
Double Machine Learning Density Estimation for Local Treatment Effects with Instruments

Yonghan Jung
Purdue University
jung222@purdue.edu

Jin Tian
Iowa State University
jtian@iastate.edu

Elias Bareinboim
Columbia University
eb@cs.columbia.edu

Abstract

Randomized Controlled Trials constitute a powerful tool to learn cause and effect relationships found throughout a wide range of applied settings. In practice, the treatment assignment’s compliance is hard to ascertain in many settings since patients may not feel compelled to take the treatment for various reasons. One typical quantity investigated in these settings is the local treatment effect (LTE, for short). The LTE measures the causal effect among compliers, which usually comes under the assumption of monotonicity (only the ones offered the treatment are allowed to take it). In this paper, we investigate the challenge of estimating the LTE density function (instead of its expected value) of a binary treatment on a continuous outcome given a binary instrumental variable in the presence of both observed and unobserved confounders. Specifically, we develop two families of methods for this task, *kernel-smoothing* and *model-based* approximations – the former smooths the density by convoluting with a smooth kernel function; the latter projects the density onto a finite-dimensional density class. For both approaches, we derive double/debiased machine learning (DML) based estimators. We study the asymptotic convergence rates of the estimators and show that they are robust to the biases in nuisance function estimation. We illustrate the proposed methods on synthetic data and a real dataset called 401(k).

1 Introduction

Controlled experimentation is one powerful tool used throughout the biological, medical, and social sciences to infer the effect of a certain treatment on a given outcome. The idea is to randomize the treatment assignment so as to neutralize the effect of the unobserved confounders. In some practical settings, however, it may be challenging to ascertain that individuals who are selected for treatment will follow their recommendations. In fact, issues of non-compliance and unmeasured confounders are quite common and lead to the non-identification of treatment effects in such cases [28, 46, 31, 52].

An approach known as instrumental variables (IVs) has been proposed to try to circumvent this issue [62]. The idea is to find a variable (or set) that is not the target of the analysis by itself, but that it will help to control for the unobserved confounding between the treatment and the outcome. In particular, IVs are special variables that (i) influence the treatment, (ii) do not directly influence the outcome, and (iii) are not affected by unmeasured confounders. For concreteness, consider a study of the effect of 401(k) participation (X) on the distribution of net financial assets (Y) [2]. This setting is represented in the causal graph in Fig. 1. Note that there exists a dashed-bidirected arrow between X and Y , which in graphical language represents unobserved confounding affecting both X and Y . The variable Z in this model represents the eligibility of 401(k). We note that Z qualifies as an instrument in this case – (i) it does affect the participation of 401(k) (X), (ii) has no direct influence on the net financial asset (Y), (iii) is not affected by unmeasured confounders between X and Y . The variable W represents observed covariates (e.g., income, gender, family size, etc.).

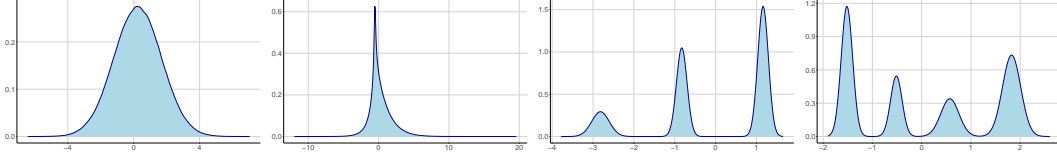


Figure 2: Densities of outcome Y among compliers under the treatment $X = 1$. All densities have a mean 0 and variance 2.

We are interested in the particular setting where only individuals who were offered the treatment may have access to it [30]. For instance, in the case of 401(k) participation ($X = 1$), only individuals who are eligible ($Z = 1$) would be allowed to join the program. This assumption is known as *monotonicity*, which rules out the possibility that any units would respond contrary to the instrument. Under monotonicity, the causal effect in the subpopulation whose actual treatment X coincides with the assigned treatment Z (called *compliers*) is identifiable [30, 2]. The average treatment effect (ATE) for the compliers is called ‘Local ATE’ (LATE) (or Complier average causal effects, CACE) [30].

The most common quantification of effects in IV settings found in practice is the average (e.g., LATE). The average is certainly an informative summary; however, it may fail to capture significant differences in the causal distributions of the outcome. For instance, consider Fig. 2 that shows the densities of outcomes Y under treatments $X = 1$ among compliers (generated from samples drawn from four synthetic data generating processes represented by the IV graph in Fig. 1, as discussed in Sec. 5). All of the distributions have the same mean 0 and variances 2. However, the difference in the distributions are self-evident.

Most of the prior work on quantifying treatment effects on outcome distributions focuses on estimating cumulative distribution functions (CDFs) or quantiles, and little attention has been given to estimating densities of treatment effects (refer to Sec. 1.1 for further comparison). As a complement to CDFs, densities have the benefits of providing more visually interpretable information of the distribution and enabling researchers/practitioners to generate counterfactual samples. One challenge with estimating densities is that while CDFs are pathwise-differentiable and enjoy \sqrt{n} -rate estimators (n is the size of data), densities are not (i.e., they are not *regular*), and therefore possess no influence functions nor \sqrt{n} -rate estimators without approximations [7, Ch. 3].

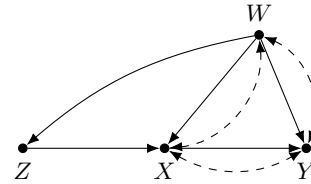


Figure 1: A causal graph for the IV setting. Bidirected arrows encode unmeasured confounders.

In this paper, our goal is to provide methods to estimate densities of local treatment effects in IV settings under the monotonicity assumption. We develop two families of methods for this task based on *kernel-smoothing* and *model-based* approximations. The former smooths the density by convolution with a smooth kernel function; the latter projects the density onto a finite-dimensional density class based on a distributional distance measure. For both approaches, we construct double/debiased machine learning (DML) style density estimators [39, 50, 48, 64, 12]. We analyze the asymptotic convergence properties of the estimators, showing that they can converge fast (i.e., \sqrt{n} -rate) even when nuisance estimates converge slowly (e.g., $n^{-1/4}$ rate) (*‘debiasedness’*¹). We illustrate the proposed methods on synthetic and real data.

1.1 Related work

Our work touches different areas, which we discuss next.

Local average treatment effect. The formal identification results for LATE under the monotonicity assumption in IV settings were developed by [30, 3]. Building on these results, semiparametric estimation for LATE has received remarkable attention [2, 55, 22, 57, 44] including robust LATE estimators that achieve debiasedness [43, 36, 34, 59]. As shown in Fig. 2, however, the average is sometimes insufficient to capture the treatment effects on distributions.

¹Also known as *‘nonparametric doubly robust* [33] or *‘rate doubly robust’* [54].

Local quantile treatment effect. A common approach to quantifying LTEs is to estimate quantiles or CDFs, which can be studied based on the LATE estimation [1, 2, 14, 23, 15, 29, 41, 17, 63]) using the fact that the expectation of $\mathbb{1}_{Y \leq y}(Y)$, an indicator that outcome Y exceeds threshold y , reduces to the quantiles (i.e., replacing Y in LATE with $\mathbb{1}_{Y \leq y}(Y)$).

Non-regular target estimand. Densities are an example of non-regular targets [7, Chap. 3]. One can approximate a non-regular target with smooth ones such that influence functions and \sqrt{n} -rate estimators can be derived. Two popular used approximation approaches are *kernel-smoothing based* (e.g., [48, 6, 38, 18, 32]) and *model-based* (e.g., [42, 48, 20, 37, 36, 35]).

Causal density estimation. There is no extensive literature on estimating the density of treatment effects. Most of the results assume that the ignorability/backdoor admissibility holds [51, 45]. [21] used the kernel-smoothing technique to estimate the density of a treatment effect, and [38] provided a kernel-smoothing based density estimator that achieves doubly robustness and debiasedness building on top of the work in [49]. Recently, [35] investigated a model-based approach and developed estimators that achieve debiasedness properties. Under the IV setting, [10] provides a local polynomial regression based density estimator for local treatment effects; we are not aware of any work studying debiased density estimators. As mentioned, this paper investigates both kernel-smoothing and model-based approaches for estimating local treatment effects under IV settings and develops DML-style density estimators for both.

2 LTE Estimation – Problem setup

Each variable is represented with a capital letter (X) and its realized value with a small letter (x). For a discrete (e.g., binary) random variable X , we use $\mathbb{1}_x(X)$ to represent the indicator function such that $\mathbb{1}_x(X) = 1$ if $X = x$; $\mathbb{1}_x(X) = 0$ otherwise. For a continuous variable X with probability density $p(x)$ and a function $f(x)$, $\mathbb{E}_P[f(X)] \equiv \int_{\mathcal{X}} f(x)p(x) d[x]$ where \mathcal{X} is the domain for X , and $\|f(X)\| \equiv \sqrt{\mathbb{E}_P[(f(X))^2]}$. \hat{f} is said to converge to f at rate r_n if $\|\hat{f}(x) - f(x)\| = O_P(1/r_n)$. For a dataset $\mathcal{D} = \{V_i\}_{i=1}^n$, we use $\mathbb{E}_{\mathcal{D}}[f(V)] \equiv (1/n) \sum_{i=1}^n f(V_i)$ to denote the empirical mean of $f(V)$ with \mathcal{D} .

Structural Causal Models (SCMs). We use the language of SCMs as our basic semantic and inferential framework [46, 4]. An SCM \mathcal{M} is a quadruple $\mathcal{M} = \langle U, V, P(U), F \rangle$ where U is a set of exogenous (latent) variables following a joint distribution $P(u)$ and V is a set of endogenous (observable) variables whose values are determined by functions $F = \{f_{V_i}\}_{V_i \in V}$ such that $V_i \leftarrow f_{V_i}(pa_i, u_i)$ where $PA_i \subseteq V$ and $U_i \subseteq U$. Each SCM \mathcal{M} induces a distribution $P(v)$ and a causal graph $G = G(\mathcal{M})$ over V in which there exists a directed edge from every variable in PA_i to V_i and dashed-bidirected arrows encode common latent variables (e.g., see Fig. 1). Within the structural semantics, performing an intervention and setting $X = x$ is represented through the do-operator, $do(X = x)$, which encodes the operation of replacing the original equations of X (i.e., $f_X(pa_x, u_x)$) by the constant x and induces a *submodel* \mathcal{M}_x and an interventional distribution $P(v|do(x))$. For any variable $Y \in V$, the *potential response* $Y_x(u)$ is defined as the solution of Y in the submodel \mathcal{M}_x given $U = u$, which induces a *counterfactual variable* Y_x .

Local Treatment Effect (LTE) with IV. We consider the IV setting represented by the causal graph G in Fig. 1², where Z is a binary instrument with domain $\{0, 1\}$, X is a binary treatment with domain $\{0, 1\}$, and Y is a (set of) continuous outcomes with bounded domain $\mathcal{Y} \subset \mathbb{R}^d$, and W is a set of covariates. G satisfies the IV assumption that Z has no direct influence on outcome Y and is not affected by unmeasured confounders between X and Y .

The causal density $p(y|do(x))$ is not identifiable from the observed distribution $p(x, y, z, w)$ due to unobserved confounders between X and Y . However, the effect may be recovered for certain subpopulation under additional assumptions. Formally, a unit in the population is an *always-taker* if $X_{Z=1} = X_{Z=0} = 1$, a *never-taker* if $X_{Z=1} = X_{Z=0} = 0$, a *complier* if $X_{Z=1} = 1, X_{Z=0} = 0$, and a *defier* if $X_{Z=1} = 0, X_{Z=0} = 1$ [3, 2]. We will make the following assumptions after the literature.

Assumption 1 (Monotonicity). *There are no defiers:* $X_{Z=1} \geq X_{Z=0}$.

Assumption 2 (Positivity). $P(x|z, w) > 0, P(z|w) > 0$ for any x, z, w .

²It is common in the literature to define IV assumptions in terms of conditional independences among counterfactuals [47, 9, 8, 2, 55, 43, 59], which connection with the causal graph in Fig. 1 is discussed in Assumption A.1

Let C denote the event that a unit is a complier. For a given constant a and a variable X , let x^a denote the event $X = a$. The LTE $p(y_x|C)$ is identifiable under monotonicity and is given by [30, 2]:

$$p(y_x|C) = \frac{\mathbb{E}_P [p(y|x, z^x, W)P(x|z^x, W) - p(y|x, z^{1-x}, W)P(x|z^{1-x}, W)]}{\mathbb{E}_P [P(x^1|z^1, W) - P(x^1|z^0, W)]}, \quad (1)$$

where the expectation is over W . In this paper, we aim to estimate the LTE density $p(y_x|C)$ in Eq. (1). We will make the following mild assumption on the target, popularly employed in density estimation (e.g., [40, 24, 26, 56, 25, 38]).

Assumption 3. $p(y|w, z, x)$ and $p(y_x|C)$ are bounded and twice differentiable for any x, z, w, y .

The proofs are provided in Appendix B in suppl. material.

3 Kernel-smoothing based approach

In this section, we develop a kernel-smoothing based approach for estimating the LTE density. The kernel-smoothing technique approximates a non-pathwise-differentiable target estimand with a differentiable estimand by convoluting the density with a kernel function $K(y)$. Properties of the kernel function includes symmetry about the origin (i.e., $\int_{\mathcal{Y}} yK(y) d[y] = 0$), non-negativity ($0 < K(y) < \infty, \forall y \in \mathcal{Y}$), and integrates to 1 (i.e., $\int_{\mathcal{Y}} K(y) d[y] = 1$) [60, Chap. 4.2].

We will consider a *product kernel* $K_{h,y}(y') \equiv h^{-d} \prod_{j=1}^d K((y_j - y'_j)/h)$ with given bandwidth $h \in \mathbb{R}$ and a fixed point $y = \{y_j\}_{j=1}^d \in \mathbb{R}^d$. We assume that the kernel of interest has a bounded second moment and norm: i.e., $\kappa_2(K) \equiv \int_{\mathcal{Y}} y^2 K(y) d[y] < \infty$ and $\|K(y)\| < \infty$ following [26, 56]. Example of kernels include Gaussian kernel: $K(u) = (1/\sqrt{2\pi}) \exp(-u^2/2)$, Epanechnikov kernel: $K(u) = (3/4)(1 - u^2) \mathbb{1}_{|u| \leq 1}(u)$, Quadratic kernel: $K(u) = (15/16)(1 - u^2)^2 \mathbb{1}_{|u| \leq 1}(u)$, Cosine kernel: $k(u) = (\pi/4) \cos(\pi u/2) \mathbb{1}_{|u| \leq 1}(u)$, etc.

For convenience, we denote the target estimand by $\psi(y) \equiv p(y_x|C)$. We will instead aim to estimate a kernel-smoothed approximation for $\psi(y)$ defined as follows:

$$\psi_h(y) \equiv \int_{\mathcal{Y}} \psi(y') K_{h,y}(y') d[y'] = \psi[K_{h,y}(Y)], \quad (2)$$

where $\psi[f(Y)]$ for any function $f(Y)$ is defined as

$$\psi[f(Y)] \equiv \frac{\mathbb{E}_P [\mathbb{E}_P [f(Y) \mathbb{1}_x(X) | z^x, W] - \mathbb{E}_P [f(Y) \mathbb{1}_x(X) | z^{1-x}, W]]}{\mathbb{E}_P [P(x^1|z^1, W) - P(x^1|z^0, W)]}. \quad (3)$$

The second equality in Eq. (2) is by Eq. (1). For a target estimand $\psi[f(Y)]$, we will denote nuisances by $\pi_z(w) \equiv P(z|w)$, $\xi_x(z, w) \equiv P(x|z, w)$, and $\theta(x, z, w)[f(Y)] \equiv \mathbb{E}_P [f(Y) \mathbb{1}_x(X) | x, w]$, shortly (π, ξ, θ) .

We aim to construct a DML estimator for the estimand ψ_h . Toward this goal, we will first derive a Neyman orthogonal score for ψ_h . Since a Neyman orthogonal score can be constructed based on *moment score functions* (a function of parameters such that its expectation is 0 at the true parameters) [13, Thm. 1], we start by defining the moment condition. Let

$$\psi^X \equiv \mathbb{E}_P [\xi_{x^1}(z^1, W) - \xi_{x^1}(z^0, W)], \quad (4)$$

$$\mathcal{V}_X(\{\pi, \xi\}) \equiv \frac{\mathbb{1}_{z^1}(Z) - \mathbb{1}_{z^0}(Z)}{\pi_Z(W)} \{ \mathbb{1}_{x^1}(X) - \xi_{x^1}(Z, W) \} + \{ \xi_{x^1}(z^1, W) - \xi_{x^1}(z^0, W) \}. \quad (5)$$

Then, the following is a moment score function for ψ_h :

$$m(\psi'; \psi_h) \equiv \frac{1}{\psi^X} (\psi_h - \psi') \mathcal{V}_X, \quad (6)$$

where ψ_h is given in Eq. (2) and ψ' is an estimate of ψ_h .

Next, we derive an influence function for $m(\psi'; \psi_h)$. Toward this, we first define the following function: for any function $f(Y) < \infty$,

$$\phi(\eta = \{\pi, \xi, \theta\}, \psi)[f(Y)] \equiv \frac{1}{\psi^X} (\mathcal{V}_X(\{\pi, \theta\})[f(Y)] - \psi[f(Y)] \mathcal{V}_X(\{\pi, \xi\})), \quad (7)$$

where \mathcal{V}_X is defined in Eq. (5), and

$$\begin{aligned} \mathcal{V}_{YX}(\{\pi, \theta\})[f(Y)] &\equiv \frac{\mathbb{1}_{z^x}(Z) - \mathbb{1}_{z^{1-x}}(Z)}{\pi_Z(W)} \{f(Y)\mathbb{1}_x(X) - \theta(x, Z, W)[f(Y)]\} \\ &\quad + \{\theta(x, z^x, W)[f(Y)] - \theta(x, z^{1-x}, W)[f(Y)]\}. \end{aligned} \quad (8)$$

Then, an influence function for the expectation of the moment score function $m(\psi'; \psi_h)$ in Eq. (6) is given as follows:

Lemma 1 (Influence function for the score $m(\psi'; \psi_h)$). *Let $m(\psi'; \psi_h)$ be the score defined in Eq. (6). Then, an influence function for $\mathbb{E}_P[m(\psi'; \psi_h)]$, denoted ϕ_m , is given by*

$$\phi_m(\eta = \{\pi, \xi, \theta\}, \psi) \equiv \phi(\eta, \psi)[K_{h,y}(Y)] \quad (9)$$

where ϕ is in Eq. (7).

For any score function (e.g., m in Eq. (6)), its addition to the influence function of the expected score (e.g., ϕ_m) is a Neyman orthogonal score³ ([13, Thm.1], [12, Sec. 2.2.5]). Specifically,

Lemma 2 (Neyman orthogonal score for ψ_h). *Let $m(\psi'; \psi_h)$ be the score function in Eq. (6), and $\phi_m(\eta = \{\pi, \xi, \theta\}, \psi_h)$ be the influence function for $\mathbb{E}_P[m(\psi'; \psi_h)]$ given in Eq. (9). Then, a Neyman orthogonal score for ψ_h is given as $\varphi(\psi'; \eta = \{\pi, \xi, \theta\}) \equiv m(\psi'; \psi_h) + \phi_m(\eta, \psi)$; Specifically,*

$$\varphi(\psi'; \eta = \{\pi, \xi, \theta\}) = \frac{1}{\psi_X} (\mathcal{V}_{YX}(\{\pi, \theta\})[K_{h,y}(Y)] - \psi' \mathcal{V}_X(\{\pi, \xi\})). \quad (10)$$

Given the Neyman orthogonal score $\varphi(\psi'; \eta)$, an estimate $\hat{\psi}_h$ satisfying $\mathbb{E}_{\mathcal{D}}[\varphi(\hat{\psi}_h; \hat{\eta} = \{\hat{\pi}, \hat{\xi}, \hat{\theta}\})] = o_P(n^{-1/2})$ gives a DML estimator. Specifically, we propose the following kernel-smoothing based estimator for the LTE density, named ‘KLTE’:

Definition 1 (KLTE estimator for ψ_h). *Let $\varphi(\psi'; \eta = \{\pi, \xi, \theta\})$ be the Neyman orthogonal score for ψ_h given in Eq. (10). Let $\{\mathcal{D}, \mathcal{D}'\}$ denote the randomly split halves of the samples, where $|\mathcal{D}| = |\mathcal{D}'| = n$. Let $\hat{\eta} = \{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ denote the estimates for the nuisance η using \mathcal{D}' . Then, the KLTE estimator for $\psi_h(y)$ for all $y \in \mathcal{Y}$, denoted $\hat{\psi}_h(y)$, is given by*

$$\hat{\psi}_h(y) \equiv \mathbb{E}_{\mathcal{D}}[\mathcal{V}_{YX}(\{\hat{\pi}, \hat{\theta}\})[K_{h,y}(Y)]] / \mathbb{E}_{\mathcal{D}}[\mathcal{V}_X(\{\hat{\pi}, \hat{\xi}\})], \quad (11)$$

where \mathcal{V}_X and \mathcal{V}_{YX} are given in Eq. (5,8), respectively.

We will show that KLTE is a DML estimator exhibiting debiasedness property. Detailed asymptotic properties are discussed next.

3.1 Asymptotic convergence

Now, we study the convergence rate of the estimator $\hat{\psi}_h(y)$. For any fixed $y \in \mathcal{Y}$, the error $\hat{\psi}_h(y) - \psi(y)$ will be analyzed in two folds: we will first analyze the error between the estimator in Eq. (11) and the smoothed estimand in Eq. (2) (i.e., $\hat{\psi}_h(y) - \psi_h(y)$), and then analyze the error between the smoothed estimand and the true estimand (i.e., $\psi_h(y) - \psi(y)$).

The following result gives the error analysis for $\hat{\psi}_h(y) - \psi_h(y)$:

Lemma 3 (Convergence rate of $\hat{\psi}_h$ to ψ_h). *For any fixed $y \in \mathcal{Y}$, suppose the estimators for nuisances are consistent; i.e., $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \eta = \{\pi, \xi, \theta\}$ for all (w, z, x) . Suppose $h < \infty$, and $nh^d \rightarrow \infty$ as $n \rightarrow \infty$. Then,*

$$\hat{\psi}_h(y) - \psi_h(y) = O_P\left(1/\sqrt{nh^d} + R_2^k + 1/\sqrt{n}\right),$$

where

$$R_2^k \equiv \sum_z \|\hat{\pi}_z - \pi_z\| \left\{ \|\hat{\theta}_z - \theta_z\| + \|\hat{\xi}_z - \xi_z\| \right\}, \quad (12)$$

where $\pi_z \equiv \pi_z(W)$, $\xi_z \equiv \xi_x(z, W)$ and $\theta_z \equiv \theta(x, z, W)[K_{h,y}(Y)]$.

³A Neyman orthogonal score is a function ϕ satisfying $\mathbb{E}_P[\phi(\psi, \eta_0)] = 0$ and $\frac{\partial}{\partial \eta} \mathbb{E}_P[\phi(V; \psi, \eta)]|_{\eta=\eta_0} = 0$, where η_0 denotes the true nuisance [12, Def.2.2]. In words, a moment condition that is not sensitive to local errors in nuisance models.

The error analysis in Lemma. 3 implies the following:

Corollary 1 (Debiasedness property of $\hat{\psi}_h$ to ψ_h). *If all nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ for any given (w, z, x, y) converge at rate $\{nh^d\}^{-1/4}$, then the target estimator $\hat{\psi}_h(y)$ achieves $\sqrt{nh^d}$ -rate convergence to ψ_h .*

We now analyze the gap between the smoothed estimand ψ_h and the true estimand ψ ; i.e., $\psi_h - \psi$:

Lemma 4 ([60, Thm. 6.28]). *The following holds:*

$$\psi_h(y) - \psi(y) = B_y \equiv 0.5h^2\kappa_2(K)(\partial^2/\partial^2y')|_{y'=y}\psi(y') + O(h^2). \quad (13)$$

Combining the results of Lemma. (3,4), we have the following result:

Theorem 1 (Convergence rate of $\hat{\psi}_h$ to ψ). *For any fixed $y \in \mathcal{Y}$, suppose the estimators for nuisances are consistent; i.e., $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \eta = \{\pi, \xi, \theta\}$ for all (w, z, x) . Suppose $h < \infty$, and $nh^d \rightarrow \infty$ as $n \rightarrow \infty$. Then*

$$\hat{\psi}_h(y) - \psi(y) = O_P\left(1/\sqrt{nh^d} + R_2^k + 1/\sqrt{n}\right) + B_y, \quad (14)$$

where B_y is defined in Eq. (13), and R_2^k is defined in Eq. (12).

Thm. 1 implies that $\hat{\psi}_h(y)$ converges fast (see Corol. 1) to $\psi(y) + B_y$. A natural question is then how to choose the bandwidth h that minimizes the gap in Eq. (14). The following provides a guideline in choosing the bandwidth h :

Lemma 5 (Data-adaptive bandwidth selection). *The bandwidth h that minimizes the error in Eq. (14) is $h = O(n^{-1/(d+4)})$. This choice of h satisfies the assumption in Lemma. 3 that $nh^d \rightarrow \infty$.*

Recall that Corol. 1 states the debiasedness property of $\hat{\psi}_h$ to ψ_h for any bandwidth h satisfying $nh^d \rightarrow \infty$. With the choice of h as in Lemma 5, $\hat{\psi}_h$ converges to ψ with the debiasedness property preserved.

Corollary 2 (Debiasedness property of $\hat{\psi}_h$ to ψ). *Let $h = O(n^{-1/(d+4)})$. If nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ converge at $\{nh^d\}^{-1/4}$ rate for any (w, z, x, y) , then the target estimator $\hat{\psi}_h(y)$ achieves $\sqrt{nh^d}$ -rate convergence to ψ .*

So far, we have analyzed the error $\hat{\psi}_h(y) - \psi(y)$ pointwise for the fixed $y \in \mathcal{Y}$. To analyze the ‘gap’ between the two densities $\hat{\psi}_h(y)$ and $\psi(y)$ for all $y \in \mathcal{Y}$, we consider the following divergence function of two densities:

Definition 2 (f -Divergence D_f [19]). Let f denote a convex function with $f(1) = 0$. $D_f(p, q) \equiv \int_{\mathcal{Y}} f(p(y), q(y))q(y) d[y]$, is a f -divergence function between two densities p, q .

f -divergence covers many well-known divergences. For example, D_f reduces to KL divergence with $f(p, q) = (p/q) \log(p/q)$. We will assume that the function $f(p, q)$ in D_f is differentiable w.r.t. p and q .

We now analyze the distance between $\hat{\psi}_h$ and ψ w.r.t. D_f . The following result provides an upper bound for D_f .

Lemma 6 (Upper bound of the divergence D_f). *Suppose D_f is a f -divergence such that $f(p, q) = 0$ if $p = q$. Then,*

$$D_f(\psi, \hat{\psi}_h) \leq \int_{\mathcal{Y}} w(y) \left(\hat{\psi}_h(y) - \psi(y) \right) d[y],$$

where $w(y) \equiv f_2'(\psi(y), \tilde{\psi}(y))\hat{\psi}_h(y)$, $f_2'(p, q) \equiv (\partial/\partial q)f(p, q)$, and $\tilde{\psi}_h(y) \equiv t\hat{\psi}_h(y) + (1-t)\psi(y)$ for some fixed $t \in [0, 1]$.

By invoking Thm. 1, we derive an upper bound for $D_f(\psi, \hat{\psi}_h)$ as follows:

Theorem 2 (Convergence rate of $\hat{\psi}_h$). *Suppose the estimators for nuisances are consistent; i.e., $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \eta = \{\pi, \xi, \theta\}$ for all (w, z, x, y) . Suppose D_f is a f -divergence such that $f(p, q) = 0$ if $p = q$. Suppose $w(y)$ in Lemma 6 is finite. Then,*

$$D_f(\psi, \hat{\psi}_h) \leq O_P\left(\sup_{y \in \mathcal{Y}} \{R_2^k + B_y\} + 1/\sqrt{nh^d} + 1/\sqrt{n}\right), \quad (15)$$

where R_2^k is defined in Eq. (12) and B_y is defined in Eq. (13).

The following result asserts that the debiasedness property is exhibited w.r.t. D_f :

Corollary 3 (Debiasedness property of $\hat{\psi}_h$ w.r.t. D_f). *Let $h = O(n^{-1/(d+4)})$. Suppose D_f satisfies $f(p, q) = 0$ if $p = q$. Suppose $w(y)$ in Lemma 6 is finite. If nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ converges at $\{nh^d\}^{-1/4}$ rate for any (w, z, x, y) , then $D_f(\psi, \hat{\psi}_h)$ converges to 0 at $\sqrt{nh^d}$ -rate.*

4 Model-based approach

In this section, we develop a *model-based approach* for estimating the LTE density $\psi(y) = p(y_x|C)$. We will approximate ψ with a class of distributions or a *density model* $\mathcal{G} = \{g(y; \beta) : \beta \in \mathbb{R}^b\}$ where $g(y; \beta) \in \mathcal{G}$ is differentiable w.r.t. β . Example density models include exponential family (e.g., Gaussian distribution), mixture of Gaussians, or more generally, mixture of exponential families. We adapt the model-based approach developed in [35] for estimating the causal density under the no unmeasured confounders assumption.

Given a density model \mathcal{G} , the best approximation for $\psi(y)$ is defined as $g(y; \beta_0) \in \mathcal{G}$ that achieves the minimum f -divergence to ψ :

$$\beta_0 \equiv \arg \min_{\beta \in \mathbb{R}^b} D_f(\psi(y), g(y; \beta)), \quad (16)$$

where D_f is the f -divergence defined in Def. 2. Our goal is estimating β_0 .

Consider $m(\beta; \psi) \equiv (\partial/\partial\beta)D_f(\psi(y), g(y; \beta))$. Definition of β_0 given in Eq. (16) implies that $m(\beta; \psi) = 0$ at $\beta = \beta_0$. We note that $m(\beta; \psi)$ serves as a *moment score function*. The closed-form expression of the score is given by [35]:

$$m(\beta; \psi) \equiv \int_{\mathcal{Y}} g'(y; \beta) \{f'_2(\psi(y), g(y; \beta))g(y; \beta) + f(\psi(y), g(y; \beta))\} d[y], \quad (17)$$

where $g'(y; \beta) = (\partial/\partial\beta)g(y; \beta)$ and $f'_2(p, q) \equiv (\partial/\partial q)f(p, q)$.

To construct a DML estimator based on the score function $m(\beta; \psi)$, we first derive an influence function for the score:

Lemma 7 (Influence Function for $m(\beta, \psi)$). *An influence function for $m(\beta; \psi)$ in Eq. (17), denoted ϕ_m , is given by*

$$\phi_m(\beta; \eta = \{\pi, \xi, \theta\}, \psi) \equiv \phi(\eta, \psi)[R_f(Y; \beta, \psi)], \quad (18)$$

where $\phi(\eta, \psi)[\cdot]$ is defined in Eq. (7), and

$$R_f(Y; \beta, \psi) \equiv g'(Y; \beta) \{f''_{21}(\psi(Y), g(Y; \beta))g(Y; \beta) + f'_1(\psi(Y), g(Y; \beta))\},$$

where $g'(y; \beta) \equiv (\partial/\partial\beta)g(y; \beta)$, $f'_1(p, q) \equiv (\partial/\partial p)f(p, q)$ and $f''_{21}(p, q) \equiv (\partial/\partial p)f'_2(p, q)$.

We derive a Neyman orthogonal score based on the moment score $m(\beta, \psi)$ and its influence function $\phi_m(\beta, \eta, \psi)$:

Lemma 8 (Neyman orthogonal score for β). *A Neyman orthogonal score for estimating β , denoted $\varphi(\beta'; (\eta = \{\pi, \xi, \theta\}, \psi))$, is given by*

$$\varphi(\beta'; (\eta = \{\pi, \xi, \theta\}, \psi)) \equiv m(\beta', \psi) + \phi_m(\beta, \eta, \psi), \quad (19)$$

where $\phi_m(\beta, \eta, \psi)$ is defined in Eq. (18).

Given the orthogonal score $\varphi(\beta'; (\eta, \psi))$ in Eq. (19), we propose the following estimator for β , named ‘MLTE’ (model-based estimator for LTE):

Definition 3 (MLTE estimator for β). *Let $\varphi(\beta'; \eta = \{\pi, \xi, \theta\}, \psi)$ be the Neyman orthogonal score for β given in Eq. (19). Let $\{\mathcal{D}, \mathcal{D}'\}$ denote the randomly split halves of the samples, where $|\mathcal{D}| = |\mathcal{D}'| = n$. Let $\hat{\eta} = \{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ denote the estimators for the nuisance η using \mathcal{D}' . Then, the MLTE estimator for β , denoted $\hat{\beta}$, is given as a solution satisfying $\mathbb{E}_{\mathcal{D}} [\varphi(\hat{\beta}; \hat{\eta}, \hat{\psi})] = o_P(n^{-1/2})$.*

To illustrate, we exemplify Eq. (17) and Lemma (7, 8) for the case where D_f is a KL-divergence and $g(y; \beta = \{\mu, \sigma^2\})$ is a normal distribution. First, $m(\beta; \psi) = \{m_\mu(\mu; \psi), m_\sigma(\sigma^2; \psi, \mu)\}$, where

$m_\mu(\mu; \psi, \sigma) = (1/\sigma^2)(\psi[Y] - \mu)$ and $m_\sigma(\sigma^2; \psi, \mu) = (0.5/\sigma^4)(\sigma^2 - \psi[(Y - \mu)^2])$. We note that $\hat{\mu}_m \equiv \hat{\psi}[Y]$ and $\hat{\sigma}_m^2 \equiv \hat{\psi}[(Y - \hat{\mu})^2]$ are estimators for $\beta_0 = \{\mu_0, \sigma_0^2\}$ for the score $m(\beta; \psi)$.

Also, $R_f(Y; \beta, \psi) \equiv -(\partial/\partial\beta) \log(g(Y; \beta)) = \{R_f(Y; \mu, \psi), R_f(Y; \sigma^2, \psi)\}$, where $R_f(Y; \mu, \psi) \equiv (\mu - Y)/\sigma^2$ and $R_f(Y; \sigma^2, \psi) \equiv 0.5\{\sigma^2 - (Y - \mu)^2\}/\sigma^4$. Then, the Neyman orthogonal score is given as $\varphi(\mu; \sigma^2, \eta, \psi) = (1/\sigma^2)\{\mu - \psi[Y] - \phi(\eta, \psi)[Y]\}$ and $\varphi(\sigma^2; \mu, \eta, \psi) = (0.5/\sigma^4)\{\sigma^2 - \psi[(Y - \mu)^2] - \phi(\eta, \psi)[(Y - \mu)^2]\}$. Finally, solutions for $\varphi(\mu; \sigma^2, \eta, \psi)$ and $\varphi(\sigma^2; \mu, \eta, \psi)$ are given by $(\hat{\mu}, \hat{\sigma}^2)$, where, for $\phi[\cdot]$ in Eq. (7), $\hat{\mu} = \hat{\psi}[Y] + \mathbb{E}_{\mathcal{D}}[\phi(\hat{\eta}, \hat{\psi})[Y]]$ and $\hat{\sigma}^2 = \psi[(Y - \hat{\mu})^2] + \mathbb{E}_{\mathcal{D}}[\phi(\hat{\eta}, \hat{\psi})[(Y - \hat{\mu})^2]]$.

The MLTE estimator in Def. 3 is consistent provided that nuisances estimates $\hat{\eta}$ are consistent [13, Thm.4]. Such $\hat{\beta}$ is known to achieve debiasedness [12], since $\hat{\beta}$ is a DML estimator. Specifically,

Theorem 3 (Convergence rate of $\hat{\beta}$). *Let $\varphi(\beta'; (\eta = \{\pi, \xi, \theta\}, \psi))$ be given in Eq. (19). Let $\phi_m(\beta, \eta, \psi)$ be given in Eq. (18). Let β_0, η_0, ψ_0 denote the true parameters. Let $\hat{\beta}$ be the MLTE estimator for β defined in Def. 3. Suppose (1) $R_f(y; \beta, \psi)$ is bounded and $R'_f(y; \beta, \psi) \equiv (\partial/\partial\psi)R_f(y; \beta, \psi) < \infty$; (2) There exists a function $H(y) < \infty$ s.t. $\sup_{\beta, \psi} \max\{R_f(y; \beta, \psi), R'_f(y; \beta, \psi)\} = O(H(y))$; (3) $\{\varphi(\beta; (\eta, \psi))\}$ is Donsker⁴ w.r.t. β for the fixed η ; (3) The estimators are consistent: $\hat{\beta} - \beta_0 = o_P(1)$ and $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \{\pi_z(w), \xi_x(z, w), \theta(x, z, w)[H(Y)]\}$ for all (w, z, x, y) ; and (4) $\mathbb{E}_P[\varphi(\beta; (\eta, \psi))]$ is differentiable w.r.t. β at $\beta = \beta_0$ with non-singular matrix $M(\beta_0, (\eta, \psi)) \equiv (\partial/\partial\beta)|_{\beta=\beta_0} \mathbb{E}_P[\varphi(\beta; (\eta, \psi))]$ for all (η, ψ) , where $M(\beta_0, (\hat{\eta}, \hat{\psi})) \xrightarrow{P} M \equiv M(\beta_0, (\eta_0, \psi_0))$. Then,*

$$\hat{\beta} - \beta_0 = -M^{-1} \mathbb{E}_{\mathcal{D}}[\phi_m(\beta_0; (\psi_0, \eta_0))] + o_P(n^{-1/2}) + O_P(R_2^m),$$

where

$$R_2^m = \sum_z \left(\|\hat{\pi}_z - \pi_z\| \left\{ \|\hat{\theta}_z - \theta_z\| + \|\hat{\xi}_z - \xi_z\| \right\} + \|\hat{\xi}_z - \xi_z\|^2 + \|\theta_z - \hat{\theta}_z\|^2 + \|\hat{\xi}_z - \xi_z\| \|\theta_z - \hat{\theta}_z\| \right),$$

where $\pi_z \equiv \pi_z(W)$, $\xi_z \equiv \xi_x(z, W)$, and $\theta_z \equiv \theta(x, z, W)[H(Y)]$.

Corollary 4 (Debiasedness property for $\hat{\beta}$). *If nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ converges at $n^{-1/4}$ rate, then the target estimator $\hat{\beta}$ converges to β_0 at \sqrt{n} -rate.*

For the above example where D_f is the KL divergence and $g(y; \beta)$ is a normal distribution, $H(Y) = Y$ for $R_f(y; \mu, \psi)$, and $H(Y) = Y^2$ for $R_f(y; \sigma^2, \psi)$.

5 Empirical applications

In this section, we apply the proposed methods to synthetic and real datasets. For the kernel-smoothing based approach, we compare KLTE with a baseline plug-in estimator ('kernel-smoothing'), where estimates of nuisances $\hat{\eta} = \{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ are plugged in the estimand Eq. (2). We use the Gaussian kernel. The bandwidth is set to $h = 0.5n^{-1/5}$. In estimating the density, we choose 200 equi-spaced points $\{y_{(i)}\}_{i=1}^{200}$ in \mathcal{Y} and evaluate both estimators at $K_{h, y_{(i)}}$ for $i = 1, \dots, 200$. For the model-based approach, we compare MLTE (e.g., $\hat{\mu}, \hat{\sigma}^2$) with a moment-score-based estimator (called 'moment'), defined as $\hat{\beta}_m$ satisfying $m(\hat{\beta}_m; \hat{\psi}) = o_P(n^{-1/2})$ (e.g., $\{\hat{\mu}_m, \hat{\sigma}_m^2\}$). We use KL divergence for D_f and the normal distribution for $g(y; \beta)$. For both approaches, nuisances are estimated through a gradient boosting model XGBoost [11], which is known to be flexible.

5.1 Synthetic datasets

We applied the proposed estimators to estimate the LTE $p(y_x|C)$ where the true densities are given as in the 4th plot in Fig. 2. As shown in the ground-truth in Fig. 3a, true densities $p(y_{x_0}|C), p(y_{x_1}|C)$ are given as a mixture of four Gaussians. Estimated densities for Moment and MLTE are given in

⁴A function class where complexities are restricted. Refer [58, Page 269] for the definition. Donsker class include Sobolev, Bounded monotone, Lipschitz class, etc.

Fig. (3b, 3c). We note that model-based approaches fail to capture important characteristics (such as the number of modes) of the true density (‘ground-truth’ in Fig. 3a) because the assumed density class is misspecified. The ‘kernel-smoothing’ (Fig. 3d) captures only one of the modes having the highest densities, and this leads to misinterpretation of the true densities. KLTE (Fig. 3e) is able to capture the number, location, and scales of modes correctly.

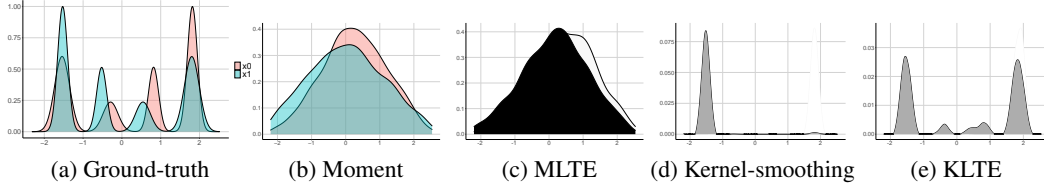


Figure 3: LTE estimation with a synthetic dataset. The ground-truth density is in (a). Red and Green for x^0 and x^1 , respectively.

5.2 Application to 401(k) data

We applied the proposed estimators (KLTE and MLTE) on 401(k) data, where the data generating processes corroborate with Fig. 1. Monotonicity assumption holds naturally, since ineligible units ($Z = 0$) cannot participate ($X = 1$) in 401(k). In our analysis, we used the dataset introduced by [2] containing 9275 individuals, which has been studied in [2, 16, 5, 43, 53, 59], to cite a few. Model-based approaches (Moment in Fig. 4a and MLTE in Fig. 4b) and kernel-smoothing based approaches (kernel-smoothing in Fig. 4c and KLTE in Fig. 4d) are implemented to analyze the data.

The model-based (Fig. (4a,4b)) and kernel-smoothing based (Fig. (4c,4d)) estimates both capture important characteristics of the distribution, such as mode, location, and scale parameters. The results of proposed estimators (MLTE and KLTE in Fig. (4b,4d)) are consistent with findings from previous analyses [2, 16, 5, 53]: The effects of the 401(k) participation (i.e., $X = 1$) on net financial assets are positive over the whole range of asset distributions. To connect to CDF method, we provide in Fig. 4e the CDF estimate induced by KLTE density estimation (Fig. 4a). We note that the CDF in Fig. 4e captures the nonconstant impact trend of the 401(k) participation on the net financial assets, which has been also described in the previous analyses [2, 16, 5, 53].

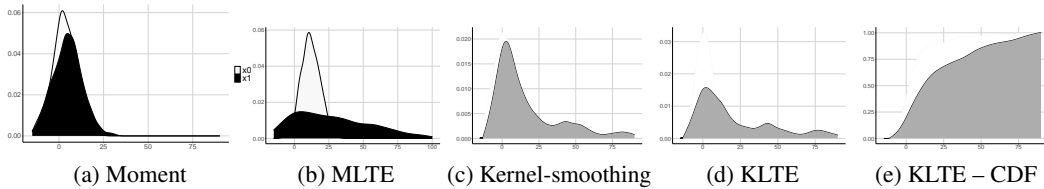


Figure 4: LTE of 401(k) participation (X) on net financial asset (Y). Red and Green for x^0 and x^1 , respectively.

6 Conclusion

In this paper, we develop *kernel-smoothing-based* and *model-based* approaches for estimating the LTE density in the presence of instruments. For each approach, we give Neyman orthogonal scores (Lemma (2,8)) and constructed corresponding DML estimators (KLTE in Def. 1 and MLTE in Def. 3), that exhibit debiasedness property (Corol. (3, 4)). We demonstrated our work through synthetic and real datasets. The performance of model-based estimators depends critically on the choice of the density class. Kernel-based estimators do not have to make assumptions about the true density class but will suffer from the curse of dimensionality.

References

- [1] A. Abadie. Bootstrap tests for distributional treatment effects in instrumental variable models. *Journal of the American statistical Association*, 97(457):284–292, 2002.

- [2] A. Abadie. Semiparametric instrumental variable estimation of treatment response models. *Journal of econometrics*, 113(2):231–263, 2003.
- [3] J. D. Angrist, G. W. Imbens, and D. B. Rubin. Identification of causal effects using instrumental variables. *Journal of the American statistical Association*, 91(434):444–455, 1996.
- [4] E. Bareinboim, J. D. Correa, D. Ibeling, and T. Icard. On pearl’s hierarchy and the foundations of causal inference. Technical Report R-60. Causal Artificial Intelligence Laboratory, Columbia University, 2020.
- [5] A. Belloni, V. Chernozhukov, and C. Hansen. Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650, 2014.
- [6] A. F. Bibaut and M. J. van der Laan. Data-adaptive smoothing for optimal-rate estimation of possibly non-regular parameters. *arXiv preprint arXiv:1706.07408*, 2017.
- [7] P. J. Bickel, C. A. Klaassen, P. J. Bickel, Y. Ritov, J. Klaassen, J. A. Wellner, and Y. Ritov. *Efficient and adaptive estimation for semiparametric models*, volume 4. Johns Hopkins University Press Baltimore, 1993.
- [8] C. Brito. *Instrumental sets*. In R. Dechter, H. Geffner, and J. Y. Halpern, editors, *Heuristics, Probability and Causality. A Tribute to Judea Pearl*. College Publications, 2010.
- [9] C. Brito and J. Pearl. Generalized instrumental variables. In *Proceedings of the Eighteenth conference on Uncertainty in artificial intelligence*, pages 85–93, 2002.
- [10] M. D. Cattaneo, M. Jansson, and X. Ma. Local regression distribution estimators. *Journal of Econometrics*, 2021.
- [11] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794, 2016.
- [12] V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins. Double/debiased machine learning for treatment and structural parameters: Double/debiased machine learning. *The Econometrics Journal*, 21(1), 2018.
- [13] V. Chernozhukov, J. C. Escanciano, H. Ichimura, W. K. Newey, and J. M. Robins. Locally robust semiparametric estimation. *arXiv preprint arXiv:1608.00033*, 2016.
- [14] V. Chernozhukov, I. Fernández-Val, and A. Galichon. Quantile and probability curves without crossing. *Econometrica*, 78(3):1093–1125, 2010.
- [15] V. Chernozhukov, I. Fernández-Val, and B. Melly. Inference on counterfactual distributions. *Econometrica*, 81(6):2205–2268, 2013.
- [16] V. Chernozhukov and C. Hansen. The effects of 401 (k) participation on the wealth distribution: an instrumental quantile regression analysis. *Review of Economics and statistics*, 86(3):735–751, 2004.
- [17] V. Chernozhukov, C. Hansen, and K. Wuthrich. Instrumental variable quantile regression. *arXiv preprint arXiv:2009.00436*, 2020.
- [18] K. Colangelo and Y.-Y. Lee. Double debiased machine learning nonparametric inference with continuous treatments. *arXiv preprint arXiv:2004.03036*, 2020.
- [19] I. Csiszár. Information-type measures of difference of probability distributions and indirect observation. *studia scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- [20] I. Díaz and M. J. van der Laan. Targeted data adaptive estimation of the causal dose–response curve. *Journal of Causal Inference*, 1(2):171–192, 2013.
- [21] J. DiNardo, N. M. Fortin, and T. Lemieux. Labor market institutions and the distribution of wages, 1973-1992: A semiparametric approach. *Econometrica: Journal of the Econometric Society*, pages 1001–1044, 1996.

- [22] M. Frölich. Nonparametric iv estimation of local average treatment effects with covariates. *Journal of Econometrics*, 139(1):35–75, 2007.
- [23] M. Frölich and B. Melly. Unconditional quantile treatment effects under endogeneity. *Journal of Business & Economic Statistics*, 31(3):346–357, 2013.
- [24] S. Ghosal et al. Convergence rates for density estimation with bernstein polynomials. *The Annals of Statistics*, 29(5):1264–1280, 2001.
- [25] E. Giné, R. Nickl, et al. Confidence bands in density estimation. *The Annals of Statistics*, 38(2):1122–1170, 2010.
- [26] L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media, 2006.
- [27] J. Hahn. On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica*, pages 315–331, 1998.
- [28] J. J. Heckman. Randomization as an instrumental variable, 1995.
- [29] Y.-C. Hsu, T.-C. Lai, and R. P. Lieli. Estimation and inference for distribution functions and quantile functions in endogenous treatment effect models. Technical report, Institute of Economics, Academia Sinica, Taipei, Taiwan, 2015.
- [30] G. W. Imbens and J. D. Angrist. Identification and estimation of local average treatment effects. *Econometrica*, 62(2):467–475, 1994.
- [31] G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- [32] N. Kallus and M. Uehara. Doubly robust off-policy value and gradient estimation for deterministic policies. *Advances in Neural Information Processing Systems*, 33, 2020.
- [33] E. H. Kennedy. Semiparametric theory and empirical processes in causal inference. In *Statistical causal inferences and their applications in public health research*, pages 141–167. Springer, 2016.
- [34] E. H. Kennedy, S. Balakrishnan, M. G’Sell, et al. Sharp instruments for classifying compliers and generalizing causal effects. *Annals of Statistics*, 48(4):2008–2030, 2020.
- [35] E. H. Kennedy, S. Balakrishnan, and L. Wasserman. Semiparametric counterfactual density estimation. *arXiv preprint arXiv:2102.12034*, 2021.
- [36] E. H. Kennedy, S. Lorch, and D. S. Small. Robust causal inference with continuous instruments using the local instrumental variable curve. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 81(1):121–143, 2019.
- [37] E. H. Kennedy, Z. Ma, M. D. McHugh, and D. S. Small. Nonparametric methods for doubly robust estimation of continuous treatment effects. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, 79(4):1229, 2017.
- [38] K. Kim, J. Kim, and E. H. Kennedy. Causal effects based on distributional distances. *arXiv preprint arXiv:1806.02935*, 2018.
- [39] C. A. Klaassen. Consistent estimation of the influence function of locally asymptotically linear estimators. *The Annals of Statistics*, pages 1548–1562, 1987.
- [40] T. Leonard. Density estimation, stochastic processes and prior information. *Journal of the Royal Statistical Society: Series B (Methodological)*, 40(2):113–132, 1978.
- [41] B. Melly and K. Wüthrich. Local quantile treatment effects. 2016.
- [42] R. Neugebauer and M. van der Laan. Nonparametric causal effects based on marginal structural models. *Journal of Statistical Planning and Inference*, 137(2):419–434, 2007.

- [43] E. L. Ogburn, A. Rotnitzky, and J. M. Robins. Doubly robust estimation of the local average treatment effect curve. *Journal of the Royal Statistical Society. Series B, Statistical methodology*, 77(2):373, 2015.
- [44] R. Okui, D. S. Small, Z. Tan, and J. M. Robins. Doubly robust instrumental variable regression. *Statistica Sinica*, pages 173–205, 2012.
- [45] J. Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–710, 1995.
- [46] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2000. 2nd edition, 2009.
- [47] J. Pearl. Parameter identification: A new perspective. Technical Report R-276, 2001.
- [48] J. Robins, L. Li, E. Tchetgen, A. van der Vaart, et al. Higher order influence functions and minimax estimation of nonlinear functionals. In *Probability and statistics: essays in honor of David A. Freedman*, pages 335–421. Institute of Mathematical Statistics, 2008.
- [49] J. Robins and A. Rotnitzky. Comment on “inference for semiparametric models: Some questions and an answer,” by pj bickel and j. kwon. *Statistica Sinica*, 11:920–936, 2001.
- [50] J. M. Robins and Y. Ritov. Toward a curse of dimensionality appropriate (coda) asymptotic theory for semi-parametric models. *Statistics in medicine*, 16(3):285–319, 1997.
- [51] P. R. Rosenbaum and D. B. Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- [52] K. J. Rothman and S. Greenland. Causation and causal inference in epidemiology. *American journal of public health*, 95(S1):S144–S150, 2005.
- [53] R. Singh and L. Sun. De-biased machine learning for compliers. *arXiv preprint arXiv:1909.05244*, 2019.
- [54] E. Smucler, A. Rotnitzky, and J. M. Robins. A unifying approach for doubly-robust ℓ_1 regularized estimation of causal contrasts. *arXiv preprint arXiv:1904.03737*, 2019.
- [55] Z. Tan. Regression and weighting methods for causal inference using instrumental variables. *Journal of the American Statistical Association*, 101(476):1607–1618, 2006.
- [56] A. B. Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.
- [57] S. D. Uysal. Doubly robust iv estimation of the local average treatment effects, 2011.
- [58] A. W. Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [59] L. Wang, Y. Zhang, T. S. Richardson, and J. M. Robins. Estimation of local treatment effects under the binary instrumental variable model. *Biometrika*, 2021.
- [60] L. Wasserman. *All of nonparametric statistics*. Springer Science & Business Media, 2006.
- [61] J. M. Wooldridge. *Econometric analysis of cross section and panel data*. MIT press, 2010.
- [62] P. G. Wright. *Tariff on animal and vegetable oils*. Macmillan Company, New York, 1928.
- [63] K. Wüthrich. A comparison of two quantile models with endogeneity. *Journal of Business & Economic Statistics*, 38(2):443–456, 2020.
- [64] W. Zheng and M. J. van der Laan. Cross-validated targeted minimum-loss-based estimation. In *Targeted Learning*, pages 459–474. Springer, 2011.

Appendix – Double Machine Learning Density Estimation for Local Treatment Effects with Instruments

A IV Settings and LTE

In this work, we consider the IV setting represented by the causal graph G in Fig. 1. It is common in the literature to define IV assumptions in terms of conditional independences among counterfactuals [2, 55, 43, 59], as given in the following:

Assumption A.1 (IV assumptions).

1. **Exclusion restriction:** $Y_{x,z} = Y_x$ almost surely for all z, x .
2. **Independence:** $Z \perp\!\!\!\perp (Y_x, X_z)|W$ for all z, x .
3. **Instruments relevance:** $P(X_{Z=1} = 1|W) \neq P(X_{Z=0} = 1|W)$ almost surely.

We show that the causal graph in Fig. 1 captures the set of IV assumptions in Assumption A.1.

Lemma A.1. *The causal graph G in Fig. 1 satisfies the set of IV assumptions in Assumption A.1.*

Proof. We will show the first item. We have $Y_{x,z} = Y_{x,z,W_x} = Y_{x,W_x} = Y_x$, where the first equality is due to the composition property [46, Property 1 (pp. 229)], the second due to exclusion restrictions [46, Eq.(7.25)], and the third by composition.

We will show the second. We have $(Z_w \perp\!\!\!\perp \{W, X_{z,w}, Y_{x,w}\})$ by independence restrictions [46, Eq.(7.26)]. Then by the weak union graphoid axiom (Refer [46, pp.11]), $(Z_w \perp\!\!\!\perp \{X_{z,w}, Y_{x,w}\}|W)$, which leads to $(Z \perp\!\!\!\perp (Y_x, X_z)|W)$ by composition.

We will show the third. By $(Z \perp\!\!\!\perp X_z|W)$, $P(x_z|w) = P(x_z|w, z) = P(x|w, z)$, where the second equality is by composition. The third assumption is reflected by that X is not independent of Z given W in G . \square

Definition A.1 (Local treatment effect (LTE) density). *The local treatment effect (LTE) density is the density of outcome Y under treatment $X = x$ among compliers (i.e., $X_{Z=1} = 1$ and $X_{Z=0} = 0$) denoted by $p(y_x|X_{Z=1} = 1, X_{Z=0} = 0)$. We will use $C = (X_{Z=1} = 1 \wedge X_{Z=0} = 0)$ to denote the event that a unit is a complier and write the LTE density as $p(y_x|C)$.*

The LTE density $p(y_x|C)$ is known to be identifiable under monotonicity in the IV settings [30, 2]. In the notations of this paper, we present the identification results as follows, where for a given constant a and a variable X , x^a denotes the event $X = a$.

Lemma A.2. *In the causal graph G in Fig. 1, $p(y_x|w, C)$ is identifiable under monotonicity and is given by*

$$p(y_x|w, C) = \frac{p(y|x, z^x, w)P(x|z^x, w) - p(y|w, x, z^{1-x})P(x|z^{1-x}, w)}{P(x^1|z^1, w) - P(x^1|z^0, w)}.$$

Theorem A.1. *In the causal graph G in Fig. 1, the LTE density $p(y_x|C)$ is identifiable under monotonicity and is given by*

$$\begin{aligned} p(y_x|C) &= \frac{\int_{\mathcal{W}} [p(y|x, z^x, w)P(x|z^x, w) - p(y|w, x, z^{1-x})P(x|z^{1-x}, w)]P(w) d[w]}{\int_{\mathcal{W}} [P(x^1|z^1, w) - P(x^1|z^0, w)]P(w) d[w]} \\ &\equiv \frac{\mathbb{E}_P [p(y|x, z^x, W)P(x|z^x, W) - p(y|x, z^{1-x}, W)P(x|z^{1-x}, W)]}{\mathbb{E}_P [P(x^1|z^1, W) - P(x^1|z^0, W)]}. \end{aligned}$$

B Proofs

Notations We will use $P_\epsilon \equiv P(1 + \epsilon g)$, where g is a mean zero bounded random function, to denote a parametric submodel for the probability measure P . Also, we note that the causal effect

$\psi[f(Y)]$ in Eq. (3) can be written as

$$\psi^{YX}[f(Y)] \equiv \mathbb{E}_P [\theta_y(x, z^x, W)[f(Y)] - \theta_y(x, z^{1-x}, W)[f(Y)]] \quad (\text{B.1})$$

$$\psi^X \equiv \mathbb{E}_P [\xi_{x^1}(z^1, W) - \xi_{x^1}(z^0, W)], \quad (\text{B.2})$$

for $\pi_z(w) \equiv P(z|w)$, $\xi_x(z, w) \equiv P(x|z, w)$ and $\theta_y(x, z, w) = \mathbb{E}_P [f(Y)\mathbb{1}_x(X)|z, w]$

Lemma S.1 ([58, Thm.5.31],[35, Lemma 3]). *Let $\phi(\mathbf{V}; \theta, \eta)$ denote a vector estimating function for target parameter $\theta \in \mathbb{R}^p$ and nuisance functions $\eta \in H$ for some function space H . Suppose $\mathbb{E}_P [\phi(\mathbf{V}; \theta_0, \eta_0)] = 0$ (where θ_0, η_0 denote true parameters) and define the estimator $\hat{\theta}$ as a solution to $\mathbb{E}_D [\phi(\mathbf{V}; \hat{\theta}, \hat{\eta})] = o_P(n^{-1/2})$, where η is estimated on a separate independent sample. Assume*

1. $\{\phi(\mathbf{V}; \theta, \eta) : \theta \in \mathbb{R}^p\}$ is Donsker for any fixed η .
2. $\hat{\theta} - \theta_0 = o_P(1)$ and $\|\hat{\eta} - \eta\|_2 = o_P(1)$.
3. The map $\theta \mapsto \mathbb{E}_P [\phi(\mathbf{V}; \theta, \eta)]$ is differentiable at θ_0 uniformly in η , with non-singular matrix $M(\theta_0, \eta) \equiv (\partial/\partial\theta)|_{\theta_0} \mathbb{E}_P [\phi(\mathbf{V}; \theta, \eta)]$, where $M(\theta_0, \hat{\eta}) \xrightarrow{P} M \equiv M(\theta_0, \eta_0)$.

Then,

$$\hat{\theta} - \theta_0 = -M^{-1} \mathbb{E}_D [\phi(\mathbf{V}; \theta_0, \eta_0)] - M^{-1} \mathbb{E}_P [\phi(\mathbf{V}; \theta_0, \hat{\eta})] + o_P(n^{-1/2}).$$

B.1 Proofs for Sec. 3

Lemma S.2 ([27, Proof of Thm. 1]). *For a target estimand $\gamma \equiv \mathbb{E}_P [\mathbb{E}_P [f(Y)|x^1, W] - \mathbb{E}_P [f(Y)|x^0, W]]$ for binary $X \in \{0, 1\}$ and $f(\cdot) < \infty$, an influence function ϕ_γ is given by*

$$\phi_\gamma \equiv \frac{\mathbb{1}_{x^1}(X) - \mathbb{1}_{x^0}(X)}{P(X|W)} (f(Y) - \mathbb{E}_P [f(Y)|X, W]) + (\mathbb{E}_P [f(Y)|x^1, W] - \mathbb{E}_P [f(Y)|x^0, W]) - \gamma.$$

Lemma S.3. *An influence function for $\psi[f(Y)]$ for $f(Y) < \infty$ is given by the mapping function in Eq. (7), which is*

$$\phi(\eta = \{\pi, \xi, \theta\}, \psi)[f(Y)] \equiv \frac{1}{\psi^X} (\mathcal{V}_{YX}(\{\pi, \theta\})[f(Y)] - \psi[f(Y)]\mathcal{V}_X(\{\pi, \xi\})).$$

Proof. We note that the estimand is given as $\psi[f(Y)] = \psi^{YX}[f(Y)]/\psi^X$, where

$$\psi^X = \mathbb{E}_P [\xi_{x^1}(z^1, W) - \xi_{x^1}(z^0, W)] \quad (\text{B.3})$$

$$\psi^{YX}[f(Y)] = \mathbb{E}_P [\theta(x, z^x, W)[f(Y)] - \theta(x, z^{1-x}, W)[f(Y)]] . \quad (\text{B.4})$$

By Lemma S.2, influence functions corresponding to ψ^X and $\psi^{YX}[f(Y)]$, denoted ϕ_X and $\phi_{YX}[f(Y)]$ respectively, are given as

$$\phi_X \equiv \frac{\mathbb{1}_{z^1}(Z) - \mathbb{1}_{z^0}(Z)}{\pi_Z(W)} (\mathbb{1}_{x^1}(X) - \xi_{x^1}(Z, W)) + (\xi_{x^1}(z^1, W) - \xi_{x^1}(z^0, W)) - \psi^X \quad (\text{B.5})$$

$$\begin{aligned} \phi_{YX}[f(Y)] &\equiv \frac{\mathbb{1}_{z^x}(Z) - \mathbb{1}_{z^{1-x}}(Z)}{\pi_Z(W)} (f(Y)\mathbb{1}_x(X) - \theta(x, Z, W)[f(Y)]) \\ &+ (\theta(x, z^x, W)[f(Y)] - \theta(x, z^{1-x}, W)[f(Y)]) - \psi^{YX}[f(Y)]. \end{aligned} \quad (\text{B.6})$$

(YJ, from R4) 1. “influence(A)/B - influence(B)/A/B2, just like the quotient rule for differentiation. Is this obvious?”;

Then, by chain rule, an influence function for $\psi[f(Y)] = \psi^{YX}[f(Y)]/\psi^X$ is given as

$$\begin{aligned} & \frac{1}{\psi^X} (\phi_{YX}[f(Y)] - \psi[f(Y)]\phi_X[f(Y)]) \\ &= \frac{1}{\psi^X} (\mathcal{V}_{YX}[f(Y)] - \psi^{YX}[f(Y)] - \psi[f(Y)] (\mathcal{V}_X[f(Y)] - \psi^X[f(Y)])) \\ &= \frac{1}{\psi^X} (\mathcal{V}_{YX}[f(Y)] - \psi[f(Y)]\mathcal{V}_X - \psi[f(Y)] + \psi[f(Y)]) \\ &= \frac{1}{\psi^X} (\mathcal{V}_{YX}[f(Y)] - \psi[f(Y)]\mathcal{V}_X). \end{aligned}$$

□

Lemma B.1 (Restated Lemma 1). *Let $m(\psi'; \psi_h)$ be the score defined in Eq. (6). Then, an influence function for $\mathbb{E}_P[m(\psi'; \psi_h)]$, denoted ϕ_m , is given by*

$$\phi_m(\eta = \{\pi, \xi, \theta\}, \psi) \equiv \phi(\eta, \psi)[K_{h,y}(Y)] \quad (\text{B.7})$$

where ϕ is given as

$$\phi(\eta = \{\pi, \xi, \theta\}, \psi)[f(Y)] \equiv \frac{1}{\psi^X} (\mathcal{V}_{YX}(\{\pi, \theta\})[f(Y)] - \psi[f(Y)]\mathcal{V}_X(\{\pi, \xi\}))$$

Proof. Let ϕ_X denote an influence function corresponding to ψ^X , given in Eq. (B.5). This implies that $\mathbb{E}_P[\mathcal{V}_X] = \psi^X$. Then,

$$\mathbb{E}_P[m(\psi'; \psi_h)] = \mathbb{E}_P \left[\frac{1}{\psi^X} (\psi_h - \psi') \mathcal{V}_X \right] = \frac{1}{\psi^X} (\psi_h - \psi') \mathbb{E}_P[\mathcal{V}_X] = \psi_h - \psi'.$$

Then, an influence function for $\mathbb{E}_P[m(\psi'; \psi_h)]$ coincides with the influence function for ψ_h , which is given by Eq. (B.7) based on Lemma S.3. □

Lemma B.2 (Restated Lemma 2). *Let $m(\psi'; \psi_h)$ be the score function in Eq. (6), and $\phi_m(\eta = \{\pi, \xi, \theta\}, \psi_h)$ be the influence function for $\mathbb{E}_P[m(\psi'; \psi_h)]$ given in Eq. (9). Then, a Neyman orthogonal score for ψ_h is given as $\varphi(\psi'; \eta = \{\pi, \xi, \theta\}) \equiv m(\psi'; \psi_h) + \phi_m(\eta, \psi)$; Specifically,*

$$\varphi(\psi'; \eta = \{\pi, \xi, \theta\}) = \frac{1}{\psi^X} (\mathcal{V}_{YX}(\{\pi, \theta\})[K_{h,y}(Y)] - \psi' \mathcal{V}_X(\{\pi, \xi\})). \quad (\text{B.8})$$

Proof. For a score function for ψ , denoted $m(\cdot)$, and an influence function for $\mathbb{E}_P[m(\cdot)]$, denoted $\phi_m(\cdot)$, a Neyman orthogonal score for ψ is given as $m + \phi_m$ [13, Thm. 1]. Applying this, $m(\psi'; \psi_h) + \phi_m(\eta, \psi_h)$ is a Neyman orthogonal score. Specifically,

$$\begin{aligned} & \varphi(\psi'; \eta = \{\pi, \xi, \theta\}) \\ &= m(\psi'; \psi_h) + \phi_m(\eta, \psi_h) \\ &= \frac{1}{\psi^X} (\psi[K_{h,y}(Y)] - \psi') \mathcal{V}_X + \frac{1}{\psi^X} (\mathcal{V}_{YX}(\{\pi, \theta\})[K_{h,y}(Y)] - \psi[K_{h,y}(Y)]\mathcal{V}_X(\{\pi, \xi\})) \\ &= \frac{1}{\psi^X} (\mathcal{V}_{YX}(\eta = \{\pi, \xi, \theta\})[K_{h,y}(Y)] - \psi' \mathcal{V}_X(\{\pi, \xi\})). \end{aligned}$$

□

Lemma B.3 (Restated Lemma 3). *For any fixed $y \in \mathcal{Y}$, suppose the estimators for nuisances are consistent; i.e., $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \eta = \{\pi, \xi, \theta\}$ for all (w, z, x) . Suppose $h < \infty$, and $nh^d \rightarrow \infty$ as $n \rightarrow \infty$. Then,*

$$\hat{\psi}_h(y) - \psi_h(y) = O_P \left(1/\sqrt{nh^d} + R_2^k + 1/\sqrt{n} \right),$$

where

$$R_2^k \equiv \sum_z \|\hat{\pi}_z - \pi_z\| \left\{ \left\| \hat{\theta}_z - \theta_z \right\| + \left\| \hat{\xi}_z - \xi_z \right\| \right\}, \quad (\text{B.9})$$

where $\pi_z \equiv \pi_z(W)$, $\xi_z \equiv \xi_x(z, W)$ and $\theta_z \equiv \theta(x, z, W)[K_{h,y}(Y)]$.

Proof. We note that the condition $nh^d \rightarrow \infty$ means that $h = O(n^{-\alpha})$ for some $\alpha < 1/d$. $h < \infty$ implies that such h is either constant or decreasing function over n . Combining, the condition implies $h = O(n^{-\alpha})$ for $\alpha \in [0, 1/d)$.

We recall that ψ^X, ψ^{YX} are defined in Eq. (B.2.B.1) and $\mathcal{V}_X, \mathcal{V}_{YX}$ are defined in Eq. (5,8).

Now, we will prove this Lemma through the master result in Lemma S.1. The KLTE estimator $\hat{\psi}_h$ in Eq. (11) satisfies $\mathbb{E}_{\mathcal{D}} [\varphi(\hat{\psi}_h, \eta)] = o_P(n^{-1/2})$, because

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\varphi(\hat{\psi}_h, \eta)] &= \frac{1}{\psi^X} \left(\mathbb{E}_{\mathcal{D}} [\mathcal{V}_{YX}[K_{h,y}(Y)]] - \hat{\psi}_h \mathbb{E}_{\mathcal{D}} [\mathcal{V}_X] \right) \\ &= \frac{1}{\psi^X} \left(\mathbb{E}_{\mathcal{D}} [\mathcal{V}_{YX}[K_{h,y}(Y)]] - \frac{\mathbb{E}_{\mathcal{D}} [\mathcal{V}_{YX}[K_{h,y}(Y)]]}{\mathbb{E}_{\mathcal{D}} [\mathcal{V}_X]} \mathbb{E}_{\mathcal{D}} [\mathcal{V}_X] \right) \\ &= 0. \end{aligned}$$

We have that the score function φ in Lemma 2 satisfies the assumptions in Lemma S.1, since φ is a linear function of ψ when nuisances are fixed. Also, M in Lemma S.1 is given as -1 . (YJ, from R4 why is $M = -1$ obvious?)

Then, by the result of Lemma S.1,

$$\hat{\psi}_h - \psi_h = \mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)] + \mathbb{E}_P [\phi_m(\psi_h, \hat{\eta})] + o_P(n^{-1/2}).$$

We will first study the convergence behavior of $\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)]$. We will show that $\mathbb{E}_P [\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)]] = O(1/\sqrt{nh^d})$. Then, $\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)]$ being $\sqrt{nh^d}$ -consistency (i.e., $\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)] = O_P(1/\sqrt{nh^d})$) can be shown immediately by Markov inequality. This implies that $\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)]$ is consistent if $nh^d \rightarrow \infty$. Let $\phi_m(V_i, \psi, \eta)$ denote the influence function evaluated at $V_i \in \mathcal{D}$. Consider the following:

$$\begin{aligned} \mathbb{E}_P [|\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)]|] &\leq \sqrt{\mathbb{E}_P [(\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)])^2]} \\ &= \sqrt{\text{var}_P (\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)])} \\ &= \sqrt{(1/n) \mathbb{E}_P [\phi_m^2(\psi_h, \eta)]}, \end{aligned}$$

where the first inequality is by Cauchy-Schwarz inequality, the second and third equality are from the iid assumption and $\mathbb{E}_P [\phi_m] = 0$.

We note that

$$\begin{aligned} \phi_m &= \frac{1}{\psi^X} (\mathcal{V}_{YX}[K_{h,y}(Y)] - \psi_h \mathcal{V}_X) \\ &= \frac{1}{\psi^X} (\mathcal{V}_{YX}[K_{h,y}(Y)] - \psi_h \mathcal{V}_X) + \underbrace{\frac{\psi^{YX}[K_{h,y}(Y)]}{\psi^X} - \frac{\psi^X}{\psi^X} \psi_h}_{=0} \\ &= \frac{1}{\psi^X} (\{\mathcal{V}_{YX}[K_{h,y}(Y)] - \psi^{YX}[K_{h,y}(Y)]\} - \psi_h \{\mathcal{V}_X - \psi^X\}) \\ &= \frac{1}{\psi^X} (\phi_{YX}[K_{h,y}(Y)] - \psi_h \phi_X). \end{aligned}$$

Next,

$$\begin{aligned} \mathbb{E}_P [\phi_m^2(\psi_h, \eta)] &= \mathbb{E}_P \left[\frac{1}{\psi_X^2} \{\phi_{XY}[K_{h,y}(Y)] - \psi_h \phi_X\}^2 \right] \\ &= \frac{1}{\psi_X^2} \mathbb{E}_P \left[\{\phi_{XY}[K_{h,y}(Y)] - \psi_h \phi_X\}^2 \right] \\ &= \frac{1}{\psi_X^2} \mathbb{E}_P [\phi_{XY}^2[K_{h,y}(Y)] + \psi_h^2 \phi_X^2 - 2\phi_{XY}[K_{h,y}(Y)]\phi_X\psi_h]. \end{aligned}$$

We first analyze $\mathbb{E}_P [\phi_{XY}^2[K_{h,y}(Y)]] = \text{var}_P[\phi_{XY}[K_{h,y}(Y)]]$. By [27, Thm. 1],

$$\begin{aligned} \text{var}_P[\phi_{XY}[K_{h,y}(Y)]] &= \mathbb{E}_P \left[\frac{\text{Var}_P(K_{h,y}(Y)\mathbb{1}_x(X)|z^x, W)}{\pi_{z^x}(W)} + \frac{\text{Var}_P(K_{h,y}(Y)\mathbb{1}_x(X)|z^{1-x}, W)}{\pi_{z^{1-x}}(W)} \right] \\ &\quad + \mathbb{E}_P \left[\left\{ \mathbb{E}_P[K_{h,y}(Y)\mathbb{1}_x(X)|z^x, W] - \mathbb{E}_P[K_{h,y}(Y)\mathbb{1}_x(X)|z^{1-x}, W] - \psi^{YX}[K_{h,y}] \right\}^2 \right]. \end{aligned}$$

First,

$$\mathbb{E}_P [\text{Var}_P(K_{h,y}(Y)\mathbb{1}_x(X)|z^x, W)] = O(\text{var}_P(K_{h,y}(Y))) \quad (\text{B.10})$$

$$\begin{aligned} &= O\left(\int_{\mathcal{Y}} K_{h,y}^2(Y) d[y]\right) \\ &= O\left(\left(1/h^{2d}\right) \int_{\mathcal{Y}} K^2((Y-y)/h) d[y]\right) \\ &= O\left(\left(1/h^d\right) \int_{\mathcal{U}} K^2(u) d[u]\right) = O(1/h^d). \quad (\text{B.11}) \end{aligned}$$

The first equality holds by Law of total variance, the second by

Also,

$$\begin{aligned} &\mathbb{E}_P \left[\left\{ \mathbb{E}_P[K_{h,y}(Y)\mathbb{1}_x(X)|z^x, W] - \mathbb{E}_P[K_{h,y}(Y)\mathbb{1}_x(X)|z^{1-x}, W] - \psi^{YX}[K_{h,y}] \right\}^2 \right] \\ &= \text{var}_P \left(\left\{ \mathbb{E}_P[K_{h,y}(Y)\mathbb{1}_x(X)|z^x, W] - \mathbb{E}_P[K_{h,y}(Y)\mathbb{1}_x(X)|z^{1-x}, W] \right\} \right) \\ &\leq 2 \sup_{z \in \{0,1\}} \text{var}_P(\mathbb{E}_P[K_{h,y}(Y)\mathbb{1}_x(X)|z^x, W]) \\ &\leq 2 \sup_{z \in \{0,1\}} \text{var}_P(K_{h,y}(Y)\mathbb{1}_x(X)|z^x, W) \\ &= O(1/h^d), \end{aligned}$$

where the first (in)equality is by the definition of the variance, the second by the linear combination of the variance, the third by the law of total variance, the fourth by Eq. (B.11). Therefore, $\text{var}_P[\phi_{XY}[K_{h,y}(Y)]] = O(1/h^d)$.

Next, we will study $\mathbb{E}_P [\psi_h^2 \phi_X^2]$. We first note that $\mathbb{E}_P [\psi_h^2 \phi_X^2] = \psi_h^2 \mathbb{E}_P [\phi_X^2] = O(\psi_h^2)$. Therefore, it suffices to show ψ_h^2 .

$$\begin{aligned} \psi_h^2 &= \left(\int_{\mathcal{Y}} K_{h,y}(y') \psi(y') d[y'] \right)^2 \\ &\leq \int_{\mathcal{Y}} K_{h,y}^2(y') \psi^2(y') d[y'] \\ &\leq \int_{\mathcal{Y}} K_{h,y}^2(y') \psi(y') d[y'] \\ &= \int_{\mathcal{Y}} \frac{1}{h^{2d}} K^2\left(\frac{y'-y}{h}\right) \psi(y') d[y'] \\ &= \int_{\mathcal{U}} \frac{1}{h^d} K^2(u) \psi(y+uh) d[u] \\ &= O(1/h^d). \end{aligned}$$

Finally, consider the term $-2\mathbb{E}_P[\phi_{YX}[K_{h,y}(Y)] \cdot \phi_X \cdot \psi_h]$. Note, $\mathbb{E}_P[\phi_{YX}[K_{h,y}(Y)] \cdot \phi_X \cdot \psi_h] = \psi_h \cdot \mathbb{E}_P[\phi_{YX}[K_{h,y}(Y)] \cdot \phi_X]$. We first consider $\mathbb{E}_P[\phi_{YX}[K_{h,y}(Y)] \cdot \phi_X]$:

$$\begin{aligned} \mathbb{E}_P[\phi_{YX}[K_{h,y}(Y)] \cdot \phi_X] &= \mathbb{E}_P[\phi_{YX}[K_{h,y}(Y)] \cdot \phi_X] \\ &\leq \sqrt{\mathbb{E}_P[\phi_{YX}^2[K_{h,y}(Y)]] \cdot \mathbb{E}_P[\phi_X^2]} \\ &= O\left(\sqrt{\mathbb{E}_P[\phi_{YX}^2[K_{h,y}(Y)]]}\right) \\ &= O\left(h^{-d/2}\right). \end{aligned}$$

Now, consider ψ_h :

$$\begin{aligned}
\psi_h &\equiv \int_{\mathcal{Y}} K_{h,y}(y') \psi(y') d[y'] \\
&= \int_{\mathcal{Y}} \frac{1}{h} K\left(\frac{y' - y}{h}\right) \psi(y') d[y'] \\
&= \int_{\mathcal{U}} K(u) \psi(hu + y) d[u] \\
&= \int_{\mathcal{U}} K(u) \left(\psi(y) + hu\psi^{(1)}(y) + h^2u^2\psi^{(2)}(y) + O(h^2u^2) \right) d[u] \\
&= O(h^2).
\end{aligned}$$

With $h = O(n^{-\alpha})$, we note $O(h^{-d/2}) = O(n^{\alpha d/2})$. Therefore, $\mathbb{E}_P [\phi_m^2(\psi_{p,h}, \eta)] = O(h^{-d} + h^{-d/2} + h^2) = O(n^{\alpha d}) = O(h^{-d})$ since $h = O(n^{-\alpha})$ for some $\alpha \in [0, 1)$. This shows that $\mathbb{E}_P [\mathbb{E}_{\mathcal{D}} [\phi_m(\psi_h, \eta)]] = O(1/\sqrt{nh^d})$.

We now consider $\mathbb{E}_P [\phi_m(\psi_h, \hat{\eta})]$.

$$\begin{aligned}
&\mathbb{E}_P [\phi(\psi_h, \hat{\eta})] \\
&= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left(\hat{\mathcal{V}}_{YX}[K_{h,y}(Y)] - \psi[K_{h,y}(Y)]\hat{\mathcal{V}}_X \right) \right] \\
&= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left(\hat{\mathcal{V}}_{YX}[K_{h,y}(Y)] - \psi[K_{h,y}(Y)]\hat{\mathcal{V}}_X \right) + \left(\frac{1}{\hat{\psi}_X} - \frac{1}{\psi_X} \right) \left(\hat{\mathcal{V}}_{YX}[K_{h,y}(Y)] - \psi[K_{h,y}(Y)]\hat{\mathcal{V}}_X \right) \right].
\end{aligned} \tag{B.12}$$

For further analysis, we consider $\mathbb{E}_P [\hat{\mathcal{V}}_{YX}[K_{h,y}(Y)] - \mathcal{V}_{YX}[K_{h,y}(Y)]]$. First, define

$$\mathcal{V}_{YX,(x,z)}(\pi, \theta)[f(Y)] \equiv \frac{\mathbb{1}_z(Z)}{\pi_Z(W)} (f(Y)\mathbb{1}_x(X) - \theta(x, Z, W)[f(Y)]) + \theta(x, z, W)[f(Y)].$$

Then, $\mathcal{V}_{YX}[f(Y)] = \mathcal{V}_{YX,(x,z^x)}[f(Y)] - \mathcal{V}_{YX,(x,z^{1-x})}[f(Y)]$. Now, consider

$\mathbb{E}_P [\hat{\mathcal{V}}_{YX,(x,z)}[K_{h,y}(Y)] - \mathcal{V}_{YX,(x,z)}[K_{h,y}(Y)]]$. We have

$$\begin{aligned}
&\mathbb{E}_P \left[\mathcal{V}_{YX,(x,z)}(\hat{\pi}, \hat{\theta})[f(Y)] - \mathcal{V}_{YX,(x,z)}(\pi, \theta)[f(Y)] \right] \\
&= \mathbb{E}_P \left[\frac{\mathbb{1}_z(Z)}{\hat{\pi}_Z(W)} \left(f(Y)\mathbb{1}_x(X) - \hat{\theta}(x, Z, W)[f(Y)] \right) + \hat{\theta}(x, z, W)[f(Y)] - \theta(x, z, W)[f(Y)] \right] \\
&= \mathbb{E}_P \left[\frac{\mathbb{1}_z(Z)}{\hat{\pi}_Z(W)} \left(\theta(x, Z, W)[f(Y)] - \hat{\theta}(x, Z, W)[f(Y)] \right) + \left\{ \hat{\theta}(x, z, W)[f(Y)] - \theta(x, z, W)[f(Y)] \right\} \right] \\
&= \mathbb{E}_P \left[\frac{\pi_z(W)}{\hat{\pi}_z(W)} \left(\theta(x, z, W)[f(Y)] - \hat{\theta}(x, z, W)[f(Y)] \right) + \left\{ \hat{\theta}(x, z, W)[f(Y)] - \theta(x, z, W)[f(Y)] \right\} \right] \\
&= \mathbb{E}_P \left[\left(\theta(x, z, W)[f(Y)] - \hat{\theta}(x, z, W)[f(Y)] \right) \left(1 - \frac{\pi_z(W)}{\hat{\pi}_z(W)} \right) \right] \\
&= \mathbb{E}_P \left[\left(\theta(x, z, W)[f(Y)] - \hat{\theta}(x, z, W)[f(Y)] \right) \left(\frac{\hat{\pi}_z(W) - \pi_z(W)}{\hat{\pi}_z(W)} \right) \right] \\
&= O_P \left(\left\| \theta(x, z, W)[f(Y)] - \hat{\theta}(x, z, W)[f(Y)] \right\| \left\| \hat{\pi}_z(W) - \pi_z(W) \right\| \right),
\end{aligned}$$

where the first and the second are by the fact that $\mathbb{E}_P [f(Y)\mathbb{1}_x(X)|W, Z, X] = \theta(x, Z, W)[f(Y)]$, the third is by taking an expectation over Z conditioned on W , the fourth and the fifth by rearrangement, and the sixth by Cauchy-Schwarz inequality and Positivity. Then,

$$\begin{aligned}
R_{YX} &\equiv \mathbb{E}_P \left[\mathcal{V}_{YX}(\hat{\pi}, \hat{\theta})[f(Y)] - \mathcal{V}_{YX}(\pi, \theta)[f(Y)] \right] \\
&= \sum_{z \in \{0,1\}} O_P \left(\left\| \theta(x, z, W)[f(Y)] - \hat{\theta}(x, z, W)[f(Y)] \right\| \left\| \hat{\pi}_z(W) - \pi_z(W) \right\| \right).
\end{aligned}$$

Also, let

$$\mathcal{V}_{X,x}(\pi, \xi) \equiv \frac{\mathbb{1}_z(Z)}{\pi_Z(W)} (\mathbb{1}_x(X) - \xi_x(Z, W)) + \xi_x(z, W).$$

Then, with the similar proof as above, we have

$$\mathbb{E}_P \left[\mathcal{V}_{X,x}(\hat{\pi}, \hat{\xi}) - \mathcal{V}_{X,x}(\pi, \xi) \right] = O_P \left(\left\| \xi_x(z, W) - \hat{\xi}_x(z, W) \right\| \left\| \hat{\pi}_z(W) - \pi_z(W) \right\| \right),$$

and

$$\mathbb{E}_P \left[\mathcal{V}_X(\hat{\pi}, \hat{\xi}) - \mathcal{V}_X(\pi, \xi) \right] = \sum_{z \in \{0,1\}} O_P \left(\left\| \xi_x(z, W) - \hat{\xi}_x(z, W) \right\| \left\| \hat{\pi}_z(W) - \pi_z(W) \right\| \right).$$

Let $R_{YX} \equiv \mathbb{E}_P \left[\hat{\mathcal{V}}_{YX} - \mathcal{V}_{YX} \right]$ and $R_X \equiv \mathbb{E}_P \left[\hat{\mathcal{V}}_X - \mathcal{V}_X \right]$. Then, continuing from Eq. (B.12),

$$\begin{aligned} \text{Eq. (B.12)} &= \mathbb{E}_P \left[\frac{1}{\psi_X} \left(\psi_{YX} + R_{YX} - \frac{\psi_{YX}}{\psi_X} (\psi_X + R_X) \right) + \left(\frac{1}{\hat{\psi}_X} - \frac{1}{\psi_X} \right) \left(\psi_{YX} + R_{YX} - \frac{\psi_{YX}}{\psi_X} (\psi_X + R_X) \right) \right] \\ &= \mathbb{E}_P \left[\frac{1}{\psi_X} (R_{YX} - \psi R_X) + \left(\frac{1}{\hat{\psi}_X} - \frac{1}{\psi_X} \right) (R_{YX} - \psi R_X) \right] \\ &= O_P(R_{YX} + R_X) \\ &= O_P(R_2^k), \end{aligned}$$

where

$$R_2^k = \sum_{z \in \{0,1\}} O_P \left(\left\| \hat{\pi}_z - \pi_z \right\| \left\{ \left\| \hat{\theta}_z - \theta_z \right\| + \left\| \hat{\xi}_z - \xi_z \right\| \right\} \right).$$

Note the first equality is by $\hat{\mathcal{V}}_{YX} = R_{YX} + \mathcal{V}_{YX}$ and $\hat{\mathcal{V}}_X = R_X + \mathcal{V}_X$, the second by rearrangement, the third by Positivity, the fourth by the definition of R_{YX} and R_X .

Summing up, we have shown that $\mathbb{E}_P [\phi_m(\psi_h, \eta)] = O(1/\sqrt{nh^d})$ and $\mathbb{E}_P [\phi_m(\psi_h, \hat{\eta})] = O_P(R_2^k)$. \square

Corollary 1 (Restated Corol. 1). *If all nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ for any given (w, z, x, y) converge at rate $\{nh^d\}^{-1/4}$, then the target estimator $\hat{\psi}_h(y)$ achieves $\sqrt{nh^d}$ -rate convergence to ψ_h .*

Proof. This result follows immediately from Lemma 3. \square

Theorem B.1 (Restated Thm. 1). *For any fixed $y \in \mathcal{Y}$, suppose the estimators for nuisances are consistent; i.e., $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \eta = \{\pi, \xi, \theta\}$ for all (w, z, x) . Suppose $h < \infty$, and $nh^d \rightarrow \infty$ as $n \rightarrow \infty$. Then*

$$\hat{\psi}_h(y) - \psi(y) = O_P \left(1/\sqrt{nh^d} + R_2^k + 1/\sqrt{n} \right) + B_y, \quad (\text{B.13})$$

where B_y is defined in Eq. (13), and R_2^k is defined in Eq. (12).

Proof. This result follows immediately from Lemmas 3 and 4. \square

Lemma B.5 (Restated Lemma 5). *The bandwidth h that minimizes the error in Eq. (14) is $h = O(n^{-1/(d+4)})$. This choice of h satisfies the assumption in Lemma. 3 that $nh^d \rightarrow \infty$.*

Proof. We note that the error in Eq. (14) w.r.t. h is $O_P(1/\sqrt{nh^d} + h^2)$. Since the function $1/\sqrt{nh^d} + h^2$ is convex w.r.t. h and the global minimum is at $h = n^{-1/(d+4)}$, the optimal h minimizing the error is $h = O(n^{-1/(d+4)})$. Then, $O(nh^d) = O(n^{4/(d+4)})$, implying that $nh^d \rightarrow \infty$. \square

Corollary 2 (Restated Corol. 2). *Let $h = O(n^{-1/(d+4)})$. If nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ converge at $\{nh^d\}^{-1/4}$ rate for any (w, z, x, y) , then the target estimator $\hat{\psi}_h(y)$ achieves $\sqrt{nh^d}$ -rate convergence to ψ .*

Proof. It suffices to show that B_y converges at $\sqrt{nh^d}$ -rate with the choice of h as in Lemma 5, since the rest is guaranteed by Corol. 1. We first note that $B_y = O(h^2)$. Since $O(nh^d) = O(n^{4/(d+4)})$, we have $O(1/\sqrt{nh^d}) = O(n^{-2/(d+4)}) = O(h^2)$. \square

Lemma B.6 (Restated Lemma 6). *Suppose D_f is a f -divergence such that $f(p, q) = 0$ if $p = q$. Then,*

$$D_f(\psi, \hat{\psi}_h) \leq \int_{\mathcal{Y}} w(y) \left(\hat{\psi}_h(y) - \psi(y) \right) d[y],$$

where $w(y) \equiv f'_2(\psi(y), \tilde{\psi}(y))\hat{\psi}_h(y)$, $f'_2(p, q) \equiv (\partial/\partial q)f(p, q)$, and $\tilde{\psi}_h(y) \equiv t\hat{\psi}_h(y) + (1-t)\psi(y)$ for some fixed $t \in [0, 1]$.

Proof. For $f(p, q)$, by applying Taylor's expansion, we have

$$f(p, q) = f(p, p) + f'_2(p, \tilde{p})(q - p),$$

for some fixed $\tilde{p} \in [p, q]$. Applying this idea,

$$\begin{aligned} D_f(\psi, \hat{\psi}_h) &= \int_{\mathcal{Y}} f(\psi(y), \hat{\psi}_h(y))\hat{\psi}_h(y) d[y] \\ &= \int_{\mathcal{Y}} \left\{ \underbrace{f(\psi(y), \psi(y))}_{=0} + f'_2(\psi(y), \tilde{\psi}(y)) \left(\hat{\psi}_h(y) - \psi(y) \right) \right\} \hat{\psi}_h(y) d[y], \\ &= \int_{\mathcal{Y}} w(y) \left(\hat{\psi}_h(y) - \psi(y) \right) d[y], \end{aligned}$$

where the second equality holds by Taylor's expansion on f , and the third equality by the given assumption that $f(p, q) = 0$ whenever $p = q$. \square

Theorem B.2 (Restated Thm. 2). *Suppose the estimators for nuisances are consistent; i.e., $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \eta = \{\pi, \xi, \theta\}$ for all (w, z, x, y) . Suppose D_f is a f -divergence such that $f(p, q) = 0$ if $p = q$. Suppose $w(y)$ in Lemma 6 is finite. Then,*

$$D_f(\psi, \hat{\psi}_h) \leq O_P \left(\sup_{y \in \mathcal{Y}} \{R_2^k + B_y\} + 1/\sqrt{nh^d} + 1/\sqrt{n} \right), \quad (\text{B.14})$$

where R_2^k is defined in Eq. (12) and B_y is defined in Eq. (13).

Proof. Under the given conditions, with Thm. 1,

$$\begin{aligned} D_f(\psi, \hat{\psi}_h) &\leq \int_{\mathcal{Y}} w(y) \left(\hat{\psi}_h(y) - \psi(y) \right) d[y] \\ &= \int_{\mathcal{Y}} w(y) \left(O_P \left(1/\sqrt{nh^d} + R_2^k + 1/\sqrt{n} \right) + B_y \right) d[y] \\ &= O_P(1/\sqrt{nh^d} + 1/\sqrt{n}) + \int_{\mathcal{Y}} (w(y)O_P(R_2^k) + B_y) d[y] \\ &= O_P(1/\sqrt{nh^d} + 1/\sqrt{n}) + O_P \left(\sup_{y \in \mathcal{Y}} \{R_2^k + B_y\} \right). \end{aligned}$$

\square

Corollary 3 (Restated Corol. 3). *Let $h = O(n^{-1/(d+4)})$. Suppose D_f satisfies $f(p, q) = 0$ if $p = q$. Suppose $w(y)$ in Lemma 6 is finite. If nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ converges at $\{nh^d\}^{-1/4}$ rate for any (w, z, x, y) , then $D_f(\psi, \hat{\psi}_h)$ converges to 0 at $\sqrt{nh^d}$ -rate.*

Proof. This result follows immediately from Thm. 2. \square

B.2 Proofs for Sec. 4

We will use ψ_p to denote ψ as a functional for p . Let p_ϵ denote a parametric submodel. We will use S_ϵ to denote a score function for p_ϵ .

Lemma B.7 (Restated Lemma 7). *An influence function for $m(\beta; \psi)$ in Eq. (17), denoted ϕ_m , is given by*

$$\phi_m(\beta; \eta = \{\pi, \xi, \theta\}, \psi) \equiv \phi(\eta, \psi)[R_f(Y; \beta, \psi)], \quad (\text{B.15})$$

where $\phi(\eta, \psi)[\cdot]$ is defined in Eq. (7), and

$$R_f(Y; \beta, \psi) \equiv g'(Y; \beta) \{f_{21}''(\psi(Y), g(Y; \beta))g(Y; \beta) + f_1'(\psi(Y), g(Y; \beta))\},$$

where $g'(y; \beta) \equiv (\partial/\partial\beta)g(y; \beta)$, $f_1'(p, q) \equiv (\partial/\partial p)f(p, q)$ and $f_{21}''(p, q) \equiv (\partial/\partial p)f_2''(p, q)$.

Proof. Let ψ_ϵ denote the estimand ψ written w.r.t. the parametric submodel $p_\epsilon = p(1 + \epsilon g)$ where g is a bounded mean-zero random function. Let $S_\epsilon \equiv ((\partial/\partial\epsilon)|_{\epsilon=0} \log p_\epsilon)$.

Let

$$\bar{m}(y; \beta, \psi) \equiv g'(y; \beta) \{f_2'(\psi(y), g(y; \beta))g(y; \beta) + f(\psi(y), g(y; \beta))\}. \quad (\text{B.16})$$

Note $m(\beta, \psi) = \int_{\mathcal{Y}} \bar{m}(y; \beta, \psi) d[y]$. Also, we note that $(\partial/\partial\psi)\bar{m}(y; \beta, \psi) = R_f(y; \beta, \psi)$.

Also, recall that an influence function for $\psi[f(Y)]$ (for $f(Y) < \infty$) is given as $\phi(\eta, \psi)[f(Y)]$ in Lemma S.3. Then, by the definition of the influence function, $\psi[f(Y)]$ satisfies the following,

$$(\partial/\partial\epsilon)|_{\epsilon=0}\psi_\epsilon[f(Y)] = \mathbb{E}_P[\phi(\psi, \eta)[f(Y)] \cdot S_\epsilon].$$

Now, we will prove that $\phi_m(\beta; \eta = \{\pi, \xi, \theta\}, \psi) \equiv \phi(\eta, \psi)[R_f(Y; \beta, \psi)]$ is a functional satisfying

$$(\partial/\partial\epsilon)|_{\epsilon=0}m(\beta, \psi) = \mathbb{E}_P[\phi(\psi, \eta)[R_f(Y; \beta, \psi)] \cdot S_\epsilon],$$

then this equation implies that $\phi(\eta, \psi)[R_f(Y; \beta, \psi)]$ is an influence function for the score $m(\beta, \psi)$.

This can be shown as follows:

$$\begin{aligned} & (\partial/\partial\epsilon)|_{\epsilon=0}m(\beta, \psi) \\ &= (\partial/\partial\epsilon)|_{\epsilon=0} \int_{\mathcal{Y}} \bar{m}(y; \beta, \psi) d[y] \\ &= \int_{\mathcal{Y}} (\partial/\partial\epsilon)|_{\epsilon=0}\bar{m}(y; \beta, \psi) d[y] \\ &= \int_{\mathcal{Y}} (\partial/\partial\epsilon)|_{\epsilon=0}\psi_\epsilon(y)(\partial/\partial\psi'(y))|_{\psi'=\psi}\bar{m}(y; \beta, \psi_\epsilon) d[y] \\ &= (\partial/\partial\epsilon)|_{\epsilon=0} \int_{\mathcal{Y}} \psi_\epsilon(y)R_f(y; \beta, \psi) d[y] \\ &= (\partial/\partial\epsilon)|_{\epsilon=0}\psi_\epsilon[R_f(Y; \beta, \psi)] \\ &= \mathbb{E}_P[\phi(\psi, \eta)[R_f(Y; \beta, \psi)] \cdot S_\epsilon], \end{aligned}$$

where the first equality is by the definition of \bar{m} , the second by the exchange of derivation/integration, the third by the chain rule, the fourth by the fact that $(\partial/\partial\psi)\bar{m}(y; \beta, \psi) = R_f(y; \beta, \psi)$ and the exchange of derivation/integration, the fifth by the definition of $\psi[f(Y)]$ in Eq. (7), the sixth by the definition of the influence function (i.e., the influence function for $\psi[f(Y)]$ is a function $\phi[f(Y)]$ satisfying $(\partial/\partial\epsilon)|_{\epsilon=0}\psi_\epsilon[f(Y)] = \mathbb{E}_P[\phi[f(Y)] \cdot S_\epsilon]$.

□

Lemma B.8 ((Restated Lemma 8)). *A Neyman orthogonal score for estimating β , denoted $\varphi(\beta'; (\eta = \{\pi, \xi, \theta\}, \psi))$, is given by*

$$\varphi(\beta'; (\eta = \{\pi, \xi, \theta\}, \psi)) \equiv m(\beta', \psi) + \phi_m(\beta, \eta, \psi), \quad (\text{B.17})$$

where $\phi_m(\beta, \eta, \psi)$ is defined in Eq. (18).

Proof. We first note that $\mathbb{E}_P [m(\beta', \psi)] = m(\beta', \psi)$, because this is not a random function. Then, the influence function for $\mathbb{E}_P [m(\beta', \psi)]$ is given by Lemma 7. For any score function which expectation is zero at the true parameter, its addition with the influence function is a Neyman orthogonal score [13, Thm.1]. That is, $m(\beta', \psi) + \phi_m(\beta, \eta, \psi)$ is a Neyman orthogonal score. \square

Theorem B.3 ((Restated Thm. 3)). *Let $\varphi(\beta'; (\eta = \{\pi, \xi, \theta\}, \psi))$ be given in Eq. (19). Let $\phi_m(\beta, \eta, \psi)$ be given in Eq. (18). Let β_0, η_0, ψ_0 denote the true parameters. Let $\hat{\beta}$ be the MLTE estimator for β defined in Def. 3. Suppose (1) $R_f(y; \beta, \psi)$ is bounded and $R'_f(y; \beta, \psi) \equiv (\partial/\partial\psi)R_f(y; \beta, \psi) < \infty$; (2) There exists a function $H(y) < \infty$ s.t. $\sup_{\beta, \psi} \max\{R_f(y; \beta, \psi), R'_f(y; \beta, \psi)\} = O(H(y))$; (3) $\{\varphi(\beta; (\eta, \psi))\}$ is Donsker⁵ w.r.t. β for the fixed η ; (3) The estimators are consistent: $\hat{\beta} - \beta_0 = o_P(1)$ and $\|\nu - \hat{\nu}\| = o_P(1)$ for $\nu \in \{\pi_z(w), \xi_x(z, w), \theta(x, z, w)[H(Y)]\}$ for all (w, z, x, y) ; and (4) $\mathbb{E}_P [\varphi(\beta; (\eta, \psi))]$ is differentiable w.r.t. β at $\beta = \beta_0$ with non-singular matrix $M(\beta_0, (\eta, \psi)) \equiv (\partial/\partial\beta)|_{\beta=\beta_0} \mathbb{E}_P [\varphi(\beta; (\eta, \psi))]$ for all (η, ψ) , where $M(\beta_0, (\hat{\eta}, \hat{\psi})) \xrightarrow{P} M \equiv M(\beta_0, (\eta_0, \psi_0))$. Then,*

$$\hat{\beta} - \beta_0 = -M^{-1} \mathbb{E}_{\mathcal{D}} [\phi_m(\beta_0; (\psi_0, \eta_0))] + o_P(n^{-1/2}) + O_P(R_2^m),$$

where

$$R_2^m = \sum_z \left(\|\hat{\pi}_z - \pi_z\| \left\{ \|\hat{\theta}_z - \theta_z\| + \|\hat{\xi}_z - \xi_z\| \right\} + \|\hat{\xi}_z - \xi_z\|^2 + \|\theta_z - \hat{\theta}_z\|^2 + \|\hat{\xi}_z - \xi_z\| \|\theta_z - \hat{\theta}_z\| \right),$$

where $\pi_z \equiv \pi_z(W)$, $\xi_z \equiv \xi_x(z, W)$, and $\theta_z \equiv \theta(x, z, W)[H(Y)]$.

Proof. We follow the proof strategy used in [35, Lemma 1, Thm.3]. First,

$$\begin{aligned} \hat{\beta} - \beta_0 &= -M^{-1} \mathbb{E}_{\mathcal{D}} [\varphi(\beta_0, (\psi_0, \eta_0))] - M^{-1} \mathbb{E}_P [\varphi(\beta_0, (\hat{\psi}, \hat{\eta}))] + o_P(n^{-1/2}) \\ &= -M^{-1} \mathbb{E}_{\mathcal{D}} [\phi_m(\beta_0, \{\psi_0, \eta_0\})] - M^{-1} \mathbb{E}_P [\varphi(\beta_0, (\hat{\psi}, \hat{\eta}))] + o_P(n^{-1/2}), \end{aligned} \quad (\text{B.18})$$

where the first equality holds by Lemma S.1, and the second holds since $m(\beta_0, \psi_0) = 0$ by the moment condition in Eq. (17). Since $\mathbb{E}_{\mathcal{D}} [\phi_m(\beta_0, \eta_0, \psi_0)]$ converges to $N(0, \text{var}(\phi_m^2))$ in distribution at \sqrt{n} -rate, the only remaining term to analyze is

$$\mathbb{E}_P [\varphi(\beta_0, (\hat{\psi}, \hat{\eta}))] = m(\beta_0, \hat{\psi}) + \mathbb{E}_P [\phi(\beta_0, (\hat{\psi}, \hat{\eta})) [R_f(Y; \beta_0, \hat{\psi})]], \quad (\text{B.19})$$

which can be analyzed as

$$\begin{aligned} &\mathbb{E}_P [\phi(\beta_0, (\hat{\psi}, \hat{\eta})) [R_f(Y; \beta_0)]] \\ &= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left\{ \hat{\mathcal{V}}_{YX} [R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi} [R_f(Y; \beta_0, \hat{\psi})] \hat{\mathcal{V}}_X \right\} \right] \\ &= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\mathcal{V}}_{YX} [R_f(Y; \beta_0, \hat{\psi})] \right] - \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi} [R_f(Y; \beta_0, \hat{\psi})] \hat{\mathcal{V}}_X \right] \\ &= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left\{ \frac{\hat{\pi}_{z^x}(W)}{\pi_{z^x}(W)} \left\{ \theta(x, z^x, W) [R_f(Y; \beta_0, \hat{\psi})] - \hat{\theta}(x, z^x, W) [R_f(Y; \beta_0, \hat{\psi})] \right\} + \hat{\theta}(x, z^x, W) R_f(Y; \beta_0, \hat{\psi}) \right\} \right] \end{aligned} \quad (\text{B.20})$$

$$- \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left\{ \frac{\hat{\pi}_{z^{1-x}}(W)}{\pi_{z^{1-x}}(W)} \left\{ \theta(x, z^{1-x}, W) [R_f(Y; \beta_0, \hat{\psi})] - \hat{\theta}(x, z^{1-x}, W) [R_f(Y; \beta_0, \hat{\psi})] \right\} + \hat{\theta}(x, z^{1-x}, W) [R_f(Y; \beta_0, \hat{\psi})] \right\} \right] \quad (\text{B.21})$$

$$- \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi} [R_f(Y; \beta_0, \hat{\psi})] \left\{ \frac{\pi_{z^x}(W)}{\hat{\pi}_{z^x}(W)} \left\{ \xi_x(z^x, W) - \hat{\xi}_x(z^x, W) \right\} + \hat{\xi}_x(z^x, W) \right\} \right] \quad (\text{B.22})$$

$$+ \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi} [R_f(Y; \beta_0, \hat{\psi})] \left\{ \frac{\pi_{z^{1-x}}(W)}{\hat{\pi}_{z^{1-x}}(W)} \left\{ \xi_x(z^{1-x}, W) - \hat{\xi}_x(z^{1-x}, W) \right\} + \hat{\xi}_x(z^{1-x}, W) \right\} \right], \quad (\text{B.23})$$

⁵A function class where complexities are restricted. Refer [58, Page 269] for the definition. Donsker class include Sobolev, Bounded monotone, Lipschitz class, etc.

where

$$(B.20) = \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \cdot \left\{ \left(\frac{\hat{\pi}_{z^x}(W)}{\pi_{z^x}(W)} - 1 \right) \left\{ \theta(x, z^x, W)[R_f(Y; \beta_0, \hat{\psi})] - \hat{\theta}(x, z^x, W)[R_f(Y; \beta_0, \hat{\psi})] \right\} \right\} \right] \quad (B.24)$$

$$+ \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \theta(x, z^x, W)[R_f(Y; \beta_0, \hat{\psi})] \right] \quad (B.25)$$

$$(B.21) = -\mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left\{ \left(\frac{\hat{\pi}_{z^{1-x}}(W)}{\pi_{z^{1-x}}(W)} - 1 \right) \left\{ \theta(x, z^{1-x}, W)[R_f(Y; \beta_0, \hat{\psi})] - \hat{\theta}(x, z^{1-x}, W)[R_f(Y; \beta_0, \hat{\psi})] \right\} \right\} \right] \quad (B.26)$$

$$- \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \theta(x, z^{1-x}, W)[R_f(Y; \beta_0, \hat{\psi})] \right] \quad (B.27)$$

$$(B.22) = -\mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \left\{ \left(\frac{\pi_{z^x}(W)}{\hat{\pi}_{z^x}(W)} - 1 \right) \left\{ \xi_x(z^x, W) - \hat{\xi}_x(z^x, W) \right\} \right\} \right] \quad (B.28)$$

$$- \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \xi_x(z^x, W) \right] \quad (B.29)$$

$$(B.23) = \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \left\{ \left(\frac{\pi_{z^{1-x}}(W)}{\hat{\pi}_{z^{1-x}}(W)} - 1 \right) \left\{ \xi_x(z^{1-x}, W) - \hat{\xi}_x(z^{1-x}, W) \right\} \right\} \right] \quad (B.30)$$

$$+ \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \xi_x(z^{1-x}, W) \right] \quad (B.31)$$

First, consider the summation of (B.25,B.27,B.29,B.31):

Eq. (B.25) + Eq. (B.27) + Eq. (B.29) + Eq. (B.31)

$$\begin{aligned} &= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left\{ \theta(x, z^x, W)[R_f(Y; \beta_0, \hat{\psi})] - \theta(x, z^{1-x}, W)[R_f(Y; \beta_0, \hat{\psi})] \right\} \right] \\ &- \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \left\{ \xi_x(z^x, W) - \xi_x(z^{1-x}, W) \right\} \right] \\ &= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left(\psi_{YX}[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \cdot \psi_X \right) \right] \\ &= \mathbb{E}_P \left[\frac{1}{\hat{\psi}_X} \left(\psi_{YX}[R_f(Y; \beta_0, \hat{\psi})] - \frac{\hat{\psi}_{YX}[[R_f(Y; \beta_0, \hat{\psi})]]}{\hat{\psi}_X} \cdot \psi_X \right) \right] \\ &= \mathbb{E}_P \left[\frac{\psi_X}{\hat{\psi}_X} \left(\frac{\psi_{YX}[R_f(Y; \beta_0, \hat{\psi})]}{\psi_X} - \frac{\hat{\psi}_{YX}[R_f(Y; \beta_0, \hat{\psi})]}{\hat{\psi}_X} \right) \right] \\ &= \mathbb{E}_P \left[\frac{\psi_X}{\hat{\psi}_X} \left(\psi[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \right) \right]. \\ &= \mathbb{E}_P \left[\left\{ \frac{\psi_X}{\hat{\psi}_X} - 1 \right\} \left(\psi[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \right) \right] + \mathbb{E}_P \left[\left(\psi[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \right) \right]. \end{aligned} \quad (B.32)$$

Then,

Eq. (B.19) = $m(\beta_0, \hat{\psi})$ + Sum of (B.25, B.27, B.29, B.31) + Sum of (B.24, B.26, B.28, B.30)

$$= m(\beta_0, \hat{\psi}) + \mathbb{E}_P \left[\left(\psi[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \right) \right] \quad (B.33)$$

$$+ \mathbb{E}_P \left[\left\{ \frac{\psi_X}{\hat{\psi}_X} - 1 \right\} \left(\psi[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \right) \right] + \text{Sum of (B.24,B.26,B.28,B.30)}. \quad (B.34)$$

To analyze (B.33), we recall that $(\partial/\partial\psi)m(\beta_0, \psi) = \int_{\mathcal{Y}} R_f(y; \beta_0, \psi) d[y]$ and $m(\beta_0, \psi) = 0$. Also, by Taylor's expansion to $\bar{m}(y; \beta, \psi)$ defined in Eq. (B.16),

$$\bar{m}(y; \beta_0, \psi) = \bar{m}(y; \beta_0, \hat{\psi}) + R_f(y; \beta, \hat{\psi})(\psi(y) - \hat{\psi}(y)) + R_f^{(1)}(y; \beta, \tilde{\psi})(\psi(y) - \hat{\psi}(y))^2,$$

where $R_f^{(1)}$ is a first derivative of R_f w.r.t. ψ . This implies that

$$0 = m(\beta_0, \psi) = m(\beta_0, \hat{\psi}) + \int_{\mathcal{Y}} R_f(y; \beta, \hat{\psi}) (\psi(y) - \hat{\psi}(y)) d[y] + \int_{\mathcal{Y}} R_f^{(1)}(y; \beta, \tilde{\psi}) (\psi(y) - \hat{\psi}(y))^2 d[y],$$

where $\tilde{\psi}$ is some unknown estimand within the interval $[\psi, \hat{\psi}]$. We obtain

$$- \int_{\mathcal{Y}} R_f^{(1)}(y; \beta, \tilde{\psi}) (\psi(y) - \hat{\psi}(y))^2 d[y] = m(\beta_0, \hat{\psi}) + \int_{\mathcal{Y}} R_f(y; \beta, \hat{\psi}) (\psi(y) - \hat{\psi}(y)) d[y].$$

By taking expectations for both sides,

$$-\mathbb{E}_P \left[\int_{\mathcal{Y}} R_f^{(1)}(y; \beta, \tilde{\psi}) (\psi(y) - \hat{\psi}(y))^2 d[y] \right] = m(\beta_0, \hat{\psi}) + \mathbb{E}_P \left[\int_{\mathcal{Y}} R_f(y; \beta, \hat{\psi}) (\psi(y) - \hat{\psi}(y)) d[y] \right]. \quad (\text{B.35})$$

We have

$$\begin{aligned} - \int_{\mathcal{Y}} R_f^{(1)}(y; \beta, \tilde{\psi}) (\psi(y) - \hat{\psi}(y))^2 d[y] &= O \left(\int_{\mathcal{Y}} R_f^{(1)}(y; \beta, \tilde{\psi}) (\psi(y) - \hat{\psi}(y))^2 d[y] \right) \\ &= O \left(\int_{\mathcal{Y}} H(y) (\psi(y) - \hat{\psi}(y))^2 d[y] \right) \\ &= O \left(\int_{\mathcal{Y}} H^2(y) (\psi(y) - \hat{\psi}(y))^2 d[y] \right) \\ &= O \left(\left\| \psi[H(Y)] - \hat{\psi}[H(Y)] \right\|^2 \right), \end{aligned}$$

where the second equality is by the definition of $H(y)$, the third by $H(y) < \infty$, and the fourth by the definition of L_2 norm.

This implies that

$$(\text{B.33}) = -\mathbb{E}_P \left[\int_{\mathcal{Y}} R_f^{(1)}(y; \beta, \tilde{\psi}) (\psi - \hat{\psi})^2 d[y] \right] = O \left(\left\| \psi[H(Y)] - \hat{\psi}[H(Y)] \right\|^2 \right),$$

where the first equality is by Eq. (B.35) and the second equality is by the above.

Also, Sum of (B.24,B.26,B.28,B.30) in (B.34) can be written as follows:

Sum of (B.24,B.26,B.28,B.30)

$$\begin{aligned} &= \sum_{z \in \{0,1\}} O_P \left(\|\hat{\pi}_z(W) - \pi_z(W)\| \left\{ \left\| \hat{\theta}(x, z, W)[R_f(Y; \beta_0, \hat{\psi})] - \theta(x, z, W)[R_f(Y; \beta_0, \hat{\psi})] \right\| + \left\| \hat{\xi}_x(z, W) - \xi_x(z, W) \right\| \right\} \right) \\ &= \sum_{z \in \{0,1\}} O_P \left(\|\hat{\pi}_z(W) - \pi_z(W)\| \left\{ \left\| \hat{\theta}(x, z, W)[H(Y)] - \theta(x, z, W)[H(Y)] \right\| + \left\| \hat{\xi}_x(z, W) - \xi_x(z, W) \right\| \right\} \right). \end{aligned}$$

For simplicity, we assume, for any x, z ,

$$\begin{aligned} O_P \left(\left\{ \xi_x(z, W) - \hat{\xi}_x(z, W) \right\} \cdot \left\{ \xi_x(1-z, W) - \hat{\xi}_x(1-z, W) \right\} \right) &= \sum_{z \in \{0,1\}} O_P \left(\left\| \xi_x(z, W) - \hat{\xi}_x(z, W) \right\|^2 \right), \text{ and} \\ O_P \left(\left\| \xi_x(z, W) - \hat{\xi}_x(z, W) \right\| \left\| \theta(x, 1-z, W)[H(Y)] - \hat{\theta}(x, 1-z, W)[H(Y)] \right\| \right) \\ &= \sum_{z \in \{0,1\}} O_P \left(\left\| \xi_x(z, W) - \hat{\xi}_x(z, W) \right\| \left\| \theta(x, z, W)[H(Y)] - \hat{\theta}(x, z, W)[H(Y)] \right\| \right). \end{aligned}$$

The other part of Eq. (B.34) is given as

$$\begin{aligned}
& \mathbb{E}_P \left[\left\{ \frac{\psi^X}{\hat{\psi}^X} - 1 \right\} \left(\psi[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \right) \right] \\
&= O_P \left(\left\| \psi^X - \hat{\psi}^X \right\| \left\| \psi[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}[R_f(Y; \beta_0, \hat{\psi})] \right\| \right) \\
&= O_P \left(\left\| \psi^X - \hat{\psi}^X \right\| \left\| \frac{\psi^{Y^X}[R_f(Y; \beta_0, \hat{\psi})]}{\psi^X} - \frac{\hat{\psi}^{Y^X}[R_f(Y; \beta_0, \hat{\psi})]}{\hat{\psi}^X} + \frac{\hat{\psi}^{Y^X}[R_f(Y; \beta_0, \hat{\psi})]}{\psi^X} - \frac{\hat{\psi}^{Y^X}[R_f(Y; \beta_0, \hat{\psi})]}{\hat{\psi}^X} \right\| \right) \\
&= O_P \left(\left\| \psi^X - \hat{\psi}^X \right\| \left(\left\| \psi^{Y^X}[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}^{Y^X}[R_f(Y; \beta_0, \hat{\psi})] \right\| + \left\| \frac{1}{\psi^X} - \frac{1}{\hat{\psi}^X} \right\| \right) \right) \\
&= O_P \left(\left\| \psi^X - \hat{\psi}^X \right\| \left(\left\| \psi^{Y^X}[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}^{Y^X}[R_f(Y; \beta_0, \hat{\psi})] \right\| + \left\| \psi^X - \hat{\psi}^X \right\| \right) \right) \\
&= O_P \left(\left\| \psi^X - \hat{\psi}^X \right\|^2 \right) + O_P \left(\left\| \psi^X - \hat{\psi}^X \right\| \left\| \psi^{Y^X}[R_f(Y; \beta_0, \hat{\psi})] - \hat{\psi}^{Y^X}[R_f(Y; \beta_0, \hat{\psi})] \right\| \right) \\
&= O_P \left(\left\| \psi^X - \hat{\psi}^X \right\|^2 \right) + O_P \left(\left\| \psi^X - \hat{\psi}^X \right\| \left\| \psi^{Y^X}[H(Y)] - \hat{\psi}^{Y^X}[H(Y)] \right\| \right) \\
&= \sum_z O_P \left(\left\| \hat{\xi}_x(z, W) - \xi_x(z, W) \right\|^2 + \left\| \hat{\xi}_x(z, W) - \xi_x(z, W) \right\| \left\| \theta(x, z, W)[H(Y)] - \hat{\theta}(x, z, W)[H(Y)] \right\| \right),
\end{aligned}$$

where the equalities can be shown using the standard computation and the positivity assumption.

Similarly we assume, for any x, z ,

$$\begin{aligned}
& O_P \left(\left\| \theta(x, z, W)[H(Y)] - \hat{\theta}(x, z, W)[H(Y)] \right\| \left\| \theta(x, 1-z, W)[H(Y)] - \hat{\theta}(x, 1-z, W)[H(Y)] \right\| \right) \\
&= \sum_{z \in \{0,1\}} O_P \left(\left\| \theta(x, z, W)[H(Y)] - \hat{\theta}(x, z, W)[H(Y)] \right\|^2 \right).
\end{aligned}$$

We have

$$\begin{aligned}
& O_P \left(\left\| \hat{\psi}[H(Y)] - \psi[H(Y)] \right\|^2 \right) \\
&= O_P \left(\left\| \psi^{\hat{Y}^X}[H(Y)] - \psi^{Y^X}[H(Y)] + \hat{\psi}^X - \psi^X \right\|^2 \right) \\
&= O_P \left(\left\| \psi^{\hat{Y}^X}[H(Y)] - \psi^{Y^X}[H(Y)] \right\|^2 + \left\| \hat{\psi}^X - \psi^X \right\|^2 + \left\| \psi^{\hat{Y}^X}[H(Y)] - \psi^{Y^X}[H(Y)] \right\| \left\| \hat{\psi}^X - \psi^X \right\| \right) \\
&= \sum_{z \in \{0,1\}} O_P \left(\left\| \theta(x, z, W)[H(Y)] - \hat{\theta}(x, z, W)[H(Y)] \right\|^2 \right) + \sum_{z \in \{0,1\}} O_P \left(\left\| \xi_x(z, W) - \hat{\xi}_x(z, W) \right\|^2 \right) \\
&+ \sum_{z \in \{0,1\}} O_P \left(\left\| \theta(x, z, W)[H(Y)] - \hat{\theta}(x, z, W)[H(Y)] \right\| \left\| \xi_x(z, W) - \hat{\xi}_x(z, W) \right\| \right).
\end{aligned}$$

Finally

$$\begin{aligned}
\text{Eq. (B.19)} &= \sum_z O_P \left(\left\| \hat{\pi}_z(W) - \pi_z(W) \right\| \left\{ \left\| \hat{\theta}(x, z, W)[H(Y)] - \theta(x, z, W)[H(Y)] \right\| + \left\| \hat{\xi}_x(z, W) - \xi_x(z, W) \right\| \right\} \right) \\
&+ \sum_z O_P \left(\left\| \hat{\xi}_x(z, W) - \xi_x(z, W) \right\|^2 + \left\| \theta(x, z, W) - \hat{\theta}(x, z, W) \right\|^2 \right) \\
&+ \sum_z O_P \left(\left\| \hat{\xi}_x(z, W) - \xi_x(z, W) \right\| \left\| \theta(x, z, W) - \hat{\theta}(x, z, W) \right\| \right). \tag{B.36}
\end{aligned}$$

Therefore, with Eq. (B.18), the following holds

$$\hat{\beta} - \beta_0 = -M^{-1} \mathbb{E}_{\mathcal{D}} [\phi_m(\mathbf{V}; \beta_0, \psi_0, \eta_0)] + \text{Eq. (B.36)} + o_P(n^{-1/2}),$$

where Eq. (B.36) = R_2^m .

□

Corollary 4 (Restated Corol. 4). *If nuisances $\{\hat{\pi}, \hat{\xi}, \hat{\theta}\}$ converges at $n^{-1/4}$ rate, then the target estimator $\hat{\beta}$ converges to β_0 at \sqrt{n} -rate.*

Proof. If all nuisances converge at $n^{-1/4}$ rate, then the R_2^m term in Thm. 3 converges at $n^{-1/2}$ rate. Also, $\mathbb{E}_{\mathcal{D}}[\phi_m(\beta_0; (\psi_0, \eta_0))]$ converges in distribution to $N(0, \text{var}(\phi_m(\beta_0, (\psi_0, \eta_0))))$ at \sqrt{n} -rate. So $\hat{\beta}$ converges to β_0 at \sqrt{n} -rate by Thm. 3. \square

C Details of empirical applications

C.1 Data generating processes for synthetic datasets

The following structural equations are used for all four data generating processes in Fig. 2:

$$\begin{aligned} U &\sim N(0, 1) \\ f_W(U) &= 2U - 1 + \epsilon_W, \text{ where } \epsilon_W \sim N(0, 1) \\ f_Z(W) &= \mathbb{1}(0.25W + \epsilon_Z > 0), \text{ where } \epsilon_Z \sim N(0, 1) \\ f_X(W, Z, U) &= \mathbb{1}(Z + 0.25 * W + 0.25 * U + \epsilon_X > 0.5) \cdot (1 - Z) + Z, \text{ where } \epsilon_X \sim N(0, 1). \end{aligned}$$

With such data generating process, $X_{Z=1} \geq X_{Z=0}$ is satisfied. We will denote four figures in Fig. 2 as Fig. 2(a,b,c,d). For Fig. 2a,

$$f_Y(W, X, U) = 0.6501(W \cdot (2X - 1) + 2U + 0.374).$$

For Fig. 2b,

$$f_Y(W, X, U) = 0.9515(2X - 1 + W) + 0.8(-2X + 1 + U) + WU + 0.082.$$

For Fig. 2c,

$$\begin{aligned} f_Y(W, X, U) &= 1.0854\mathbb{1}(W < 0)(2X - 1 + 0.1U) + \mathbb{1}(0 \leq W < 1)(-2X + 1 + 0.1U) \\ &\quad + 1.0854 \cdot 0.9163(\mathbb{1}(W \geq 1)(-3(2X - 1) + 0.2U + 0.3) - 0.122) \end{aligned}$$

For Fig. 2d,

$$\begin{aligned} f_Y(W, X, U) &= 0.7865 \cdot 1.0628 \cdot \mathbb{1}(W < -1)(-0.8(2X - 1) + 0.1U) + \mathbb{1}(-1 \leq W < 0)(-2(2X - 1) + 0.1U) \\ &\quad + 0.7865 \cdot 1.0628 \cdot (\mathbb{1}(0 \leq W < 1)(2(2X - 1) + 0.2U) + \mathbb{1}(W > 1)(0.5(2X - 1) + 0.2U) + 0.0525) \\ &\quad + 1.0628 \cdot 0.104 \end{aligned}$$

C.2 Application to 401(k) data

We use the 401(k) dataset that is initially introduced by [2]. Specifically, we used the version of the data named ‘The Woodridge Data Set [61]’ originally entitled ‘401ksu.dta’ in STATA format (available in <https://www.stata.com/texts/eacsap/>). In the dataset, we used `nettfa` (net financial asset in \$1000) as Y , `p401k` (participation in 401(k), participation = 1) as X , `e401k` (eligibility for 401(k), eligible = 1) as Z , and $W = \{W_1, W_2, W_3, W_4, W_5\} = \{\text{agesq}, \text{fsize}, \text{male}, \text{marr}, \text{incsq}\}$, where `agesq` means the square of the age, `fsize` the family size, `male` the gender (male = 1), `marr` the marital status (married = 1) and `incsq` the square of the income.