordance with the definition of  $\int_{\mathcal{Y}} L'(y - g(t; \beta_0)) f_{Y|X,T}(y|\mathbf{x}, t) dy =: \varepsilon(t, \mathbf{x}; \beta_0)$ .

For the term (6), we have

$$\begin{aligned}
(6) &= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T} \times \mathcal{Y}} \left\{ \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot L'(y - g(t; \beta_0)) - \frac{f_T(t)}{f_{T|X}(t|\mathbf{x})} m(t; \beta_0) \cdot \varepsilon(t, \mathbf{x}; \beta_0) \right. \\
&\quad \left. + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X} = \mathbf{x}] + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | T = t] \right\} \\
&\quad \times f_{Y|X,T}(y|\mathbf{x}, t) f_{T|X}(t|\mathbf{x}) \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_X^\alpha(\mathbf{x}) dy d\mathbf{x} dt \\
&= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \left\{ \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X} = \mathbf{x}] + \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | T = t] \right\} \\
&\quad \times f_{T|X}(t|\mathbf{x}) \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_X^\alpha(\mathbf{x}) d\mathbf{x} dt \\
&= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X} = \mathbf{x}] \cdot f_{T|X}(t|\mathbf{x}) \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_X^\alpha(\mathbf{x}) d\mathbf{x} dt \\
&= H_0^{-1} \int_{\mathcal{X}} \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0) \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X} = \mathbf{x}] \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_X^\alpha(\mathbf{x}) d\mathbf{x} \\
&= H_0^{-1} \int_{\mathcal{X} \times \mathcal{T}} \varepsilon(t, \mathbf{x}; \beta_0) m(t; \beta_0) \cdot f_T(t) \cdot \frac{\partial}{\partial \alpha} \Big|_{\alpha=0} f_X^\alpha(\mathbf{x}) d\mathbf{x} dt.
\end{aligned}$$

We have proved (3) holds, hence  $S_{eff}$  is the efficient influence function of  $\beta_0$ .

## 2.2 Particular Case I: Binary Average Treatment Effects

In this section, we show that when  $T \in \{0, 1\}$ ,  $g(t; \beta) = \beta_0 + \beta_1 \cdot t$  and  $L(v) = v^2$ , our general efficiency bound derived in Theorem 3.1 reduces to the well-known efficiency bound for average treatment effects in [Robins, Rotnitzky, and Zhao \(1994\)](#) and [Hahn \(1998\)](#). In accordance with

our identification condition,  $\beta_0$  and  $\beta_1$  are identified by minimizing the following loss function

$$\sum_{t \in \{0,1\}} \mathbb{E}[(Y^*(t) - \beta_0 - \beta_1 \cdot t)^2] \cdot \mathbb{P}(T = t).$$

The solutions are given by

$$\beta_0 = \mathbb{E}[Y^*(0)], \quad \beta_1 = \mathbb{E}[Y^*(1) - Y^*(0)].$$

Here  $\beta_1$  is the average treatment effects.

**Corollary 2.1** *Suppose  $T \in \{0, 1\}$ ,  $L(v) = v^2$ ,  $g(t; \boldsymbol{\beta}) = \beta_0 + \beta_1 \cdot t$  and the conditions in Theorem 3.1 hold, the efficient influence functions of  $\beta_0$  and  $\beta_1$  given by Theorem 3.1 reduce to*

$$\begin{aligned} S_{eff}(T, \mathbf{X}, Y; \beta_0) &= \phi_2(T, \mathbf{X}, Y; \beta_0), \\ S_{eff}(T, \mathbf{X}, Y; \beta_1, \beta_0) &= \phi_2(T, \mathbf{X}, Y; \beta_0) - \phi_1(T, \mathbf{X}, Y; \beta_1, \beta_0), \end{aligned}$$

where

$$\begin{aligned} \phi_1(T, \mathbf{X}, Y; \boldsymbol{\beta}) &= \frac{T}{\mathbb{P}(T = 1|\mathbf{X})} \cdot Y^*(1) - \left\{ \frac{T}{\mathbb{P}(T = 1|\mathbf{X})} - 1 \right\} \cdot \mathbb{E}[Y^*(1)|\mathbf{X}] - \beta_0 - \beta_1, \\ \phi_2(T, \mathbf{X}, Y; \boldsymbol{\beta}) &= \frac{1 - T}{\mathbb{P}(T = 0|\mathbf{X})} \cdot Y^*(0) - \left\{ \frac{1 - T}{\mathbb{P}(T = 0|\mathbf{X})} - 1 \right\} \cdot \mathbb{E}[Y^*(0)|\mathbf{X}] - \beta_0, \end{aligned}$$

and they are the same as the efficient influence functions given in [Robins, Rotnitzky, and Zhao \(1994\)](#) and [Hahn \(1998\)](#).

**Proof.** Using our notation, we have

$$\begin{aligned} \boldsymbol{\beta}_0 &= (\beta_0, \beta_1)^\top, \quad g(t; \boldsymbol{\beta}_0) = \beta_0 + \beta_1 \cdot t, \quad m(t; \boldsymbol{\beta}_0) = \begin{bmatrix} 1 \\ t \end{bmatrix}, \quad H_0 = \mathbb{E} [m(T; \boldsymbol{\beta}_0)m(T; \boldsymbol{\beta}_0)^\top], \\ \varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) &= T \cdot \{\mathbb{E}[Y^*(1) - Y^*(0)|\mathbf{X}] - \beta_1\} + \mathbb{E}[Y^*(0)|\mathbf{X}] - \beta_0, \\ \pi_0(T, \mathbf{X}) &= \frac{T \cdot p + (1 - T) \cdot q}{T \cdot \mathbb{P}(T = 1|\mathbf{X}) + T \cdot \mathbb{P}(T = 0|\mathbf{X})} = \frac{T}{\mathbb{P}(T = 1|\mathbf{X})} \cdot p + \frac{1 - T}{\mathbb{P}(T = 0|\mathbf{X})} \cdot q, \end{aligned}$$

where  $p = \mathbb{P}(T = 1)$  and  $q = \mathbb{P}(T = 0)$ . In accordance with our Theorem 3.1, the efficient influence function of  $(\beta_0, \beta_1)$  is

$$H_0^{-1} \left\{ \pi_0(T, \mathbf{X})m(T; \boldsymbol{\beta}_0) \{Y - \mathbb{E}[Y|\mathbf{X}, T]\} + \mathbb{E} [\varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0)\pi_0(T, \mathbf{X})m(T; \boldsymbol{\beta}_0)|\mathbf{X}] \right\}.$$

With some computation, we have

$$H_0^{-1} = \begin{bmatrix} 1 & p \\ p & p \end{bmatrix}^{-1} = \frac{1}{pq} \cdot \begin{bmatrix} p & -p \\ -p & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{q} & -\frac{1}{q} \\ -\frac{1}{q} & \frac{1}{pq} \end{bmatrix}. \quad (7)$$

and

$$\begin{aligned} & \pi_0(T, \mathbf{X})m(T; \beta_0) \{Y - \mathbb{E}[Y|\mathbf{X}, T]\} \\ &= \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p \cdot \begin{bmatrix} 1 \\ T \end{bmatrix} \cdot \left\{ Y - T \cdot \mathbb{E}[Y^*(1)|\mathbf{X}] - (1-T) \cdot \mathbb{E}[Y^*(0)|\mathbf{X}] \right\} \\ & \quad + \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \cdot \begin{bmatrix} 1 \\ T \end{bmatrix} \cdot \left\{ Y - T \cdot \mathbb{E}[Y^*(1)|\mathbf{X}] - (1-T) \cdot \mathbb{E}[Y^*(0)|\mathbf{X}] \right\} \\ &= \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} \cdot \left\{ Y^*(1) - \mathbb{E}[Y^*(1)|\mathbf{X}] \right\} + \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} \cdot \left\{ Y^*(0) - \mathbb{E}[Y^*(0)|\mathbf{X}] \right\} \\ &= \begin{bmatrix} \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot \{Y^*(1) - \mathbb{E}[Y^*(1)|\mathbf{X}]\} \cdot p + \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot \{Y^*(0) - \mathbb{E}[Y^*(0)|\mathbf{X}]\} \cdot q \\ \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot \{Y^*(1) - \mathbb{E}[Y^*(1)|\mathbf{X}]\} \cdot p \end{bmatrix} \end{aligned} \quad (8)$$

and

$$\begin{aligned} & \mathbb{E}[\varepsilon(T, \mathbf{X}; \beta_0)\pi_0(T, \mathbf{X})m(T; \beta_0)|\mathbf{X}] \\ &= \mathbb{E} \left[ \left( T \cdot \{\mathbb{E}[Y^*(1) - Y^*(0)|\mathbf{X}] - \beta_1\} + \mathbb{E}[Y^*(0)|\mathbf{X}] - \beta_0 \right) \cdot \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p \cdot \begin{bmatrix} 1 \\ T \end{bmatrix} \middle| \mathbf{X} \right] \\ & \quad + \mathbb{E} \left[ \left( T \cdot \{\mathbb{E}[Y^*(1) - Y^*(0)|\mathbf{X}] - \beta_1\} + \mathbb{E}[Y^*(0)|\mathbf{X}] - \beta_0 \right) \cdot \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \cdot \begin{bmatrix} 1 \\ T \end{bmatrix} \middle| \mathbf{X} \right] \\ &= \mathbb{E} \left[ \left( \mathbb{E}[Y^*(1)|\mathbf{X}] - \beta_1 - \beta_0 \right) \cdot \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} \middle| \mathbf{X} \right] \\ & \quad + \mathbb{E} \left[ \left( \mathbb{E}[Y^*(0)|\mathbf{X}] - \beta_0 \right) \cdot \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} \middle| \mathbf{X} \right] \\ &= \begin{bmatrix} \left( \mathbb{E}[Y^*(1)|\mathbf{X}] - \beta_1 - \beta_0 \right) \cdot p + \left( \mathbb{E}[Y^*(0)|\mathbf{X}] - \beta_1 \right) \cdot q \\ \left( \mathbb{E}[Y^*(1)|\mathbf{X}] - \beta_1 - \beta_0 \right) \cdot p \end{bmatrix}. \end{aligned} \quad (9)$$

Therefore, with (7), (8), and (9) we can obtain that

$$\begin{aligned} & \pi_0(T, \mathbf{X})m(T; \beta_0) \{Y - \mathbb{E}[Y|\mathbf{X}, T]\} + \mathbb{E}[\varepsilon(T, \mathbf{X}; \beta_0)\pi_0(T, \mathbf{X})m(T; \beta_0)|\mathbf{X}] \\ &= \begin{pmatrix} p \cdot \phi_1(T, \mathbf{X}, Y; \beta_0) + q \cdot \phi_2(T, \mathbf{X}, Y; \beta_0) \\ p \cdot \phi_1(T, \mathbf{X}, Y; \beta_0) \end{pmatrix}, \end{aligned}$$

and the efficient influence functions of  $\beta_1$  and  $\beta_2$  are given by

$$\begin{bmatrix} \frac{1}{q} & -\frac{1}{q} \\ -\frac{1}{q} & \frac{1}{pq} \end{bmatrix} \cdot \begin{pmatrix} p \cdot \phi_1(T, \mathbf{X}, Y; \beta) + q \cdot \phi_2(T, \mathbf{X}, Y; \beta) \\ p \cdot \phi_1(T, \mathbf{X}, Y; \beta) \end{pmatrix} = \begin{pmatrix} \phi_2(T, \mathbf{X}, Y; \beta) \\ \phi_1(T, \mathbf{X}, Y; \beta) - \phi_2(T, \mathbf{X}, Y; \beta) \end{pmatrix}.$$

■

### 2.3 Particular Case II: Multiple Average Treatment Effects

In this section, we show that when  $T \in \{0, 1, \dots, J\}$ ,  $J \in \mathbb{N}$ ,  $g(t; \beta) = \sum_{j=0}^J \beta_j \cdot I(t = j)$  and  $L(v) = v^2$ , our general efficiency bound derived in Theorem 3.1 reduces to the efficiency bound of multi-level treatment effects given in Cattaneo (2010). In accordance with our proposed identification condition,  $\{\beta_j\}_{j=0}^J$  are identified by minimizing the following loss function

$$\sum_{j=0}^J \mathbb{E}[(Y^*(j) - \beta_j)^2] \cdot \mathbb{P}(T = j).$$

The solutions are  $\beta_j = \mathbb{E}[Y^*(j)]$  for  $j \in \{0, \dots, J\}$ .

**Corollary 2.2** *Suppose  $T \in \{0, 1, \dots, J\}$ ,  $J \in \mathbb{N}$ ,  $g(t; \beta) = \sum_{j=0}^J \beta_j \cdot I(t = j)$ ,  $L(v) = v^2$ , and the conditions in Theorem 3.1 hold, the efficient influence functions of  $\{\beta_j\}_{j=0}^J$  given by Theorem 3.1 reduce to*

$$S_{eff}(T, \mathbf{X}, Y; \beta_j) = \frac{I(T = j)}{\mathbb{P}(T = j|\mathbf{X})} \cdot \{Y^*(j) - \mathbb{E}[Y^*(j)|\mathbf{X}]\} + \mathbb{E}[Y^*(j)|\mathbf{X}] - \beta_j, \quad j \in \{0, \dots, J\},$$

and they are the same as the efficient influence functions given in Cattaneo (2010).

**Proof.** Using our notation, we have

$$\boldsymbol{\beta}_0 = (\beta_0, \dots, \beta_J)^\top, \quad g(t; \boldsymbol{\beta}_0) = \sum_{j=0}^J \beta_j \cdot I(t = j), \quad m(t; \boldsymbol{\beta}_0) = \begin{bmatrix} I(t = 0) \\ I(t = 1) \\ \vdots \\ I(t = J) \end{bmatrix}, \quad H_0 = \mathbb{E} [m(T; \boldsymbol{\beta}_0) m(T; \boldsymbol{\beta}_0)^\top].$$

Then

$$\begin{aligned} \varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) &= \mathbb{E}[Y|T, X] - g(T; \boldsymbol{\beta}_0) \\ &= \sum_{j=0}^J \mathbb{E}[Y^*(j)|X] \cdot I(T = j) - \sum_{j=0}^J \beta_j \cdot I(T = j) \\ &= \sum_{j=0}^J (\mathbb{E}[Y^*(j)|X] - \beta_j) \cdot I(T = j) \end{aligned}$$

and

$$\pi_0(T, X) = \sum_{j=0}^J \frac{I(T = j)}{\mathbb{P}(T = j|\mathbf{X})} \cdot p_j, \quad \text{where } p_j = \mathbb{P}(T = j).$$

Then we have

$$H_0^{-1} = \mathbb{E} [m(T; \boldsymbol{\beta}_0) m(T; \boldsymbol{\beta}_0)^\top]^{-1} = \begin{bmatrix} p_0^{-1} & & & \\ & p_1^{-1} & & \\ & \dots & & \\ & & & p_J^{-1} \end{bmatrix},$$

and

$$\begin{aligned} &\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \{Y - \mathbb{E}[Y|\mathbf{X}, T]\} \\ &= \left\{ \sum_{j=0}^J \frac{I(T = j)}{\mathbb{P}(T = j|\mathbf{X})} \cdot p_j \right\} \cdot \begin{bmatrix} I(T = 0) \\ I(T = 1) \\ \vdots \\ I(T = J) \end{bmatrix} \cdot \left\{ Y - \sum_{j=0}^J I(T = j) \cdot \mathbb{E}[Y^*(j)|\mathbf{X}] \right\} \end{aligned}$$



$$\begin{aligned}
&= \begin{bmatrix} I(T=0) \\ I(T=1) \\ \vdots \\ I(T=J) \end{bmatrix} \left\{ \sum_{j=0}^J \frac{I(T=j)}{\mathbb{P}(T=j|\mathbf{X})} \cdot p_j \cdot Y^*(j) - \sum_{j=0}^J \frac{I(T=j)}{\mathbb{P}(T=j|\mathbf{X})} \cdot p_j \cdot \mathbb{E}[Y^*(j)|\mathbf{X}] \right\} \\
&= \begin{bmatrix} \frac{I(T=0)}{\mathbb{P}(T=0|\mathbf{X})} \cdot p_0 \cdot \{Y^*(0) - \mathbb{E}[Y^*(0)|\mathbf{X}]\} \\ \frac{I(T=1)}{\mathbb{P}(T=1|\mathbf{X})} \cdot p_1 \cdot \{Y^*(1) - \mathbb{E}[Y^*(1)|\mathbf{X}]\} \\ \vdots \\ \frac{I(T=J)}{\mathbb{P}(T=J|\mathbf{X})} \cdot p_J \cdot \{Y^*(j) - \mathbb{E}[Y^*(j)|\mathbf{X}]\} \end{bmatrix} \tag{10}
\end{aligned}$$

and

$$\begin{aligned}
&\varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \\
&= \left\{ \sum_{j=0}^J (\mathbb{E}[Y^*(j)|X] - \beta_j) \cdot I(T=j) \right\} \left\{ \sum_{j=0}^J \frac{I(T=j)}{\mathbb{P}(T=j|\mathbf{X})} \cdot p_j \right\} \begin{bmatrix} I(T=0) \\ I(T=1) \\ \vdots \\ I(T=J) \end{bmatrix} \\
&= \begin{bmatrix} \frac{I(T=0)}{\mathbb{P}(T=0|\mathbf{X})} \cdot p_0 \cdot \{\mathbb{E}[Y^*(0)|X] - \beta_0\} \\ \frac{I(T=1)}{\mathbb{P}(T=1|\mathbf{X})} \cdot p_1 \cdot \{\mathbb{E}[Y^*(1)|X] - \beta_1\} \\ \vdots \\ \frac{I(T=J)}{\mathbb{P}(T=J|\mathbf{X})} \cdot p_J \cdot \{\mathbb{E}[Y^*(j)|X] - \beta_J\} \end{bmatrix}
\end{aligned}$$

and

$$\mathbb{E} [\varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) | \mathbf{X}] = \begin{bmatrix} p_0 \cdot \{\mathbb{E}[Y^*(0)|X] - \beta_0\} \\ p_1 \cdot \{\mathbb{E}[Y^*(1)|X] - \beta_1\} \\ \vdots \\ p_J \cdot \{\mathbb{E}[Y^*(j)|X] - \beta_J\} \end{bmatrix}. \tag{11}$$

From Theorem 3.1, the efficient influence function of  $\boldsymbol{\beta}_0 = (\beta_0, \dots, \beta_J)$  is given by

$$H_0^{-1} \{ \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \{Y - \mathbb{E}[Y|\mathbf{X}, T]\} + \mathbb{E} [\varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) | \mathbf{X}] \}$$

$$= \begin{bmatrix} \frac{I(T=0)}{\mathbb{P}(T=0|\mathbf{X})} \cdot \{Y^*(0) - \mathbb{E}[Y^*(0)|\mathbf{X}]\} + \mathbb{E}[Y^*(0)|X] - \beta_0 \\ \frac{I(T=1)}{\mathbb{P}(T=1|\mathbf{X})} \cdot \{Y^*(1) - \mathbb{E}[Y^*(1)|\mathbf{X}]\} + \mathbb{E}[Y^*(1)|X] - \beta_1 \\ \vdots \\ \frac{I(T=J)}{\mathbb{P}(T=J|\mathbf{X})} \cdot \{Y^*(j) - \mathbb{E}[Y^*(j)|\mathbf{X}]\} + \mathbb{E}[Y^*(j)|X] - \beta_J \end{bmatrix},$$

which is the same as the efficient influence function developed in Corollary 1 of [Cattaneo \(2010\)](#).

■

## 2.4 Particular Case III: Binary Quantile Treatment Effects

In this section, we show that when  $T \in \{0, 1\}$  is a binary treatment variable,  $L(v) = v(\tau - I(v \leq 0))$  is the check function with  $\tau \in (0, 1)$ , and  $g(t; \beta_0) = \beta_0 \cdot (1 - t) + \beta_1 \cdot t$ , where  $\beta_0 = (\beta_0, \beta_1)$ , our general efficiency bound derived in Theorem 3.1 reduces to the efficiency bound of quantile treatment effects given in [Firpo \(2007\)](#). In accordance with our identification condition,  $\beta_0$  and  $\beta_1$  are identified by minimizing the following loss function

$$\sum_{j \in \{0, 1\}} \mathbb{P}(T = j) \cdot \mathbb{E}[(Y^*(j) - \beta_j) \{\tau - I(Y^*(j) \leq \beta_j)\}].$$

The solutions are  $\beta_0 = \inf\{q : \mathbb{P}(Y^*(0) \leq q) \geq \tau\}$  and  $\beta_1 = \inf\{q : \mathbb{P}(Y^*(1) \leq q) \geq \tau\}$ , which are the  $\tau^{th}$  quantiles of potential outcomes.

**Corollary 2.3** *Let  $T \in \{0, 1\}$ ,  $f_{Y^*(1)}$  and  $f_{Y^*(0)}$  be the probability densities of the potential outcomes  $Y^*(1)$  and  $Y^*(0)$  respectively,  $g(t; \beta_0) = \beta_0 \cdot (1 - t) + \beta_1 \cdot t$ ,  $L(v) = v(\tau - I(v \leq 0))$ , and the conditions in Theorem 3.1 hold, then the efficient influence function of  $\beta_0$  given by Theorem 3.1 reduces to*

$$S_{eff}(Y, T, \mathbf{X}; \beta_0) = \left[ \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot \left\{ \frac{\tau - I(Y^*(0) \leq \beta_0)}{f_{Y^*(0)}(\beta_0)} \right\} - \left( \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} - 1 \right) \cdot \mathbb{E} \left[ \frac{\tau - I(Y^*(0) \leq \beta_0)}{f_{Y^*(0)}(\beta_0)} \middle| \mathbf{X} \right] \right] \\ - \left[ \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot \left\{ \frac{\tau - I(Y^*(1) \leq \beta_1)}{f_{Y^*(1)}(\beta_1)} \right\} - \left( \frac{T}{\mathbb{P}(T=1|\mathbf{X})} - 1 \right) \cdot \mathbb{E} \left[ \frac{\tau - I(Y^*(1) \leq \beta_1)}{f_{Y^*(1)}(\beta_1)} \middle| \mathbf{X} \right] \right],$$

which is the same as the efficient influence function given in [Firpo \(2007\)](#).

**Proof.** Using our notation, we have

$$\beta_0 = (\beta_0, \beta_1)^\top, \quad g(t; \beta_0) = \beta_0 \cdot (1 - t) + \beta_1 \cdot t, \quad m(t; \beta_0) = \begin{bmatrix} 1 - t \\ t \end{bmatrix},$$

$$L(v) = v(\tau - I(v \leq 0)), \quad L'(v) = \tau - I(v \leq 0) \text{ a.s.},$$

$$\begin{aligned}\varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) &= T \cdot \mathbb{E}[\tau - I(Y^*(1) \leq \beta_1) | \mathbf{X}] + (1 - T) \cdot \mathbb{E}[\tau - I(Y^*(0) \leq \beta_0) | \mathbf{X}], \\ \pi_0(T, \mathbf{X}) &= \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p + \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q, \quad p = \mathbb{P}(T=1), \quad q = \mathbb{P}(T=0).\end{aligned}$$

Direct computation yields

$$\begin{aligned}& \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) L'(Y - g(T; \boldsymbol{\beta}_0)) \\ &= \left\{ \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p + \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \right\} \cdot \begin{bmatrix} 1-T \\ T \end{bmatrix} \cdot \left\{ \tau - I(Y \leq \beta_0 \cdot (1-T) + \beta_1 \cdot T) \right\} \\ &= \begin{bmatrix} \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \cdot \{\tau - I(Y^*(0) \leq \beta_0)\} \\ \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p \cdot \{\tau - I(Y^*(1) \leq \beta_1)\} \end{bmatrix}\end{aligned}$$

and

$$\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) = \begin{bmatrix} \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \cdot \mathbb{E}[\tau - I(Y^*(0) \leq \beta_0) | \mathbf{X}] \\ \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p \cdot \mathbb{E}[\tau - I(Y^*(1) \leq \beta_1) | \mathbf{X}] \end{bmatrix}$$

and

$$\mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) | \mathbf{X}] = \begin{bmatrix} q \cdot \mathbb{E}[\tau - I(Y^*(0) \leq \beta_0) | \mathbf{X}] \\ p \cdot \mathbb{E}[\tau - I(Y^*(1) \leq \beta_1) | \mathbf{X}] \end{bmatrix}$$

and

$$H_0 = \nabla_{\boldsymbol{\beta}} \mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}) L'(Y - g(T; \boldsymbol{\beta}))] = \begin{bmatrix} -q \cdot f_{Y^*(0)}(\beta_0) & 0 \\ 0 & -p \cdot f_{Y^*(1)}(\beta_1) \end{bmatrix}.$$

Therefore, by Theorem 3.1, the efficient influence function of  $\boldsymbol{\beta}_0$  is

$$\begin{aligned}& S_{eff}(Y, T, \mathbf{X}; \boldsymbol{\beta}_0) \\ &= H_0^{-1} \cdot \left\{ \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) L'(Y - g(T; \boldsymbol{\beta}_0)) - \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) \right. \\ &\quad \left. + \mathbb{E}[\varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) \pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) | \mathbf{X}] \right\} \\ &= \begin{bmatrix} q^{-1} \cdot \frac{1}{f_{Y^*(0)}(\beta_0)} & 0 \\ 0 & p^{-1} \cdot \frac{1}{f_{Y^*(1)}(\beta_1)} \end{bmatrix}\end{aligned}$$

$$\begin{aligned}
& \times \left[ \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot q \cdot \{\tau - I(Y^*(0) \leq \beta_0)\} - q \cdot \left( \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} - 1 \right) \cdot \mathbb{E}[\tau - I(Y^*(0) \leq \beta_0)|\mathbf{X}] \right] \\
& \left[ \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot p \cdot \{\tau - I(Y^*(1) \leq \beta_1)\} - p \cdot \left( \frac{T}{\mathbb{P}(T=1|\mathbf{X})} - 1 \right) \cdot \mathbb{E}[\tau - I(Y^*(1) \leq \beta_1)|\mathbf{X}] \right] \\
& = \left[ \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} \cdot \left\{ \frac{\tau - I(Y^*(0) \leq \beta_0)}{f_{Y^*(0)}(\beta_0)} \right\} - \left( \frac{1-T}{\mathbb{P}(T=0|\mathbf{X})} - 1 \right) \cdot \mathbb{E} \left[ \frac{\tau - I(Y^*(0) \leq \beta_0)}{f_{Y^*(0)}(\beta_0)} \middle| \mathbf{X} \right] \right] \\
& \left[ \frac{T}{\mathbb{P}(T=1|\mathbf{X})} \cdot \left\{ \frac{\tau - I(Y^*(1) \leq \beta_1)}{f_{Y^*(1)}(\beta_1)} \right\} - \left( \frac{T}{\mathbb{P}(T=1|\mathbf{X})} - 1 \right) \cdot \mathbb{E} \left[ \frac{\tau - I(Y^*(1) \leq \beta_1)}{f_{Y^*(1)}(\beta_1)} \middle| \mathbf{X} \right] \right],
\end{aligned}$$

which coincides with efficiency bound derived in [Firpo \(2007\)](#). ■

### 3 Convergence Rate of Estimated Stabilized Weights

In this section, we establish the convergence rate of estimated stabilized weights  $\hat{\pi}_K(T, \mathbf{X})$ . Let  $G_{K_1 \times K_2}^*$ ,  $\Lambda_{K_1 \times K_2}^*$  and  $\pi_K^*(t, \mathbf{x})$  be the theoretical counterparts of  $\hat{G}_{K_1 \times K_2}$ ,  $\hat{\Lambda}_{K_1 \times K_2}$  and  $\hat{\pi}_K(t, \mathbf{x})$  respectively:

$$\begin{aligned}
G_{K_1 \times K_2}^*(\Lambda) &:= \mathbb{E}[\hat{G}_{K_1 \times K_2}(\Lambda)] = \mathbb{E}[\rho(u_{K_1}(T)^\top \Lambda v_{K_2}(\mathbf{X}))] - \mathbb{E}[u_{K_1}(T)^\top] \cdot \Lambda \cdot \mathbb{E}[v_{K_2}(\mathbf{X})], \\
\Lambda_{K_1 \times K_2}^* &:= \arg \max G_{K_1 \times K_2}^*(\Lambda), \\
\pi_K^*(t, \mathbf{x}) &:= \rho'(u_{K_1}(t)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})).
\end{aligned}$$

As discussed in Appendix A.3, we assume the sieve bases  $u_{K_1}(T)$  and  $v_{K_2}(\mathbf{X})$  are orthonormalized, i.e.,

$$\mathbb{E}[u_{K_1}(T)u_{K_1}^\top(T)] = I_{K_1 \times K_1}, \quad \mathbb{E}[v_{K_2}(\mathbf{X})v_{K_2}^\top(\mathbf{X})] = I_{K_2 \times K_2}. \quad (12)$$

Let

$$\zeta_1(K_1) := \sup_{t \in \mathcal{T}} \|u_{K_1}(t)\|, \quad \zeta_2(K_2) := \sup_{\mathbf{x} \in \mathcal{X}} \|v_{K_2}(\mathbf{x})\|, \quad K = K_1 \cdot K_2, \quad \zeta(K) = \zeta_1(K_1)\zeta_2(K_2).$$

We also recall the following property satisfied by  $\pi_0(T, \mathbf{X})$ : for any integrable functions  $u(t)$  and  $v(\mathbf{X})$ ,

$$\mathbb{E}[\pi_0(T, \mathbf{X})u(T)v(\mathbf{X})] = \mathbb{E}[u(T)] \cdot \mathbb{E}[v(\mathbf{X})]. \quad (13)$$

#### 3.1 Lemma 3.1

The first lemma states that  $\pi_K^*(t, \mathbf{x})$  is arbitrarily close to the true stabilized weights  $\pi_0(t, \mathbf{x})$ .

**Lemma 3.1** Under Assumption 1.2-1.6, we have

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_0(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})| = O(K^{-\alpha} \zeta(K)),$$

and

$$\mathbb{E} [|\pi_0(T, \mathbf{X}) - \pi_K^*(T, \mathbf{X})|^2] = O(K^{-2\alpha}),$$

and

$$\frac{1}{N} \sum_{i=1}^N |\pi_0(T_i, \mathbf{X}_i) - \pi_K^*(T_i, \mathbf{X}_i)|^2 = O_p(K^{-2\alpha}).$$

**Proof.** By Assumption 1.3,  $\pi_0(t, \mathbf{x}) \in [\eta_1, \eta_2]$ ,  $\forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}$  and  $(\rho')^{-1}$  is strictly decreasing. Define

$$\bar{\gamma} := \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} (\rho')^{-1}(\pi_0(t, \mathbf{x})) \leq (\rho')^{-1}(\eta_1) \quad \text{and} \quad \underline{\gamma} := \inf_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} (\rho')^{-1}(\pi_0(t, \mathbf{x})) \geq (\rho')^{-1}(\eta_2),$$

which are two finite constants. By Assumptions 1.4, there exist a constant  $C > 0$  and a  $K_1 \times K_2$  matrix  $\Lambda_{K_1 \times K_2} \in \mathbb{R}^{K_1 \times K_2}$  such that

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |(\rho')^{-1}(\pi_0(t, \mathbf{x})) - u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})| < CK^{-\alpha},$$

which implies

$$\begin{aligned} u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) &\in ((\rho')^{-1}(\pi_0(t, \mathbf{x})) - CK^{-\alpha}, (\rho')^{-1}(\pi_0(t, \mathbf{x})) + CK^{-\alpha}) \\ &\subset [\underline{\gamma} - CK^{-\alpha}, \bar{\gamma} + CK^{-\alpha}], \quad \forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}, \end{aligned} \tag{14}$$

and

$$\begin{aligned} &\rho'(u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) + CK^{-\alpha}) - \rho'(u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) \\ &< \pi_0(t, \mathbf{x}) - \rho'(u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})) \\ &< \rho'(u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) - CK^{-\alpha}) - \rho'(u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})), \quad \forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}. \end{aligned}$$

Let  $\Gamma_1 := [\underline{\gamma} - 1, \bar{\gamma} + 1]$ , by Mean Value Theorem, for large enough  $K$ , there exist

$$\begin{aligned} \xi_1(t, \mathbf{x}) &\in (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}), u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) + CK^{-\alpha}) \\ &\subset [\underline{\gamma} - CK^{-\alpha}, \bar{\gamma} + 2CK^{-\alpha}] \subset \Gamma_1, \end{aligned}$$

$$\begin{aligned}\xi_2(t, \mathbf{x}) &\in (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) - CK^{-\alpha}, u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})) \\ &\subset [\underline{\gamma} - 2CK^{-\alpha}, \bar{\gamma} + CK^{-\alpha}] \subset \Gamma_1,\end{aligned}$$

such that

$$\rho' (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) + CK^{-\alpha}) - \rho' (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})) = \rho''(\xi_1(t, x))CK^{-\alpha} \geq -a_1CK^{-\alpha}$$

and

$$\rho' (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) - CK^{-\alpha}) - \rho' (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})) = -\rho''(\xi_2(t, \mathbf{x}))CK^{-\alpha} \leq a_2CK^{-\alpha},$$

where  $-a_1 := \inf_{\gamma \in \Gamma_1} \rho''(\gamma)$  and  $a_2 := \sup_{\gamma \in \Gamma_1} (-\rho''(\gamma))$ . Let  $a := \max\{a_1, a_2\}$ , we have

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_0(t, \mathbf{x}) - \rho' (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}))| < aCK^{-\alpha}. \quad (15)$$

For some fixed  $C_2 > 0$  (to be chosen later), define

$$\Upsilon_{K_1 \times K_2} := \{\Lambda \in \mathbb{R}^{K_1 \times K_2} : \|\Lambda - \Lambda_{K_1 \times K_2}\| \leq C_2K^{-\alpha}\}.$$

For sufficiently large  $K_1$  and  $K_2$ , we have that  $\forall \Lambda \in \Upsilon_{K_1 \times K_2}$ ,  $\forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}$ ,

$$\begin{aligned}&|u_{K_1}(t)^\top \Lambda v_{K_2}(\mathbf{x}) - u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})| \\ &\leq \|\Lambda - \Lambda_{K_1 \times K_2}\| \cdot \sup_{\mathbf{x} \in \mathcal{X}} \|v_{K_2}(\mathbf{x})\| \cdot \sup_{t \in \mathcal{T}} \|u_{K_1}(t)\| \leq C_2K^{-\alpha} \zeta_1(K_1) \zeta_2(K_2).\end{aligned}$$

Then in light of (14) and Assumption 1.15, for large enough  $K_1$  and  $K_2$ ,  $\forall \Lambda \in \Upsilon_{K_1 \times K_2}$  and  $\forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}$ , we can deduce that

$$\begin{aligned}u_{K_1}(t)^\top \Lambda v_{K_2}(\mathbf{x}) &\in (u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) - C_2K^{-\alpha} \zeta_1(K_1) \zeta_2(K_2), \\ &\quad u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}) + C_2K^{-\alpha} \zeta_1(K_1) \zeta_2(K_2)) \\ &\subset [\underline{\gamma} - CK^{-\alpha} - C_2K^{-\alpha} \zeta_1(K_1) \zeta_2(K_2), \\ &\quad \bar{\gamma} + CK^{-\alpha} + C_2K^{-\alpha} \zeta_1(K_1) \zeta_2(K_2)] \subset \Gamma_1.\end{aligned} \quad (16)$$

By definition

$$G_{K_1 \times K_2}^*(\Lambda) = \mathbb{E} [\rho (u_{K_1}(T)^\top \Lambda v_{K_2}(\mathbf{X}))] - \mathbb{E}[u_{K_1}(T)]^\top \Lambda \mathbb{E}[v_{K_2}(\mathbf{X})],$$

is a strictly concave function of  $\Lambda$ . By (13), the formula  $\text{tr}(AB) = \text{tr}(BA)$  for matrices  $A$  and  $B$ , the facts  $\mathbb{E}[v_{K_2}(\mathbf{X})v_{K_2}(\mathbf{X})^\top] = I_{K_2 \times K_2}$  and  $\mathbb{E}[u_{K_1}(T)u_{K_1}(T)^\top] = I_{K_1 \times K_1}$ , we can deduce that

$$\begin{aligned}
& \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\|^2 \\
&= \|\mathbb{E}[\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) u_{K_1}(T)v_{K_2}(\mathbf{X})^\top] - \mathbb{E}[u_{K_1}(T)]\mathbb{E}[v_{K_2}(\mathbf{X})]^\top]\|^2 \\
&= \|\mathbb{E}[\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) u_{K_1}(T)v_{K_2}(\mathbf{X})^\top] - \mathbb{E}[\pi_0(T, \mathbf{X})u_{K_1}(T)v_{K_2}(\mathbf{X})]^\top]\|^2 \quad (\text{by (13)}) \\
&= \left\| \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \right\|^2 \\
&= \text{tr} \left\{ \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \right. \\
&\quad \times \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} v_{K_2}(\mathbf{X})u_{K_1}(T)^\top \right] \Big\} \\
&= \text{tr} \left\{ \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \cdot \mathbb{E}[u_{K_2}(\mathbf{X})u_{K_2}(\mathbf{X})^\top] \right. \\
&\quad \times \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} v_{K_2}(\mathbf{X})u_{K_1}(T)^\top \right] \cdot \mathbb{E}[u_{K_1}(T)u_{K_1}(T)^\top] \Big\} \\
&= \mathbb{E} \left[ \text{tr} \left\{ u_{K_1}(T)^\top \cdot \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \cdot \mathbb{E}[u_{K_2}(\mathbf{X})u_{K_2}(\mathbf{X})^\top] \right. \right. \\
&\quad \times \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} v_{K_2}(\mathbf{X})u_{K_1}(T)^\top \right] \cdot u_{K_1}(T) \Big\} \Big] \\
&= \mathbb{E} \left[ \pi_0(T, \mathbf{X}) \cdot u_{K_1}(T)^\top \cdot \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \cdot u_{K_2}(\mathbf{X}) \right. \\
&\quad \times \cdot u_{K_2}(\mathbf{X})^\top \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} v_{K_2}(\mathbf{X})u_{K_1}(T)^\top \right] \cdot u_{K_1}(T) \Big] \quad (\text{by (13)}) \\
&= \mathbb{E} \left[ \left| \pi_0(T, \mathbf{X})^{\frac{1}{4}} u_{K_1}(T) \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \pi_0(T, \mathbf{X})^{\frac{1}{4}} v_{K_2}(\mathbf{X}) \right|^2 \right]. \tag{17}
\end{aligned}$$

Note that the term in the last expression

$$\pi_0(T, \mathbf{X})^{\frac{1}{4}} u_{K_1}(T) \cdot \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \pi_0(T, \mathbf{X})^{\frac{1}{4}} v_{K_2}(\mathbf{X})$$

is the  $L^2(dF_{T, \mathbf{X}})$ -projection of  $\frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}}$  on the space spanned by  $\{\pi_0(T, \mathbf{X})^{\frac{1}{4}} u_{K_1}(T), \pi_0(T, \mathbf{X})^{\frac{1}{4}} v_{K_2}(\mathbf{X})\}$ , which implies that

$$\mathbb{E} \left[ \left| \pi_0(T, \mathbf{X})^{\frac{1}{4}} u_{K_1}(T) \mathbb{E} \left[ \sqrt{\pi_0(T, \mathbf{X})} \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} u_{K_1}(T)v_{K_2}(\mathbf{X})^\top \right] \pi_0(T, \mathbf{X})^{\frac{1}{4}} v_{K_2}(\mathbf{X}) \right|^2 \right]$$

$$\leq \mathbb{E} \left[ \left| \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} \right|^2 \right]. \quad (18)$$

Now, with (17), (18), we can obtain that

$$\begin{aligned} & \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| \\ & \leq \mathbb{E} \left[ \left| \frac{\{\rho'(u_{K_1}(T)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X})) - \pi_0(T, \mathbf{X})\}}{\sqrt{\pi_0(T, \mathbf{X})}} \right|^2 \right]^{\frac{1}{2}} \\ & \leq \frac{1}{\sqrt{\eta_1}} \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'(u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})) - \pi_0(t, \mathbf{x})| \quad (\text{by Assumption 1.3}) \\ & \leq \frac{aC}{\sqrt{\eta_1}} \cdot K^{-\alpha} \quad (\text{by (15)}). \end{aligned} \quad (19)$$

Note that for any  $\Lambda \in \partial \Upsilon_{K_1 \times K_2}$ , i.e.  $\|\Lambda - \Lambda_{K_1 \times K_2}\| = C_2 K^{-\alpha}$ , by Mean Value Theorem and the fact  $\rho''(y) = -\rho'(y)$ , we can deduce that

$$\begin{aligned} & G_{K_1 \times K_2}^*(\Lambda) - G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2}) \\ & = \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^K)^\top \frac{\partial}{\partial \lambda_i} G_{K_1 \times K_2}^*(\lambda_1^K, \dots, \lambda_{K_2}^K) \\ & \quad + \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} \frac{1}{2} (\lambda_j - \lambda_j^K)^\top \frac{\partial^2}{\partial \lambda_i \partial \lambda_l} G_{K_1 \times K_2}^*(\bar{\lambda}_1^K, \dots, \bar{\lambda}_{K_2}^K) (\lambda_l - \lambda_l^K) \\ & \leq \|\Lambda - \Lambda_{K_1 \times K_2}\| \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| \\ & \quad + \frac{1}{2} \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^K)^\top \mathbb{E} \left[ \rho''(u_{K_1}^\top(T) \bar{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X})) u_{K_1}(T) u_{K_1}(T)^\top v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right] (\lambda_l - \lambda_l^K) \\ & = \|\Lambda - \Lambda_{K_1 \times K_2}\| \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| \\ & \quad - \frac{1}{2} \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^K)^\top \mathbb{E} \left[ \frac{\rho'(u_{K_1}^\top(T) \bar{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X}))}{\pi_0(T, \mathbf{X})} \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right] (\lambda_l - \lambda_l^K) \\ & \leq \|\Lambda - \Lambda_{K_1 \times K_2}\| \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| \\ & \quad - \frac{a_3}{2\eta_2} \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^K)^\top \mathbb{E} \left[ \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right] (\lambda_l - \lambda_l^K) \quad (\text{by } a_3 = \inf_{y \in \Gamma_1} \{\rho'(y)\}) \\ & = \|\Lambda - \Lambda_{K_1 \times K_2}\| \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| \\ & \quad - \frac{a_3}{2\eta_2} \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^K)^\top \mathbb{E} \left[ u_{K_1}(T) u_{K_1}(T)^\top \right] \mathbb{E} [v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X})] (\lambda_l - \lambda_l^K) \quad (\text{by (13)}) \end{aligned}$$



$$\begin{aligned}
&= \|\Lambda - \Lambda_{K_1 \times K_2}\| \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| - \frac{a_3}{2\eta_2} \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^K)^\top \mathbb{E}[v_{K_2,j}(\mathbf{X})v_{K_2,l}(\mathbf{X})] (\lambda_l - \lambda_l^K) \\
&= \|\Lambda - \Lambda_{K_1 \times K_2}\| \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| - \frac{a_3}{2\eta_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^K)^\top (\lambda_j - \lambda_j^K) \\
&= \|\Lambda - \Lambda_{K_1 \times K_2}\| \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| - \frac{a_3}{2\eta_2} \|\Lambda - \Lambda_{K_1 \times K_2}\|^2 \\
&= \|\Lambda - \Lambda_{K_1 \times K_2}\| \left( \|\nabla G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2})\| - \frac{a_3}{2\eta_2} \|\Lambda - \Lambda_{K_1 \times K_2}\| \right) \\
&\leq \|\Lambda - \Lambda_{K_1 \times K_2}\| \left( \frac{aC}{\sqrt{\eta_1}} K^{-\alpha} - \frac{a_3}{2\eta_2} \cdot C_2 K^{-\alpha} \right), \quad (\text{by (19)})
\end{aligned}$$

where  $\bar{\Lambda}_{K_1 \times K_2} = (\bar{\lambda}_1^K, \dots, \bar{\lambda}_{K_2}^K)$  lies on the line joining  $\Lambda = (\lambda_1, \dots, \lambda_{K_2})$  and  $\Lambda_{K_1 \times K_2} = (\lambda_1^K, \dots, \lambda_{K_2}^K)$ , which implies  $u_{K_1}^\top(t) \bar{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \in \Gamma_1$  by (16);  $a_3 = \inf_{y \in \Gamma_1} \{\rho'(y)\} > 0$  is a finite positive constant; the fourth and fifth equalities follow from  $\mathbb{E}[u_{K_1}(T)u_{K_1}(T)^\top] = I_{K_1 \times K_1}$  and  $\mathbb{E}[v_{K_2}(\mathbf{X})v_{K_2}(\mathbf{X})^\top] = I_{K_2 \times K_2}$  respectively. Therefore, by choosing

$$C_2 > \frac{2\eta_2}{a_3} \cdot \frac{aC}{\sqrt{\eta_1}},$$

we can obtain the following conclusion:

$$G_{K_1 \times K_2}^*(\Lambda_{K_1 \times K_2}) > G_{K_1 \times K_2}^*(\Lambda), \quad \forall \Lambda \in \partial \Upsilon_{K_1 \times K_2}. \quad (20)$$

Since  $G_{K_1 \times K_2}^*$  is continuous, (20) implies that there exists a local maximum of  $G_{K_1 \times K_2}^*$  in the interior of  $\Upsilon_{K_1 \times K_2}$ . Note that  $G_{K_1 \times K_2}^*$  is strictly concave with a unique global maximum point  $\Lambda_{K_1 \times K_2}^*$ , therefore we can claim that

$$\Lambda_{K_1 \times K_2}^* \in \Upsilon_{K_1 \times K_2}^\circ, \quad \text{i.e. } \|\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\| = O(K^{-\alpha}). \quad (21)$$

By Mean Value Theorem, (16) and (21), we can deduce that

$$\begin{aligned}
&|\rho'(u_{K_1}(t)\Lambda_{K_1 \times K_2}v_{K_2}(\mathbf{x})) - \rho'(u_{K_1}(t)\Lambda_{K_1 \times K_2}^*v_{K_2}(\mathbf{x}))| \\
&= |\rho''(\xi^*(t, \mathbf{x}))| |u_{K_1}(t)\Lambda_{K_1 \times K_2}v_{K_2}(\mathbf{x}) - u_{K_1}(t)\Lambda_{K_1 \times K_2}^*v_{K_2}(\mathbf{x})| \\
&\leq -\rho''(\xi^*(t, \mathbf{x})) \times \|\Lambda_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^*\| \times \sup_{t \in \mathcal{T}} \|u_{K_1}(t)\| \times \sup_{\mathbf{x} \in \mathcal{X}} \|v_{K_2}(\mathbf{x})\| \\
&\leq a_2 C_2 K^{-\alpha} \zeta_1(K_1) \zeta_2(K_2),
\end{aligned}$$

where  $a_2 = \sup_{\gamma \in \Gamma_1} \{-\rho''(\gamma)\} < \infty$  is a finite positive constant, and  $\xi^*(t, \mathbf{x})$  lies between the point  $u_{K_1}(t)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})$  and  $u_{K_1}(t)^\top \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})$  (note (16) implies  $\xi^*(t, \mathbf{x}) \in \Gamma_1$  for all  $(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}$  and large enough  $K$ ). Therefore, using the triangle inequality, and Assumption 1.15, we can have

$$\begin{aligned}
& \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_0(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})| \\
& \leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_0(t, \mathbf{x}) - \rho'(u_{K_1}(t) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}))| \\
& \quad + \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'(u_{K_1}(t) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x})) - \rho'(u_{K_1}(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}))| \\
& \leq aCK^{-\alpha} + a_2 C_2 K^{-\alpha} \zeta_1(K_1) \zeta_2(K_2) = O(K^{-\alpha} \zeta(K)),
\end{aligned}$$

where  $\zeta(K) = \zeta_1(K_1) \zeta_2(K_2)$ .

We next prove  $\mathbb{E} [|\pi_0(T, \mathbf{X}) - \pi_K^*(T, \mathbf{X})|^2] = O(K^{-2\alpha})$ . By Assumption 1.15, we can deduce that

$$\begin{aligned}
& \mathbb{E} [|\pi_0(T, \mathbf{X}) - \pi_K^*(T, \mathbf{X})|^2] \\
& \leq 2 \cdot \mathbb{E} [|\pi_0(T, \mathbf{X}) - \rho'(u_{K_1}(T) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X}))|^2] + 2 \cdot \mathbb{E} [|\rho'(u_{K_1}(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X})) - \rho'(u_{K_1}(T) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X}))|^2] \\
& \leq 2 \cdot \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_0(t, \mathbf{x}) - \rho'(u_{K_1}(t) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}))|^2 + 2 \sup_{\gamma \in \Gamma_1} |\rho''(\gamma)|^2 \cdot \mathbb{E} [|u_{K_1}^\top(T) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} v_{K_2}(\mathbf{X})|^2] \\
& \leq O(K^{-2\alpha}) + O(1) \cdot \mathbb{E} [|u_{K_1}^\top(T) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} v_{K_2}(\mathbf{X})|^2].
\end{aligned}$$

We next compute the order of  $\mathbb{E} [|u_{K_1}^\top(T) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} v_{K_2}(\mathbf{X})|^2]$ . Note that  $\mathbb{E}[u_{K_1}(T) u_{K_1}(T)^\top] = I_{K_1 \times K_1}$ ,  $\mathbb{E}[v_{K_2}(\mathbf{X}) v_{K_2}(\mathbf{X})^\top] = I_{K_2 \times K_2}$ , (13), (21) and Assumption 1.3, we can deduce that

$$\begin{aligned}
& \mathbb{E} [|u_{K_1}^\top(T) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} v_{K_2}(\mathbf{X})|^2] \\
& = \mathbb{E} [u_{K_1}^\top(T) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} v_{K_2}(\mathbf{X}) v_{K_2}(\mathbf{X})^\top \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\}^\top u_{K_1}(T)] \\
& = \mathbb{E} \left[ \frac{1}{\pi_0(T, \mathbf{X})} \pi_0(T, \mathbf{X}) u_{K_1}^\top(T) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} v_{K_2}(\mathbf{X}) v_{K_2}(\mathbf{X})^\top \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\}^\top u_{K_1}(T) \right] \\
& \leq \frac{1}{\eta_1} \cdot \mathbb{E} [\pi_0(T, \mathbf{X}) u_{K_1}^\top(T) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} v_{K_2}(\mathbf{X}) v_{K_2}(\mathbf{X})^\top \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\}^\top u_{K_1}(T)] \\
& = \frac{1}{\eta_1} \cdot \int_{\mathcal{T}} u_{K_1}^\top(t) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} \mathbb{E} [v_{K_2}(\mathbf{X}) v_{K_2}(\mathbf{X})^\top] \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\}^\top u_{K_1}(t) dF_T(t) \quad (\text{by (13)}) \\
& = \frac{1}{\eta_1} \cdot \int_{\mathcal{T}} u_{K_1}^\top(t) \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} \cdot \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\}^\top u_{K_1}(t) dF_T(t) \\
& = \frac{1}{\eta_1} \cdot \int_{\mathcal{T}} \text{tr} \left( \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\} \cdot \{\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\}^\top u_{K_1}(t) u_{K_1}^\top(t) \right) dF_T(t)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\eta_1} \cdot \text{tr} \left( \left\{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \right\} \cdot \left\{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \right\}^\top \right) \\
&\leq \frac{1}{\eta_1} \cdot \|\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\|^2 = O(K^{-2\alpha}). \quad (\text{by (21)})
\end{aligned} \tag{22}$$

Therefore, we can obtain

$$\mathbb{E} [|\pi_0(T, \mathbf{X}) - \pi_K^*(T, \mathbf{X})|^2] = O(K^{-2\alpha}).$$

We finally prove  $N^{-1} \sum_{i=1}^N |\pi_0(T_i, \mathbf{X}_i) - \pi_K^*(T_i, \mathbf{X}_i)|^2 = O_p(K^{-2\alpha})$ . Note that by (22), we can have

$$\begin{aligned}
&\mathbb{E} \left[ \left\{ \frac{1}{N} \sum_{i=1}^N |u_{K_1}^\top(T_i) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X}_i)|^2 - \mathbb{E} [ |u_{K_1}^\top(T) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X})|^2 ] \right\}^2 \right] \\
&\leq \frac{1}{N} \cdot \mathbb{E} [ |u_{K_1}^\top(T) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X})|^4 ] \\
&\leq \frac{1}{N} \cdot \mathbb{E} [ |u_{K_1}^\top(T) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X})|^2 ] \cdot \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |u_{K_1}^\top(t) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{x})|^2 \\
&\leq \frac{1}{N} \cdot O(K^{-2\alpha}) \cdot \zeta_1(K_1)^2 \zeta_2(K_2)^2 \cdot \|\Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2}\|^2 \leq \frac{1}{N} \cdot \zeta(K)^2 \cdot O(K^{-4\alpha}),
\end{aligned}$$

then in light of Chebyshev's inequality and Assumption 1.6, we have

$$\begin{aligned}
&\frac{1}{N} \sum_{i=1}^N |u_{K_1}^\top(T_i) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X}_i)|^2 - \mathbb{E} [ |u_{K_1}^\top(T) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X})|^2 ] \\
&= O_p \left( \frac{\zeta(K)}{\sqrt{N}} K^{-2\alpha} \right) = o_p(K^{-2\alpha}).
\end{aligned} \tag{23}$$

With (21), (22), (23), and Assumption 1.3, we can deduce that

$$\begin{aligned}
&\frac{1}{N} \sum_{i=1}^N |\pi_0(T_i, \mathbf{X}_i) - \pi_K^*(T_i, \mathbf{X}_i)|^2 \\
&\leq \frac{2}{N} \sum_{i=1}^N |\pi_0(T_i, \mathbf{X}_i) - \rho'(u_{K_1}(T_i) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i))|^2 \\
&\quad + \frac{2}{N} \sum_{i=1}^N |\rho'(u_{K_1}(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) - \rho'(u_{K_1}(T_i) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i))|^2 \\
&\leq 2 \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_0(t, \mathbf{x}) - \rho'(u_{K_1}(t) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}))|^2 \\
&\quad + \sup_{\gamma \in \Gamma_1} |\rho''(\gamma)|^2 \cdot \frac{2}{N} \sum_{i=1}^N |u_{K_1}^\top(T_i) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X}_i)|^2
\end{aligned}$$

$$\begin{aligned}
&\leq 2 \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_0(t, \mathbf{x}) - \rho'(u_{K_1}(t) \Lambda_{K_1 \times K_2} v_{K_2}(\mathbf{x}))|^2 \\
&\quad + 2 \cdot \sup_{\gamma \in \Gamma_1} |\rho''(\gamma)|^2 \cdot \mathbb{E} \left[ \left| u_{K_1}^\top(T) \{ \Lambda_{K_1 \times K_2}^* - \Lambda_{K_1 \times K_2} \} v_{K_2}(\mathbf{X}) \right|^2 \right] + o_p(K^{-2\alpha}) \\
&= O(K^{-2\alpha}) + O(K^{-2\alpha}) + o_p(K^{-2\alpha}) = O_p(K^{-2\alpha}). \quad (\text{by (21)})
\end{aligned}$$

■

### 3.2 Lemma 3.2

**Lemma 3.2** *Under Assumption 1.2-1.6, we have*

$$\left\| \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\| = O_p \left( \sqrt{\frac{K}{N}} \right).$$

**Proof.** Define

$$\hat{S}_N := \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \pi_0(T_i, \mathbf{X}_i) u_{K_1}(T_i) u_{K_1}(T_i)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}_i) v_{K_2,l}(\mathbf{X}_i),$$

where  $\lambda_j$  and  $\lambda_j^*$  are the  $j$ -th column of  $\Lambda$  and  $\Lambda_{K_1 \times K_2}^*$  respectively. Since  $\hat{S}_N$  is symmetric, using (13) and the facts that  $\mathbb{E}[u_{K_1}(T) u_{K_1}(T)^\top] = I_{K_1 \times K_1}$  and  $\mathbb{E}[v_{K_2}(\mathbf{X}) v_{K_2}(\mathbf{X})^\top] = I_{K_2 \times K_2}$ , we can have

$$\begin{aligned}
\mathbb{E}[\hat{S}_N] &= \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \mathbb{E}[\pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X})] (\lambda_l - \lambda_l^*) \\
&= \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \mathbb{E}[u_{K_1}(T) u_{K_1}(T)^\top] \mathbb{E}[v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X})] (\lambda_l - \lambda_l^*) \\
&= \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top (\lambda_j - \lambda_j^*) = \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\|.
\end{aligned}$$

Then we can further deduce that

$$\begin{aligned}
&\mathbb{E} \left[ \left\| \hat{S}_N - \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\| \right\|^2 \right] \\
&= \mathbb{E}[\hat{S}_N^2] - 2\mathbb{E}[\hat{S}_N] \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\| + \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\|^2 \\
&= \frac{N}{N^2} \cdot \mathbb{E} \left[ \left( \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right)^2 \right]
\end{aligned}$$

$$\begin{aligned}
& + 2 \cdot \frac{C_N^2}{N^2} \cdot \mathbb{E} \left[ \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right]^2 \\
& - \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\|^2 \\
& = \frac{1}{N} \mathbb{E} \left[ \left( \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right)^2 \right] \\
& + \frac{N(N-1)}{N^2} \cdot \mathbb{E} \left[ \hat{S}_N \right]^2 - \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\|^2 \\
& = \frac{1}{N} \mathbb{E} \left[ \left( \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right)^2 \right] \\
& - \frac{1}{N} \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\|^2 \\
& < \frac{1}{N} \mathbb{E} \left[ \left( \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \right)^2 \right].
\end{aligned}$$

In light of the fact that

$$0 \leq y^\top \left\{ \pi_0(t, \mathbf{x}) u_{K_1}(t) u_{K_1}(t)^\top \right\} y \leq \eta_2 \zeta_1(K_1)^2 y^\top y, \quad \forall y \in \mathbb{R}^{K_1}, \quad \forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X},$$

we can deduce that

$$\begin{aligned}
& \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \left\{ \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top \right\} (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}) v_{K_2,l}(\mathbf{X}) \\
& = \left[ \sum_{j=1}^{K_2} v_{K_2,j}(\mathbf{X}) (\lambda_j - \lambda_j^*)^\top \right] \cdot \left\{ \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top \right\} \cdot \left[ \sum_{l=1}^{K_2} (\lambda_l - \lambda_l^*) v_{K_2,l}(\mathbf{X}) \right] \\
& \leq \eta_2 \cdot \|u_{K_1}(T)\|^2 \cdot \left\| \sum_{i=1}^{K_2} (\lambda_i - \lambda_i^*)^\top v_{K_2,i}(\mathbf{X}) \right\|^2 \\
& \leq \eta_2 \cdot \|u_{K_1}(T)\|^2 \cdot \left( \sum_{i=1}^{K_2} \|\lambda_i - \lambda_i^*\|^2 \right) \left( \sum_{i=1}^{K_2} v_{K_2,i}(\mathbf{X})^2 \right) \\
& = \eta_2 \cdot \|u_{K_1}(T)\|^2 \cdot \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\|^2 \|v_{K_2}(\mathbf{X})\|^2.
\end{aligned}$$

Therefore, we can obtain that

$$\mathbb{E} \left[ \left| \hat{S}_N - \left\| \Lambda - \Lambda_{K_1 \times K_2}^* \right\|^2 \right| \right]$$

$$\begin{aligned}
&\leq \frac{1}{N} \eta_2^2 \cdot \mathbb{E} [\|u_{K_1}(T)\|^4 \cdot \|v_{K_2}(\mathbf{X})\|^4] \cdot \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4 \\
&\leq \frac{1}{N} \eta_2^2 \cdot \zeta_1(K_1)^2 \cdot \zeta_2(K_2)^2 \cdot \mathbb{E} [\|u_{K_1}(T)\|^2 \cdot \|v_{K_2}(\mathbf{X})\|^2] \cdot \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4 \\
&= \frac{1}{N} \eta_2^2 \cdot \zeta_1(K_1)^2 \cdot \zeta_2(K_2)^2 \cdot \mathbb{E} \left[ \frac{1}{\pi_0(T, \mathbf{X})} \cdot \pi_0(T, \mathbf{X}) \|u_{K_1}(T)\|^2 \cdot \|v_{K_2}(\mathbf{X})\|^2 \right] \cdot \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4 \\
&\leq \frac{1}{N} \frac{\eta_2^2}{\eta_1} \cdot \zeta_1(K_1)^2 \cdot \zeta_2(K_2)^2 \cdot \mathbb{E} [\pi_0(T, \mathbf{X}) \|u_{K_1}(T)\|^2 \cdot \|v_{K_2}(\mathbf{X})\|^2] \cdot \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4 \quad (\text{by Assumption 1.3}) \\
&= \frac{1}{N} \frac{\eta_2^2}{\eta_1} \cdot \zeta_1(K_1)^2 \cdot \zeta_2(K_2)^2 \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \cdot \mathbb{E} [\|v_{K_2}(\mathbf{X})\|^2] \cdot \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4 \quad (\text{by (13)}) \\
&= \frac{1}{N} \frac{\eta_2^2}{\eta_1} \cdot \zeta_1(K_1)^2 \cdot \zeta_2(K_2)^2 \cdot K_1 \cdot K_2 \cdot \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4 \quad (\text{since } \mathbb{E}[\|u_{K_1}(T)\|^2] = K_1 \text{ and } \mathbb{E}[\|v_{K_2}(\mathbf{X})\|^2] = K_2) \\
&= \frac{1}{N} \frac{\eta_2^2}{\eta_1} \cdot \zeta(K)^2 \cdot K \cdot \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4 \quad (\text{since } \zeta(K) = \zeta_1(K_1)\zeta_2(K_2) \text{ and } K = K_1 \cdot K_2) \tag{24}
\end{aligned}$$

Considering the event set

$$E_N := \left\{ \hat{S}_N > \frac{1}{2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2, \Lambda \neq \Lambda_{K_1 \times K_2}^* \right\},$$

by Chebyshev's inequality, (24), and Assumption 1.15 we can get

$$\begin{aligned}
&\mathbb{P} \left( \left| \hat{S}_N - \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2 \right| \geq \frac{1}{2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2, \Lambda \neq \Lambda_{K_1 \times K_2}^* \right) \\
&\leq \frac{4\mathbb{E} \left[ \left| \hat{S}_N - \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2 \right|^2 \right]}{\|\Lambda - \Lambda_{K_1 \times K_2}^*\|^4} \\
&\leq \frac{4}{N} \frac{\eta_2^2}{\eta_1} \cdot \zeta(K)^2 \cdot K \leq O \left( \frac{\zeta(K)^2 K}{N} \right) = o(1),
\end{aligned}$$

which implies that for any  $\epsilon > 0$ , there exists  $N_0(\epsilon) \in \mathbb{N}$  such that  $N > N_0(\epsilon)$  large enough

$$\mathbb{P}((E_N)^c) < \mathbb{P} \left( \left| \hat{S}_N - \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2 \right| \geq \frac{1}{2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2, \Lambda \neq \Lambda_{K_1 \times K_2}^* \right) < \frac{\epsilon}{2}. \tag{25}$$

Note that

$$\begin{aligned}
&\frac{\partial}{\partial \lambda_j} \hat{G}_{K_1 \times K_2}(\lambda_1, \dots, \lambda_{K_2}) \\
&= \frac{1}{N} \sum_{i=1}^N \rho' \left( u_{K_1}^\top(T_i) \Lambda v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_i) - \frac{1}{N^2} \sum_{i=1}^N \sum_{l=1}^N u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_l)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{i=1}^N \left\{ \rho' \left( u_{K_1}^\top(T_i) \Lambda v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_i) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] u_{K_1}(T_i) \right\} \\
&\quad - \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \left\{ \frac{1}{N} \sum_{l=1}^N v_{K_2,j}(\mathbf{X}_l) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] \right\}
\end{aligned}$$

and

$$\frac{\partial}{\partial \lambda_j} G_{K_1 \times K_2}^*(\lambda_1, \dots, \lambda_{K_2}) = \mathbb{E} \left[ \rho' \left( u_{K_1}^\top(T_i) \Lambda v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_i) \right] - \mathbb{E}[u_{K_1}(T)] \cdot \mathbb{E}[v_{K_2,j}(\mathbf{X})].$$

Since  $\Lambda_{K_1 \times K_2}^*$  is the unique maximizer of  $G_{K_1 \times K_2}^*(\cdot)$ , then for each  $j \in \{1, \dots, K_2\}$ ,

$$\begin{aligned}
&\frac{\partial}{\partial \lambda_j} G_{K_1 \times K_2}^*(\lambda_1^*, \dots, \lambda_{K_2}^*) \\
&= \mathbb{E} \left[ \rho' \left( u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}) \right) u_{K_1}(T) v_{K_2,j}(\mathbf{X}) \right] - \mathbb{E}[u_{K_1}(T)] \mathbb{E}[v_{K_2,j}(\mathbf{X})] = 0.
\end{aligned}$$

Therefore, for large enough  $K$ , we can deduce that

$$\begin{aligned}
&\mathbb{E} \left[ \|\nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*)\|^2 \right] = \sum_{j=1}^{K_2} \mathbb{E} \left[ \left\| \frac{\partial}{\partial \lambda_j} \hat{G}_{K_1 \times K_2}(\lambda_1^*, \dots, \lambda_{K_2}^*) \right\|^2 \right] \tag{26} \\
&\leq 2 \sum_{j=1}^{K_2} \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N \left\{ \rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_i) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] u_{K_1}(T_i) \right\} \right\|^2 \right] \\
&\quad + 2 \sum_{j=1}^{K_2} \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \left\{ \frac{1}{N} \sum_{l=1}^N v_{K_2,j}(\mathbf{X}_l) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] \right\} \right\|^2 \right] \\
&= \frac{2}{N^2} \sum_{j=1}^{K_2} \sum_{i=1}^N \mathbb{E} \left[ \left\| \rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_i) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] u_{K_1}(T_i) \right\|^2 \right] \\
&\quad + 2 \sum_{j=1}^{K_2} \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \left\{ \frac{1}{N} \sum_{l=1}^N v_{K_2,j}(\mathbf{X}_l) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] \right\} \right\|^2 \right] \\
&\leq \frac{2}{N^2} \sum_{j=1}^{K_2} \sum_{i=1}^N \mathbb{E} \left[ \left\| \rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_i) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] u_{K_1}(T_i) \right\|^2 \right] \\
&\quad + 2 \sum_{j=1}^{K_2} \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{l=1}^N v_{K_2,j}(\mathbf{X}_l) - \mathbb{E}[v_{K_2,j}(\mathbf{X})] \right\|^2 \right] \cdot \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \right\|^2 \right] \\
&\leq \frac{4}{N} \sum_{j=1}^{K_2} \left\{ \mathbb{E} \left[ \left\| \rho' \left( u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}) \right) u_{K_1}(T) v_{K_2,j}(\mathbf{X}) \right\|^2 \right] + \mathbb{E}[v_{K_2,j}(\mathbf{X})^2] \mathbb{E}[\|u_{K_1}(T)\|^2] \right\}
\end{aligned}$$

$$\begin{aligned}
& + \frac{2}{N} \sum_{j=1}^{K_2} \mathbb{E} [v_{K_2,j}(\mathbf{X})^2] \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \\
& = \frac{4}{N} \sum_{j=1}^{K_2} \left\{ \mathbb{E} \left[ \frac{|\rho'(u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}))|^2}{\pi_0(T, \mathbf{X})} \cdot \pi_0(T, \mathbf{X}) \cdot \|u_{K_1}(T) v_{K_2,j}(\mathbf{X})\|^2 \right] + \mathbb{E}[v_{K_2,j}(\mathbf{X})^2] \mathbb{E} [\|u_{K_1}(T)\|^2] \right\} \\
& + \frac{2}{N} \sum_{j=1}^{K_2} \mathbb{E} [v_{K_2,j}(\mathbf{X})^2] \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \\
& \leq \frac{4}{N} \sum_{j=1}^{K_2} \left\{ \frac{(\sup_{\gamma \in \Gamma_1} \rho'(\gamma))^2}{\eta_1} \cdot \mathbb{E} [\pi_0(T, \mathbf{X}) \cdot \|u_{K_1}(T) v_{K_2,j}(\mathbf{X})\|^2] + \mathbb{E}[v_{K_2,j}(\mathbf{X})^2] \mathbb{E} [\|u_{K_1}(T)\|^2] \right\} \\
& + \frac{2}{N} \sum_{j=1}^{K_2} \mathbb{E} [v_{K_2,j}(\mathbf{X})^2] \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \\
& = \frac{4}{N} \sum_{j=1}^{K_2} \left\{ \frac{(\sup_{\gamma \in \Gamma_1} \rho'(\gamma))^2}{\eta_1} \cdot \mathbb{E} [v_{K_2,j}(\mathbf{X})^2] \mathbb{E} [\|u_{K_1}(T)\|^2] + \mathbb{E}[v_{K_2,j}(\mathbf{X})^2] \mathbb{E} [\|u_{K_1}(T)\|^2] \right\} \\
& + \frac{2}{N} \sum_{j=1}^{K_2} \mathbb{E} [v_{K_2,j}(\mathbf{X})^2] \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \\
& \leq \frac{1}{N} \left\{ \frac{4}{\eta_1} \left( \sup_{\gamma \in \Gamma_1} \rho'(\gamma) \right)^2 + 4 + 2 \right\} \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \sum_{j=1}^{K_2} \mathbb{E} [v_{K_2,j}(\mathbf{X})^2] \\
& = \frac{1}{N} \left\{ \frac{4}{\eta_1} \left( \sup_{\gamma \in \Gamma_1} \rho'(\gamma) \right)^2 + 6 \right\} K_1 K_2 = C_4^2 \frac{K}{N},
\end{aligned}$$

where the last inequality follows by Assumption 1.15 and  $C_4 := \sqrt{\frac{4}{\eta_1} (\sup_{\gamma \in \Gamma_1} \rho'(\gamma))^2 + 6}$  is a finite universal constant.

Let  $\epsilon > 0$ , fix  $C_5(\epsilon) > 0$  (to be chosen later) and define

$$\hat{\Upsilon}_{K_1 \times K_2}(\epsilon) := \left\{ \Lambda \in \mathbb{R}^{K_1 \times K_2} : \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \leq C_5(\epsilon) C_4 \sqrt{\frac{K}{N}} \right\}.$$

For  $\forall \Lambda \in \hat{\Upsilon}_{K_1 \times K_2}(\epsilon)$ ,  $\forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}$ , we can have

$$\begin{aligned}
& |u_{K_1}(t)^\top \Lambda v_{K_2}(\mathbf{x}) - u_{K_1}(t)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})| \\
& \leq \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \sup_{t \in \mathcal{T}} \|u_{K_1}(t)\| \sup_{\mathbf{x} \in \mathcal{X}} \|v_{K_2}(\mathbf{x})\| \leq C_5(\epsilon) C_4 \sqrt{\frac{K}{N}} \zeta_1(K_1) \zeta_2(K_2),
\end{aligned}$$



thus for large enough  $N$ , in accordance with Assumption 1.15 and (14), we have

$$\begin{aligned}
u_{K_1}(t)^\top \Lambda v_{K_2}(\mathbf{x}) &\in \left[ u_{K_1}(t)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) - C_5(\epsilon) C_4 \zeta_1(K_1) \zeta_2(K_2) \sqrt{\frac{K}{N}}, \right. \\
&\quad \left. u_{K_1}(t)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) + C_5(\epsilon) C_4 \zeta_1(K_1) \zeta_2(K_2) \sqrt{\frac{K}{N}} \right] \\
&\subset \left[ \underline{\gamma} - CK^{-\alpha} - C_5(\epsilon) C_4 \zeta_1(K_1) \zeta_2(K_2) \sqrt{\frac{K}{N}}, \right. \\
&\quad \left. \bar{\gamma} + CK^{-\alpha} + C_5(\epsilon) C_4 \zeta_1(K_1) \zeta_2(K_2) \sqrt{\frac{K}{N}} \right] \subset \Gamma_2(\epsilon), \quad (27)
\end{aligned}$$

where  $\Gamma_2(\epsilon) := [\underline{\gamma} - 1 - C_5(\epsilon), \bar{\gamma} + 1 + C_5(\epsilon)]$  is a compact set and independent of  $(t, \mathbf{x})$ .

For any  $\Lambda \in \partial \hat{\Upsilon}_{K_1 \times K_2}(\epsilon)$ , there exists  $\bar{\Lambda}$  on the line joining  $\Lambda$  and  $\Lambda_{K_1 \times K_2}^*$  such that

$$\begin{aligned}
\hat{G}_{K_1 \times K_2}(\Lambda) &= \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*) + \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \frac{\partial}{\partial \lambda_i} \hat{G}_{K_1 \times K_2}(\lambda_1^*, \dots, \lambda_{K_2}^*) \\
&\quad + \frac{1}{2} \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \frac{\partial^2}{\partial \lambda_i \partial \lambda_l} \hat{G}_{K_1 \times K_2}(\bar{\lambda}_1, \dots, \bar{\lambda}_{K_2}) (\lambda_l - \lambda_l^*),
\end{aligned}$$

where  $\bar{\lambda}_j$  denotes the  $j$ -th column of  $\bar{\Lambda}$ . For the second order term in above equality, note that  $u_{K_1}^\top(t) \bar{\Lambda} v_{K_2}(\mathbf{x}) \in \Gamma_2(\epsilon)$  for all  $(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}$ , we can further deduce that

$$\begin{aligned}
&\sum_{l=1}^{K_2} \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \frac{\partial^2}{\partial \lambda_i \partial \lambda_l} \hat{G}_{K_1 \times K_2}(\bar{\lambda}_1, \dots, \bar{\lambda}_{K_2}) (\lambda_l - \lambda_l^*) \quad (28) \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{K_2} \sum_{l=1}^{K_2} (\lambda_j - \lambda_j^*)^\top u_{K_1}(T_i) \rho''(u_{K_1}^\top(T_i) \bar{\Lambda} v_{K_2}(\mathbf{X}_i)) (\lambda_l - \lambda_l^*)^\top u_{K_1}(T_i) v_{K_2,j}(\mathbf{X}_i) v_{K_2,l}(\mathbf{X}_i) \\
&\leq -\frac{\bar{b}(\epsilon)}{N} \sum_{i=1}^N \sum_{j=1}^{K_2} \sum_{l=1}^{K_2} (\lambda_j - \lambda_j^*)^\top u_{K_1}(T_i) u_{K_1}(T_i)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}_i) v_{K_2,l}(\mathbf{X}_i) \\
&= -\frac{\bar{b}(\epsilon)}{N} \sum_{i=1}^N \sum_{j=1}^{K_2} \sum_{l=1}^{K_2} \frac{1}{\pi_0(T_i, \mathbf{X}_i)} (\lambda_j - \lambda_j^*)^\top \pi_0(T_i, \mathbf{X}_i) u_{K_1}(T_i) u_{K_1}(T_i)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}_i) v_{K_2,l}(\mathbf{X}_i) \\
&\leq -\frac{\bar{b}(\epsilon)}{N \eta_2} \sum_{i=1}^N \sum_{j=1}^{K_2} \sum_{l=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \pi_0(T_i, \mathbf{X}_i) u_{K_1}(T_i) u_{K_1}(T_i)^\top (\lambda_l - \lambda_l^*) v_{K_2,j}(\mathbf{X}_i) v_{K_2,l}(\mathbf{X}_i)
\end{aligned}$$

$$= -\frac{\bar{b}(\epsilon)}{\eta_2} \hat{S}_N ,$$

where  $-\bar{b}(\epsilon) := \sup_{\gamma \in \Gamma_2(\epsilon)} \rho''(\gamma) < \infty$  for each fixed  $\epsilon$ . Therefore, on the event  $E_N$  and for large enough  $N$ , we can deduce that for any  $\Lambda \in \partial \hat{\Upsilon}_{K_1 \times K_2}(\epsilon)$ ,

$$\begin{aligned} \text{on the event } E_N : \quad & \hat{G}_{K_1 \times K_2}(\Lambda) - \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*) \\ &= \sum_{j=1}^{K_2} (\lambda_j - \lambda_j^*)^\top \frac{\partial}{\partial \lambda_i} \hat{G}_{K_1 \times K_2}(\lambda_1^*, \dots, \lambda_{K_2}^*) \\ &\quad + \sum_{l=1}^{K_2} \sum_{j=1}^{K_2} \frac{1}{2} (\lambda_j - \lambda_j^*)^\top \frac{\partial^2}{\partial \lambda_i \partial \lambda_l} \hat{G}_{K_1 \times K_2}(\bar{\lambda}_1, \dots, \bar{\lambda}_{K_2})(\lambda_l - \lambda_l^*) \\ &\leq \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \|\nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*)\| - \frac{\bar{b}(\epsilon)}{2\eta_2} \hat{S}_N \text{ (by (28))} \\ &\leq \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \|\nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*)\| - \frac{\bar{b}(\epsilon)}{4\eta_2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2 \\ &= \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \left( \|\nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*)\| - \frac{\bar{b}(\epsilon)}{4\eta_2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \right) , \end{aligned} \tag{29}$$

where the second inequality follows from definition of the event  $E_N$ .

Note that for sufficiently large  $N$ , by Chebyshev's inequality and (26) we have

$$\begin{aligned} & \mathbb{P} \left\{ \|\nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*)\| \geq \frac{\bar{b}(\epsilon)}{4\eta_2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \right\} \\ & \leq \frac{16\eta_2^2}{\bar{b}(\epsilon)^2} \cdot \frac{\mathbb{E} \left[ \left\| \nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*) \right\|^2 \right]}{\|\Lambda - \Lambda_{K_1 \times K_2}^*\|^2} \leq \frac{16\eta_2^2}{\bar{b}(\epsilon)^2 C_5^2(\epsilon)} \leq \frac{\epsilon}{2} , \end{aligned} \tag{30}$$

where the last inequality holds by choosing

$$C_5(\epsilon) \geq \sqrt{\frac{32\eta_2^2}{\bar{b}(\epsilon)^2 \epsilon}} .$$

Therefore, for sufficiently large  $N$ , by (25) and (30) we can derive

$$\mathbb{P} \left( (E_N)^c \text{ or } \|\nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*)\| \geq \frac{\bar{b}(\epsilon)}{2\eta_2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \right) \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon$$

$$\Rightarrow \mathbb{P} \left( E_N \text{ and } \|\nabla \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*)\| < \frac{\bar{b}(\epsilon)}{2\eta_2} \|\Lambda - \Lambda_{K_1 \times K_2}^*\| \right) > 1 - \epsilon. \quad (31)$$

With (29) and (31), we can obtain that

$$\mathbb{P} \left\{ \hat{G}_{K_1 \times K_2}(\Lambda) - \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*) < 0, \quad \forall \Lambda \in \partial \hat{\Upsilon}_{K_1 \times K_2}(\epsilon) \right\} \geq 1 - \epsilon.$$

Note that the event  $\left\{ \hat{G}_{K_1 \times K_2}(\Lambda_{K_1 \times K_2}^*) > \hat{G}_{K_1 \times K_2}(\Lambda), \quad \forall \Lambda \in \partial \hat{\Upsilon}_{K_1 \times K_2}(\epsilon) \right\}$  implies that there exists a local maximizer in the interior of  $\hat{\Upsilon}_{K_1 \times K_2}(\epsilon)$ . Since  $\hat{G}_{K_1 \times K_2}(\cdot)$  is strictly concave and  $\hat{\Lambda}_{K_1 \times K_2}$  is the unique global maximizer of  $\hat{G}_{K_1 \times K_2}$ , then

$$\mathbb{P} \left( \hat{\Lambda}_{K_1 \times K_2} \in \hat{\Upsilon}_{K_1 \times K_2}(\epsilon) \right) > 1 - \epsilon, \quad (32)$$

i.e.  $\left\| \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\| = O_p \left( \sqrt{\frac{K}{N}} \right).$  ■

### 3.3 Corollary 3.3

The next corollary states that  $\hat{\pi}_K(t, \mathbf{x})$  is arbitrarily close to  $\pi_K^*(t, \mathbf{x})$ .

**Corollary 3.3** *Under Assumption 1.2-1.6, we have*

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})| = O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right),$$

and

$$\int_{\mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})|^2 dF_{T, X}(t, \mathbf{x}) = O_p \left( \frac{K}{N} \right),$$

and

$$\frac{1}{N} \sum_{i=1}^N |\hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_K^*(T_i, \mathbf{X}_i)|^2 = O_p \left( \frac{K}{N} \right).$$

**Proof.** From the proof of Lemma 3.2, we know the facts  $\mathbb{P} \left( \hat{\Lambda}_{K_1 \times K_2} \in \hat{\Upsilon}_{K_1 \times K_2}(\epsilon) \right) > 1 - \epsilon$  and (27). Then for any element  $\tilde{\Lambda}_{K_1 \times K_2}$  lying on the line joining  $\hat{\Lambda}_{K_1 \times K_2}$  and  $\Lambda_{K_1 \times K_2}^*$ , we can have that  $\mathbb{P}(u_{K_1}(t)^\top \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \in \Gamma_2(\epsilon) \text{ for all } (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}) \geq 1 - \epsilon$ , which implies

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho''(u_{K_1}(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}))| = O_p(1). \quad (33)$$

Using Mean Value Theorem, Lemma 3.1, and (33), we can obtain that

$$\begin{aligned}
& \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})| \\
&= \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho' \left( u_{K_1}(t) \hat{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) - \rho' \left( u_{K_1}(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right)| \\
&\leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho''(u_{K_1}(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}))| \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} \left| u_{K_1}(t) \hat{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) - u_{K_1}(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right| \\
&\leq O_p(1) \cdot \|\hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^*\| \cdot \sup_{t \in \mathcal{T}} \|u_{K_1}(t)\| \cdot \sup_{\mathbf{x} \in \mathcal{X}} \|v_{K_2}(\mathbf{x})\| \\
&\leq O_p(1) \cdot O_p \left( \sqrt{\frac{K}{N}} \right) \zeta_1(K_1) \cdot \zeta_2(K_2) = O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right).
\end{aligned}$$

Note that by Mean Value Theorem and (33), we can deduce that

$$\begin{aligned}
& \int_{\mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})|^2 dF_{T, X}(t, \mathbf{x}) \\
&\leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho''(u_{K_1}(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}))|^2 \int_{\mathcal{T} \times \mathcal{X}} \left| u_{K_1}(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right|^2 dF_{T, X}(t, \mathbf{x}) \\
&\leq O_p(1) \cdot \int_{\mathcal{T} \times \mathcal{X}} \left| u_{K_1}(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right|^2 dF_{T, X}(t, \mathbf{x}).
\end{aligned}$$

We estimate  $\int_{\mathcal{T} \times \mathcal{X}} \left| u_{K_1}(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right|^2 dF_{T, X}(t, \mathbf{x})$ . Note that  $\mathbb{E}[u_{K_1}(T)u_{K_1}(T)^\top] = I_{K_1 \times K_1}$ ,  $\mathbb{E}[v_{K_2}(\mathbf{X})v_{K_2}(\mathbf{X})^\top] = I_{K_2 \times K_2}$ , (13) and Assumption 1.3, we can deduce that

$$\begin{aligned}
& \int_{\mathcal{T} \times \mathcal{X}} \left| u_{K_1}(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right|^2 dF_{T, X}(t, \mathbf{x}) \\
&\leq \int_{\mathcal{T} \times \mathcal{X}} u_{K_1}^\top(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) v_{K_2}(\mathbf{x})^\top \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\}^\top u_{K_1}(t) dF_{T, X}(t, \mathbf{x}) \\
&= \int_{\mathcal{T} \times \mathcal{X}} \frac{1}{\pi_0(t, \mathbf{x})} \pi_0(t, \mathbf{x}) u_{K_1}^\top(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) v_{K_2}(\mathbf{x})^\top \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\}^\top u_{K_1}(t) dF_{T, X}(t, \mathbf{x}) \\
&\leq \frac{1}{\eta_1} \int_{\mathcal{T} \times \mathcal{X}} \pi_0(t, \mathbf{x}) \cdot u_{K_1}^\top(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) v_{K_2}(\mathbf{x})^\top \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\}^\top u_{K_1}(t) dF_{T, X}(t, \mathbf{x}) \\
&= \frac{1}{\eta_1} \int_{\mathcal{T}} u_{K_1}^\top(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} \left( \int_{\mathcal{X}} v_{K_2}(\mathbf{x}) v_{K_2}(\mathbf{x})^\top dF_X(\mathbf{x}) \right) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\}^\top u_{K_1}(t) dF_T(t) \\
&= \frac{1}{\eta_1} \int_{\mathcal{T}} u_{K_1}^\top(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\}^\top u_{K_1}(t) dF_T(t) \\
&= \frac{1}{\eta_1} \text{tr} \left( \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\}^\top \int_{\mathcal{T}} u_{K_1}(t) u_{K_1}^\top(t) dF_T(t) \right)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\eta_1} \text{tr} \left( \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\}^\top \right) \\
&= \frac{1}{\eta_1} \cdot \left\| \hat{\Lambda}_{K_2 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\|^2 = O_p \left( \frac{K}{N} \right).
\end{aligned} \tag{34}$$

Then we obtain

$$\int_{\mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})|^2 dF_{T,X}(t, \mathbf{x}) = O_p \left( \frac{K}{N} \right).$$

Similar to (23), we have

$$\begin{aligned}
&\frac{1}{N} \sum_{i=1}^N \left| u_{K_1}^\top(T_i) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{X}_i) \right|^2 - \int_{\mathcal{T} \times \mathcal{X}} \left| u_{K_1}(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right|^2 dF_{T,X}(t, \mathbf{x}) \\
&= O_p \left( \frac{\zeta(K)}{\sqrt{N}} \cdot \left\| \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\|^2 \right) = O_p \left( \frac{\zeta(K)}{\sqrt{N}} \cdot \frac{K}{N} \right) = o_p \left( \frac{K}{N} \right).
\end{aligned} \tag{35}$$

where the last equality holds in light of Assumption 1.6. Hence, with (34) and (35), we have

$$\begin{aligned}
&\frac{1}{N} \sum_{i=1}^N |\hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_K^*(T_i, \mathbf{X}_i)|^2 \\
&\leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho''(u_{K_1}(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}))|^2 \cdot \frac{1}{N} \sum_{i=1}^N \left| u_{K_1}(T_i) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{X}_i) \right|^2 \\
&\leq O_p(1) \cdot \int_{\mathcal{T} \times \mathcal{X}} \left| u_{K_1}(t) \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right|^2 dF_{T,X}(t, \mathbf{x}) + o_p \left( \frac{K}{N} \right) \\
&\leq O_p \left( \frac{K}{N} \right) + o_p \left( \frac{K}{N} \right) = O_p \left( \frac{K}{N} \right).
\end{aligned}$$

■

## 4 Efficient Estimation

### 4.1 Proof of Theorem 4

The proposed estimator  $\hat{\beta}$  is a special case of [Chen, Linton, and Van Keilegom \(2003\)](#), where the authors establish the consistency and asymptotic normality of a class of semiparametric optimization estimators under that the criterion function does not satisfy standard smoothness conditions. The asymptotic distribution of the proposed estimator can be derived by applying

Theorem 2 of [Chen, Linton, and Van Keilegom \(2003\)](#). Using their notation, we denote

$$M_N(\boldsymbol{\beta}, \pi(\cdot)) := \frac{1}{N} \sum_{i=1}^N \pi(T_i, \mathbf{X}_i) L'(Y_i - g(T_i; \boldsymbol{\beta})) m(T_i; \boldsymbol{\beta}),$$

$$M(\boldsymbol{\beta}, \pi(\cdot)) := \mathbb{E}[M_N(\boldsymbol{\beta}, \pi(\cdot))] = \mathbb{E}[\pi(T, \mathbf{X}) L'(Y - g(T; \boldsymbol{\beta})) m(T; \boldsymbol{\beta})].$$

The ordinary derivative  $\Gamma_1(\boldsymbol{\beta}, \pi(\cdot))$  in  $\boldsymbol{\beta}$  of  $M(\boldsymbol{\beta}, \pi(\cdot))$  is

$$\begin{aligned} \Gamma_1(\boldsymbol{\beta}, \pi(\cdot))(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta}) &= : \lim_{\tau \rightarrow 0} \frac{M(\boldsymbol{\beta} + \tau(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta}), \pi(\cdot)) - M(\boldsymbol{\beta}, \pi(\cdot))}{\tau} \\ &= \nabla_{\boldsymbol{\beta}} \mathbb{E}[\pi(T, \mathbf{X}) L'(Y - g(T; \boldsymbol{\beta})) m(T; \boldsymbol{\beta})], \end{aligned}$$

and the functional derivative  $\Gamma_2(\boldsymbol{\beta}, \pi_0(\cdot))[\pi(\cdot) - \pi_0(\cdot)]$  of  $M(\boldsymbol{\beta}, \pi_0(\cdot))$  along the direction  $\pi(\cdot) - \pi_0(\cdot)$  is

$$\begin{aligned} \Gamma_2(\boldsymbol{\beta}, \pi_0(\cdot))[\pi(\cdot) - \pi_0(\cdot)] &:= \lim_{\tau \rightarrow 0} \frac{M(\boldsymbol{\beta}, \pi_0(\cdot) + \tau(\pi(\cdot) - \pi_0(\cdot))) - M(\boldsymbol{\beta}, \pi_0(\cdot))}{\tau} \\ &= \mathbb{E}[(\pi(T, \mathbf{X}) - \pi_0(T, \mathbf{X})) L'(Y - g(T; \boldsymbol{\beta})) m(T; \boldsymbol{\beta})]. \end{aligned}$$

In order to apply Theorem 2 of [Chen, Linton, and Van Keilegom \(2003\)](#), we need to verify their Conditions (2.1)-(2.6) hold. Conditions (2.1)-(2.5) of [Chen, Linton, and Van Keilegom \(2003\)](#) can be easily verified by using following facts:

- Theorem 10 ensures  $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\| \xrightarrow{p} 0$ ;
- Assumption 1.11 implies  $\|M_N(\hat{\boldsymbol{\beta}}, \hat{\pi}_K(\cdot))\| = \left\| N^{-1} \sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i) m(T_i; \hat{\boldsymbol{\beta}}) L'\{Y_i - g(T_i; \hat{\boldsymbol{\beta}})\} \right\| = o_P(1/\sqrt{N})$ ;
- Assumption 1.15 implies  $K = o_p(N^{1/4})$  and  $K^{-\alpha} = o_p(N^{-1/2})$ , then by Theorem 5.6 we have  $\int_{\mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t, \mathbf{x}) - \pi_0(t, \mathbf{x})|^2 dF_{T, \mathbf{X}}(t, \mathbf{x}) = O_p(K^{-\alpha}) + O_p(\sqrt{K/N}) = o_P(N^{-1/2}) + o_P(N^{-3/8}) \leq o_p(N^{-1/4})$ .

The most important step toward the application of Theorem 2 of [Chen, Linton, and Van Keilegom \(2003\)](#) is to check their Condition (2.6) holds, which states that there exists some finite matrix  $V_1$  such that

$$\sqrt{N} \{M_N(\boldsymbol{\beta}_0, \pi_0(\cdot)) + \Gamma_2(\boldsymbol{\beta}_0, \pi_0(\cdot))[\hat{\pi}_K(\cdot) - \pi_0(\cdot)]\} \xrightarrow{d} N(0, V_1). \quad (36)$$

If Conditions (2.1)-(2.6) hold, Theorem 2 of [Chen, Linton, and Van Keilegom \(2003\)](#) ensures that

$$\sqrt{N} \left( \hat{\beta} - \beta_0 \right) \xrightarrow{d} \mathcal{N}(0, \Omega),$$

where  $\Omega := \Gamma_1(\beta_0, \pi_0(\cdot))^{-1} V_1 (\Gamma_1(\beta_0, \pi_0(\cdot))^{-1})^\top = H_0^{-1} V_1 (H_0^{-1})^\top$ . However, [Chen, Linton, and Van Keilegom \(2003\)](#) do not give the expression of  $V_1$  and the verification of (36) is difficult which is also admitted by the authors themselves (see the first paragraph in Section 3.3 of [Chen, Linton, and Van Keilegom \(2003\)](#)). In Section 4.2, we prove (36) holds and give

$$V_1 = \mathbb{E}[\psi(Y, T, \mathbf{X}; \beta_0) \psi(Y, T, \mathbf{X}; \beta_0)^\top].$$

Therefore, we can have  $\Omega = V_{eff}$  which justifies Theorem 3.17.

## 4.2 Proof of (36)

Before proving (36), we prepare some preliminary notation and results that will be used later. Since  $\hat{\Lambda}_{K_1 \times K_2}$  is a unique maximizer of the concave function  $\hat{G}_{K_1 \times K_2}$ , then

$$\frac{1}{N} \sum_{i=1}^N \rho' \left( u_{K_1}(T_i)^\top \hat{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2}(\mathbf{X}_i)^\top - \frac{1}{N^2} \sum_{i=1}^N \sum_{l=1}^N u_{K_1}(T_l) v_{K_2}(\mathbf{X}_i)^\top = 0.$$

Using Mean Value Theorem, we can have

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N \rho' \left( u_{K_1}(T_i)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2}(\mathbf{X}_i)^\top \\ & + \frac{1}{N} \sum_{i=1}^N \rho'' \left( u_{K_1}(T_i)^\top \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) u_{K_1}(T_i)^\top \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{X}_i) v_{K_2}(\mathbf{X}_i)^\top \\ & = \frac{1}{N^2} \sum_{i=1}^N \sum_{l=1}^N u_{K_1}(T_l) v_{K_2}(\mathbf{X}_i)^\top, \end{aligned} \tag{37}$$

where  $\tilde{\Lambda}_{K_1 \times K_2}$  lies on the line joining from  $\hat{\Lambda}_{K_1 \times K_2}$  to  $\Lambda_{K_1 \times K_2}^*$ . We define the following notation:

$$\hat{A}_{K_1 \times K_2} := \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^*, \tag{38}$$

$$\tilde{A}_{K_1 \times K_2} := \tilde{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^*, \tag{39}$$

and

$$\begin{aligned}
A_{K_1 \times K_2}^* &:= \nabla \hat{G}_{K_1 \times K_2} (\Lambda_{K_1 \times K_2}^*) \\
&= \frac{1}{N} \sum_{i=1}^N \rho' \left( u_{K_1}(T_i)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) v_{K_2}(\mathbf{X}_i)^\top - \left( \frac{1}{N} \sum_{l=1}^N u_{K_1}(T_l) \right) \left( \frac{1}{N} \sum_{i=1}^N v_{K_2}(\mathbf{X}_i)^\top \right).
\end{aligned} \tag{40}$$

In light of (26) we have

$$\|A_{K_1 \times K_2}^*\| = O_p \left( \sqrt{\frac{K}{N}} \right).$$

From (37),  $A_{K_1 \times K_2}^*$  can also be written as

$$A_{K_1 \times K_2}^* = -\frac{1}{N} \sum_{i=1}^N \rho'' \left( u_{K_1}(T_i)^\top \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) u_{K_1}(T_i)^\top \left\{ \hat{\Lambda}_{K_1 \times K_2} - \Lambda_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{X}_i) v_{K_2}(\mathbf{X}_i)^\top. \tag{41}$$

We now start to (36). We decompose  $\sqrt{N} \{M_N(\boldsymbol{\beta}_0, \pi_0(\cdot)) + \Gamma_2(\boldsymbol{\beta}_0, \pi_0(\cdot))[\hat{\pi}_K(\cdot) - \pi_0(\cdot)]\}$  as follows:

$$\begin{aligned}
&\sqrt{N} \{M_N(\boldsymbol{\beta}_0, \pi_0(\cdot)) + \Gamma_2(\boldsymbol{\beta}_0, \pi_0(\cdot))[\hat{\pi}_K(\cdot) - \pi_0(\cdot)]\} \\
&= \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \pi_0(T_i, \mathbf{X}_i) L' \{Y_i - g(T_i; \boldsymbol{\beta}_0)\} m(T_i; \boldsymbol{\beta}_0) + \int_{\mathcal{T}} \int_{\mathcal{X}} (\hat{\pi}_K(t, \mathbf{x}) - \pi_0(t, \mathbf{x})) \varepsilon(\mathbf{x}, t; \boldsymbol{\beta}_0) m(t; \boldsymbol{\beta}_0) dF_{X,T}(\mathbf{x}, t) \right\} \\
&= \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \boldsymbol{\beta}_0) \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) (\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})) dF_{X,T}(\mathbf{x}, t) \\
&\quad + \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} (\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})) \varepsilon(\mathbf{x}, t; \boldsymbol{\beta}_0) m(t; \boldsymbol{\beta}_0) dF_{X,T}(\mathbf{x}, t) \\
&\quad + \frac{1}{\sqrt{N}} \sum_{i=1}^N \pi_0(T_i, \mathbf{X}_i) L' \{Y_i - g(T_i; \boldsymbol{\beta}_0)\} m(T_i; \boldsymbol{\beta}_0) \\
&= \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \boldsymbol{\beta}_0) \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) (\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})) dF_{X,T}(\mathbf{x}, t) \tag{42}
\end{aligned}$$

$$\begin{aligned}
&+ \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} (\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})) \varepsilon(\mathbf{x}, t; \boldsymbol{\beta}_0) m(t; \boldsymbol{\beta}_0) dF_{X,T}(\mathbf{x}, t) \tag{43} \\
&\quad - \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \rho'' \left( u_{K_1}^\top(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) m(t; \boldsymbol{\beta}_0) dF_{X,T}(\mathbf{x}, t)
\end{aligned}$$

$$+ \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \rho'' \left( u_{K_1}^\top(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) m(t; \boldsymbol{\beta}_0) dF_{X,T}(\mathbf{x}, t) \tag{44}$$



$$\begin{aligned}
& -\sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) A_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) m(t; \boldsymbol{\beta}_0) dF_{X,T}(\mathbf{x}, t) \\
& + \sqrt{N} \int_{\mathcal{X}} \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) A_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) m(t; \boldsymbol{\beta}_0) dF_{X,T}(\mathbf{x}, t) \quad (45) \\
& + \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \pi_0(T_i, \mathbf{X}_i) m(T_i; \boldsymbol{\beta}_0) \varepsilon(T_i, \mathbf{X}_i; \boldsymbol{\beta}_0) - \mathbb{E} [\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{x}; \boldsymbol{\beta}_0) | \mathbf{X} = \mathbf{X}_i] \right. \\
& \quad \left. - \mathbb{E} [\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{x}; \boldsymbol{\beta}_0) | T = T_i] \right\} \\
& + \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \pi_0(T_i, \mathbf{X}_i) L' \{Y_i - g(T_i; \boldsymbol{\beta}_0)\} m(T_i; \boldsymbol{\beta}_0) - \pi_0(T_i, \mathbf{X}_i) m(T_i; \boldsymbol{\beta}_0) \varepsilon(T_i, \mathbf{X}_i; \boldsymbol{\beta}_0) \right. \quad (46) \\
& \quad \left. + \mathbb{E} [\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) | \mathbf{X} = \mathbf{X}_i] + \mathbb{E} [\pi_0(T, \mathbf{X}) m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{x}; \boldsymbol{\beta}_0) | T = T_i] \right\},
\end{aligned}$$

where  $\hat{A}_{K_1 \times K_2}$  and  $A_{K_1 \times K_2}^*$  are defined in (38) and (41). We show that the terms (42)-(45) are all of  $o_p(1)$ , while the term (46) is asymptotically normal.

**For term (42):** By Lemma 3.1 and Assumption 1.15, we can deduce that

$$\begin{aligned}
& \left\| \sqrt{N} \cdot \mathbb{E} [m(T; \boldsymbol{\beta}_0) \varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0) (\pi_K^*(T, \mathbf{X}) - \pi_0(T, \mathbf{X}))] \right\| \\
& \leq \sqrt{N} \sup_{t \in \mathcal{T}} \|m(t; \boldsymbol{\beta}_0)\| \cdot \mathbb{E} [|\varepsilon(T, \mathbf{X}; \boldsymbol{\beta}_0)|^2]^{\frac{1}{2}} \cdot \mathbb{E} [|\pi_K^*(T, \mathbf{X}) - \pi_0(T, \mathbf{X})|^2]^{\frac{1}{2}} = O \left( \sqrt{N} K^{-\alpha} \right).
\end{aligned}$$

**For term (43):** By Mean Value Theorem and the definition of  $\hat{A}_{K_1 \times K_2}$  in (38), the term (43) is exactly equal to zero.

**For term (44):** We can telescope (44) as follows:

$$\begin{aligned}
& \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \boldsymbol{\beta}_0) \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \rho'' \left( u_{K_1}^\top(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \\
& - \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \boldsymbol{\beta}_0) \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) A_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \\
& = \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \boldsymbol{\beta}_0) \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \left\{ \rho'' \left( u_{K_1}^\top(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) - \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) \right\} \\
& \quad \times u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \quad (47)
\end{aligned}$$

$$\begin{aligned}
& + \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \boldsymbol{\beta}_0) \varepsilon(t, \mathbf{x}; \boldsymbol{\beta}_0) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \\
& \quad \times \left\{ \hat{A}_{K_1 \times K_2} - A_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t). \quad (48)
\end{aligned}$$

For the term (47), by Mean Value Theorem,

$$(47) = \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho'''(\xi_3(t, \mathbf{x})) \left\{ u_{K_1}^\top(t) \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right\} \left\{ u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right\} dF_{X,T}(\mathbf{x}, t).$$

Since  $\xi_3(t, \mathbf{x})$  lies between  $u_{K_1}(t)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})$  and  $u_{K_1}(t)^\top \tilde{\Lambda}_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})$ , which implies  $\xi_3(t, \mathbf{x})$  lies between  $u_{K_1}(t)^\top \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})$  and  $u_{K_1}(t)^\top \hat{\Lambda}_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})$ . Then in light of (27) and (32), we have  $\mathbb{P}(\xi_3(t, \mathbf{x}) \in \Gamma_2(\epsilon), \forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}) > 1 - \epsilon$ , therefore,

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| = O_p(1). \quad (49)$$

With (34), (49), the fact  $\|\tilde{A}_{K_1 \times K_2}\| \leq \|\hat{A}_{K_1 \times K_2}\|$ , Lemma 3.2, and Assumption 1.15, we can derive that

$$\begin{aligned} \|(47)\| &\leq \sqrt{N} \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| \sup_{t \in \mathcal{T}} \|m(t; \beta_0)\| \cdot \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\varepsilon(t, \mathbf{x}; \beta_0)| \\ &\quad \cdot \int_{\mathcal{T}} \int_{\mathcal{X}} \left| u_{K_1}(t)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right| \cdot \left| u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right| dF_{X,T}(\mathbf{x}, t) \\ &\leq \sqrt{N} \cdot O_p(1) \cdot O(1) \cdot O(1) \cdot \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} \left| u_{K_1}(t)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right|^2 dF_{X,T}(\mathbf{x}, t) \right\}^{\frac{1}{2}} \\ &\quad \cdot \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} \left| u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right|^2 dF_{X,T}(\mathbf{x}, t) \right\}^{\frac{1}{2}} \\ &= \sqrt{N} \cdot O_p(1) \cdot O(1) \cdot O(1) \cdot O_p\left(\sqrt{\frac{K}{N}}\right) \cdot O_p\left(\sqrt{\frac{K}{N}}\right) = O_p\left(\sqrt{\frac{K^2}{N}}\right) \quad (\text{by } ((34))). \end{aligned} \quad (50)$$

For the term (48), we first compute the probability order of  $\|A_{K_1 \times K_2}^* - \hat{A}_{K_1 \times K_2}\|$ . Using (41), the fact  $\rho''(v) = -\rho'(v)$  and Mean Value Theorem, we have

$$\begin{aligned} &A_{K_1 \times K_2}^* - \hat{A}_{K_1 \times K_2} \\ &= -\frac{1}{N} \sum_{i=1}^N \rho''(u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \\ &\quad - \frac{1}{N} \sum_{i=1}^N \rho'''(\xi_3(T_i, \mathbf{X}_i)) \left\{ u_{K_1}(T_i)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right\} u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \\ &\quad - \hat{A}_{K_1 \times K_2} \\ &= \frac{1}{N} \sum_{i=1}^N \left\{ \rho'(u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) - \hat{A}_{K_1 \times K_2} \right\} \end{aligned} \quad (51)$$

$$-\frac{1}{N} \sum_{i=1}^N \rho'''(\xi_3(T_i, \mathbf{X}_i)) \left\{ u_{K_1}(T_i)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right\} u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i). \quad (52)$$

For the term (51), by (13) we can write  $\hat{A}_{K_1 \times K_2}$  as

$$\hat{A}_{K_1 \times K_2} = \mathbb{E}_{T, \mathbf{X}} \left[ \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}) v_{K_2}^\top(\mathbf{X}) \right],$$

where  $\mathbb{E}_{T, \mathbf{X}}[\cdot]$  denotes taking expectation with respect to  $(T, \mathbf{X})$ . We telescope (51) as follows:

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N \left\{ \rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) - \hat{A}_{K_1 \times K_2} \right\} \\ &= \frac{1}{N} \sum_{i=1}^N \left\{ \left\{ \rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) - \pi_0(T_i, \mathbf{X}_i) \right\} u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \right\} \end{aligned} \quad (53)$$

$$\begin{aligned} & - \frac{1}{N} \sum_{i=1}^N \left\{ \pi_0(T_i, \mathbf{X}_i) u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \right. \\ & \quad \left. - \mathbb{E}_{T, \mathbf{X}} \left[ \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}) v_{K_2}^\top(\mathbf{X}) \right] \right\}. \end{aligned} \quad (54)$$

For the term (53), by Lemmas 3.1 and 3.2, we have that

$$\begin{aligned} & \left\| \frac{1}{N} \sum_{i=1}^N \left\{ \rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) - \pi_0(T_i, \mathbf{X}_i) \right\} u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \right\| \\ & \leq \frac{1}{N} \sum_{i=1}^N \left| \rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) - \pi_0(T_i, \mathbf{X}_i) \right| \cdot \|u_{K_1}(T_i)\|^2 \cdot \|v_{K_2}(\mathbf{X}_i)\|^2 \cdot \|\hat{A}_{K_1 \times K_2}\| \\ & = \left\{ \mathbb{E} \left[ \left| \rho' \left( u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}) \right) - \pi_0(T, \mathbf{X}) \right| \cdot \|u_{K_1}(T)\|^2 \cdot \|v_{K_2}(\mathbf{X})\|^2 \right] + O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) \right\} \cdot \|\hat{A}_{K_1 \times K_2}\| \\ & \leq \left\{ \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})| \cdot \frac{1}{\eta_1} \cdot \mathbb{E} [\pi_0(T, \mathbf{X}) \|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2] + O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) \right\} \cdot \|\hat{A}_{K_1 \times K_2}\| \\ & = \left\{ \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})| \cdot \frac{1}{\eta_1} \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \cdot \mathbb{E} [\|v_{K_2}(\mathbf{X})\|^2] + O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) \right\} \cdot \|\hat{A}_{K_1 \times K_2}\| \\ & \leq \left\{ O(K^{-\alpha} \zeta(K)) \cdot O(K_1) \cdot O(K_2) + O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) \right\} \cdot O_p \left( \sqrt{\frac{K}{N}} \right) \\ & \leq O_p \left( N^{-\frac{1}{2}} \zeta(K) \cdot K^{\frac{3}{2}-\alpha} \right). \end{aligned}$$

For the term (54), define the linear map  $\mathcal{J}(\cdot) : \mathbb{R}^{K_1 \times K_2} \rightarrow \mathbb{R}$  by

$$\mathcal{J}(M) := \frac{1}{N} \sum_{i=1}^N \left\{ \pi_0(T_i, \mathbf{X}_i) u_{K_1}(T_i) u_{K_1}(T_i)^\top M v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) - \mathbb{E}_{T, \mathbf{X}} [\pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top M v_{K_2}(\mathbf{X}) v_{K_2}^\top(\mathbf{X})] \right\},$$

then (54) =  $\mathcal{J}(\hat{A}_{K_1 \times K_2})$ . For any fixed  $M \in \mathbb{R}^{K_1 \times K_2}$ , by (13) and  $M = \mathbb{E}[\pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top M v_{K_2}(\mathbf{X}) v_{K_2}^\top(\mathbf{X})]$ , then we have

$$\begin{aligned} & \mathbb{E} [\mathcal{J}(M)^2] \\ &= \frac{1}{N} \cdot \mathbb{E} \left[ \left\| \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top M v_{K_2}(\mathbf{X}) v_{K_2}^\top(\mathbf{X}) - \mathbb{E} [\pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top M v_{K_2}(\mathbf{X}) v_{K_2}^\top(\mathbf{X})] \right\|^2 \right] \\ &\leq \frac{1}{N} \cdot \mathbb{E} \left[ \left\| \pi_0(T, \mathbf{X}) u_{K_1}(T) u_{K_1}(T)^\top M v_{K_2}(\mathbf{X}) v_{K_2}^\top(\mathbf{X}) \right\|^2 \right] \\ &\leq \frac{1}{N} \cdot \eta_2 \cdot \mathbb{E} \left[ \pi_0(T, \mathbf{X}) \cdot \|u_{K_1}(T)\|^4 \|v_{K_2}(\mathbf{X})\|^4 \right] \cdot \|M\|^2 \\ &= \frac{1}{N} \cdot \eta_2 \cdot \mathbb{E}[\|u_{K_1}(T)\|^4] \cdot \mathbb{E}[\|v_{K_2}(\mathbf{X})\|^4] \cdot \|M\|^2 \\ &\leq \frac{1}{N} \cdot \eta_2 \cdot \zeta_1(K)^2 \cdot \zeta_2(K)^2 \cdot \mathbb{E}[\|u_{K_1}(T)\|^2] \cdot \mathbb{E}[\|v_{K_2}(\mathbf{X})\|^2] \cdot \|M\|^2 \\ &= \|M\|^2 \cdot O\left(\zeta(K)^2 \frac{K}{N}\right). \end{aligned}$$

Using Chebyshev's inequality we have

$$|\mathcal{J}(M)| = \|M\| O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right),$$

then in light of Lemma 3.2,

$$(54) = \mathcal{J}(\hat{A}_{K_1 \times K_2}) = \|\hat{A}_{K_1 \times K_2}\| O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) = O_p \left( \zeta(K) \frac{K}{N} \right).$$

Therefore,

$$(51) = (53) + (54) = O_p \left( N^{-\frac{1}{2}} \zeta(K) \cdot K^{\frac{3}{2}-\alpha} \right) + O_p \left( \zeta(K) \frac{K}{N} \right).$$

For the term (52), we can deduce that

$$\left\| \frac{1}{N} \sum_{i=1}^N \rho'''(\xi_3(T_i, \mathbf{X}_i)) \left\{ u_{K_1}(T_i)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right\} \left\{ u_{K_1}(T_i) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \right\} \right\|$$

$$\begin{aligned}
&\leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| \cdot \|\hat{A}_{K_1 \times K_2}\|^2 \cdot \frac{1}{N} \sum_{i=1}^N \|u_{K_1}(T_i)\|^3 \cdot \|v_{K_2}(\mathbf{X}_i)\|^3 \\
&\leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| \cdot \|\hat{A}_{K_1 \times K_2}\|^2 \cdot \zeta_1(K_1) \cdot \zeta_2(K_2) \cdot \frac{1}{N} \sum_{i=1}^N \|u_{K_1}(T_i)\|^2 \cdot \|v_{K_2}(\mathbf{X}_i)\|^2 \\
&\leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| \cdot \|\hat{A}_{K_1 \times K_2}\|^2 \cdot \zeta(K) \left\{ \mathbb{E} [\|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2] + O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) \right\} \\
&\leq O_p(1) \cdot O_p \left( \frac{K}{N} \right) \cdot \zeta(K) \cdot O(K) = O_p \left( \zeta(K) \frac{K^2}{N} \right),
\end{aligned}$$

where the fourth inequality follows from (49) and Lemma 3.2. Now, we can obtain

$$\begin{aligned}
\|\hat{A}_{K_1 \times K_2} - A_{K_1 \times K_2}^*\| &= (51) + (52) = O_p \left( N^{-\frac{1}{2}} \zeta(K) K^{\frac{3}{2}-\alpha} \right) + O_p \left( \zeta(K) \frac{K}{N} \right) + O_p \left( \zeta(K) \frac{K^2}{N} \right) \\
&= O_p \left( N^{-\frac{1}{2}} \zeta(K) \cdot K^{\frac{3}{2}-\alpha} \right) + O_p \left( \zeta(K) \frac{K^2}{N} \right). \tag{55}
\end{aligned}$$

Using (55), Assumptions 1.9 and 1.15, for large enough  $N$ , we have

$$\begin{aligned}
(48) &= \left\| \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho''(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) u_{K_1}^\top(t) \left\{ \hat{A}_{K_1 \times K_2} - A_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right\| \\
&\leq \sqrt{N} \sup_{t \in \mathcal{T}} \|m(t; \beta_0)\| \sup_{\gamma \in \Gamma_1} |\rho''(\gamma)| \cdot \mathbb{E} [\varepsilon(T, \mathbf{X}; \beta_0)^2]^{\frac{1}{2}} \cdot \left[ \int_{\mathcal{T} \times \mathcal{X}} \left( u_{K_1}(t) \left\{ \hat{A}_{K_1 \times K_2} - A_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right)^2 dF_{T,X}(t, \mathbf{x}) \right]^{\frac{1}{2}} \\
&\leq \sqrt{N} \cdot O(1) \cdot O(1) \cdot O(1) \cdot O(1) \cdot O(\|\hat{A}_{K_1 \times K_2} - A_{K_1 \times K_2}^*\|) \\
&\leq O_p \left( \zeta(K) \cdot K^{\frac{3}{2}-\alpha} \right) + O_p \left( \zeta(K) \frac{K^2}{\sqrt{N}} \right), \tag{56}
\end{aligned}$$

where the second inequality holds since by using the same argument of establishing (34), we have

$$\int_{\mathcal{T} \times \mathcal{X}} \left( u_{K_1}(t) \left\{ \hat{A}_{K_1 \times K_2} - A_{K_1 \times K_2}^* \right\} v_{K_2}(\mathbf{x}) \right)^2 dF_{T,X}(t, \mathbf{x}) = O(\|\hat{A}_{K_1 \times K_2} - A_{K_1 \times K_2}^*\|).$$

Therefore, by combining (50) and (56), we can obtain that

$$\begin{aligned}
(44) &= (47) + (48) = O_p \left( \sqrt{\frac{K^2}{N}} \right) + O_p \left( \zeta(K) \cdot K^{\frac{3}{2}-\alpha} \right) + O_p \left( \zeta(K) \frac{K^2}{\sqrt{N}} \right) \\
&= O_p \left( \zeta(K) \cdot K^{\frac{3}{2}-\alpha} \right) + O_p \left( \zeta(K) \frac{K^2}{\sqrt{N}} \right).
\end{aligned}$$

**For term (45):** By the definition of  $A_{K_1 \times K_2}^*$  in (40), we have

$$(45) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho''(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) u_{K_1}^\top(t) \right. \quad (57)$$

$$\begin{aligned} & \times \left\{ u_{K_1}(T_i) \rho'(u_{K_1}(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) v_{K_2}^\top(\mathbf{X}_i) \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) + m(T_i; \beta_0) \varepsilon(T_i, \mathbf{X}_i; \beta_0) \pi_0(T_i, \mathbf{X}_i) \Big\} \\ & - \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta^*) \rho''(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}^\top(\mathbf{x})) u_{K_1}^\top(t) \left( \frac{1}{N} \sum_{l=1}^N u_{K_1}(T_l) \right) \right. \quad (58) \\ & \times \left( \frac{1}{N} \sum_{j=1}^N v_{K_2}^\top(\mathbf{X}_j) \right) v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) + \mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \beta_0) \varepsilon(T, \mathbf{X}; \beta_0) | \mathbf{X} = \mathbf{X}_i] \\ & \left. + \mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \beta_0) \varepsilon(T, \mathbf{X}; \beta_0) | T = T_i] \right\}. \end{aligned}$$

We shall show that both (57) and (58) are of  $o_p(1)$ . Noting  $\rho'' = -\rho'$ , we can telescope (57) as follows:

$$(57) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho'(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) u_{K_1}^\top(t) \right. \quad (59)$$

$$\begin{aligned} & \times \left\{ u_{K_1}(T_i) \left[ -\rho'(u_{K_1}(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) + \pi_0(T_i, \mathbf{X}_i) \right] v_{K_2}^\top(\mathbf{X}_i) \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \Big\} \\ & - \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \left\{ \rho'(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) - \pi_0(t, \mathbf{x}) \right\} u_{K_1}^\top(t) \right. \quad (60) \\ & \times \left\{ u_{K_1}(T_i) \pi_0(T_i, \mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \Big\} \end{aligned}$$

$$\begin{aligned} & - \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \pi_0(t, \mathbf{x}) u_{K_1}^\top(t) \left\{ u_{K_1}(T_i) \pi_0(T_i, \mathbf{X}_i) v_{K_2}^\top(\mathbf{X}_i) \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right. \\ & \left. + m(T_i; \beta_0) \varepsilon(T_i, \mathbf{X}_i; \beta_0) \pi_0(T_i, \mathbf{X}_i) \right\}. \quad (61) \end{aligned}$$

We shall show that (59), (60) and (61) are all of  $o_p(1)$ . Note that second moment of (59) is

$$\begin{aligned} \mathbb{E}[(59)^2] &= \mathbb{E} \left[ \left| \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho'(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) u_{K_1}^\top(t) \right. \right. \\ & \quad \times \left. \left\{ u_{K_1}(T_i) \left[ -\rho'(u_{K_1}(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) + \pi_0(T_i, \mathbf{X}_i) \right] v_{K_2}^\top(\mathbf{X}_i) \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right|^2 \Big] \\ &= \mathbb{E} \left[ \left| \int_{\mathcal{T}} \int_{\mathcal{X}} \pi_0(t, \mathbf{x}) \cdot m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \left[ \frac{\rho'(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}))}{\pi_0(t, \mathbf{x})} \right] u_{K_1}^\top(t) \right. \right. \\ & \quad \times \left. \left\{ u_{K_1}(T_i) \left[ -\rho'(u_{K_1}(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) + \pi_0(T_i, \mathbf{X}_i) \right] v_{K_2}^\top(\mathbf{X}_i) \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right|^2 \Big] \\ &\leq \mathbb{E} \left[ \left| \int_{\mathcal{T}} \int_{\mathcal{X}} \pi_0(t, \mathbf{x}) \cdot m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \left[ \frac{\rho'(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}))}{\pi_0(t, \mathbf{x})} \right] u_{K_1}^\top(t) \left\{ u_{K_1}(T_i) v_{K_2}^\top(\mathbf{X}_i) \right\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right|^2 \right] \end{aligned}$$

$$\begin{aligned}
& \times \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} \left\{ -\rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) + \pi_0(t, \mathbf{x}) \right\}^2 \\
& = \left\{ \mathbb{E} \left[ \left| m(T_i; \beta_0) \varepsilon(T_i, \mathbf{X}_i; \beta_0) \left[ \frac{\rho' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right)}{\pi_0(T_i, \mathbf{X}_i)} \right] \right|^2 \right] + o(1) \right\} \times \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} \left\{ -\pi_K^*(t, \mathbf{x}) + \pi_0(t, \mathbf{x}) \right\}^2 \\
& = O(1) \cdot O(K^{-2\alpha} \zeta(K)^2) = O(K^{-2\alpha} \zeta(K)^2) \rightarrow 0, \text{ (by Assumption 1.6)}
\end{aligned}$$

where the third equality holds because

$$\int_{\mathcal{T}} \int_{\mathcal{X}} \pi_0(t, \mathbf{x}) \cdot m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \left[ \frac{\rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right)}{\pi_0(t, \mathbf{x})} \right] u_{K_1}^\top(t) \{u_{K_1}(T) v_{K_2}^\top(\mathbf{X})\} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t)$$

is the weighted  $L^2$ -projection of  $m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \left[ \frac{\rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right)}{\pi_0(t, \mathbf{x})} \right]$  on the space linearly spanned by  $\{u_{K_1}(t), v_{K_2}(\mathbf{x})\}$  with the weighted measure  $\pi_0(t, \mathbf{x}) dF_{T,X}(t, \mathbf{x})$ . Similarly, we can also show (60) and (61) are of  $o_p(1)$ . Therefore, (57) is of  $o_p(1)$ .

For the term (58), since  $\rho''(v) = -\rho'(v)$  and the fact  $\mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \beta_0) \varepsilon(T, \mathbf{X}; \beta_0)] = 0$ , we telescope it as follows:

$$\begin{aligned}
(58) &= \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \left( \frac{1}{N} \sum_{l=1}^N u_{K_1}(T_l) - \mathbb{E}[u_{K_1}(T)] \right) \\
&\quad \times \left( \frac{1}{N} \sum_{j=1}^N v_{K_2}^\top(\mathbf{X}_j) - \mathbb{E}[v_{K_2}^\top(\mathbf{X})] \right) v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \tag{62}
\end{aligned}$$

$$\begin{aligned}
&+ \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \mathbb{E}[u_{K_1}(T)] v_{K_2}^\top(\mathbf{X}_i) v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right. \\
&\quad \left. - \mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \beta_0) \varepsilon(T, \mathbf{X}; \beta_0) | \mathbf{X} = \mathbf{X}_i] \right\} \tag{63}
\end{aligned}$$

$$\begin{aligned}
&+ \frac{1}{\sqrt{N}} \sum_{l=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) u_{K_1}(T_l) \mathbb{E}[v_{K_2}^\top(\mathbf{X})] v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right. \\
&\quad \left. - \mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \beta_0) \varepsilon(T, \mathbf{X}; \beta_0) | T = T_l] \right\} \tag{64}
\end{aligned}$$

$$\begin{aligned}
&- \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \mathbb{E}[u_{K_1}(T)] \mathbb{E}[v_{K_2}^\top(\mathbf{X})] v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right. \\
&\quad \left. - \mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \beta_0) \varepsilon(T, \mathbf{X}; \beta_0)] \right\}. \tag{65}
\end{aligned}$$

For the term (62), since

$$\left\| \frac{1}{N} \sum_{l=1}^N u_{K_1}(T_l) - \mathbb{E}[u_{K_1}(T)] \right\| = O_p \left( \sqrt{\frac{K_1}{N}} \right),$$

$$\left\| \frac{1}{N} \sum_{j=1}^N v_{K_2}(\mathbf{X}_j) - \mathbb{E}[v_{K_2}(\mathbf{X})] \right\| = O_p \left( \sqrt{\frac{K_2}{N}} \right),$$

$$\sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} \left| \rho' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) \right| = O(1),$$

and by Assumptions 1.9, 1.10, and 1.15, we can deduce that

$$(62) = \sqrt{N} \cdot O(\zeta(K)) O_p \left( \sqrt{\frac{K_1}{N}} \right) O_p \left( \sqrt{\frac{K_2}{N}} \right) = O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) = o_p(1).$$

For the term (63), noting the fact that  $\mathbb{E}[\pi_0(T, \mathbf{X}) m(T; \beta_0) \varepsilon(T, \mathbf{X}; \beta_0) | \mathbf{X}] = \int_{\mathcal{T}} m(t; \beta_0) \varepsilon(t, \mathbf{X}; \beta_0) dF_T(t)$ , we can rewrite (63) as follows:

$$(63) = \frac{1}{\sqrt{N}} \sum_{j=1}^N \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \frac{\pi_K^*(t, \mathbf{x})}{\pi_0(t, \mathbf{x})} u_{K_1}^\top(t) \mathbb{E}[u_{K_1}(T)] v_{K_2}^\top(\mathbf{X}_j) v_{K_2}(\mathbf{x}) dF_X(\mathbf{x}) dF_T(t) \right. \\ \left. - \int_{\mathcal{T}} m(t; \beta_0) \varepsilon(t, \mathbf{X}_j; \beta_0) dF_T(t) \right\}.$$

By computing the second moment of (63), we can obtain that

$$\mathbb{E} \left[ \left\| \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \frac{\pi_K^*(t, \mathbf{x})}{\pi_0(t, \mathbf{x})} u_{K_1}^\top(t) \mathbb{E}[u_{K_1}(T)] v_{K_2}^\top(\mathbf{X}) v_{K_2}(\mathbf{x}) dF_X(\mathbf{x}) dF_T(t) - \int_{\mathcal{T}} m(t; \beta_0) \varepsilon(t, \mathbf{X}; \beta_0) dF_T(t) \right\|^2 \right] \\ \leq \mathbb{E} \left[ \left\| \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \frac{\pi_K^*(t, \mathbf{x})}{\pi_0(t, \mathbf{x})} u_{K_1}^\top(t) u_{K_1}(T^*) v_{K_2}^\top(\mathbf{X}^*) v_{K_2}(\mathbf{x}) dF_X(\mathbf{x}) dF_T(t) - m(T^*; \beta_0) \varepsilon(T^*, \mathbf{X}^*; \beta_0) \right\|^2 \right] \\ \leq 2 \cdot \mathbb{E} \left[ \left\| \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) u_{K_1}^\top(t) u_{K_1}(T^*) v_{K_2}^\top(\mathbf{X}^*) v_{K_2}(\mathbf{x}) dF_X(\mathbf{x}) dF_T(t) - m(T^*; \beta_0) \varepsilon(T^*, \mathbf{X}^*; \beta_0) \right\|^2 \right] \\ + 2 \cdot \mathbb{E} \left[ \left\| \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \frac{\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})}{\pi_0(t, \mathbf{x})} u_{K_1}^\top(t) u_{K_1}(T^*) v_{K_2}^\top(\mathbf{X}^*) v_{K_2}(\mathbf{x}) dF_X(\mathbf{x}) dF_T(t) \right\|^2 \right] \rightarrow 0,$$

where  $T^* \sim F_T$ ,  $\mathbf{X}^* \sim F_X$ , and  $T^*$  is independent of  $\mathbf{X}^*$ ; the first inequality holds by Jensen's inequality; the last convergence result follows from Lemma 3.1 and the fact that

$$\int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) u_{K_1}^\top(t) u_{K_1}(T^*) v_{K_2}^\top(\mathbf{X}^*) v_{K_2}(\mathbf{x}) dF_X(\mathbf{x}) dF_T(t)$$

is the  $L^2$ -projection of  $m(T^*; \beta_0) \varepsilon(T^*, \mathbf{X}^*; \beta_0)$  on the space spanned by  $\{u_{K_1}(T^*), v_{K_2}(\mathbf{X}^*)\}$ . Thus (63) is of  $o_p(1)$  by Chebyshev's inequality. Similar argument can be applied to show that both (64) and (65) are of  $o_p(1)$ . Therefore, we can have that

$$|(58)| \leq |(62)| + |(63)| + |(64)| = o_p(1).$$



Then, we can obtain that

$$|(45)| \leq |(57)| + |(58)| = o_p(1) .$$

Summing up all orders (42)-(45) and using Assumption 1.15, we have

$$\begin{aligned} & (42) + (43) + (44) + (45) \\ &= O(\sqrt{N}K^{-\alpha}) + 0 + \left\{ O_p\left(\zeta(K) \cdot K^{\frac{3}{2}-\alpha}\right) + O_p\left(\zeta(K)\frac{K^2}{\sqrt{N}}\right) \right\} + o_p(1) = o_p(1). \end{aligned}$$

## 5 Some Extensions

### 5.1 Proof of Theorem 6

(Proof of Consistency). Let

$$\hat{\gamma} = \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \sum_{i=1}^N u_{K_1}(T_i) \hat{\pi}_K(T_i, \mathbf{X}_i) Y_i \right]$$

then  $\hat{\theta}_K(t) = \hat{\gamma}^\top u_{K_1}(t)$ . By assumption, there exists  $\gamma^* \in \mathbb{R}^{K_1}$  such that

$$\sup_{t \in \mathcal{T}} |\mathbb{E}[\pi_0(T, X)Y | T = t] - (\gamma^*)^\top u_{K_1}(t)| = O(K_1^{-\bar{\alpha}}). \quad (66)$$

We first claim that

$$\|\hat{\gamma} - \gamma^*\| = O_p\left(\zeta(K)\sqrt{\frac{K}{N}} + \zeta(K)K^{-\alpha} + K_1^{-\bar{\alpha}}\right), \quad (67)$$

and the proof will be established later. With the claim (67), we first show that  $\int_{\mathcal{T}} |\hat{\theta}_K(t) - \theta(t)|^2 dF_T(t) = O_p\left(\frac{\zeta(K)^2 K}{N} + \zeta(K)^2 K^{-2\alpha} + K_1^{-2\bar{\alpha}}\right)$ . Note that

$$\begin{aligned} & \int_{\mathcal{T}} [\hat{\theta}_K(t) - \theta(t)]^2 dF_T(t) \\ &= \int_{\mathcal{T}} [\hat{\gamma}^\top u_{K_1}(t) - (\gamma^*)^\top u_{K_1}(t) + (\gamma^*)^\top u_{K_1}(t) - \theta(t)]^2 dF_T(t) \\ &\leq 2(\hat{\gamma} - \gamma^*)^\top \left[ \int_{\mathcal{T}} u_{K_1}(t) u_{K_1}(t)^\top dF_T(t) \right] (\hat{\gamma} - \gamma^*) + 2 \int_{\mathcal{T}} [(\gamma^*)^\top u_{K_1}(t) - \theta(t)]^2 dF_T(t) \end{aligned}$$

$$\begin{aligned}
&\leq 2\|\hat{\gamma} - \gamma^*\|^2 \cdot \lambda_{\max}(\mathbb{E}[u_{K_1}(T)u_{K_1}(T)^\top]) + 2\sup_{t \in \mathcal{T}} |(\gamma^*)^\top u_{K_1}(t) - \theta(t)|^2 \\
&= O_p\left(\zeta(K)^2 \frac{K}{N} + \zeta(K)^2 K^{-2\alpha} + K_1^{-2\tilde{\alpha}}\right).
\end{aligned}$$

With the claim (67), we next show that  $\sup_{t \in \mathcal{T}} |\hat{\theta}_K(t) - \theta(t)| = O_p[\zeta_1(K_1)(\zeta(K)\sqrt{K/N} + \zeta(K)K^{-\alpha} + K_1^{-\alpha})]$ . Note that

$$\begin{aligned}
&\sup_{t \in \mathcal{T}} |\hat{\theta}_K(t) - \theta(t)| \\
&= \sup_{t \in \mathcal{T}} |\hat{\gamma}^\top u_{K_1}(t) - (\gamma^*)^\top u_{K_1}(t) + (\gamma^*)^\top u_{K_1}(t) - \theta(t)| \\
&\leq \sup_{t \in \mathcal{T}} \|u_{K_1}(t)\| \cdot \|\hat{\gamma} - \gamma^*\| + \sup_{t \in \mathcal{T}} |(\gamma^*)^\top u_{K_1}(t) - \theta(t)| \\
&\leq \zeta_1(K_1) \cdot \left\{ O_p\left(\sup_{(t,x) \in \mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t,x) - \pi_0(t,x)|\right) + O(K_1^{-\tilde{\alpha}}) \right\} + O(K_1^{-\tilde{\alpha}}) \\
&\leq O_p\left[\zeta_1(K_1) \left(\zeta(K)\sqrt{\frac{K}{N}} + \zeta(K)K^{-\alpha} + K_1^{-\alpha}\right)\right] + O(K_1^{-\tilde{\alpha}}) \\
&= O_p\left[\zeta_1(K_1) \left(\zeta(K)\sqrt{\frac{K}{N}} + \zeta(K)K^{-\alpha} + K_1^{-\tilde{\alpha}}\right)\right].
\end{aligned}$$

Finally, we turn back to prove the claim (67). Note that

$$\begin{aligned}
\hat{\gamma} - \gamma^* &= \left[\sum_{i=1}^N u_{K_1}(T_i)u_{K_1}(T_i)^\top\right]^{-1} \left[\sum_{i=1}^N \hat{\pi}_K(T_i, \mathbf{X}_i)u_{K_1}(T_i)Y_i\right] - \gamma^* \\
&= \left[\sum_{i=1}^N u_{K_1}(T_i)u_{K_1}(T_i)^\top\right]^{-1} \left[\sum_{i=1}^N u_{K_1}(T_i) \{\hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_0(T_i, \mathbf{X}_i)\} Y_i\right] \\
&\quad + \left[\sum_{i=1}^N u_{K_1}(T_i)u_{K_1}(T_i)^\top\right]^{-1} \left[\sum_{i=1}^N u_{K_1}(T_i) \{\pi_0(T_i, \mathbf{X}_i)Y_i - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i)Y_i|T_i]\}\right] \\
&\quad + \left[\sum_{i=1}^N u_{K_1}(T_i)u_{K_1}(T_i)^\top\right]^{-1} \left[\sum_{i=1}^N u_{K_1}(T_i) \{\mathbb{E}[\pi_0(T_i, \mathbf{X}_i)Y_i|T_i] - (\gamma^*)^\top u_{K_1}(T_i)\}\right] \\
&= A_{1N} + A_{2N} + A_{3N}
\end{aligned}$$

where

$$\begin{aligned}
A_{1N} &= \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \sum_{i=1}^N u_{K_1}(T_i) \{ \hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_0(T_i, \mathbf{X}_i) \} Y_i \right], \\
A_{2N} &= \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \sum_{i=1}^N u_{K_1}(T_i) \{ \pi_0(T_i, \mathbf{X}_i) Y_i - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] \} \right], \\
A_{3N} &= \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \sum_{i=1}^N u_{K_1}(T_i) \{ \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] - (\gamma^*)^\top u_{K_1}(T_i) \} \right].
\end{aligned}$$

We first compute the probability order of  $A_{1N}$ . We use the following notation:

$$\begin{aligned}
\hat{H}_N &:= (\{ \hat{\pi}_K(T_1, X_1) - \pi_0(T_1, X_1) \} Y_1, \dots, \{ \hat{\pi}_K(T_N, X_N) - \pi_0(T_N, X_N) \} Y_N)^\top, \\
U_{N \times K_1} &:= (u_{K_1}(T_1), \dots, u_{K_1}(T_N))^\top, \\
\hat{\Phi}_{K_1 \times K_1} &:= \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top.
\end{aligned}$$

Then we can obtain that

$$\begin{aligned}
\|A_{1N}\|^2 &= \left\| \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \sum_{i=1}^N u_{K_1}(T_i) \{ \hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_0(T_i, \mathbf{X}_i) \} Y_i \right] \right\|^2 \\
&= N^{-2} \text{tr} \left( \hat{\Phi}_{K_1 \times K_1}^{-1} U_{N \times K_1}^\top \hat{H}_N \hat{H}_N^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \\
&= N^{-2} \text{tr} \left( U_{N \times K_1}^\top \hat{H}_N \hat{H}_N^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \\
&= N^{-2} \text{tr} \left( \hat{\Phi}_{K_1 \times K_1}^{-1/2} U_{N \times K_1}^\top \hat{H}_N \hat{H}_N^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1/2} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \\
&\leq \lambda_{\max}(\hat{\Phi}_{K_1 \times K_1}^{-1}) N^{-2} \text{tr} \left( \hat{\Phi}_{K_1 \times K_1}^{-1/2} U_{N \times K_1}^\top \hat{H}_N \hat{H}_N^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1/2} \right) \\
&= \lambda_{\max}(\hat{\Phi}_{K_1 \times K_1}^{-1}) N^{-1} \text{tr} \left( \hat{H}_N \hat{H}_N^\top U_{N \times K_1} (U_{N \times K_1}^\top U_{N \times K_1})^{-1} U_{N \times K_1}^\top \right) \\
&\leq [\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})]^{-1} N^{-1} \|\hat{H}_N\|^2 \\
&= [\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})]^{-1} \cdot \frac{1}{N} \sum_{i=1}^N \{ \hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_0(T_i, \mathbf{X}_i) \}^2 Y_i^2 \\
&\leq [\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})]^{-1} \sup_{(t,x) \in \mathcal{T} \times \mathcal{X}} | \hat{\pi}_K(t, x) - \pi_0(t, x) |^2 \cdot \frac{1}{N} \sum_{i=1}^N Y_i^2 \\
&\leq O_p(1) \cdot O_p \left( \zeta(K)^2 K^{-2\alpha} + \frac{\zeta(K)^2 K}{N} \right) \cdot O_p(1)
\end{aligned}$$

$$=O_p\left(\zeta(K)^2K^{-2\alpha} + \frac{\zeta(K)^2K}{N}\right),$$

where the first inequality follows from the fact that  $\text{tr}(AB) \leq \lambda_{\max}(B)\text{tr}(A)$  for any symmetric matrix  $B$  and positive semidefinite matrix  $A$ , the second inequality follows from the same fact and the fact that  $U_{N \times K_1}(U_{N \times K_1}^\top U_{N \times K_1})^{-1}U_{N \times K_1}^\top$  is a projection matrix with maximum eigenvalue 1, and the fourth inequality follows from the facts that  $|\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})|^{-1} = O_p(1)$ ,  $\sup_{(t,x) \in \mathcal{T} \times \mathcal{X}} |\hat{\pi}_K(t,x) - \pi_0(t,x)| = O_p\left(\zeta(K)K^{-\alpha} + \zeta(K)\sqrt{K/N}\right)$  and  $N^{-1} \sum_{i=1}^N Y_i^2 = O_p(1)$ .

Next, we compute the probability order of  $A_{2N}$ . We can deduce that

$$\begin{aligned} \|A_{2N}\|^2 &= \left\| \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \sum_{i=1}^N u_{K_1}(T_i) \varepsilon_i \right] \right\|^2 \\ &= N^{-2} \text{tr} \left( \hat{\Phi}_{K_1 \times K_1}^{-1} U_{N \times K_1}^\top \mathcal{E}_N \mathcal{E}_N^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \\ &= N^{-2} \text{tr} \left( U_{N \times K_1}^\top \mathcal{E}_N \mathcal{E}_N^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \\ &\leq [\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})]^{-2} N^{-2} \|U_{N \times K_1}^\top \mathcal{E}_N\|^2 = O_p\left(\frac{K_1}{N}\right), \end{aligned}$$

where the last equality follows that  $|\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})|^{-1} = O_p(1)$  and  $N^{-2} \|U_{N \times K_1}^\top \mathcal{E}_N\|^2 = O_p(K_1/N)$  by Markov's inequality.

We finally compute the probability order of  $A_{3N}$ . We define the notation

$$R_N(\gamma^*) = \left( \left\{ \mathbb{E}[\pi_0(T_1, X_1) Y_1 | T_1] - (\gamma^*)^\top u_{K_1}(T_1) \right\}, \dots, \left\{ \mathbb{E}[\pi_0(T_N, X_N) Y_N | T_N] - (\gamma^*)^\top u_{K_1}(T_N) \right\} \right)^\top,$$

then it follows that with probability approaching to 1,

$$\begin{aligned} \|A_{3N}\|^2 &= \left\| \left[ \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \sum_{i=1}^N u_{K_1}(T_i) \left\{ \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] - (\gamma^*)^\top u_{K_1}(T_i) \right\} \right] \right\|^2 \\ &= N^{-2} \left\| \hat{\Phi}_{K_1 \times K_1}^{-1} U_{N \times K_1}^\top R_N(\gamma^*) \right\|^2 \\ &= N^{-2} \text{tr} \left( \hat{\Phi}_{K_1 \times K_1}^{-1} U_{N \times K_1}^\top R_N(\gamma^*) R_N(\gamma^*)^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \\ &= N^{-2} \text{tr} \left( U_{N \times K_1}^\top R_N(\gamma^*) R_N(\gamma^*)^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \\ &= N^{-2} \text{tr} \left( \hat{\Phi}_{K_1 \times K_1}^{-1/2} U_{N \times K_1}^\top R_N(\gamma^*) R_N(\gamma^*)^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1/2} \hat{\Phi}_{K_1 \times K_1}^{-1} \right) \end{aligned}$$

$$\begin{aligned}
&\leq \lambda_{\max}(\hat{\Phi}_{K_1 \times K_1}^{-1}) N^{-2} \text{tr} \left( \hat{\Phi}_{K_1 \times K_1}^{-1/2} U_{N \times K_1}^\top R_N(\gamma^*) R_N(\gamma^*)^\top U_{N \times K_1} \hat{\Phi}_{K_1 \times K_1}^{-1/2} \right) \\
&= \lambda_{\max}(\hat{\Phi}_{K_1 \times K_1}^{-1}) N^{-1} \text{tr} \left( R_N(\gamma^*) R_N(\gamma^*)^\top U_{N \times K_1} (U_{N \times K_1}^\top U_{N \times K_1})^{-1} U_{N \times K_1}^\top \right) \\
&\leq [\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})]^{-1} N^{-1} \|R_N(\gamma^*)\|^2 \\
&= [\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})]^{-1} \cdot \frac{1}{N} \sum_{i=1}^N \left\{ \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] - (\gamma^*)^\top u_{K_1}(T_i) \right\}^2 = O_p(K_1^{-2\tilde{\alpha}}),
\end{aligned}$$

where the first inequality follows from the fact that  $\text{tr}(AB) \leq \lambda_{\max}(B) \text{tr}(A)$  for any symmetric matrix  $B$  and positive semidefinite matrix  $A$ , the second inequality follows from the same fact and the fact that  $U_{N \times K_1} (U_{N \times K_1}^\top U_{N \times K_1})^{-1} U_{N \times K_1}^\top$  is a projection matrix with maximum eigenvalue 1, and the last equality follows from the fact that  $|\lambda_{\min}(\hat{\Phi}_{K_1 \times K_1})|^{-1} = O_p(1)$  and the fact that  $\frac{1}{N} \sum_{i=1}^N \left\{ \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] - (\gamma^*)^\top u_{K_1}(T_i) \right\}^2 \leq \sup_{t \in \mathcal{T}} |\mathbb{E}[\pi_0(T, X) Y | T] - (\gamma^*)^\top u_{K_1}(t)|^2 = O(K_1^{-2\tilde{\alpha}})$ . Thus we complete the proof of (67).

**(Proof of Asymptotic Normality).** We have the following decomposition for  $\hat{\theta}(t) - \theta(t)$ :

$$\begin{aligned}
&\hat{\theta}_K(t) - \theta(t) \\
&= u_{K_1}(t)^\top (\hat{\gamma} - \gamma^*) + [(\gamma^*)^\top u_{K_1}(t) - \theta(t)] \\
&= u_{K_1}(t)^\top \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \left\{ \hat{\pi}_K(T_i, \mathbf{X}_i) Y_i - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] \right\} \right] \\
&\quad + u_{K_1}(t)^\top \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \cdot \left\{ \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] - (\gamma^*)^\top u_{K_1}(T_i) \right\} \right] \\
&\quad + \left[ (\gamma^*)^\top u_{K_1}(t) - \theta(t) \right] \\
&= b_{1N}(t) + b_{2N}(t) + b_{3N}(t),
\end{aligned}$$

where

$$\begin{aligned}
b_{1N}(t) &= u_{K_1}(t)^\top \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \left\{ \hat{\pi}_K(T_i, \mathbf{X}_i) Y_i - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] \right\} \right] \\
b_{2N}(t) &= u_{K_1}(t)^\top \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \cdot \left\{ \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i] - (\gamma^*)^\top u_{K_1}(T_i) \right\} \right] \\
b_{3N}(t) &= (\gamma^*)^\top u_{K_1}(t) - \theta(t).
\end{aligned}$$

Then we have that

$$\begin{aligned} & \sqrt{N}V_K(t)^{-1/2} \left[ \hat{\theta}_K(t) - \theta(t) \right] \\ &= \sqrt{N}V_K(t)^{-1/2}b_{1N}(t) + \sqrt{N}V_K(t)^{-1/2}b_{2N}(t) + \sqrt{N}V_K(t)^{-1/2}b_{3N}(t). \end{aligned}$$

We shall show that  $b_{1N}(t)$  contributes to the asymptotic variance; and  $b_{2N}(t)+b_{3N}(t)$  contributes to the asymptotic bias which is asymptotically negligible. Thus to complete the proof of asymptotic normality, it is sufficient to prove the following results:

- (i)  $V_K \geq c\|u_{K_1}(t)\|^2$  for some  $c > 0$ ;
- (ii)  $\sqrt{N}V_K^{-1/2}b_{1N}(t) \xrightarrow{d} N(0, 1)$ ;
- (iii)  $\sqrt{N}V_K^{-1/2}b_{2N}(t) = o_p(1)$ ;
- (iv)  $\sqrt{N}V_K^{-1/2}b_{3N}(t) = o_p(1)$ .

We first prove Result (i). Note that  $\lambda_{\min}(\Sigma_{K_1 \times K_1}) \geq \underline{c}_{\sigma^2} \lambda_{\min}(\Phi_{K_1 \times K_1}) = \underline{c}_{\sigma^2} \lambda_{\min}(I_{K_1 \times K_1}) \geq \underline{c}_{\sigma^2}$ , we can have

$$\begin{aligned} V_K &= u_{K_1}^\top(t) \Phi_{K_1 \times K_1}^{-1} \Sigma_{K_1 \times K_1} \Phi_{K_1 \times K_1}^{-1} u_{K_1}(t) \\ &\geq \lambda_{\min}(\Sigma_{K_1 \times K_1}) u_{K_1}^\top(t) \Phi_{K_1 \times K_1}^{-1} \Phi_{K_1 \times K_1}^{-1} u_{K_1}(t) \\ &\geq \underline{c}_{\sigma^2} \|u_{K_1}(t)\|^2. \end{aligned}$$

For the claim (ii). Let

$$\begin{aligned} \tilde{b}_{1N}(t) &= u_{K_1}(t)^\top \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \right]^{-1} \\ &\times \left[ \frac{1}{N} \sum_{i=1}^N u_{K_1}(T_i) \left\{ \pi_0(T_i, \mathbf{X}_i) Y_i - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | T_i, \mathbf{X}_i] + \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i | \mathbf{X}_i] - \mathbb{E}[\pi_0(T_i, \mathbf{X}_i) Y_i] \right\} \right]. \end{aligned}$$

Similar to the proof of (36), we can have that

$$\sqrt{N}V_K(t)^{-1/2} \cdot (b_{1N}(t) - \tilde{b}_{1N}(t)) = o_p(1).$$

Then

$$\sqrt{N}V_K(t)^{-1/2}b_{1N}(t)$$

$$\begin{aligned}
&= \sqrt{N} V_K(t)^{-1/2} \tilde{b}_{1N}(t) + o_p(1) \\
&= \sqrt{N} V_K(t)^{-1/2} u_{K_1}(t)^\top \hat{\Phi}_{K_1 \times K_1}^{-1} N^{-1} U_{N \times K_1}^\top \mathcal{E}_N \\
&= \sqrt{N} V_K(t)^{-1/2} u_{K_1}(t)^\top \Phi_{K_1 \times K_1}^{-1} N^{-1} U_{N \times K_1}^\top \mathcal{E}_N + \sqrt{N} V_K(t)^{-1/2} u_{K_1}(t)^\top [\hat{\Phi}_{K_1 \times K_1}^{-1} - \Phi_{K_1 \times K_1}^{-1}] N^{-1} U_{N \times K_1}^\top \mathcal{E}_N \\
&= B_{1N,1}(t) + B_{1N,2}(t) .
\end{aligned}$$

For  $B_{1N,1}(t)$ , we can simply apply the Liapounov CLT and show that  $B_{1N,1}(t) \xrightarrow{d} N(0, 1)$ . For  $B_{1N,2}(t)$ , let  $\mathbf{T} = (T_1, \dots, T_N)$ , we can obtain that

$$\begin{aligned}
&\mathbb{E} [B_{1N,2}(t)^2 | \mathbf{T}] \\
&= N^{-1} V_K(t) u_{K_1}(t)^\top [\hat{\Phi}_{K_1 \times K_1}^{-1} - \Phi_{K_1 \times K_1}^{-1}] \cdot \mathbb{E} [U_{N \times K_1}^\top \mathcal{E}_N \mathcal{E}_N^\top U_{N \times K_1} | \mathbf{T}] \cdot [\hat{\Phi}_{K_1 \times K_1}^{-1} - \Phi_{K_1 \times K_1}^{-1}] u_{K_1}(t) \\
&= \lambda_{\max} (N^{-1} \mathbb{E} [U_{N \times K_1}^\top \mathcal{E}_N \mathcal{E}_N^\top U_{N \times K_1} | \mathbf{T}]) V_K(t) u_{K_1}(t)^\top [\hat{\Phi}_{K_1 \times K_1}^{-1} - \Phi_{K_1 \times K_1}^{-1}] \cdot [\hat{\Phi}_{K_1 \times K_1}^{-1} - \Phi_{K_1 \times K_1}^{-1}] u_{K_1}(t) \\
&\leq \lambda_{\max} (N^{-1} \mathbb{E} [U_{N \times K_1}^\top \mathcal{E}_N \mathcal{E}_N^\top U_{N \times K_1} | \mathbf{T}]) \{V_K^{-1} \|u_{K_1}(t)\|^2\} \cdot \|\hat{\Phi}_{K_1 \times K_1}^{-1} - \Phi_{K_1 \times K_1}^{-1}\|^2 \\
&= O_p(1) O_p(1) o_p(1) = o_p(1)
\end{aligned}$$

where we use the fact that  $N^{-1} \mathbb{E} [U_{N \times K_1}^\top \mathcal{E}_N \mathcal{E}_N^\top U_{N \times K_1} | \mathbf{T}] = N^{-1} \sum_{i=1}^N u_{K_1}(T_i) u_{K_1}(T_i)^\top \sigma^2(T_i)$  has bounded maximum eigenvalue. Therefore,  $B_{1N,1}(t) = o_p(1)$  by the conditional Chebyshev's inequality. Thus (ii) holds.

For (iii), by Cauchy-Schwarz's inequality, we can obtain that

$$\begin{aligned}
&\sqrt{N} V_K^{-1/2} |b_{2N}(t)| \\
&= N^{-1/2} V_K^{-1/2} \left| u_{K_1}(t)^\top \hat{\Phi}_{K_1 \times K_1}^{-1} U_{N \times K_1}^\top R_N(\gamma^*) \right| \\
&\leq V_K^{-1/2} \left\{ u_{K_1}(t)^\top \hat{\Phi}_{K_1 \times K_1}^{-1} (N^{-1} U_{N \times K_1}^\top U_{N \times K_1}) \hat{\Phi}_{K_1 \times K_1}^{-1} u_{K_1}(t) \right\}^{\frac{1}{2}} \{R_N(\gamma^*)^\top R_N(\gamma^*)\}^{\frac{1}{2}} \\
&\leq V_K^{-1/2} \left\{ u_{K_1}(t)^\top \hat{\Phi}_{K_1 \times K_1}^{-1} u_{K_1}(t) \right\}^{\frac{1}{2}} \{R_N(\gamma^*)^\top R_N(\gamma^*)\}^{\frac{1}{2}} \\
&\leq \{V_K^{-1/2} \|u_{K_1}(t)\|\} \cdot |\lambda_{\max}(\Phi_{K_1 \times K_1}^{-1})|^{\frac{1}{2}} \cdot O(N^{\frac{1}{2}} \cdot K_1^{-\tilde{\alpha}}) \\
&= O(1) \cdot O_p(1) \cdot o_p(1) = o_p(1) .
\end{aligned}$$

Similarly, we can show that  $\sqrt{N} V_K^{-1/2} |b_{3N}(t)| = o_p(1)$ . This completes the proof of the Theorem.

## 5.2 Proof of Theorem 7

Let  $\mu(t, \mathbf{x}) = \mathbb{E}[Y|T = t, \mathbf{X} = \mathbf{x}]$ , we decompose  $\sqrt{N}(\hat{\psi}_K - \psi)$  as follows:

$$\begin{aligned} \sqrt{N}(\hat{\psi}_K - \psi) &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \{ \hat{\pi}_K(T_i, \mathbf{X}_i) Y_i - \psi \} \\ &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ (\hat{\pi}_K(T_i, \mathbf{X}_i) - \pi_K^*(T_i, \mathbf{X}_i)) Y_i - \int_{\mathcal{T}} \int_{\mathcal{X}} (\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})) \mu(\mathbf{x}, t) dF_{X,T}(\mathbf{x}, t) \right\} \end{aligned} \quad (68)$$

$$+ \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ (\pi_K^*(T_i, \mathbf{X}_i) - \pi_0(T_i, \mathbf{X}_i)) Y_i - \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t, \mathbf{x}) (\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})) dF_{X,T}(\mathbf{x}, t) \right\} \quad (69)$$

$$+ \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t, \mathbf{x}) (\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})) dF_{X,T}(\mathbf{x}, t) \quad (70)$$

$$\begin{aligned} &+ \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} (\hat{\pi}_K(t, \mathbf{x}) - \pi_K^*(t, \mathbf{x})) \mu(\mathbf{x}, t) dF_{X,T}(\mathbf{x}, t) \\ &\quad - \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t, \mathbf{x}) \rho'' \left( u_{K_1}^\top(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) m(t; \beta_0) dF_{X,T}(\mathbf{x}, t) \end{aligned} \quad (71)$$

$$\begin{aligned} &+ \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \varepsilon(t, \mathbf{x}; \beta_0) \rho'' \left( u_{K_1}^\top(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \\ &\quad - \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t, \mathbf{x}) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) A_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \end{aligned} \quad (72)$$

$$+ \sqrt{N} \int_{\mathcal{X}} \mu(t, \mathbf{x}) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) A_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \quad (73)$$

$$\begin{aligned} &+ \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \pi_0(T_i, \mathbf{X}_i) \mu(T_i, \mathbf{X}_i) - \mathbb{E}[\pi_0(T, \mathbf{X}) \mu(T, \mathbf{X}) | \mathbf{X} = \mathbf{X}_i] \right. \\ &\quad \left. - \mathbb{E}[\pi_0(T, \mathbf{X}) \mu(T, \mathbf{X}) | T = T_i] + \mathbb{E}[\pi_0(T, \mathbf{X}) \mu(T, \mathbf{X})] \right\} \\ &+ \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ \pi_0(T_i, \mathbf{X}_i) Y_i - \pi_0(T_i, \mathbf{X}_i) \mu(T_i, \mathbf{X}_i) + \mathbb{E}[\pi_0(T, \mathbf{X}) \mu(T, \mathbf{X}) | T = T_i] \right. \\ &\quad \left. + \mathbb{E}[\pi_0(T, \mathbf{X}) \mu(T, \mathbf{X}) | \mathbf{X} = \mathbf{X}_i] - \mathbb{E}[\pi_0(T, \mathbf{X}) \mu(T, \mathbf{X})] - \psi \right\}. \end{aligned} \quad (74)$$

Using the similar argument for showing that (42)-(45) are all  $o_p(1)$ , we can obtain that the terms (70)-(73) are all  $o_p(1)$ . The term (74) is asymptotically normal and attains the efficiency bound.



Hence to complete the proof, it suffices to show (68) and (69) are of  $o_P(1)$ .

**For term (68):**

Denoting (68) by  $W_K$  and applying Mean Value Theorem twice, we can obtain

$$\begin{aligned}
W_K &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ Y_i \rho'' \left( u_{K_1}^\top(T_i) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right. \\
&\quad \left. - \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t; \mathbf{x}) \rho'' \left( u_{K_1}^\top(t) \tilde{\Lambda}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right] \\
&= \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ Y_i \rho'' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right. \\
&\quad \left. - \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t; \mathbf{x}) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right] \\
&\quad + \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ Y_i \rho''' \left( \xi_3(T_i, \mathbf{X}_i) \right) \left\{ u_{K_1}(T_i)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right\} u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right] \\
&\quad - \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t; \mathbf{x}) \rho''' \left( \xi_3(t, \mathbf{x}) \right) \left\{ u_{K_1}(t)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right\} u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \\
&= W_{1K} + W_{2K} + W_{3K},
\end{aligned}$$

where  $\tilde{A}_{K_1 \times K_2}$  is defined in (39), and

$$\begin{aligned}
W_{1K} &:= \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ L' \{Y_i - g(T_i; \beta_0)\} m(T_i; \beta_0) \rho'' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right. \\
&\quad \left. - \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho'' \left( u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \right) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right], \\
W_{2K} &:= \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ L' \{Y_i - g(T_i; \beta_0)\} m(T_i; \beta_0) \rho''' \left( \xi_3(T_i, \mathbf{X}_i) \right) \left\{ u_{K_1}(T_i)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right\} u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right], \\
W_{3K} &:= -\sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} m(t; \beta_0) \varepsilon(t, \mathbf{x}; \beta_0) \rho''' \left( \xi_3(t, \mathbf{x}) \right) \left\{ u_{K_1}(t)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right\} u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t),
\end{aligned}$$

and  $\xi_3(t, \mathbf{x})$  lies between  $u_{K_1}(t) \tilde{\Lambda}_{K_1 \times K_2}^\top v_{K_2}(\mathbf{x})$  and  $u_{K_1}(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})$ .

For the term  $W_{1K}$ , we have

$$W_{1K} := \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ Y_i \rho'' \left( u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i) \right) u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right]$$

$$\begin{aligned}
& - \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t; \mathbf{x}) \rho''(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \Big] \\
& = \text{tr} \left\{ \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ Y_i \rho''(u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) v_{K_2}(\mathbf{X}_i) u_{K_1}(T_i)^\top \right. \right. \\
& \quad \left. \left. - \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t; \mathbf{x}) \rho''(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) v_{K_2}(\mathbf{x}) u_{K_1}^\top(t) dF_{X,T}(\mathbf{x}, t) \right] \hat{A}_{K_1 \times K_2} \right\} \\
& = \text{tr} \left\{ U_{K_2 \times K_1} \hat{A}_{K_1 \times K_2} \right\},
\end{aligned}$$

where

$$\begin{aligned}
U_{K_2 \times K_1} := & \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ Y_i \rho''(u_{K_1}^\top(T_i) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}_i)) v_{K_2}(\mathbf{X}_i) u_{K_1}(T_i)^\top \right. \\
& \left. - \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t; \mathbf{x}) \rho''(u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x})) v_{K_2}(\mathbf{x}) u_{K_1}^\top(t) dF_{X,T}(\mathbf{x}, t) \right].
\end{aligned}$$

We compute the second moment of  $U_{K_2 \times K_1}$  to get that

$$\begin{aligned}
& \mathbb{E} [\|U_{K_2 \times K_1}\|^2] = \mathbb{E} [\text{tr}\{(U_{K_2 \times K_1})^\top U_{K_2 \times K_1}\}] \\
& = \mathbb{E} \left[ Y^2 \rho''(u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}))^2 \|v_{K_2}(\mathbf{X})\|^2 \|u_{K_1}(T)\|^2 \right] \\
& \quad - \text{tr} \left\{ \mathbb{E} [\mu(T; \mathbf{X}) \rho''(u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X})) u_{K_1}(T) v_{K_2}^\top(\mathbf{X})] \cdot \mathbb{E} [\mu(T; \mathbf{X}) \rho''(u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X})) v_{K_2}(\mathbf{X}) u_{K_1}(T)^\top] \right\} \\
& \leq \mathbb{E} \left[ Y^2 \rho''(u_{K_1}^\top(T) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{X}))^2 \|v_{K_2}(\mathbf{X})\|^2 \|u_{K_1}(T)\|^2 \right] \\
& \leq \mathbb{E}[Y^2] \cdot a_3 \cdot O(\zeta(K)^2) = O(\zeta(K)^2),
\end{aligned}$$

where  $a_3 := \sup_{\gamma \in \Gamma_1} |\rho''(\gamma)|^2 < +\infty$ , the second inequality follows from this definition and the fact that  $u_{K_1}^\top(t) \Lambda_{K_1 \times K_2}^* v_{K_2}(\mathbf{x}) \in \Gamma_1$ ,  $\forall (t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}$  when  $K$  is large enough. Then in light of Chebyshev's inequality, Lemma 3.2 and Assumption 1.15, we have

$$|W_{1K}| \leq \|U_{K_2 \times K_1}\| \|\hat{A}_{K_1 \times K_2}\| = O_p(\zeta(K)) O_p\left(\sqrt{\frac{K}{N}}\right) = O_p\left(\zeta(K) \sqrt{\frac{K}{N}}\right).$$

For the term  $W_{3K}$ . With (34), (49), the fact  $\|\tilde{A}_{K_1 \times K_2}\| \leq \|\hat{A}_{K_1 \times K_2}\|$ , Lemma 3.2, and Assumption 1.15, we can derive that

$$\|W_{3K}\| = \left\| \sqrt{N} \int_{\mathcal{T}} \int_{\mathcal{X}} \mu(t, \mathbf{x}) \rho'''(\xi_3(t, \mathbf{x})) \left\{ u_{K_1}(t)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right\} u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) dF_{X,T}(\mathbf{x}, t) \right\|$$

$$\begin{aligned}
&\leq \sqrt{N} \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| \cdot \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\mu(t, \mathbf{x})| \\
&\quad \cdot \int_{\mathcal{T}} \int_{\mathcal{X}} \left| u_{K_1}(t)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right| \cdot \left| u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right| dF_{X,T}(\mathbf{x}, t) \\
&\leq \sqrt{N} \cdot O_p(1) \cdot O(1) \cdot \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} \left| u_{K_1}(t)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right|^2 dF_{X,T}(\mathbf{x}, t) \right\}^{\frac{1}{2}} \\
&\quad \cdot \left\{ \int_{\mathcal{T}} \int_{\mathcal{X}} \left| u_{K_1}^\top(t) \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{x}) \right|^2 dF_{X,T}(\mathbf{x}, t) \right\}^{\frac{1}{2}} \\
&= \sqrt{N} \cdot O_p(1) \cdot O(1) \cdot O_p\left(\sqrt{\frac{K}{N}}\right) \cdot O_p\left(\sqrt{\frac{K}{N}}\right) = O_p\left(\sqrt{\frac{K^2}{N}}\right) \quad (\text{by } ((34))). \tag{75}
\end{aligned}$$

For the term  $W_{2K}$ , we can deduce that

$$\begin{aligned}
&\left\| \frac{1}{\sqrt{N}} \sum_{i=1}^N \left[ Y_i \rho'''(\xi_3(T_i, \mathbf{X}_i)) \left\{ u_{K_1}(T_i)^\top \tilde{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right\} u_{K_1}(T_i)^\top \hat{A}_{K_1 \times K_2} v_{K_2}(\mathbf{X}_i) \right] \right\| \\
&\leq \left\{ \frac{1}{\sqrt{N}} \sum_{i=1}^N |Y_i| \cdot \|u_{K_1}(T_i)\|^2 \|v_{K_2}(\mathbf{X}_i)\|^2 \right\} \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| \cdot \|\hat{A}_{K_1 \times K_2}\|^2 \\
&\leq \sqrt{N} \left\{ \frac{1}{N} \sum_{i=1}^N Y_i^2 \right\}^{\frac{1}{2}} \left\{ \frac{1}{N} \sum_{i=1}^N \|u_{K_1}(T_i)\|^4 \|v_{K_2}(\mathbf{X}_i)\|^4 \right\}^{\frac{1}{2}} \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\rho'''(\xi_3(t, \mathbf{x}))| \cdot \|\hat{A}_{K_1 \times K_2}\|^2 \\
&\leq \sqrt{N} \cdot O_p(1) \cdot \{\zeta_1(K_1) \zeta_2(K_2)\} \cdot \left\{ \frac{1}{N} \sum_{i=1}^N \|u_{K_1}(T_i)\|^2 \|v_{K_2}(\mathbf{X}_i)\|^2 \right\}^{\frac{1}{2}} \cdot O_p(1) \cdot O_p\left(\frac{K}{N}\right) \\
&\leq \sqrt{N} \cdot O_p(1) \cdot \zeta(K) \cdot \left\{ \mathbb{E} [\|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2] + O_p\left(\zeta(K) \sqrt{\frac{K}{N}}\right) \right\}^{\frac{1}{2}} \cdot O_p(1) \cdot O_p\left(\frac{K}{N}\right) \\
&\leq \sqrt{N} \cdot O_p(1) \cdot \zeta(K) \cdot O_p(\sqrt{K}) \cdot O_p(1) \cdot O_p\left(\frac{K}{N}\right) = O_p\left(\zeta(K) \sqrt{\frac{K^3}{N}}\right),
\end{aligned}$$

where the fourth inequality follows from the fact that

$$\begin{aligned}
&\mathbb{E} \left[ \left( \frac{1}{N} \sum_{i=1}^N \|u_{K_1}(T_i)\|^2 \|v_{K_2}(\mathbf{X}_i)\|^2 - \mathbb{E} [\|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2] \right)^2 \right] \\
&= \frac{1}{N} \cdot \mathbb{E} \left[ (\|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2 - \mathbb{E} [\|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2])^2 \right] \\
&\leq \frac{1}{N} \cdot \mathbb{E} [\|u_{K_1}(T)\|^4 \|v_{K_2}(\mathbf{X})\|^4] \leq \frac{1}{N} \cdot \zeta_1(K_1)^2 \zeta_2(K_2)^2 \cdot \mathbb{E} [\|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2] \\
&= \frac{1}{N} \cdot \zeta_1(K_1)^2 \zeta_2(K_2)^2 \cdot \mathbb{E} \left[ \frac{1}{\pi_0(T, \mathbf{X})} \cdot \pi_0(T, \mathbf{X}) \|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2 \right] \\
&\leq \frac{1}{N} \cdot \frac{1}{\eta_1} \cdot \zeta_1(K_1)^2 \zeta_2(K_2)^2 \cdot \mathbb{E} [\pi_0(T, \mathbf{X}) \|u_{K_1}(T)\|^2 \|v_{K_2}(\mathbf{X})\|^2] \\
&= \frac{1}{N} \cdot \frac{1}{\eta_1} \cdot \zeta_1(K_1)^2 \zeta_2(K_2)^2 \cdot \mathbb{E} [\|u_{K_1}(T)\|^2] \cdot \mathbb{E} [\|v_{K_2}(\mathbf{X})\|^2] = O\left(\frac{K}{N} \zeta(K)^2\right).
\end{aligned}$$

Therefore, we can obtain that

$$(68) = W_{1K} + W_{2K} + W_{3K} = O_p \left( \zeta(K) \sqrt{\frac{K}{N}} \right) + O_p \left( \zeta(K) \sqrt{\frac{K^3}{N}} \right) + O_p \left( \sqrt{\frac{K^2}{N}} \right) = O_p \left( \zeta(K) \sqrt{\frac{K^3}{N}} \right).$$

Finally, it follows that the term (68) is of  $o_p(1)$  in light of Assumption 1.15.

**For term (69):** Note that

$$\begin{aligned} & \mathbb{E} \left[ \left\| \frac{1}{\sqrt{N}} \sum_{i=1}^N \left\{ (\pi_K^*(T_i, \mathbf{X}_i) - \pi_0(T_i, \mathbf{X}_i)) Y_i - \mathbb{E}[\mu(T, \mathbf{X}) (\pi_K^*(T, \mathbf{X}) - \pi_0(T, \mathbf{X}))] \right\} \right\|^2 \right] \\ & \leq \mathbb{E} [(\pi_K^*(T, \mathbf{X}) - \pi_0(T, \mathbf{X}))^2 Y^2] \leq \sup_{(t, \mathbf{x}) \in \mathcal{T} \times \mathcal{X}} |\pi_K^*(t, \mathbf{x}) - \pi_0(t, \mathbf{x})|^2 \cdot \mathbb{E} [Y^2] \leq O(\zeta(K)^2 K^{-2\alpha}), \end{aligned}$$

where the last equality follows from Lemma 3.1. Then by Chebyshev's inequality, we can claim that the term (69) is of  $O_p(\zeta(K)K^{-\alpha})$ .

## 6 Variance Estimation in Monte Carlo Simulations

### 6.1 Proposed Variance Estimator

In Monte Carlo simulations, the estimated parameter is the average treatment effects, which corresponds to a differentiable loss function  $L(v) = v^2$ . The variance estimator can be simply defined as follows:

$$\begin{aligned} \hat{V}_{eff} &= \left[ \frac{1}{N} \sum_{i=1}^N m(T_i; \hat{\beta}) m(T_i; \hat{\beta})^\top \right]^{-1} \\ &\quad \times \left\{ \frac{1}{N} \sum_{j=1}^N \left[ \hat{\psi}(Y_j, T_j, \mathbf{X}_j; \hat{\beta}) - \text{mean}(\hat{\psi}) \right] \cdot \left[ \hat{\psi}(Y_j, T_j, \mathbf{X}_j; \hat{\beta}) - \text{mean}(\hat{\psi}) \right]^\top \right\} \\ &\quad \times \left[ \frac{1}{N} \sum_{i=1}^N m(T_i; \hat{\beta}) m(T_i; \hat{\beta})^\top \right]^{-1}, \end{aligned} \tag{76}$$

where

$$\hat{\psi}(Y, T, \mathbf{X}; \hat{\beta}) = \hat{\pi}_{K'}(T, \mathbf{X}) m(T; \hat{\beta}) \{Y - \hat{\mathbb{E}}[Y|T, \mathbf{X}]\} + \hat{\mathbb{E}}[\{Y - g(T; \beta_0)\} \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X}],$$

and

$$\widehat{\mathbb{E}}[Y|T, \mathbf{X}] = \left[ \sum_{i=1}^N Y_i w_{K_0}(T_i, \mathbf{X}_i)^\top \right] \left[ \sum_{i=1}^N w_{K_0}(T_i, \mathbf{X}_i) w_{K_0}(T_i, \mathbf{X}_i)^\top \right]^{-1} w_{K_0}(T, \mathbf{X})$$

and

$$\begin{aligned} \widehat{\mathbb{E}}[\{Y - g(T; \beta_0)\} \pi_0(T, \mathbf{X}) m(T; \beta_0) | \mathbf{X}] &= \left[ \sum_{i=1}^N \hat{\pi}_{K'}(T_i, \mathbf{X}_i) (Y_i - g(T_i; \hat{\beta})) m(T_i; \hat{\beta}) v_{M_0}(\mathbf{X}_i)^\top \right] \\ &\times \left[ \sum_{i=1}^N v_{M_0}(\mathbf{X}_i) v_{M_0}(\mathbf{X}_i)^\top \right]^{-1} v_{M_0}(\mathbf{X}), \end{aligned}$$

and

$$\text{mean}(\hat{\psi}) := \frac{1}{N} \sum_{j=1}^N \hat{\psi}(Y_j, T_j, \mathbf{X}_j; \hat{\beta}).$$

## 6.2 True Values of $V_{eff}$ in Monte Carlo Simulations

In Section 9 of the main paper, Monte Carlo simulations on variance estimation are performed. Computing the bias, standard deviation, and RMSE of the variance estimator  $\hat{V}_{eff}$  requires to compute the true variance  $V_{eff}$ . We describe how to compute  $V_{eff}$  under DGP-L1 in Section 6.2.1 and DGP-NL1 in Section 6.2.2. DGP-L2 and DGP-NL2 are omitted since they can be handled in the same way as DGP-L1 and DGP-NL1.

To reduce notation, the single covariate  $X_1$  is redefined as  $X$ . Note that the influence function is written as

$$\begin{aligned} \psi(Y, T, X; \beta_0) &= \pi_0(T, X) m(T; \beta_0) \{Y - g(T; \beta_0)\} - \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | T, X] \\ &\quad + \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | T] + \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | X] \\ &\quad - \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0)] \\ &= \pi_0(T, X) m(T; \beta_0) \{Y - g(T; \beta_0)\} - \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | T, X] \\ &\quad + \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | X]. \end{aligned} \tag{77}$$

### 6.2.1 DGP-L1

Recall that DGP-L1 is

$$T = 1 + \rho_{T,X} \cdot X + \xi, \quad Y = 1 + X + T + \varepsilon. \quad (\rho_{T,X} = 0.2)$$

We have

$$\mathbb{E}[Y|T, X] = 1 + X + T, \quad \mathbb{E}[Y(t)] = g(t; \beta_0) = 1 + t.$$

We directly compute

$$\begin{aligned} \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | T = t, X = x] &= \int \{y - g(t; \beta_0)\} \cdot \frac{f_T(t)}{f_{T|X}(t|x)} \cdot m(t) \cdot f_{Y|T,X}(y|t, x) dy \\ &= m(t) \cdot \pi_0(t, x) \cdot \{\mathbb{E}[Y|X = x, T = t] - g(t; \beta_0)\} = m(t) \cdot \pi_0(t, x) \cdot \{1 + x + t - (1 + t)\} \\ &= m(t) \cdot \pi_0(t, x) \cdot x \end{aligned} \quad (78)$$

and

$$\begin{aligned} \mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | X = x] &= \int m(t) \cdot \pi_0(t, x) \cdot x \cdot f_{T|X}(t|x) dt \\ &= x \cdot \int m(t) f(t) dt = x \cdot \mathbb{E}[m(T)]. \end{aligned} \quad (79)$$

Substitute (78) and (79) into (77) to get

$$\begin{aligned} \psi(Y, T, X; \beta_0) &= \pi_0(T, X) m(T) \{Y - 1 - T\} - m(T) \cdot \pi_0(T, X) \cdot X + X \cdot \mathbb{E}[m(T)] \\ &= \pi_0(T, X) \cdot m(T) \cdot \{Y - 1 - X - T\} + X \cdot \mathbb{E}[m(T)] \\ &= \pi_0(T, X) \cdot \{Y - 1 - X - T\} \cdot \begin{bmatrix} 1 \\ T \end{bmatrix} + \begin{bmatrix} X \\ X \end{bmatrix}. \end{aligned} \quad (80)$$

To compute  $\pi_0(t, x)$ , note that

$$\begin{aligned} f_{T|X}(t|x) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(t - 1 - \rho_{T,X} \cdot x)^2}{2}\right), \\ f_T(t) &= \frac{1}{\sqrt{2\pi \cdot (1 + \rho_{T,X}^2)}} \exp\left(-\frac{(t - 1)^2}{2 \cdot (1 + \rho_{T,X}^2)}\right), \quad (T \sim N(1, 1 + \rho_{T,X}^2)). \end{aligned}$$

Hence,

$$\pi_0(t, x) = \frac{f_T(t)}{f_{T|X}(t|x)} = \frac{1}{\sqrt{1 + \rho_{T,X}^2}} \exp \left\{ \frac{\rho_{T,X}^2 \cdot (t-1)^2 + (1 + \rho_{T,X}^2) \cdot \rho_{T,X}^2 \cdot x^2 - 2 \cdot (1 + \rho_{T,X}^2) \cdot \rho_{T,X} \cdot x(t-1)}{2(1 + \rho_{T,X}^2)} \right\}. \quad (81)$$

Using (80) and (81), the true variance can be computed as

$$V_{eff} = \left[ \frac{1}{N} \sum_{i=1}^N m(T_i) m(T_i)^\top \right]^{-1} \left\{ \frac{1}{N} \sum_{i=1}^N \hat{\psi}(Y_i, T_i, X_i; \beta_0) \psi(Y_i, T_i, X_i; \beta_0)^\top \right\} \left[ \frac{1}{N} \sum_{i=1}^N m(T_i) m(T_i)^\top \right]^{-1}. \quad (82)$$

Based on a simulated sample with large enough size  $N = 10^8$ , it follows that  $V_{11} = 3.142$ ,  $V_{12} = -1.097$ , and  $V_{22} = 1.097$ .

### 6.2.2 DGP-NL1

Recall that DGP-NL1 is

$$T = \rho_{T,X} \cdot X^2 + \xi, \quad Y = X^2 + T + \epsilon. \quad (\rho_{T,X} = 0.1)$$

We have

$$\mathbb{E}[Y|T, X] = X^2 + T, \quad \mathbb{E}[Y(t)] = g(t; \beta_0) = 1 + t.$$

Eq. (78) is now rewritten as

$$\mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | T = t, X = x] = m(t) \cdot \pi_0(t, x) \cdot \{x^2 - 1\}.$$

Eq. (79) is now rewritten as

$$\mathbb{E}[\{Y - g(T; \beta_0)\} \pi_0(T, X) m(T; \beta_0) | X = x] = \{x^2 - 1\} \cdot \mathbb{E}[m(T)].$$

Substitute those equations into (77) to get

$$\psi(Y, T, X; \beta_0) = \pi_0(T, X) \cdot \{Y - X^2 - T\} \cdot \begin{bmatrix} 1 \\ T \end{bmatrix} + \begin{bmatrix} X^2 - 1 \\ \{X^2 - 1\} \cdot \rho_{T,X} \end{bmatrix}. \quad (83)$$

To compute  $\pi_0(t, x)$ , note that

$$f_{T|X}(t|x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(t - \rho_{T,X} \cdot x^2)^2}{2}\right),$$

$$\hat{f}_T(t) = \frac{1}{Nh} \sum_{i=1}^N k\left(\frac{T_i - t}{h}\right) = \frac{1}{Nh} \sum_{i=1}^N \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(T_i - t)^2}{2h^2}\right),$$

where  $k(\cdot)$  is the kernel function which can be taken as  $k(z) = (2\pi)^{-1/2} \exp(-z^2/2)$ , and the bandwidth can be taken as, say,  $h = 0.1$ . Then

$$\pi_0(t, x) = \frac{\hat{f}_T(t)}{f_{T|X}(t|x)} = \frac{1}{Nh} \sum_{i=1}^N \exp\left(-\frac{(T_i - t)^2}{2h^2} + \frac{(t - \rho_{T,X} \cdot x^2)^2}{2}\right). \quad (84)$$

Using (83) and (84), the true variance can be computed from (82). Based on a simulated sample with large enough size  $N = 50000$ , it follows that  $V_{11} = 3.043$ ,  $V_{12} = -0.118$ , and  $V_{22} = 1.074$ .

## References

- BICKEL, P. J., C. A. J. KLAASSEN, Y. RITOV, AND J. A. WELLNER (1993): *Efficient and Adaptive Estimation for Semiparametric Models*. The Johns Hopkins University Press.
- CATTANEO, M. D. (2010): “Efficient semiparametric estimation of multi-valued treatment effects under ignorability,” *Journal of Econometrics*, 155(2), 138–154.
- CHEN, X., O. LINTON, AND I. VAN KEILEGOM (2003): “Estimation of Semiparametric Models When the Criterion Function Is Not Smooth,” *Econometrica*, 71(5), 1591–1608.
- FIRPO, S. (2007): “Efficient Semiparametric Estimation of Quantile Treatment Effects,” *Econometrica*, 75(1), 259–276.
- HAHN, J. (1998): “On the role of the propensity score in efficient semiparametric estimation of average treatment effects,” *Econometrica*, 66(2), 315–331.
- ROBINS, J. M., A. ROTNITZKY, AND L. P. ZHAO (1994): “Estimation of regression coefficients when some regressors are not always observed,” *Journal of the American Statistical Association*, 89(427), 846–866.
- TCHETGEN TCHETGEN, E. J., AND I. SHPITSER (2012): “Semiparametric theory for causal mediation analysis: Efficiency bounds, multiple robustness, and sensitivity analysis,” *The Annals of Statistics*, 40(3), 1816–1845.