
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

1 It is common to quantify causal effects with mean values, which, however, may
2 fail to capture significant distribution differences of the outcome under different
3 treatments. We study the problem of estimating the density of the causal effect
4 of a binary treatment on a continuous outcome given a binary instrumental vari-
5 able in the presence of covariates. Specifically, we consider the local treatment
6 effect, which measures the effect of treatment among those who comply with the
7 assignment under the assumption of monotonicity (only the ones who were offered
8 the treatment take it). We develop two families of methods for this task, *kernel-*
9 *smoothing* and *model-based* approximations – the former smoothes the density by
10 convoluting with a smooth kernel function; the latter projects the density onto a
11 finite-dimensional density class. For both approaches, we derive double/debiased
12 machine learning (DML) based estimators. We study the asymptotic convergence
13 rates of the estimators and show that they are robust to the biases in nuisance
14 function estimation. We illustrate the proposed methods on synthetic data and a
15 real dataset called 401(k).

16 1 Introduction

17 Controlled experimentation is one powerful tool used throughout the biological, medical, and social
18 sciences to infer the effect of a certain treatment on a given outcome. The idea is to randomize the
19 treatment assignment so as to neutralize the effect of the unobserved confounders. In some practical
20 settings, however, it may be challenging to ascertain that individuals who are selected for treatment
21 will follow their recommendations. In fact, issues of non-compliance and unmeasured confounders
22 are quite common and lead to the non-identification of treatment effects in such cases [29, 50, 32, 56].

23 An approach known as instrumental variables (IVs) has been proposed to try to circumvent this issue
24 [68]. The idea is to find a variable (or set) that is not the target of the analysis by itself, but that will
25 help to control for the unobserved confounding between the treatment and the outcome. In particular,
26 IVs are special variables that (i) influence the treatment, (ii) do not directly influence the outcome,
27 and (iii) are not affected by unmeasured confounders. For concreteness, consider a study of the
28 effect of 401(k) participation (X) on the distribution of net financial assets (Y) [2]. This setting is
29 represented in the causal graph in Fig. 1. Note that there exists a dashed-bidirected arrow between X
30 and Y , which in graphical language represents unobserved confounding affecting both X and Y . The
31 variable Z in this model represents the eligibility of 401(k). We note that Z qualifies as an instrument
32 in this case – (i) it does affect the participation of 401(k) (X), (ii) has no direct influence on the net
33 financial asset (Y), (iii) is not affected by unmeasured confounders between X and Y . The variable
34 W represents observed covariates (e.g., income, gender, family size, etc.).

35 We are interested in the particular setting where only individuals who were offered the treatment
36 may have access to it [31]. For instance, in the case of 401(k) participation ($X = 1$), only eligible
37 individuals ($Z = 1$) would be allowed to join the program. This assumption is known as *monotonicity*,

Figure 2: Densities of outcomes among compliers under the treatment $x=1$. All densities have a mean μ and a variance σ^2 .

38 which rules out the possibility that any units would respond contrary to the instrument. Under
 39 monotonicity, the causal effect in the subpopulation whose actual treatment coincides with the
 40 assigned treatment (called compliers) is identifiable [31, 2]. The average treatment effect (ATE)
 41 for the compliers is called ‘Local ATE’ (LATE) (or Complier average causal effects, CACE) [31].
 42 The most common quantification of effects in IV settings found in practice is the average (e.g.,
 43 LATE). The average is certainly an informative summary; however, it may fail to capture significant
 44 differences in the causal distributions of the outcome. For instance, consider Fig. 2 that shows the
 45 densities of outcomes Y under treatment $x = 1$ among compliers (generated from samples drawn
 46 from four synthetic data generating processes represented by the IV graph in Fig. 1, as discussed in
 47 Sec. 5). All of the distributions have the same mean and variance σ^2 . However, the difference in the
 48 distributions is self-evident.

49 Most of the prior work on quantifying treatment effects on
 50 outcome distributions focuses on estimating cumulative distribution functions (CDFs) or quantiles, and little attention
 51 has been given to estimating densities of treatment effects (refer to Sec. 1.1 for further comparison). As a comple-
 52 ment to CDFs, densities have the benefits of providing more visually interpretable information of the distribution and en-
 53 abling researchers/practitioners to generate counterfactual samples. One challenge with estimating densities is that
 54 while CDFs are pathwise-differentiable and enjoy rate $n^{-1/2}$ estimators (if n is the size of data), densities are not (i.e., they
 55 are not regular), and therefore possess no influence functions nor $n^{-1/2}$ -rate estimators without approximations [7, Ch. 3].

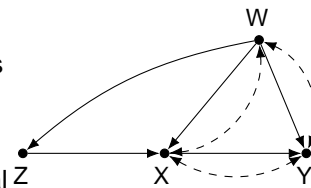


Figure 1: A causal graph for the IV setting. Bidirected arrows encode unmeasured confounders.

62 In this paper, our goal is to provide methods to estimate densities of local treatment effects in IV
 63 settings under the monotonicity assumption. We develop two families of methods for this task based
 64 on kernel-smoothing and model-based approximations. The former smooths the density by convolu-
 65 tion with a smooth kernel function; the latter projects the density onto a finite-dimensional density
 66 class based on a distributional distance measure. For both approaches, we construct double/debiased
 67 machine learning (DML) style density estimators [54, 52, 70, 13]. We analyze the asymptotic
 68 convergence properties of the estimators, showing that they can converge fast (rate) even
 69 when nuisance estimates converge slowly (rate) (debiasedness). We illustrate the
 70 proposed methods on synthetic and real data.

71 1.1 Related work

72 Our work touches different areas, which we will discuss next.

73 DML-based causal effect estimators. The DML framework has been adapted for estimating the
 74 average causal effect under the setting where the exclusion criterion [50, Sec. 3.3.1] (also known
 75 as ignorability [57]) holds (e.g., [2, 19]). Recently, DML-based causal effect estimators have been
 76 developed for any identifiable causal functionals in a given causal graph [33, 34].

77 Local average & quantile treatment effect. The formal identification results for LATE under
 78 the monotonicity assumption in IV settings were developed by [3]. Building on these results,
 79 semiparametric estimation for LATE has received remarkable attention [23, 62, 48], including
 80 robust LATE estimators that achieve debiasedness [40, 38, 64]. As shown in Fig. 2, however,

¹Also known as nonparametric doubly robust [87] or ‘rate doubly robust’ [59].

81 the average is sometimes insufficient to capture the effects of the treatment on the distributions of
 82 outcomes. To address this issue, a problem of estimating quantiles or CDFs has taken attention.

83 A common approach to quantifying LTEs is to estimate quantiles or CDFs, which can be studied
 84 based on the LATE estimation ([2, 15, 24, 16, 30, 45, 18, 69]) using the fact that the expectation
 85 of $1_{Y > y}(Y)$, an indicator that outcome Y falls short of threshold y , reduces to the quantiles (i.e.,
 86 replacing Y in LATE with $1_{Y > y}(Y)$).

87 Non-regular target estimation. Densities are an example of non-regular targets [Chap. 3]. One
 88 can approximate a non-regular target with smooth ones such that inference functions made
 89 estimators can be derived. Two popular used approximation approaches are kernel-smoothing-based
 90 (e.g., [52, 6, 42, 19, 35]) and model-based (e.g., [46, 52, 21, 41, 40, 39]).

91 Causal density estimation. There is limited literature on estimating the density of treatment effects.
 92 Most of the results assume that the ignorability/backdoor admissibility holds [55, 49]. [22] used the
 93 kernel-smoothing technique to estimate the density of a treatment effect [2] and provided a kernel-
 94 smoothing-based density estimator that achieves doubly robustness and debiasedness building on top
 95 of the work in [53]. Recently, [39] investigated a model-based approach and developed estimators that
 96 achieve debiasedness properties. Under the IV setting [10] provides a local polynomial regression-
 97 based density estimator for local treatment effects, and we are not aware of any work studying
 98 debiased density estimators. As mentioned, this paper investigates both kernel-smoothing and model-
 99 based approaches for estimating local treatment effects under IV settings and develops DML-style
 100 density estimators for both.

101 2 LTE Estimation – Problem setup

102 Each variable is represented with a capital letter X and its realized value with a small letter x .
 103 For a discrete (e.g., binary) random variable X , we use $1_x(X)$ to represent the indicator function
 104 such that $1_x(X) = 1$ if $X = x$; $1_x(X) = 0$ otherwise. For a continuous variable X with a
 105 probability density $p(x)$ of a distribution P and a function $f(x)$, $E_P[f(X)] = \int_X f(x)p(x) dx$
 106 where X is the domain for X , and $k_f(X) = \sqrt{E_P[(f(X))^2]}$. f is said to converge to f^* at rate r_n if
 107 $k_f(X) - f^*(x) = O_P(1/r_n)$. For a dataset $D = \{V_i\}_{i=1}^n$, we use $E_D[f(V)] = (1/n) \sum_{i=1}^n f(V_i)$
 108 to denote the empirical mean $\bar{f}(V)$ with D .

109 Structural Causal Models (SCMs). We use the language of SCMs as our basic semantic and
 110 inferential framework [50, 4]. An SCMM is a quadruple $M = \{U; V; P(U); F\}$ where U is a set
 111 of exogenous (latent) variables following a joint distribution $P(u)$, and V is a set of endogenous
 112 (observable) variables whose values are determined by functions $f_{V_i, G_{V_i, 2V}}$ such that V_i
 113 $f_{V_i}(p_{u_i}; u_i)$ where $P_{A_i} \subseteq V$ and $U_i \subseteq U$. Each SCMM induces a distribution $P(v)$ and a causal
 114 graph $G = G(M)$ over V in which there exists a directed edge from every variable U_i to V_i and
 115 dashed-bidirected arrows encode common latent variables (e.g., see Fig. 1). Within the structural
 116 semantics, performing an intervention and setting x is represented through the do-operator,
 117 $do(X = x)$, which encodes the operation of replacing the original equations (i.e., $f_X(p_{u_x}; u_x)$)
 118 by the constant x and induces a submodel M_x and an interventional distribution $P(v|do(x))$. For
 119 any variable $Y \in V$, the potential response $y_x(u)$ is defined as the solution of Y in the submodel
 120 M_x given $U = u$, which induces a counterfactual variable Y_x .

121 Local Treatment Effect (LTE) with IV. We consider the IV setting represented by the causal graph
 122 G in Fig. 1², where Z is a binary instrument with domain $\{0, 1\}$, X is a binary treatment with domain
 123 $\{0, 1\}$, and Y is a (set of) continuous outcomes with bounded domain \mathbb{R}^d , and W is a set of
 124 covariates. G satisfies the IV assumption that Z has no direct influence on outcome Y and is not
 125 affected by unmeasured confounders between X and Y .

126 The causal density $p(y|do(x))$ is not identifiable from the observed density $p(y; x; z; w)$ due to
 127 unobserved confounders between X and Y . However, the effect may be recovered for certain
 128 subpopulation under additional assumptions. Formally, a unit in the population is an always-taker if

²It is common in the literature to define IV assumptions in terms of conditional independences among counterfactuals [1, 9, 8, 2, 60, 47, 64], whose connection with the causal graph in Fig. 1 is discussed in Assumption A.1

129 $X_{Z=1} = X_{Z=0} = 1$, a never-taker if $X_{Z=1} = X_{Z=0} = 0$, a complier if $X_{Z=1} = 1; X_{Z=0} = 0$, and
 130 a defier if $X_{Z=1} = 0; X_{Z=0} = 1$ [3, 2]. We will make the following assumptions after the literature.
 131 Assumption 1 (Monotonicity). There are no defiers $X_{Z=1} = 0; X_{Z=0} = 1$.
 132 Assumption 2 (Positivity). $P(x|z; w) > 0; P(z|w) > 0$ for any $x; z; w$.

Let C denote the event that a unit is a complier. For a given constant a , let $\mathbb{I}(X = a)$ denote the event $X = a$. The LTE $p(y_x|C)$ is identifiable under monotonicity and is given by [1, 2]:

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

133 where the expectation is over W . In this paper, we aim to estimate the LTE density $p(y_x|C)$ in Eq. (1).
 134 We will make the following mild assumption on the target, popularly employed in density estimation
 135 (e.g., [44, 25, 27, 61, 26, 42]).

136 Assumption 3. For any $x; z; w; y$, $p(y|w; z; x)$, $p(y|z; x)$ and $p(y_x|C)$ are bounded, and $p(y_x|C)$ is
 137 twice differentiable.

138 Double/Debiased Machine Learning (DML) method [13] Let \mathbb{P}_0 denote a functional of
 139 an arbitrary distribution P^0 . We use P to denote the true distribution such that $P = P_0$. Let $\theta_0 := \theta_P$
 140 denote the true parameter to be estimated. To estimate θ_0 , DML-based estimators use Neyman
 141 Orthogonal score $(V; \theta_0)$ (where θ_0 is a set of nuisance parameters and θ_P denotes
 142 the true nuisances), a function such that $\mathbb{E}_P [V(\theta_0)] = 0$, $\mathbb{E}_P [V(\theta_P)] = 0$.
 143 Given θ_0 , an DML estimator is constructed using cross-fitting techniques as follows: For randomly
 144 split halves of D denoted $D_0; D_1$, let b_p for $p \in \{0, 1\}$ denote the estimates for D_p . Let
 145 T_p denote a solution such that $\mathbb{E}_{D_p} [V(\theta; T_p; b_p)] = o_P(N^{-1/2})$. Then, $T := (T_0 + T_1)/2$ is an
 146 DML estimator [13, Def. 3.1]. In addition to being consistent, the estimator exhibits a robustness
 147 property called debiasness T converges to θ_0 in the root N rate even when b_p converges to θ_0
 148 in slower $N^{-1/4}$ rate [13, Thm. 3.1]. A Neyman Orthogonal Score can be derived by adding
 149 its influence function [14, Thm. 1]. An influence function of the functional \mathbb{P} is defined as a
 150 solution satisfying $\mathbb{E}_P [I] = 0$, $\mathbb{E}_P [I^2] < \infty$, and $\mathbb{E}_P [I(\theta_P; T_p; b_p)] = 0$
 151 where $T_p = P(v)(1 + t g(v))$ for $t \in \mathbb{R}$ and any bounded mean-zero function g over V , and
 152 $S_t(v; t) = \mathbb{E}_P [I(\theta_P; T_p; b_p)]_{t=0}$ [63, Chap. 25].

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

154 3 Kernel-smoothing-based approach

155 In this section, we develop a kernel-smoothing-based approach for estimating the LTE density.
 156 The kernel-smoothing technique approximates a non-pathwise-differentiable target estimand with
 157 a differentiable estimand by convoluting the density with a kernel function $K(y)$. Properties of
 158 the kernel function includes symmetry about the origin (i.e. $\int_{\mathbb{R}} y K(y) dy = 0$), non-negativity
 159 ($0 < K(y) < \infty; \int_{\mathbb{R}} K(y) dy = 1$) [66, Chap. 4.2].

160 We consider product kernel $K_{h;y}(y^0) = \prod_{j=1}^d K((y_j - y_j^0)/h)$ with given bandwidth $h \in \mathbb{R}$
 161 and a fixed point $y = (y_j)_{j=1}^d \in \mathbb{R}^d$. We assume that the kernel of interest has a bounded
 162 second moment and norm: i.e. $\int_{\mathbb{R}} y^2 K(y) dy < \infty$ and $\int_{\mathbb{R}} K(y) dy < \infty$ following
 163 [27, 61]. Example of kernels include Gaussian kernel $K(u) = (1/\sqrt{2\pi}) \exp(-u^2/2)$ for $u \in \mathbb{R}$,
 164 Epanechnikov kernel $K(u) = (3/4)(1 - u^2)1_{|u| \leq 1}(u)$, Quadratic kernel $K(u) = (15/16)(1 - u^2)^2 1_{|u| \leq 1}(u)$,
 165 Cosine kernel $K(u) = (1/4) \cos(u/2) 1_{|u| \leq 2}(u)$, etc.

For convenience, we denote the target estimand $p(y_x|C)$. In the kernel-smoothing-based approach, we will aim to estimate a kernel-smoothed approximation $\hat{p}_h(y)$ defined as follows:

$$\hat{p}_h(y) = \int_{\mathbb{R}} p(y^0|C) K_{h;y}(y^0) dy^0 = \mathbb{E}[K_{h;y}(Y)]; \quad (2)$$

where $\mathbb{E}[f(Y)]$ is an expectation of a function $f(Y)$ w.r.t. $p(y_x|C)$, which is specified as

$$\mathbb{E}[f(Y)] = \frac{E_P [E_P [f(Y)1_x(X)|z^x; W)] - E_P [f(Y)1_x(X)|z^{1-x}; W]}{E_P [P(x^1|z^1; W) - P(x^1|z^0; W)]}; \quad (3)$$

166 The second equality in Eq. (2) is by Eq. (1). For a target estimator $\hat{\eta}$, we will denote nuisances
 167 by $\eta(z; w) = P(z|w)$, $\eta_x(z; w) = P(x|z; w)$, and $\eta_{f(Y)}(x; z; w) = E_P[f(Y)1_x(X)|z; w]$, shortly
 168 $(\eta; \eta_x; \eta_{f(Y)})$.

We aim to construct a DML estimator for the estimator $\hat{\eta}$. Toward this goal, we will first derive a Neyman orthogonal score for η . Since a Neyman orthogonal score can be constructed based on moment score function (a function of parameters such that its expectation is 0 at the true parameters) [14, Thm. 1], we start by defining the moment score function. Let

$$\eta^X = E_P[\eta^1(z^1; W) - \eta^1(z^0; W)]; \quad (4)$$

$$V_X(f; g) = \frac{1_{z^1}(Z) - 1_{z^0}(Z)}{z(W)} f 1_{x^1}(X) - \eta^1(Z; W)g + \eta^1(z^1; W) - \eta^1(z^0; W); \quad (5)$$

Then, the following is a moment score function for:

$$m(\eta^0; \eta) = \frac{1}{X} (\eta^0 - \eta^1) V_X; \quad (6)$$

169 where η^0 is given in Eq. (2) and η^1 is an estimate of η .

Next, we derive an influence function for the moment score function $m(\eta^0; \eta)$. We first define the following function: for a bounded function $f(Y) < 1$, let

$$\eta^{YX}[f(Y)] = E_P[\eta^X(x; z^X; W)[f(Y)] - \eta^X(x; z^{1-X}; W)[f(Y)]]; \quad (7)$$

$$V_{YX}(f; g)[f(Y)] = \frac{1_{z^X}(Z) - 1_{z^{1-X}}(Z)}{z(W)} f f(Y) 1_x(X) - \eta^X(x; Z; W)[f(Y)]g + \eta^X(x; z^X; W)[f(Y)] - \eta^X(x; z^{1-X}; W)[f(Y)]; \quad (8)$$

and

$$(\eta = f; \eta; g)[f(Y)] = \frac{1}{X} (V_{YX}(f; g)[f(Y)] - [f(Y)]V_X(f; g)); \quad (9)$$

170 where V_X is defined in Eq. (5). Then, the influence function for the expectation of the moment score
 171 function $m(\eta^0; \eta)$ in Eq. (6) is given as follows:

Lemma 1 (Influence function for $m(\eta^0; \eta)$). Let $m(\eta^0; \eta)$ be the score defined in Eq. (6). Then, the influence function for $E_P[m(\eta^0; \eta)]$, denoted m , is given by

$$m(\eta = f; \eta; g) = (\eta; \eta)[K_{\eta, Y}(Y)] \quad (10)$$

172 where η is in Eq. (9).

173 For any score function (e.g., η in Eq. (6)), its addition to the influence function of the expected score
 174 (e.g., m) is a Neyman orthogonal score [14, Thm.1], [13, Sec. 2.2.5]. Specifically,

Lemma 2 (Neyman orthogonal score for η). Let $m(\eta^0; \eta)$ be the score function in Eq. (6), and $m(\eta = f; \eta; g; \eta)$ be the influence function for $E_P[m(\eta^0; \eta)]$ given in Eq. (10). Then, a Neyman orthogonal score for η is given as $(\eta^0; \eta = f; \eta; g) = m(\eta^0; \eta) + m(\eta; \eta)$. Specifically,

$$(\eta^0; \eta = f; \eta; g) = \frac{1}{X} (V_{YX}(f; g)[K_{\eta, Y}(Y)] - \eta^0 V_X(f; g)); \quad (11)$$

175 Given the Neyman orthogonal score $(\eta^0; \eta)$, an estimate $\hat{\eta}_n$ satisfying
 176 $E_D[(\hat{\eta}_n; \eta = f; \hat{\eta}; \hat{g})] = o_P(n^{-1/2})$ gives a DML estimator. Specifically, we propose
 177 the following kernel-smoothing based estimator for the LTE density, named 'KLTE' (kernel-based
 178 estimator for LTE):

³A Neyman orthogonal score is a function satisfying $E_P[(\eta; \eta^0)] = 0$ and $E_P[(V(\eta; \eta^0))] = 0$, where η^0 denotes the true nuisance [Def.2.2]. In words, a score function that is not sensitive to local errors in nuisance models.

Definition 1 (KLTE estimator for η_h). Let $(\theta_0, \eta_0) = (f, g)$ be the Neyman orthogonal score for η_h given in Eq. (11). Let D^0, D^1 denote the randomly split halves of the samples, where $|D^j| = |D^0| = n$. Let \hat{f}, \hat{g} denote the estimates for the nuisances using D^0 . Then, the KLTE estimator for $\eta_h(y)$ for all $y \in Y$, denoted $\hat{\eta}_h(y)$, is given by

$$\hat{\eta}_h(y) = E_{D^1} V_{Y \times X}(\hat{f}, \hat{g})[K_{h,y}(Y)] - E_{D^0} V_X(\hat{f}, \hat{g}); \quad (12)$$

where V_X and $V_{Y \times X}$ are given in Eqs. (5,8), respectively.

We will show that the 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 $\hat{\eta}_h(y)$. For any fixed $y \in Y$, the error $\hat{\eta}_h(y) - \eta_h(y)$ will be analyzed in two folds: we will first analyze the error between the estimator in Eq. (12) and the smoothed estimand in Eq. (2) ($\hat{\eta}_h(y) - \eta_h(y)$), and then analyze the error between the smoothed estimand and the true estimand ($\eta_h(y) - \eta(y)$).

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

Lemma 3 (Convergence rate of $\hat{\eta}_h$ to η_h). For any fixed $y \in Y$, suppose the estimators for nuisances are consistent; i.e., $\hat{\theta}_k = o_p(1)$ for $\theta_0 = f, g$ for all $(w; z; x)$. Suppose $h < 1$, and $nh^d \rightarrow 1$ as $n \rightarrow \infty$. Then,

$$\hat{\eta}_h(y) - \eta_h(y) = O_p\left(\frac{1}{nh^d} + R_2^k + 1\right) = O_p\left(\frac{1}{n}\right);$$

where

$$R_2^k = \sum_z \sum_x k_{z,x} \hat{\theta}_z^k + \sum_z \hat{\theta}_z^0; \quad (13)$$

where $\hat{\theta}_z = \hat{\theta}_z(W)$, $\hat{\theta}_z = \hat{\theta}_z(x; W)$ and $\hat{\theta}_z = \hat{\theta}_z(x; z; W)[K_{h,y}(Y)]$.

The error analysis in Lemma. 3 implies the following:

Corollary 1 (Debiasedness property of $\hat{\eta}_h$ to η_h). If all nuisances \hat{f}, \hat{g} for any given $(w; z; x; y)$ converge at rate $nh^d \rightarrow 1$, then the target estimator $\hat{\eta}_h(y)$ achieves nh^d -rate convergence to η_h .

We now analyze the gap between the smoothed estimand and the true estimand; i.e., $\eta_h - \eta$:

Lemma 4 ([66, Thm. 6.28]) The following holds:

$$\eta_h(y) - \eta(y) = B_y + O(h^2) + o_p(h^2); \quad (14)$$

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

Theorem 1 (Convergence rate of $\hat{\eta}_h$ to η). For any fixed $y \in Y$, suppose the estimators for nuisances are consistent; i.e., $\hat{\theta}_k = o_p(1)$ for $\theta_0 = f, g$ for all $(w; z; x)$. Suppose $h < 1$, and $nh^d \rightarrow 1$ as $n \rightarrow \infty$. Then

$$\hat{\eta}_h(y) - \eta(y) = O_p\left(\frac{1}{nh^d} + R_2^k + 1\right) + B_y; \quad (15)$$

where B_y is defined in Eq. (14), and R_2^k is defined in Eq. (13).

Thm. 1 implies that $\hat{\eta}_h(y)$ converges fast (see Corol. 1) to $\eta(y) + B_y$. A natural question is then how to choose the bandwidth h that minimizes the gap in Eq. (15). The following provides a guideline in choosing the bandwidth:

Lemma 5 (Data-adaptive bandwidth selection). The bandwidth h that minimizes the error in Eq. (15) is $h = O(n^{-1/(d+4)})$. This choice of h satisfies the assumption in Lemma 3 ($nh^d \rightarrow 1$).

201 Recall that Corol. 1 states the debiasedness property of $\hat{\mu}_h$ for any bandwidth satisfying
 202 $nh^d \rightarrow 1$. With the choice of h as in Lemma 5, $\hat{\mu}_h$ converges to μ with the debiasedness property
 203 preserved.

204 Corollary 2 (Debiasedness property of $\hat{\mu}_h$ to μ). Let $h = O(n^{-1/(d+4)})$. If nuisances $\hat{\beta}, \hat{\gamma}, \hat{g}$
 205 converge at $nh^d \rightarrow 1$ rate for any $(w; z; x; y)$, then the target estimator $\hat{\mu}_h(y)$ achieves nh^d -rate
 206 convergence to μ .

207 So far, we have analyzed the error $\hat{\mu}_h(y) - \mu(y)$ pointwise for the fixed $y \in \mathcal{Y}$. To analyze the
 208 difference between the two densities $\hat{p}_h(y)$ and $p(y)$ for all $y \in \mathcal{Y}$, we consider the following
 209 divergence function of two densities:

210 Definition 2 (f-Divergence D_f [20]). Let f denote a convex function with $f(1) = 0$. $D_f(p; q)$
 211 $\int_{\mathcal{Y}} f(p(y); q(y)) q(y) d[y]$, is an f-divergence function between two densities.

212 f-divergence covers many well-known divergences. For example, it reduces to KL divergence with
 213 $f(p; q) = (p - q) \log(p/q)$. We will assume that the function $f(p; q)$ in D_f is differentiable w.r.t.
 214 p and q .

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

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

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

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

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

Theorem 2 (Convergence rate of \hat{p}_h). Suppose the estimators for nuisances are consistent; i.e.,
 $\| \hat{\beta} - \beta \|_p = o_p(1)$ for $2 \leq p \leq \infty$; $\hat{\gamma} \rightarrow \gamma$ for all $(w; z; x; y)$. Suppose D_f is an f-divergence such that
 $f(p; q) = 0$ if $p = q$. Suppose $w(y)$ in Lemma 6 is finite. Then,

$$D_f(\hat{p}_h; p) = O_p \left(\sup_{y \in \mathcal{Y}} R_2^k + B_y + 1 \right) = \frac{p}{nh^d} + 1 = \frac{p}{n}; \quad (16)$$

220 where R_2^k is defined in Eq. (13) and B_y is defined in Eq. (14).

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

222 Corollary 3 (Debiasedness property of \hat{p}_h w.r.t. D_f). Let $h = O(n^{-1/(d+4)})$. Suppose D_f
 223 satisfies $f(p; q) = 0$ if $p = q$. Suppose $w(y)$ in Lemma 6 is finite. If nuisances $\hat{\beta}, \hat{\gamma}, \hat{g}$ converges at
 224 $nh^d \rightarrow 1$ rate for any $(w; z; x; y)$, then $D_f(\hat{p}_h; p)$ converges to 0 at nh^d -rate.

225 4 Model-based approach

226 In this section, we develop a model-based approach for estimating the LTE density $(y) = p(y; \mathcal{C})$.
 227 We will approximate (y) with a class of distributions or density models $\mathcal{G} = \{g(y; \theta) : \theta \in \mathcal{R}^b\}$
 228 where $g(y; \theta) \geq 0$ is differentiable w.r.t. θ . Example density models include exponential family (e.g.,
 229 Gaussian distribution), mixture of Gaussians, or more generally, mixture of exponential families.
 230 The choice of the density model may depend on domain knowledge. Alternatively, one may choose
 231 among a set of candidate density families using separate validation data or applying cross-validation.
 232 We adapt the model-based approach developed in [16] for estimating the causal density under the no
 233 unmeasured confounders assumption.

Given a density model \mathcal{G} , the best approximation for (y) is defined as $g(y; \theta_0) \in \mathcal{G}$ that achieves
 the minimum f-divergence to :

$$\theta_0 = \arg \min_{\theta \in \mathcal{R}^b} D_f((y); g(y; \theta)); \quad (17)$$

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

Consider $m(\theta; \eta) = \int_{\mathcal{Y}} g^0(y; \eta) f f_2^0(y; g(y; \eta)) g(y; \eta) + f(y; g(y; \eta)) g^d(y; \eta) d[y]$. Definition of η_0 given in Eq. (17) implies that $m(\theta; \eta) = 0$ at $\eta = \eta_0$. We note that $m(\theta; \eta)$ serves as a moment score function. The closed-form expression of the score is given by [39]:

$$m(\theta; \eta) = \int_{\mathcal{Y}} g^0(y; \eta) f f_2^0(y; g(y; \eta)) g(y; \eta) + f(y; g(y; \eta)) g^d(y; \eta) d[y]; \quad (18)$$

235 where $g^0(y; \eta) = \int_{\mathcal{Y}} g(y; \eta) d[\eta]$ and $f f_2^0(p; q) = \int_{\mathcal{Y}} f(p; q) d[\eta]$.

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

237 Lemma 7 (Influence Function for $m(\theta; \eta)$). An influence function for $m(\theta; \eta)$ in Eq. (18), denoted m_η , is given by

$$m_\eta(\theta; \eta) = f(\theta; \eta) [R_f(Y; \eta)]; \quad (19)$$

where $(\theta; \eta) [R_f(Y; \eta)]$ is defined in Eq. (9), and

$$R_f(Y; \eta) = g^0(Y; \eta) f f_{21}^0(Y; g(Y; \eta)) g(Y; \eta) + f_1^0(Y; g(Y; \eta)) g;$$

238 where $g^0(y; \eta) = \int_{\mathcal{Y}} g(y; \eta) d[\eta]$, $f_1^0(p; q) = \int_{\mathcal{Y}} f(p; q) d[\eta]$ and $f f_{21}^0(p; q) = \int_{\mathcal{Y}} f(p; q) d[\eta]$.

239 We derive a Neyman orthogonal score based on the moment score $m(\theta; \eta)$ and its influence function $m_\eta(\theta; \eta)$:

240 Lemma 8 (Neyman orthogonal score for $m(\theta; \eta)$). A Neyman orthogonal score for estimating $m(\theta; \eta)$, denoted $\psi^0(\theta; \eta)$, is given by

$$\psi^0(\theta; \eta) = m(\theta; \eta) + m_\eta(\theta; \eta); \quad (20)$$

241 where $m_\eta(\theta; \eta)$ is defined in Eq. (19).

242 Given the orthogonal score $\psi^0(\theta; \eta)$ in Eq. (20), we propose the following estimator for $m(\theta; \eta)$ named 'MLTE' (model-based estimator for LTE):

244 Definition 3 (MLTE estimator for $m(\theta; \eta)$). Let $\psi^0(\theta; \eta) = f(\theta; \eta) [R_f(Y; \eta)]$ be the Neyman orthogonal score for $m(\theta; \eta)$ given in Eq. (20). Let D^0 and D^1 denote the randomly split halves of the samples, where $|D^j| = |D^0| = |D^1| = n$. Let $\hat{f} = f^{\wedge}$, $\hat{g} = g^{\wedge}$ denote the estimators for the nuisance parameters f and g . Then, the MLTE estimator for $m(\theta; \eta)$, denoted $\hat{m}(\theta; \hat{f}, \hat{g})$, is given as a solution satisfying $E_{D^0} \psi^0(\hat{m}(\theta; \hat{f}, \hat{g}); \hat{f}, \hat{g}) = o_p(n^{-1/2})$.

248 To illustrate, we exemplify Eq. (18) and Lemma (7, 8) for the case where D_f is a KL-divergence and $g(y; \eta) = f(y; \eta) g$ is a normal distribution. First, $m(\theta; \eta) = f(\theta; \eta) [m(\theta; \eta); g]$, where $m(\theta; \eta) = (1 - \eta^2) \int_{\mathcal{Y}} [Y] d[\eta]$ and $m(\theta; \eta) = (0.5 - \eta^2) \int_{\mathcal{Y}} [(Y - \eta)^2] d[\eta]$. We note that $\hat{m} = \hat{f} [Y]$ and $\hat{m}^2 = \hat{f} [(Y - \hat{\eta})^2]$ are estimators for $m_0 = f_0$; $m_0^2 = g$ for the score $m(\theta; \eta)$.

252 Also, $R_f(Y; \eta) = \int_{\mathcal{Y}} \log(g(Y; \eta)) = f R_f(Y; \eta); R_f(Y; \eta) = f R_f(Y; \eta); R_f(Y; \eta) = f R_f(Y; \eta)$, where $R_f(Y; \eta) = \int_{\mathcal{Y}} (Y - \eta)^2$ and $R_f(Y; \eta) = 0.5 \int_{\mathcal{Y}} (Y - \eta)^2 g = 0.5 \int_{\mathcal{Y}} (Y - \eta)^2 g$. Then, the Neyman orthogonal score is given as $\psi^0(\theta; \eta) = (1 - \eta^2) f [Y] (\theta; \eta) [Y] g$ and $\psi^0(\theta; \eta) = (0.5 - \eta^2) \int_{\mathcal{Y}} [(Y - \eta)^2] (\theta; \eta) [(Y - \eta)^2]$. Finally, solutions for $\hat{m}(\theta; \hat{f}, \hat{g})$ and $\psi^0(\theta; \hat{f}, \hat{g})$ are given by $\hat{m}(\theta; \hat{f}, \hat{g})$, where, for $[Y]$ in Eq. (9), $\hat{f} = \hat{f} [Y] + E_{D^0} (\hat{f}^{\wedge}) [Y]$ and $\hat{m}^2 = \int_{\mathcal{Y}} [(Y - \hat{\eta})^2] + E_{D^0} (\hat{f}^{\wedge}) [(Y - \hat{\eta})^2]$.

258 The MLTE estimator in Def. 3 is consistent provided that nuisance estimators are consistent [4, Thm.4]. Such $\hat{m}(\theta; \hat{f}, \hat{g})$ is known to achieve debiasedness [13], since a DML estimator. Specifically, Theorem 3 (Convergence rate of $\hat{m}(\theta; \hat{f}, \hat{g})$). Let $\psi^0(\theta; \eta) = f(\theta; \eta) [R_f(Y; \eta)]$ be given in Eq. (20). Let $m(\theta; \eta)$ be given in Eq. (19). Let $\theta_0; \eta_0; \eta_0$ denote the true parameters. Let $\hat{m}(\theta; \hat{f}, \hat{g})$ be the MLTE estimator for $m(\theta; \eta)$ defined in Def. 3. Suppose (1) $R_f(y; \eta)$ is bounded and $R_f^0(y; \eta) = \int_{\mathcal{Y}} R_f(y; \eta) d[\eta] < 1$; (2) There exists a function $H(y) < 1$ s.t. $\sup_y \max_f R_f(y; \eta); R_f^0(y; \eta) g = O(H(y))$; (3) $f'(\theta; \eta) g$ is Donsker w.r.t. for

⁴A function class where complexities are restricted. See Def. S.1 in the Appendix for the definition. Donsker class include Sobolev, Bounded monotone, Lipschitz class, etc.

the xed ; (3) The estimators are consistent: $\hat{\theta}_0 = o_P(1)$ and $k^{\wedge}k = o_P(1)$ for $2 f_z(w); x(z; w); (x; z; w)[H(Y)]g$ for all $(w; z; x; y)$; and (4) $E_P [f'(\cdot; \cdot; \cdot)]$ is differentiable w.r.t. θ at θ_0 with non-singular matrix $M(\theta_0; \cdot; \cdot)$ ($\partial = \partial_j = \partial_{\theta_j} E_P [f'(\cdot; \cdot; \cdot)]$) for all $(\cdot; \cdot)$, where $M(\theta_0; \cdot; \cdot) = M(\theta_0; \theta_0)$. Then,

$$\hat{\theta}_0 = M^{-1} E_D [m(\theta_0; \theta_0)] + o_P(n^{-1/2}) + O_P(R_2^m);$$

where

$$R_2^m = \int k_{\hat{z}}^n + \int k_{\hat{z}}^n + \int k_{\hat{z}}^n + \int k_{\hat{z}}^n + \int k_{\hat{z}}^n + \int k_{\hat{z}}^n + \int k_{\hat{z}}^n + \int k_{\hat{z}}^n;$$

260 where $z_z(W)$, $z_x(z; W)$, and $z(x; z; W)[H(Y)]$.

261 Corollary 4 (Debiasedness property for $\hat{\theta}$). If nuisances f^{\wedge}, g^{\wedge} converges at $n^{-1/4}$ rate, then
 262 the target estimator $\hat{\theta}$ converges to θ_0 at $n^{-1/2}$ -rate.

263 For the above example where D is the KL divergence and $(y; \cdot)$ is a normal distribution $H(Y) =$
 264 Y for $R_f(y; \cdot)$, and $H(Y) = Y^2$ for $R_g(y; \cdot)$.

265 5 Empirical applications

266 In this section, we apply the proposed methods to synthetic and real datasets. For the kernel-smoothing
 267 based approach, we compare KLTE with a baseline plug-in estimator ('kernel-smoothing'), where
 268 estimates of nuisances f^{\wedge}, g^{\wedge} are plugged in the estimand Eq. (2). We use the Gaussian kernel.
 269 The bandwidth is set to $h = 0.5n^{-1/5}$. In estimating the density, we choose 200 equi-spaced points
 270 $y_{(i)}^j$ in Y and evaluate both estimators $\hat{f}_{k, y_{(i)}}^j$ for $i = 1; \dots; 200$. For the model-based
 271 approach, we compare MLTE (e.g., $\hat{\theta}^2$) with a moment-score-based estimator (called 'moment'),
 272 defined as $\hat{\theta}_m$ satisfying $m(\hat{\theta}_m; \hat{\theta}) = o_P(n^{-1/2})$ (e.g., $f_{\hat{\theta}_m}; \hat{\theta}_m^2 g$). We use KL divergence d_{θ_f}
 273 and the normal distribution $\phi(y; \cdot)$. For both approaches, nuisances are estimated through a
 274 gradient boosting model XGBoost [11], which is known to be flexible.

275 5.1 Synthetic dataset

276 We applied the proposed estimators to estimate the density $p(y; C)$ where the true densities are given as
 277 in the 4th plot in Fig. 2. As shown in the ground-truth in Fig. 3a, true densities $p(y_0; C); p(y_{x^1}; C)$
 278 are given as a mixture of four Gaussians. Estimated densities for Moment and MLTE are given in
 279 Fig. (3b, 3c). We note that model-based approaches fail to capture important characteristics (such as
 280 the number of modes) of the true density ('ground-truth' in Fig. 3a) because the assumed density
 281 class is misspecified. The 'kernel-smoothing' (Fig. 3d) captures only one of the modes having the
 282 highest densities, and this leads to misinterpretation of the true densities. KLTE (Fig. 3e) is able to
 283 capture the number, location, and scales of modes correctly.

(a) Ground-truth (b) Moment (c) MLTE (d) Kernel-smoothing (e) KLTE

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.

284 5.2 Application to 401(k) data

285 We applied the proposed estimators (KLTE and MLTE) on 401(k) data, where the data generating
 286 processes corroborate with Fig. 1. Monotonicity assumption holds naturally, since ineligible units
 287 ($Z = 0$) cannot participate ($X = 1$) in 401(k). In our analysis, we used the dataset introduced

288 by [2] containing 9275 individuals, which has been studied in [7, 5, 47, 58, 64], to cite a few.
 289 Model-based approaches (Moment in Fig. 4a and MLTE in Fig. 4b) and kernel-smoothing based
 290 approaches (kernel-smoothing in Fig. 4c and KLTE in Fig. 4d) are implemented to analyze the data.
 291 The model-based (Fig. (4a,4b)) and kernel-smoothing based (Fig. (4c,4d)) estimates both capture
 292 important characteristics of the distribution, such as mode, location, and scale parameters. The results
 293 of proposed estimators (MLTE and KLTE in Fig. (4b,4d)) are consistent with findings from previous
 294 analyses [2, 17, 5, 58]: The effects of the 401(k) participation ($\alpha = 1$) on net financial assets
 295 are positive over the whole range of asset distributions. To connect to CDF method, we provide in
 296 Fig. 4e the CDF estimate induced by KLTE density estimation (Fig. 4a). We note that the CDF in
 297 Fig. 4e captures the nonconstant impact trend of the 401(k) participation on the net financial assets,
 298 which has been also described in the previous analyses [2, 17, 5, 58].

(a) Moment (b) MLTE (c) Kernel-smoothing (d) KLTE (e) KLTE – CDF

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

299 6 Conclusion

300 In this paper, we develop kernel-smoothing-based and model-based approaches for estimating the
 301 LTE density in the presence of instruments. For each approach, we give Neyman orthogonal scores
 302 (Lemma (2,8)) and constructed corresponding DML estimators (KLTE in Def. 1 and MLTE in Def. 3),
 303 that exhibit debiasedness property (Corol. (3, 4)). We demonstrated our work through synthetic
 304 and real datasets. The performance of model-based estimators depends critically on the choice of
 305 the density class. Kernel-based estimators do not have to make assumptions about the true density
 306 class but will suffer from the curse of dimensionality. This work is limited to settings where the
 307 monotonicity assumption holds, i.e., there are no de ders. One could perform sensitivity analyses on
 308 the impact of potential de ders to the estimates as conducted in [65, 36].

309 Acknowledgements

310 We thank the reviewers for their feedback helping to improve this manuscript. Elias Bareinboim and
 311 Yonghan Jung were partially supported by grants from NSF IIS-1750807 (CAREER). Jin Tian was
 312 partially supported by ONR grant N000141712140.

313 References

- 314 [1] A. Abadie. Bootstrap tests for distributional treatment effects in instrumental variable models.
 315 *Journal of the American statistical Association*, 97(457):284–292, 2002.
- 316 [2] A. Abadie. Semiparametric instrumental variable estimation of treatment response models.
 317 *Journal of econometrics*, 13(2):231–263, 2003.
- 318 [3] J. D. Angrist, G. W. Imbens, and D. B. Rubin. Identification of causal effects using instrumental
 319 variables. *Journal of the American statistical Association*, 91(434):444–455, 1996.
- 320 [4] E. Bareinboim, J. D. Correa, D. Ibeling, and T. Icard. On pearl's hierarchy and the foundations
 321 of causal inference. Technical Report R-60. Causal Artificial Intelligence Laboratory, Columbia
 322 University, 2020.
- 323 [5] A. Belloni, V. Chernozhukov, and C. Hansen. Inference on treatment effects after selection
 324 among high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650, 2014.

- 325 [6] A. F. Bibaut and M. J. van der Laan. Data-adaptive smoothing for optimal-rate estimation of
326 possibly non-regular parameters. *arXiv preprint arXiv:1706.07408*, 2017.
- 327 [7] P. J. Bickel, C. A. Klaassen, P. J. Bickel, Y. Ritov, J. Klaassen, J. A. Wellner, and Y. Ritov.
328 cient and adaptive estimation for semiparametric models. volume 4. Johns Hopkins University
329 Press Baltimore, 1993.
- 330 [8] C. Brito. Instrumental sets. In R. Dechter, H. Geffner, and J. Y. Halpern, editors, *Heuristics,
331 Probability and Causality. A Tribute to Judea Pearl*. College Publications, 2010.
- 332 [9] C. Brito and J. Pearl. Generalized instrumental variables. *Proceedings of the Eighteenth
333 conference on Uncertainty in artificial intelligence*, pages 85–93, 2002.
- 334 [10] M. D. Cattaneo, M. Jansson, and X. Ma. Local regression distribution estimation. *Journal of
335 Econometrics*, 2021.
- 336 [11] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. *Proceedings of the 22nd
337 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages
338 785–794, 2016.
- 339 [12] V. Chernozhukov, D. Chetverikov, M. Demirer, E. Du o, C. Hansen, and W. Newey. Dou-
340 ble/debiased/neyman machine learning of treatment effects. *American Economic Review*
341 107(5):261–65, 2017.
- 342 [13] V. Chernozhukov, D. Chetverikov, M. Demirer, E. Du o, C. Hansen, W. Newey, and J. Robins.
343 Double/debiased machine learning for treatment and structural parameters: Double/debiased
344 machine learning. *The Econometrics Journal*, 21(1), 2018.
- 345 [14] V. Chernozhukov, J. C. Escanciano, H. Ichimura, W. K. Newey, and J. M. Robins. Locally
346 robust semiparametric estimation. *arXiv preprint arXiv:1608.00033*, 2016.
- 347 [15] V. Chernozhukov, I. Fernández-Val, and A. Galichon. Quantile and probability curves without
348 crossing. *Econometrica* 78(3):1093–1125, 2010.
- 349 [16] V. Chernozhukov, I. Fernández-Val, and B. Melly. Inference on counterfactual distributions.
350 *Econometrica* 81(6):2205–2268, 2013.
- 351 [17] V. Chernozhukov and C. Hansen. The effects of 401 (k) participation on the wealth distribution:
352 an instrumental quantile regression analysis. *Review of Economics and Statistics*, 86(3):735–751,
353 2004.
- 354 [18] V. Chernozhukov, C. Hansen, and K. Wuthrich. Instrumental variable quantile regression.
355 preprint arXiv:2009.00436, 2020.
- 356 [19] K. Colangelo and Y.-Y. Lee. Double debiased machine learning nonparametric inference with
357 continuous treatments. *arXiv preprint arXiv:2004.03036*, 2020.
- 358 [20] I. Csiszár. Information-type measures of difference of probability distributions and indirect
359 observation. *studia scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- 360 [21] I. Díaz and M. J. van der Laan. Targeted data adaptive estimation of the causal dose–response
361 curve. *Journal of Causal Inference*, 1(2):171–192, 2013.
- 362 [22] J. DiNardo, N. M. Fortin, and T. Lemieux. Labor market institutions and the distribution of
363 wages, 1973-1992: A semiparametric approach. *Econometrica: Journal of the Econometric
364 Society*, pages 1001–1044, 1996.
- 365 [23] M. Frölich. Nonparametric iv estimation of local average treatment effects with covariates.
366 *Journal of Econometrics*, 139(1):35–75, 2007.
- 367 [24] M. Frölich and B. Melly. Unconditional quantile treatment effects under endogeneity. *Journal
368 of Business & Economic Statistics*, 31(3):346–357, 2013.
- 369 [25] S. Ghosal et al. Convergence rates for density estimation with bernstein polynomials.
370 *Annals of Statistics*, 29(5):1264–1280, 2001.

- 371 [26] E. Giné, R. Nickl, et al. Confidence bands in density estimation. *The Annals of Statistics*
372 38(2):1122–1170, 2010.
- 373 [27] L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A distribution-free theory of nonparametric*
374 *regression*. Springer Science & Business Media, 2006.
- 375 [28] J. Hahn. On the role of the propensity score in efficient semiparametric estimation of average
376 treatment effects. *Econometrica* pages 315–331, 1998.
- 377 [29] J. J. Heckman. *Randomization as an instrumental variable*, 1995.
- 378 [30] Y.-C. Hsu, T.-C. Lai, and R. P. Lieli. Estimation and inference for distribution functions
379 and quantile functions in endogenous treatment effect models. Technical report, Institute of
380 Economics, Academia Sinica, Taipei, Taiwan, 2015.
- 381 [31] G. W. Imbens and J. D. Angrist. Identification and estimation of local average treatment effects.
382 *Econometrica* 62(2):467–475, 1994.
- 383 [32] G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*
384 Cambridge University Press, 2015.
- 385 [33] Y. Jung, J. Tian, and E. Bareinboim. Estimating identifiable causal effects on markov equivalence
386 class through double machine learning. *Proceedings of the 38th International Conference on*
387 *Machine Learning (ICML)* 2021.
- 388 [34] Y. Jung, J. Tian, and E. Bareinboim. Estimating identifiable causal effects through double
389 machine learning. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, 2021.
- 390 [35] N. Kallus and M. Uehara. Doubly robust off-policy value and gradient estimation for determin-
391 istic policies. *Advances in Neural Information Processing Systems* 33, 2020.
- 392 [36] H. Kang, Y. Jiang, Q. Zhao, and D. S. Small. *lvmodel: an r package for inference and sensitivity*
393 *analysis of instrumental variables models with one endogenous variable*. *Observational Studies*
394 7(2):1–24, 2021.
- 395 [37] E. H. Kennedy. *Semiparametric theory and empirical processes in causal inference*. *Statistical*
396 *causal inferences and their applications in public health research* pages 141–167. Springer,
397 2016.
- 398 [38] E. H. Kennedy, S. Balakrishnan, M. G'Sell, et al. Sharp instruments for classifying compliers
399 and generalizing causal effects. *Annals of Statistics* 48(4):2008–2030, 2020.
- 400 [39] E. H. Kennedy, S. Balakrishnan, and L. Wasserman. Semiparametric counterfactual density
401 estimation. *arXiv preprint arXiv:2102.12034*, 2021.
- 402 [40] E. H. Kennedy, S. Lorch, and D. S. Small. Robust causal inference with continuous instruments
403 using the local instrumental variable curve. *Journal of the Royal Statistical Society: Series B*
404 *(Statistical Methodology)* 81(1):121–143, 2019.
- 405 [41] E. H. Kennedy, Z. Ma, M. D. McHugh, and D. S. Small. Nonparametric methods for doubly
406 robust estimation of continuous treatment effects. *Journal of the Royal Statistical Society. Series*
407 *B, Statistical Methodology* 79(4):1229, 2017.
- 408 [42] K. Kim, J. Kim, and E. H. Kennedy. Causal effects based on distributional distances.
409 preprint *arXiv:1806.02935*, 2018.
- 410 [43] C. A. Klaassen. Consistent estimation of the influence function of locally asymptotically linear
411 estimators. *The Annals of Statistics* pages 1548–1562, 1987.
- 412 [44] T. Leonard. Density estimation, stochastic processes and prior information. *Journal of the*
413 *Royal Statistical Society: Series B (Methodological)* 40(2):113–132, 1978.
- 414 [45] B. Melly and K. Wüthrich. *Local quantile treatment effects*. 2016.

- 415 [46] R. Neugebauer and M. van der Laan. Nonparametric causal effects based on marginal structural
416 models. *Journal of Statistical Planning and Inference* 37(2):419–434, 2007.
- 417 [47] E. L. Ogburn, A. Rotnitzky, and J. M. Robins. Doubly robust estimation of the local average
418 treatment effect curve. *Journal of the Royal Statistical Society. Series B, Statistical methodology*
419 77(2):373, 2015.
- 420 [48] R. Okui, D. S. Small, Z. Tan, and J. M. Robins. Doubly robust instrumental variable regression.
421 *Statistica Sinica* pages 173–205, 2012.
- 422 [49] J. Pearl. Causal diagrams for empirical research. *Biometrika* 82(4):669–710, 1995.
- 423 [50] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York,
424 2000. 2nd edition, 2009.
- 425 [51] J. Pearl. Parameter identification: A new perspective. Technical Report R-276, 2001.
- 426 [52] J. Robins, L. Li, E. Tchetgen, A. van der Vaart, et al. Higher order influence functions and
427 minimax estimation of nonlinear functionals. *Probability and statistics: essays in honor of*
428 *David A. Freedman* pages 335–421. Institute of Mathematical Statistics, 2008.
- 429 [53] J. Robins and A. Rotnitzky. Comment on “inference for semiparametric models: Some questions
430 and an answer,” by pj bickel and j. kwon. *Statistica Sinica* 11:920–936, 2001.
- 431 [54] J. M. Robins and Y. Ritov. Toward a curse of dimensionality appropriate (coda) asymptotic
432 theory for semi-parametric models. *Statistics in medicine* 16(3):285–319, 1997.
- 433 [55] P. R. Rosenbaum and D. B. Rubin. The central role of the propensity score in observational
434 studies for causal effects. *Biometrika* 70(1):41–55, 1983.
- 435 [56] K. J. Rothman and S. Greenland. Causation and causal inference in epidemiology. *American*
436 *journal of public health* 95(S1):S144–S150, 2005.
- 437 [57] D. B. Rubin. Bayesian inference for causal effects: The role of randomization. *The Annals of*
438 *statistics* pages 34–58, 1978.
- 439 [58] R. Singh and L. Sun. De-biased machine learning for compliers. arXiv preprint
440 arXiv:1909.05244, 2019.
- 441 [59] E. Smucler, A. Rotnitzky, and J. M. Robins. A unifying approach for doubly-robust regular-
442 ized estimation of causal contrasts. arXiv preprint arXiv:1904.03737, 2019.
- 443 [60] Z. Tan. Regression and weighting methods for causal inference using instrumental variables.
444 *Journal of the American Statistical Association* 101(476):1607–1618, 2006.
- 445 [61] A. B. Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media,
446 2008.
- 447 [62] S. D. Uysal. Doubly robust iv estimation of the local average treatment effects, 2011.
- 448 [63] A. W. Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- 449 [64] L. Wang, Y. Zhang, T. S. Richardson, and J. M. Robins. Estimation of local treatment effects
450 under the binary instrumental variable model. *Biometrika* 2021.
- 451 [65] X. Wang, Y. Jiang, N. R. Zhang, and D. S. Small. Sensitivity analysis and power for instrumental
452 variable studies. *Biometrics* 74(4):1150–1160, 2018.
- 453 [66] L. Wasserman. *All of nonparametric statistics*. Springer Science & Business Media, 2006.
- 454 [67] J. M. Wooldridge. *Econometric analysis of cross section and panel data*. MIT press, 2010.
- 455 [68] P. G. Wright. *Tariff on animal and vegetable oils*. Macmillan Company, New York, 1928.
- 456 [69] K. Wüthrich. A comparison of two quantile models with endogeneity. *Journal of Business &*
457 *Economic Statistics* 38(2):443–456, 2020.
- 458 [70] W. Zheng and M. J. van der Laan. Cross-validated targeted minimum-loss-based estimation. In
459 *Targeted Learning* pages 459–474. Springer, 2011.

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

460 A IV Settings and LTE

461 In this work, we consider the IV setting represented by the causal graph Fig. 1. It is common in
 462 the literature to define IV assumptions in terms of conditional independences among counterfactuals
 463 [2, 60, 47, 64], as given in the following:

464 Assumption A.1 (IV assumptions).

- 465 1. Exclusion restriction $Y_{x;z} = Y_x$ almost surely for all $z; x$.
- 466 2. Independence $Z \perp (Y_x; X_z) | W$ for all $z; x$.
- 467 3. Instruments relevance $P(X_{Z=1} = 1 | W) \notin P(X_{Z=0} = 1 | W)$ almost surely.

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

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

470 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
 471 is due to the composition property [50, Property 1 (pp. 229)], the second due to exclusion restrictions
 472 [50, Eq.(7.25)], and the third by composition.

473 We will show the second. We have $(Z \perp (Y_x; X_z) | W)$ by independence restriction [50,
 474 Eq.(7.26)]. Then by the weak union graphoid axiom (Rete, pp.11), $(Z \perp (Y_x; X_z) | W)$,
 475 which leads to $(Z \perp (Y_x; X_z) | W)$ by composition.

476 We will show the third. By $(Z \perp (Y_x; X_z) | W)$, $P(X_z | W) = P(X_z | W; Z) = P(X_z | W; Z=1)$, where the second
 477 equality is by composition. The third assumption is reflected by $X_{Z=1}$ not independent of $X_{Z=0}$ given
 478 W in G . □

479 Definition A.1 (Local treatment effect (LTE) density). The local treatment effect (LTE) density
 480 the density of outcome Y under treatment $X = x$ among compliers (i.e. $X_{Z=1} = 1$ and $X_{Z=0} = 0$)
 481 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
 482 event that a unit is a complier and write the LTE density as $p(y_x | C)$.

483 The LTE density $p(y_x | C)$ is known to be identifiable under monotonicity in the IV settings [2]. In
 484 the notations of this paper, we present the identification results as follows, where for a given constant
 485 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

$$p(y_x | W; C) = \frac{p(y | X; Z^x; W)P(XZ^x; W) - p(y | W; X; Z^{1-x})P(XZ^{1-x}; W)}{P(X^1 | Z^1; W) - P(X^1 | Z^0; W)}.$$

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

$$p(y_x | C) = \frac{\int_W [p(y | X; Z^x; W)P(XZ^x; W) - p(y | W; X; Z^{1-x})P(XZ^{1-x}; W)]P(W) d[W]}{\int_W [P(X^1 | Z^1; W) - P(X^1 | Z^0; W)]P(W) d[W]} \\ = \frac{E_P [p(y | X; Z^x; W)P(XZ^x; W) - p(y | X; Z^{1-x}; W)P(XZ^{1-x}; W)]}{E_P [P(X^1 | Z^1; W) - P(X^1 | Z^0; W)]}.$$

486 B Proofs

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

489 $[f(Y)]$ in Eq. (3) can be written as $Y^X [f(Y)] = X$, where X and $Y^X [f(Y)]$ are defined in
 490 Eqs. (4,7).

491 We provide a formal definition of a function class called Donsker class, which is used throughout the
 492 proof.

493 Definition S.1 (Donsker Class [63, page. 269]) Let $G_n(f) = \frac{1}{n} \sum_{i=1}^n f(v_{(i)}) - E_P [f(V)]$
 494 denote the empirical process evaluated at a measurable function f . A class of measurable functions
 495 F is called P -Donsker class if the sequence of processes $\{G_n(f); f \in F\}$ converges in distribution
 496 to a limit process G in the space $\mathcal{C}^1(F)$, where G is the process such that, for all $\epsilon > 0$, there is a
 497 compact set S such that $P(G \in S) > 1 - \epsilon$.

498 Lemma S.1 ([63, Thm.5.31], [9, Lemma 3]) Let $(V; \theta; \eta)$ denote a vector estimating function
 499 for target parameter $\theta \in \mathbb{R}^p$ and nuisance functions $\eta \in H$ for some function space H . Suppose
 500 $E_P [V(\theta; \eta; \theta_0)] = 0$ (where $\theta_0; \eta_0$ denote true parameters) and define the estimator $\hat{\theta}_n$ as a solution
 501 to $E_{D_n} (V; \hat{\theta}_n; \hat{\eta}_n) = 0$, where $\hat{\eta}_n$ is estimated on a separate independent sample. Assume

- 502 1. $f(V; \theta; \eta): \mathbb{R}^p \times H \rightarrow \mathbb{R}^p$ is Donsker for any θ, η .
- 503 2. $\hat{\theta}_n - \theta_0 = o_P(1)$ and $\|\hat{\eta}_n - \eta_0\|_2 = o_P(1)$.
- 504 3. The map $\theta \mapsto E_P [V(\theta; \eta; \theta_0)]$ is differentiable at θ_0 uniformly in η , with non-singular
 505 matrix $M(\theta_0; \eta) = \frac{\partial}{\partial \theta} E_P [V(\theta; \eta; \theta_0)]$, where $M(\theta_0; \eta) = \frac{\partial}{\partial \theta} M(\theta_0; \eta)$.

Then,

$$\hat{\theta}_n - \theta_0 = M^{-1} E_{D_n} [V(\theta_0; \eta_0; \theta_0)] - M^{-1} E_P [V(\theta_0; \eta_0; \theta_0)] + o_P(n^{-1/2});$$

506 B.1 Proofs for Sec. 3

Lemma S.2 ([28, Proof of Thm. 1]) For a target estimand $E_P [f(Y)]$ for binary $X \in \{0, 1\}$ and $f(\cdot) \in [0, 1]$, an
 influence function is given by

$$\frac{1_{X=1}(X) - 1_{X=0}(X)}{P(X|W)} (f(Y) - E_P [f(Y)|X; W]) + E_P [f(Y)|X=1; W] - E_P [f(Y)|X=0; W];$$

Lemma S.3. An influence function for $E_P [f(Y)]$ for $f(Y) \in [0, 1]$ is given by the mapping function in
 Eq. (9), which is

$$IF(f; \theta; \eta; g; \gamma) [f(Y)] = \frac{1}{X} (V_{YX}(f; \theta; \eta; g) [f(Y)] - [f(Y)] V_X(f; \theta; \eta; g));$$

507 Proof. We note that the estimand is given as $E_P [f(Y)] = Y^X [f(Y)] = X$, where X and
 508 $Y^X [f(Y)]$ are defined in Eqs. (4,7).

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

$$\gamma_X = \frac{1_{Z=1}(Z) - 1_{Z=0}(Z)}{Z(W)} (1_{X=1}(X) - \gamma_X(Z; W)) + \gamma_X(Z=1; W) - \gamma_X(Z=0; W); \quad (B.1)$$

$$\gamma_{Y^X [f(Y)]} = \frac{1_{Z^X=1}(Z) - 1_{Z^X=0}(Z)}{Z(W)} (f(Y) - \gamma_{Y^X [f(Y)]}(X; Z; W)) + \gamma_{Y^X [f(Y)]}(X; Z^X=1; W) - \gamma_{Y^X [f(Y)]}(X; Z^X=0; W); \quad (B.2)$$

Then, by applying the chain rule for the Gateaux derivative (since the influence function is a Gateaux derivative), an influence function for $[f(Y)] = \int Y^X [f(Y)] = \int X$ is given as

$$\begin{aligned} & \frac{1}{X} (\int Y^X [f(Y)] - [f(Y)] \int X) \\ &= \frac{1}{X} \int V_{YX} [f(Y)] - \int Y^X [f(Y)] - [f(Y)] \int V_X - \int X \\ &= \frac{1}{X} \int V_{YX} [f(Y)] - [f(Y)] \int V_X - \int [f(Y)] + \int [f(Y)] \\ &= \frac{1}{X} (\int V_{YX} [f(Y)] - [f(Y)] \int V_X) : \end{aligned}$$

509

□

Lemma B.1 (Restated Lemma 1) Let $m(\theta; h)$ be the score defined in Eq. (6). Then, an influence function for $E_P [m(\theta; h)]$, denoted m , is given by

$$m(\theta; h) = \int (f; g; \theta) (\theta; h) [K_{h,y}(Y)] \quad (B.3)$$

where \int is given as

$$\int (f; g; \theta) [f(Y)] = \frac{1}{X} (\int V_{YX} (f; g) [f(Y)] - [f(Y)] \int V_X (f; g))$$

Proof. Let \int_X denote the influence function corresponding to \int , given in Eq. (B.1). This implies that $E_P [\int_X] = \int_X$. Then, equipped with the true nuisance \int_X ,

$$E_P [m(\theta; h)] = E_P \left[\frac{1}{X} (\int_h \int \int V_X = \frac{1}{X} (\int_h \int \int) E_P [\int_X] = \int_h \int \int \right]$$

510 Then, the influence function for $E_P [m(\theta; h)]$ coincides with the influence function for \int_h , which
511 is given by Eq. (B.3) based on Lemma S.3. □

Lemma B.2 (Restated Lemma 2) Let $m(\theta; h)$ be the score function in Eq. (6), and $\int_h(\theta; h) = \int (f; g; \theta; h)$ be the influence function for $E_P [m(\theta; h)]$ given in Eq. (10). Then, a Neyman orthogonal score for \int_h is given as $\int'(\theta; h) = \int (f; g; \theta; h) + \int_m(\theta; h)$; Specifically,

$$\int'(\theta; h) = \frac{1}{X} (\int V_{YX} (f; g) [K_{h,y}(Y)] - \int V_X (f; g)) : \quad (B.4)$$

Proof. For a score function for \int , denoted $m(\theta; h)$, and the influence function for $E_P [m(\theta; h)]$, denoted $\int_m(\theta; h)$, a Neyman orthogonal score for \int_h is given as $\int_m + \int_m$ [14, Thm. 1]. Applying this, $\int_m(\theta; h) + \int_m(\theta; h)$ is a Neyman orthogonal score. Specifically,

$$\begin{aligned} & \int'(\theta; h) = \int (f; g; \theta; h) \\ &= m(\theta; h) + \int_m(\theta; h) \\ &= \frac{1}{X} (\int [K_{h,y}(Y)] - \int V_X + \frac{1}{X} (\int V_{YX} (f; g) [K_{h,y}(Y)] - [K_{h,y}(Y)] \int V_X (f; g)) \\ &= \frac{1}{X} (\int V_{YX} (f; g) [K_{h,y}(Y)] - \int V_X (f; g)) : \end{aligned}$$

512

□

Lemma B.3 (Restated Lemma 3) For any fixed $2 \leq Y$, suppose the estimators for nuisances are consistent; i.e. $\hat{k} = o_p(1)$ for $2 \leq f; g$ for all $(w; z; x)$. Suppose $\phi < 1$, and $nh^d \rightarrow 1$ as $n \rightarrow \infty$. Then,

$$\hat{h}_h(y) - h(y) = O_p \left(\frac{1}{nh^d} + R_2^k + 1 = \frac{1}{n} \right);$$

where

$$R_2^k = \sum_z^X k_z^k \sum_z^n \hat{z}_z^k + \sum_z^0 \hat{z}_z^0; \quad (B.5)$$

513 where $z = z(W)$, $z = x(z; W)$ and $z = (x; z; W) [K_{h,y}(Y)]$.

514 Proof. We note that the condition $nh^d \rightarrow 1$ means that $h = O(n^{-1/d})$ for some $0 < d < 1$.
 515 implies that such h is either constant or decreasing function over n . Combining, the condition implies
 516 $h = O(n^{-1/d})$ for $d \in [0, 1]$.

517 We recall that V_X ; V_{YX} are defined in Eq. (4.7) and V_X ; 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
 Eq. (12) satisfies $E_D[\psi(\hat{\theta}_h; b)] = o_P(n^{-1/2})$, because

$$\begin{aligned} E_D[\psi(\hat{\theta}_h; b)] &= \frac{1}{X} E_D[V_{YX}(f; b; g)[K_{h,y}(Y)] - \hat{h} E_D[V_X(f; b; g)] \\ &= \frac{1}{X} \{ E_D[V_{YX}(f; b; g)[K_{h,y}(Y)] - E_D[V_X(f; b; g)] \} \\ &= 0. \end{aligned}$$

The Neyman orthogonal score function in Lemma 2 satisfies the assumptions in Lemma S.1, since ψ is a linear function of θ when nuisances are fixed. Also ψ in Lemma S.1 is given as 1, which can be witnessed by the following:

$$M(\theta_0; \theta_0) = \frac{1}{X} E_P[fV_{YX} - V_X g] = \frac{1}{X} E_P[V_X];$$

518 and, with the true nuisance $\theta_0 = M(\theta_0; \theta_0) = 1$ since $E_P[V_X] = X$.

Then, by the result of Lemma S.1,

$$\hat{\theta}_h - \theta_0 = E_D[\psi(\hat{\theta}_h; \theta_0)] + E_P[\psi(\hat{\theta}_h; \theta_0)] + o_P(n^{-1/2});$$

519 We will first study the convergence behavior of $E_D[\psi(\hat{\theta}_h; \theta_0)]$. We will show that
 520 $E_P[E_D[\psi(\hat{\theta}_h; \theta_0)]] = O(n^{-1/2})$. Then, the $n^{-1/2}$ -consistency of $E_D[\psi(\hat{\theta}_h; \theta_0)]$ (i.e.,
 521 $E_D[\psi(\hat{\theta}_h; \theta_0)] = o_P(n^{-1/2})$) can be shown immediately by the Markov inequality. This
 522 implies that $E_D[\psi(\hat{\theta}_h; \theta_0)]$ converges in probability to 0.

523 Let $\psi_m(V_i; \theta; \theta_0)$ denote the influence function evaluated at θ_0 .

Consider the following:

$$\begin{aligned} E_P[\psi_m(\hat{\theta}_h; \theta_0)] &= \frac{1}{X} E_P[\psi_m(\hat{\theta}_h; \theta_0)] \\ &= \frac{1}{X} \text{var}_P(\psi_m(\hat{\theta}_h; \theta_0)) \\ &= \frac{1}{X} E_P[\psi_m^2(\hat{\theta}_h; \theta_0)]; \end{aligned}$$

524 where the first inequality is by Cauchy-Schwarz inequality, the second and third equality are from the
 525 iid assumption and $E_P[\psi_m] = 0$.

We note that

$$\begin{aligned} \psi_m &= \frac{1}{X} (V_{YX}[K_{h,y}(Y)] - V_X) \\ &= \frac{1}{X} (V_{YX}[K_{h,y}(Y)] - V_X) + \frac{YX[K_{h,y}(Y)] - X}{X} \\ &= \frac{1}{X} (V_{YX}[K_{h,y}(Y)] - V_X) + \frac{YX[K_{h,y}(Y)] - X}{X} \\ &= \frac{1}{X} (V_{YX}[K_{h,y}(Y)] - V_X): \end{aligned}$$

Next,

$$\begin{aligned} E_P \left[\frac{1}{h} \int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx \right]^2 &= E_P \left[\frac{1}{h} \int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx \right]^2 \\ &= \frac{1}{h^2} E_P \left[\int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx \right]^2 \\ &= \frac{1}{h^2} E_P \left[\int_{X-h}^X f_{XY}^2 [K_{h,Y}(Y)] dx + 2 \int_{X-h}^X \int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx dx \right] : \end{aligned}$$

We first analyze $E_P \left[\int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx \right]^2 = \text{var}_P \left[\int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx \right]$. By [28, Thm. 1],

$$\begin{aligned} \text{var}_P \left[\int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx \right] &= E_P \left[\frac{\text{var}_P (K_{h,Y}(Y) 1_X(X) jz^X; W)}{z^X(W)} + \frac{\text{var}_P (K_{h,Y}(Y) 1_X(X) jz^{1-X}; W)}{z^{1-X}(W)} \right] \\ &+ E_P \left[E_P [K_{h,Y}(Y) 1_X(X) jz^X; W] E_P [K_{h,Y}(Y) 1_X(X) jz^{1-X}; W] \right] : \end{aligned}$$

First,

$$\begin{aligned} E_P [\text{var}_P (K_{h,Y}(Y) 1_X(X) jz^X; W)] &= \text{var}_P (K_{h,Y}(Y) 1_X(X) jz^X) \\ &= E_P [K_{h,Y}^2(Y) 1_X(X) jz^X] \\ &= E_P [K_{h,Y}^2(Y) jx; z^X] \\ &= \int_{Z^Y} K_{h,Y}^2(y^0) p(y|x; z^X) d[y^0] \\ &= \int_{Z^Y} K_{h,Y}^2(y^0) d[y^0] \\ &= \frac{1}{h^{2d}} \int_{Z^Y} K^2 \left(\frac{y^0 - y}{h} \right) d[y^0] \\ &= \frac{1}{h^d} \int_{Z^Y} K^2(u) d[u] \\ &= O(1=h^d) : \end{aligned} \tag{B.6}$$

526 The 1st equality holds by Law of total variance, the 2nd and 3rd by the standard algebra, the 5th by
527 the assumption that $p(y|x; z^X)$ is bounded, and the remaining parts from the change of a variable in
528 the integral computation.

Also,

$$\begin{aligned} E_P \left[E_P [K_{h,Y}(Y) 1_X(X) jz^X; W] E_P [K_{h,Y}(Y) 1_X(X) jz^{1-X}; W] \right] &= O(1=h^d) \\ &= \text{var}_P \left[E_P [K_{h,Y}(Y) 1_X(X) jz^X; W] E_P [K_{h,Y}(Y) 1_X(X) jz^{1-X}; W] \right] \\ &= \frac{1}{h^{2d}} \sup_{z^f, 0; 1g} \text{var}_P (E_P [K_{h,Y}(Y) 1_X(X) jz; W]) \\ &= O(1=h^d); \end{aligned}$$

529 where the 1st (in)equality is by the definition of the variance, the second by the linear combination
530 of the variance, and the last by Eq. (B.6). Therefore $E_P \left[\int_{X-h}^X f_{XY} [K_{h,Y}(Y)] dx \right]^2 = O(1=h^d)$.

Next, we will study $E_P \frac{2}{h} \frac{2}{X}$. We first note that $E_P \frac{2}{h} \frac{2}{X} = \frac{2}{h} E_P \frac{2}{X} = O(\frac{2}{h})$. Therefore, it suffices to analyze $O(\frac{2}{h})$.

$$\begin{aligned} \frac{2}{h} &= \int_Z^Y K_{h,y}(y^0) (y^0)^2 d[y^0] \\ &\leq \int_Z^Y K_{h,y}^2(y^0) (y^0)^2 d[y^0] \\ &\leq \int_Z^Y K_{h,y}^2(y^0) d[y^0] \\ &= \int_Z^Y \frac{1}{h^{2d}} K^2 \left(\frac{y^0 - y}{h} \right) d[y^0] \\ &= \int_U \frac{1}{h^d} K^2(u) d[u] \\ &= O(1/h^d); \end{aligned}$$

531 where the 2nd line inequality by the Cauchy-Schwarz inequality, the 3rd by the assumption that
532 is bounded, the 4th by the change of variables.

Finally, consider the term $2E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right]$. Note, $E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right] = \frac{2}{h} E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right]$. We first consider $E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right]$:

$$\begin{aligned} E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right] &= \frac{E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right]}{\sqrt{E_P \left[\left(\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right)^2 \right] E_P \left[\frac{2}{X} \right]}} \\ &= O \left(\frac{1}{E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right]} \right) = O(h^{d-2}); \end{aligned}$$

533 where the last equality holds by Eq. (B.6).

Next, consider $\frac{2}{h}$:

$$\begin{aligned} \frac{2}{h} &= \int_Z^Y K_{h,y}(y^0) (y^0)^2 d[y^0] \\ &= \int_Z^Y \frac{1}{h} K \left(\frac{y^0 - y}{h} \right) (y^0)^2 d[y^0] \\ &= \int_U K(u) (hu + y) d[u] \\ &= \int_U K(u) (y + hu) d[u] \\ &= C + O(h^2); \end{aligned}$$

534 for some constant C , The 4th line equality holds by the differentiability assumption and the last
535 equality holds since (y) is bounded and twice differentiable. Combining, we can rewrite the term
536 $2E_P \left[\int_Y^X [K_{h,y}(Y)] \frac{2}{X} \right]$ as $O(h^{d-2} + h^{d-2}h^2)$.

Therefore,

$$E_P \left(\frac{2}{h} \right) = O(h^{d-2} + h^{d-2} + h^{d-2}h^2):$$

With $h = O(n^{-1/d})$ with $d \in [0; 1=d]$, we can rewrite

$$E_P \left(\frac{2}{h} \right) = O(h^{d-2} + h^{d-2} + h^{d-2}h^2) = O(n^{-d}) = O(h^d):$$

This shows that

$$E_P [E_D \left(\frac{2}{h} \right)] = O \left(\frac{1}{n} \right) E_P \left(\frac{2}{h} \right) = O \left(\frac{1}{n} \right) = O \left(\frac{1}{nh^d} \right):$$

We now consider $E_P [m(h; \hat{\alpha})]$.

$$\begin{aligned}
 & E_P [m(h; \hat{\alpha})] \\
 &= E_P \frac{1}{\hat{\alpha}_X} \hat{V}_{YX} [K_{h;Y}(Y)] [K_{h;Y}(Y)] \hat{V}_X \\
 &= E_P \frac{1}{\hat{\alpha}_X} \hat{V}_{YX} [K_{h;Y}(Y)] [K_{h;Y}(Y)] \hat{V}_X + \frac{1}{\hat{\alpha}_X} \frac{1}{\hat{\alpha}_X} \hat{V}_{YX} [K_{h;Y}(Y)] [K_{h;Y}(Y)] \hat{V}_X : \\
 & \hspace{15em} (B.7)
 \end{aligned}$$

For further analysis, we consider $E_P \hat{V}_{YX}^h [K_{h;Y}(Y)] V_{YX}^i [K_{h;Y}(Y)]$. First, define

$$V_{YX; (x; z)}(\cdot; \cdot) [f(Y)] = \frac{1_z(Z)}{z(W)} (f(Y) 1_x(X) (x; Z; W) [f(Y)] + (x; z; W) [f(Y)]):$$

Then, $V_{YX} [f(Y)] = V_{YX; (x; z^*)} [f(Y)] + V_{YX; (x; z^{1-x})} [f(Y)]$. Now, consider

$E_P \hat{V}_{YX; (x; z)} [K_{h;Y}(Y)] V_{YX; (x; z)} [K_{h;Y}(Y)]$. We have

$$\begin{aligned}
 & E_P V_{YX; (x; z)}^h (\hat{\alpha}; \hat{\alpha}) [f(Y)] V_{YX; (x; z)}^i (\cdot; \cdot) [f(Y)] \\
 &= E_P \frac{1_z(Z)}{\hat{\alpha}_Z(W)} f(Y) 1_x(X) \hat{\alpha}(x; Z; W) [f(Y)] + \hat{\alpha}(x; z; W) [f(Y)] (x; z; W) [f(Y)] \\
 &= E_P \frac{1_z(Z)}{\hat{\alpha}_Z(W)} (x; Z; W) [f(Y)] \hat{\alpha}(x; Z; W) [f(Y)] + \hat{\alpha}(x; z; W) [f(Y)] (x; z; W) [f(Y)]^0 \\
 &= E_P \frac{z(W)}{\hat{\alpha}_Z(W)} (x; z; W) [f(Y)] \hat{\alpha}(x; z; W) [f(Y)] + \hat{\alpha}(x; z; W) [f(Y)] (x; z; W) [f(Y)]^0 \\
 &= E_P (x; z; W) [f(Y)] \hat{\alpha}(x; z; W) [f(Y)] \frac{1}{\hat{\alpha}_Z(W)} \frac{z(W)}{z(W)} \\
 &= E_P (x; z; W) [f(Y)] \hat{\alpha}(x; z; W) [f(Y)] \frac{\hat{\alpha}_Z(W)}{\hat{\alpha}_Z(W)} \frac{z(W)}{z(W)} \\
 &= O_P (x; z; W) [f(Y)] \hat{\alpha}(x; z; W) [f(Y)] k_{\hat{\alpha}_Z(W)} z(W) k ;
 \end{aligned}$$

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

$$\begin{aligned}
 R_{YX} &= E_P V_{YX}^h (\hat{\alpha}; \hat{\alpha}) [f(Y)] V_{YX}^i (\cdot; \cdot) [f(Y)] \\
 &= \sum_{z^2 f 0; 1g} O_P (x; z; W) [f(Y)] \hat{\alpha}(x; z; W) [f(Y)] k_{\hat{\alpha}_Z(W)} z(W) k :
 \end{aligned}$$

Also, let

$$V_{X;x}(\cdot; \cdot) = \frac{1_z(Z)}{z(W)} (1_x(X) (x; Z; W)) + (x; z; W):$$

Then, with the similar proof as above, we have

$$E_P V_{X;x}^h (\hat{\alpha}; \hat{\alpha}) V_{X;x}^i (\cdot; \cdot) = O_P (x(z; W) \hat{\alpha}_x(z; W) k_{\hat{\alpha}_Z(W)} z(W) k) ;$$

and

$$E_P V_X^h (\hat{\alpha}; \hat{\alpha}) V_X^i (\cdot; \cdot) = \sum_{z^2 f 0; 1g} O_P (x(z; W) \hat{\alpha}_x(z; W) k_{\hat{\alpha}_Z(W)} z(W) k) :$$

Recall $R_{YX} = E_P \hat{V}_{YX} - V_{YX}$ and let $R_X = E_P \hat{V}_X - V_X$. Then, continuing from Eq. (B.7),

$$\begin{aligned} \text{Eq. (B.7)} &= E_P \frac{1}{X} YX + R_{YX} - \frac{YX}{X} (X + R_X) + \frac{1}{\hat{\Lambda}_X} \frac{1}{X} YX + R_{YX} - \frac{YX}{X} (X + R_X) \\ &= E_P \frac{1}{X} (R_{YX} - R_X) + \frac{1}{\hat{\Lambda}_X} \frac{1}{X} (R_{YX} - R_X) \\ &= O_P(R_{YX} + R_X) \\ &= O_P(R_2^k); \end{aligned}$$

where

$$R_2^k = O_P \left(k^{\wedge_z} z^k \wedge_z + \wedge_z z^o \right);$$

537 Note the first equality is by $E_P \hat{V}_{YX} = R_{YX} + E_P [V_{YX}]$ and $E_P \hat{V}_X = R_X + E_P [V_X]$, the
538 second by rearrangement, the third by Positivity, the fourth by the definition of R_{YX} and R_X .

539 Summing up, we have shown that $E_P [m(h; \wedge)] = O(1 = \overline{nh^d})$ and $E_P [m(h; \hat{\wedge})] = O_P(R_2^k)$.
540 □

541 Corollary 1 (Restated Corol. 1) If all nuisances $\wedge; \hat{\wedge}; \hat{g}$ for any given $(w; z; x; y)$ converge at
542 rate $f nh^d g^{-1=4}$, then the target estimator $\hat{h}(y)$ achieves $\overline{nh^d}$ -rate convergence to h .

543 Proof. This result follows immediately from Lemma 3. □

Theorem B.1 (Restated Thm. 1) For any fixed Y , suppose the estimators for nuisances are consistent; i.e. $k^{\wedge} = o_P(1)$ for $z = f; g$ for all $(w; z; x)$. Suppose $\phi < 1$, and $nh^d \rightarrow 1$ as $n \rightarrow \infty$. Then

$$\hat{h}(y) - h(y) = O_P \left(1 = \overline{nh^d} + R_2^k + 1 = \overline{n} + B_y \right); \quad (\text{B.8})$$

544 where B_y is defined in Eq. (14), and R_2^k is defined in Eq. (13).

545 Proof. This result follows immediately from Lemmas 3 and 4. □

546 Lemma B.5 (Restated Lemma 5) The bandwidth that minimizes the error in Eq. (15) is
547 $O(n^{-1=(d+4)})$. This choice of h satisfies the assumption in Lemma 3 that $nh^d \rightarrow 1$.

548 Proof. We note that the error in Eq. (15) w.r.t. is $O_P(1 = \overline{nh^d} + h^2)$. Since the function $1 = \overline{nh^d} +$
549 h^2 is convex w.r.t. h and the global minimum is at $h = n^{-1=(d+4)}$, the optimal h minimizing the
550 error is $h = O(n^{-1=(d+4)})$. Then, $O(nh^d) = O(n^{4=(d+4)})$, implying that $nh^d \rightarrow 1$. □

551 Corollary 2 (Restated Corol. 2) Let $h = O(n^{-1=(d+4)})$. If nuisances $\wedge; \hat{\wedge}; \hat{g}$ converge at
552 rate $f nh^d g^{-1=4}$ for any $(w; z; x; y)$, then the target estimator $\hat{h}(y)$ achieves $\overline{nh^d}$ -rate convergence
553 to h .

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

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

$$D_f(\cdot; b_h) = \int_Y w(y) \hat{h}(y) - h(y) d[y];$$

557 where $w(y) = f_2^0(y; \sim(y)) \hat{h}(y)$, $f_2^0(p; q) = \int_0^1 \phi(t) dF(p; q)$, and $\sim_h(y) = t \hat{h}(y) + (1 - t) h(y)$
558 for some fixed $t \in [0; 1]$.

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

$$f(p; q) = f(p; p) + f_2^0(p; p)(q - p);$$

for some $x \in [p, q]$. Applying this idea,

$$\begin{aligned} D_f(\cdot; b_h) &= \int_Y^Z f(\cdot(y); \hat{h}_h(y)) \hat{h}_h(y) d[y] \\ &= \int_Y^Z \left[\underbrace{f(\cdot(y); \cdot(y))}_{=0} + f_2^0(\cdot(y); \cdot(y)) \hat{h}_h(y) \right] d[y]; \\ &= \int_Y^Z w(y) \hat{h}_h(y) d[y]; \end{aligned}$$

559 where the second equality holds by Taylor's expansion and the third equality by the given
560 assumption that $f(p; q) = 0$ whenever $p = q$.

561

□

Theorem B.2 (Restated Thm. 2) Suppose the estimators for nuisances are consistent, i.e., $\hat{\alpha} = o_P(1)$ for $\alpha = f, g$ 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,

$$D_f(\cdot; b_h) = O_P \left(\sup_{y \in Y} R_2^k + B_y + 1 = \frac{p}{nh^d} + 1 = \frac{p}{\bar{n}} \right); \quad (\text{B.9})$$

562 where R_2^k is defined in Eq. (13) and B_y is defined in Eq. (14).

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

$$\begin{aligned} D_f(\cdot; b_h) &= \int_Y^Z w(y) \hat{h}_h(y) d[y] \\ &= \int_Y^Z w(y) O_P \left(\frac{p}{nh^d} + R_2^k + 1 = \frac{p}{\bar{n}} + B_y \right) d[y] \\ &= O_P \left(\frac{p}{nh^d} + 1 = \frac{p}{\bar{n}} \right) + \int_Y^Z (w(y) O_P(R_2^k) + B_y) d[y] \\ &= O_P \left(\frac{p}{nh^d} + 1 = \frac{p}{\bar{n}} \right) + O_P \left(\sup_{y \in Y} R_2^k + B_y \right); \end{aligned}$$

563

□

564 Corollary 3 (Restated Corol. 3) Let $h = O(n^{-(d+4)})$. Suppose D_f satisfies $f(p; q) = 0$ if $p = q$.
565 Suppose $w(y)$ in Lemma 6 is finite. If nuisances $\alpha, \hat{\alpha}, \hat{g}$ converges at $nh^d g^{-1=4}$ rate for any
566 $(w; z; x; y)$, then $D_f(\cdot; b_h)$ converges to 0 at nh^d -rate.

567 Proof. This result follows immediately from Thm. 2.

□

568 B.2 Proofs for Sec. 4

569 We will use ρ to denote ρ as a functional ρ . Let p denote a parametric submodel. We will use
570 S to denote a score function f .

Lemma B.7 (Restated Lemma 7) An influence function $m(\cdot; \cdot)$ in Eq. (18), denoted m , is given by

$$m(\cdot; \cdot) = f(\cdot; \cdot; \cdot; \cdot) - \int (f(\cdot; \cdot; \cdot; \cdot)) [R_f(Y; \cdot; \cdot)]; \quad (\text{B.10})$$

where $(\cdot; \cdot)[\cdot]$ is defined in Eq. (9), and

$$R_f(Y; \cdot; \cdot) = g^0(Y; \cdot) f f_{21}^{00}(\cdot(Y); g(Y; \cdot)) g(Y; \cdot) + f_1^0(\cdot(Y); g(Y; \cdot)) g;$$

571 where $g^0(y; \cdot) = g(y; \cdot)$, $f_1^0(p; q) = f_1(p; q)$ and $f_{21}^{00}(p; q) = f_{21}^0(p; q)$.

572 Proof. Let η denote the estimand written w.r.t. the parametric submodel $\eta = p(1 + g)$ where g
 573 is a bounded mean-zero random function. Set $\tau(\theta) = \int_{\mathcal{Y}} \log p$.

Let

$$m(y; \theta) = g^0(y; \theta) f_2^0(y; g(y; \theta)) g(y; \theta) + f(y; g(y; \theta)) g(y; \theta) \quad (\text{B.11})$$

574 Note $m(\theta) = \int_{\mathcal{Y}} m(y; \theta) d[y]$. Also, we note that $\tau(\theta) = \int_{\mathcal{Y}} m(y; \theta) d[y]$.

Also, recall that an influence function for $f(Y)$ (for $f(Y) < 1$) is given as $\psi_f(Y) = f(Y) - \int_{\mathcal{Y}} f(Y) d[Y]$ in Lemma S.3. Then, by the definition of the influence function $\psi_f(Y)$ satisfies the following,

$$\int_{\mathcal{Y}} \psi_f(Y) d[Y] = E_P[\psi_f(Y)] = 0$$

Now, we will prove that $m(\theta) = \int_{\mathcal{Y}} m(y; \theta) d[y]$ is a functional satisfying

$$\int_{\mathcal{Y}} m(y; \theta) d[y] = E_P[\psi_f(Y)] = 0$$

575 then this equation implies that $\psi_f(Y)$ is an influence function for the score $\tau(\theta)$.

This can be shown as follows:

$$\begin{aligned} & \int_{\mathcal{Y}} m(y; \theta) d[y] \\ &= \int_{\mathcal{Y}} \int_{\mathcal{Y}} m(y; \theta) d[y] d[\theta] \\ &= \int_{\mathcal{Y}} \int_{\mathcal{Y}} m(y; \theta) d[\theta] d[y] \\ &= \int_{\mathcal{Y}} \int_{\mathcal{Y}} (y) \psi_f(y) d[\theta] d[y] \\ &= \int_{\mathcal{Y}} \int_{\mathcal{Y}} (y) R_f(y; \theta) d[\theta] d[y] \\ &= \int_{\mathcal{Y}} R_f(Y; \theta) d[Y] \\ &= E_P[\psi_f(Y)] = 0 \end{aligned}$$

576 where the first equality is by the definition of m , the second by the exchange of derivation/integration,
 577 the third by the chain rule, the fourth by the fact that $\tau(\theta) = \int_{\mathcal{Y}} m(y; \theta) d[y]$ and the
 578 exchange of derivation/integration, the fifth by the definition of $\psi_f(Y)$ in Eq. (9), the sixth by the
 579 definition of the influence function (i.e., the influence function for $f(Y)$ is a function $\psi_f(Y)$
 580 satisfying $\int_{\mathcal{Y}} \psi_f(Y) d[Y] = E_P[\psi_f(Y)] = 0$).

581

□

Lemma B.8 (Restated Lemma 8) A Neyman orthogonal score for estimating η denoted $\psi(\theta; \eta)$ is given by

$$\psi(\theta; \eta) = m(\theta) + m(\theta; \eta) \quad (\text{B.12})$$

582 where $m(\theta; \eta)$ is defined in Eq. (19).

583 Proof. We first note that $E_P[m(\theta; \eta)] = m(\theta; \eta)$, because this is not a random function. Then, the
 584 influence function for $E_P[m(\theta; \eta)]$ is given by Lemma 7. For any score function which expectation
 585 is zero at the true parameter, its addition with the influence function is a Neyman orthogonal score
 586 [14, Thm.1]. That is $m(\theta; \eta) + m(\theta)$ is a Neyman orthogonal score. □

Theorem B.3 (Restated Thm. 3) Let $\psi(\theta; \eta)$ be given in Eq. (20). Let $m(\theta; \eta)$ be given in Eq. (19). Let θ_0 denote the true parameters. Let $\hat{\theta}_n$ be the MLTE estimator for θ_0 defined in Def. 3. Suppose (1) $R_f(y; \theta)$ is bounded and $R_f^0(y; \theta) = \int_{\mathcal{Y}} R_f(y; \theta) d[\theta] < 1$; (2) There exists a function $H(y) < 1$ s.t. $\sup_y \max_{\theta} R_f(y; \theta); R_f^0(y; \theta) g = O(H(y))$; (3) $f'(\theta; \eta) g$ is Donsker w.r.t. for the fixed θ ; (3) The estimators are consistent: $\hat{\theta}_n \rightarrow \theta_0 = o_P(1)$ and $k_n^{-1} \hat{\theta}_n \rightarrow 0 = o_P(1)$ for $\int_{\mathcal{Y}} f(z; w); x(z; w); (x; z; w)[H(Y)]g$ for all $(w; z; x; y)$; and (4)

$E_P[\cdot; (\cdot; \cdot)]$ is differentiable w.r.t. at $\theta = \theta_0$ with non-singular matrix $M(\theta_0; (\cdot; \cdot))$ ($\partial = \partial_j = \partial_{\theta_j}$) for all $(\cdot; \cdot)$, where $M(\theta_0; (\hat{\theta}; \hat{\theta})) = M(\theta_0; (\theta_0; \theta_0))$. Then,

$$b_{\theta_0} = M^{-1} E_D[m(\theta_0; (\theta_0; \theta_0))] + o_P(n^{-1/2}) + O_P(R_2^m);$$

where

$$R_2^m = \sum_z k_z^{\wedge} z k_z^{\wedge n} z + \sum_z \hat{z} z^0 + \sum_z \hat{z} z^2 + \sum_z \hat{z} z^2 + \sum_z \hat{z} z z \hat{z};$$

587 where $z = z(W)$, $z = x(z; W)$, and $z = (x; z; W)[H(Y)]$.

Proof. We follow the proof strategy used in [39, Lemma 1, Thm.3]. First,

$$\begin{aligned} b_{\theta_0} &= M^{-1} E_D[m(\theta_0; (\theta_0; \theta_0))] - M^{-1} E_P[h(\theta_0; (\hat{\theta}; \hat{\theta}))] + o_P(n^{-1/2}) \\ &= M^{-1} E_D[m(\theta_0; f_{\theta_0}; g_{\theta_0})] - M^{-1} E_P[h(\theta_0; (\hat{\theta}; \hat{\theta}))] + o_P(n^{-1/2}); \end{aligned} \quad (B.13)$$

where the first equality holds by Lemma S.1, and the second holds since $m(\theta_0; \theta_0) = 0$ by the moment condition in Eq. (18). Since $E_D[m(\theta_0; \theta_0; \theta_0)]$ converges to $N(0; \text{var}(\frac{2}{m}))$ in distribution at \sqrt{n} -rate, the only remaining term to analyze is

$$E_P[h(\theta_0; (\hat{\theta}; \hat{\theta}))] = m(\theta_0; \hat{\theta}) + E_P[h(\theta_0; (\hat{\theta}; \hat{\theta}))][R_f(Y; \theta_0; \hat{\theta})]; \quad (B.14)$$

which can be analyzed as

$$\begin{aligned} &E_P[h(\theta_0; (\hat{\theta}; \hat{\theta}))][R_f(Y; \theta_0; \hat{\theta})] \\ &= E_P\left[\frac{1}{\hat{\Lambda}_X} \sum_x \hat{V}_{YX}[R_f(Y; \theta_0; \hat{\theta})] \hat{\Lambda}[R_f(Y; \theta_0; \hat{\theta})] \hat{V}_X^0\right] \\ &= E_P\left[\frac{1}{\hat{\Lambda}_X} \hat{V}_{YX}[R_f(Y; \theta_0; \hat{\theta})] - E_P\left[\frac{1}{\hat{\Lambda}_X} \hat{\Lambda}[R_f(Y; \theta_0; \hat{\theta})] \hat{V}_X\right]\right] \\ &= E_P\left[\frac{1}{\hat{\Lambda}_X} \frac{\hat{\Lambda}_{z^x}(W)}{z^x(W)} \sum_x (x; z^x; W)[R_f(Y; \theta_0; \hat{\theta})] \hat{\Lambda}(x; z^x; W)[R_f(Y; \theta_0; \hat{\theta})]^0 + \hat{\Lambda}(x; z^x; W) R_f(Y; \theta_0; \hat{\theta})\right] \end{aligned} \quad (B.15)$$

$$E_P\left[\frac{1}{\hat{\Lambda}_X} \frac{\hat{\Lambda}_{z^{1-x}}(W)}{z^{1-x}(W)} \sum_x (x; z^{1-x}; W)[R_f(Y; \theta_0; \hat{\theta})] \hat{\Lambda}(x; z^{1-x}; W)[R_f(Y; \theta_0; \hat{\theta})]^0 + \hat{\Lambda}(x; z^{1-x}; W)[R_f(Y; \theta_0; \hat{\theta})]\right] \quad (B.16)$$

$$E_P\left[\frac{1}{\hat{\Lambda}_X} \hat{\Lambda}[R_f(Y; \theta_0; \hat{\theta})] \frac{z^x(W)}{\hat{\Lambda}_{z^x}(W)} \sum_x (z^x; W) \hat{\Lambda}_x(z^x; W)^0 + \hat{\Lambda}_x(z^x; W)\right] \quad (B.17)$$

$$+ E_P\left[\frac{1}{\hat{\Lambda}_X} \hat{\Lambda}[R_f(Y; \theta_0; \hat{\theta})] \frac{z^{1-x}(W)}{\hat{\Lambda}_{z^{1-x}}(W)} \sum_x (z^{1-x}; W) \hat{\Lambda}_x(z^{1-x}; W)^0 + \hat{\Lambda}_x(z^{1-x}; W)\right]; \quad (B.18)$$

588

where

$$(B.15) = E_P\left[\frac{1}{\hat{\Lambda}_X} \frac{\hat{\Lambda}_{z^x}(W)}{z^x(W)} \sum_x (x; z^x; W)[R_f(Y; \theta_0; \hat{\theta})] \hat{\Lambda}(x; z^x; W)[R_f(Y; \theta_0; \hat{\theta})]^0\right] \quad (B.19)$$

$$+ E_P\left[\frac{1}{\hat{\Lambda}_X} \sum_x (x; z^x; W)[R_f(Y; \theta_0; \hat{\theta})]\right] \quad (B.20)$$

$$(B.16) = E_P\left[\frac{1}{\hat{\Lambda}_X} \frac{\hat{\Lambda}_{z^{1-x}}(W)}{z^{1-x}(W)} \sum_x (x; z^{1-x}; W)[R_f(Y; \theta_0; \hat{\theta})] \hat{\Lambda}(x; z^{1-x}; W)[R_f(Y; \theta_0; \hat{\theta})]^0\right] \quad (B.21)$$

$$E_P\left[\frac{1}{\hat{\Lambda}_X} \sum_x (x; z^{1-x}; W)[R_f(Y; \theta_0; \hat{\theta})]\right] \quad (B.22)$$

$$(B.17) = E_P \frac{1}{\hat{\Lambda}_X} \hat{[R_f(Y; \theta; \hat{\Lambda})]} \frac{z^x(W)}{\hat{\Lambda}_{z^x}(W)} 1^n x(z^x; W) \hat{\Lambda}_x(z^x; W)^0 \quad (B.23)$$

$$E_P \frac{1}{\hat{\Lambda}_X} \hat{[R_f(Y; \theta; \hat{\Lambda})]} x(z^x; W) \quad (B.24)$$

$$(B.18) = E_P \frac{1}{\hat{\Lambda}_X} \hat{[R_f(Y; \theta; \hat{\Lambda})]} \frac{z^{1-x}(W)}{\hat{\Lambda}_{z^{1-x}}(W)} 1^n x(z^{1-x}; W) \hat{\Lambda}_x(z^{1-x}; W)^0 \quad (B.25)$$

$$+ E_P \frac{1}{\hat{\Lambda}_X} \hat{[R_f(Y; \theta; \hat{\Lambda})]} x(z^{1-x}; W) \quad (B.26)$$

First, consider the summation of (B.20,B.22,B.24,B.26):

Eq. (B.20) + Eq. (B.22) + Eq. (B.24) + Eq. (B.26)

$$= E_P \frac{1}{\hat{\Lambda}_X} \sum_x (x; z^x; W) [R_f(Y; \theta; \hat{\Lambda})] \sum_x (x; z^{1-x}; W) [R_f(Y; \theta; \hat{\Lambda})]^0$$

$$E_P \frac{1}{\hat{\Lambda}_X} \hat{[R_f(Y; \theta; \hat{\Lambda})]} \sum_x (x; z^x; W) \sum_x (x; z^{1-x}; W)$$

$$= E_P \frac{1}{\hat{\Lambda}_X} \sum_x [R_f(Y; \theta; \hat{\Lambda})] \sum_x \hat{[R_f(Y; \theta; \hat{\Lambda})]} \quad \#$$

$$= E_P \frac{1}{\hat{\Lambda}_X} \sum_x [R_f(Y; \theta; \hat{\Lambda})] \frac{\sum_x \hat{[R_f(Y; \theta; \hat{\Lambda})]}}{\hat{\Lambda}_X} \quad \#$$

$$= E_P \frac{X}{\hat{\Lambda}_X} \frac{\sum_x [R_f(Y; \theta; \hat{\Lambda})]}{X} \frac{\sum_x \hat{[R_f(Y; \theta; \hat{\Lambda})]}}{\hat{\Lambda}_X} \quad \#$$

$$= E_P \frac{X}{\hat{\Lambda}_X} [R_f(Y; \theta; \hat{\Lambda})] \hat{[R_f(Y; \theta; \hat{\Lambda})]} :$$

$$= E_P \frac{X}{\hat{\Lambda}_X} 1 [R_f(Y; \theta; \hat{\Lambda})] \hat{[R_f(Y; \theta; \hat{\Lambda})]} + E_P^h [R_f(Y; \theta; \hat{\Lambda})] \hat{[R_f(Y; \theta; \hat{\Lambda})]}^i : \quad (B.27)$$

Then,

$$\text{Eq. (B.14)} = m(\theta; \hat{\Lambda}) + \text{Sum of (B.20, B.22, B.24, B.26)} + \text{Sum of (B.19, B.21, B.23, B.25)}$$

$$= m(\theta; \hat{\Lambda}) + E_P [R_f(Y; \theta; \hat{\Lambda})] \hat{[R_f(Y; \theta; \hat{\Lambda})]} \quad (B.28)$$

$$+ E_P \frac{X}{\hat{\Lambda}_X} 1 [R_f(Y; \theta; \hat{\Lambda})] \hat{[R_f(Y; \theta; \hat{\Lambda})]} + \text{Sum of (B.19, B.21, B.23, B.25)} \quad (B.29)$$

To analyze (B.28), we recall that $m(\theta; \hat{\Lambda}) = \int_Y R_f(y; \theta; \hat{\Lambda}) d[y]$ and $m(\theta; \hat{\Lambda}) = 0$. Also, by Taylor's expansion on $m(y; \hat{\Lambda})$ defined in Eq. (B.11),

$$m(y; \theta; \hat{\Lambda}) = \bar{m}(y; \theta; \hat{\Lambda}) + R_f(y; \theta; \hat{\Lambda}) (y - \hat{y}) + R_f^{(1)}(y; \theta; \hat{\Lambda}) (y - \hat{y})^2;$$

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

$$0 = m(\theta; \hat{\Lambda}) = m(\theta; \hat{b}) + \int_Y R_f(y; \theta; \hat{\Lambda}) (y - \hat{y}) d[y] + \int_Y R_f^{(1)}(y; \theta; \hat{\Lambda}) (y - \hat{y})^2 d[y];$$

where \hat{b} is some unknown estimand within the interval $[\hat{a}, \hat{c}]$. We obtain

$$\int_Y R_f^{(1)}(y; \theta; \hat{\Lambda}) (y - \hat{y})^2 d[y] = m(\theta; \hat{b}) + \int_Y R_f(y; \theta; \hat{\Lambda}) (y - \hat{y}) d[y];$$

By taking expectations for both sides,

$$E_P \int_Y^Z R_f^{(1)}(y; \cdot; \cdot) (y) \hat{(y)}^2 d[y] = m(\cdot; \cdot) + E_P \int_Y^Z R_f(y; \cdot; \cdot) (y) \hat{(y)} d[y] : \quad (B.30)$$

We have

$$\begin{aligned} \int_Y^Z R_f^{(1)}(y; \cdot; \cdot) (y) \hat{(y)}^2 d[y] &= O \int_Y^Z R_f^{(1)}(y; \cdot; \cdot) (y) \hat{(y)}^2 d[y] \\ &= O \int_Y^Z H(y) (y) \hat{(y)}^2 d[y] \\ &= O \int_Y^Z H^2(y) (y) \hat{(y)}^2 d[y] \\ &= O [H(Y)] \hat{[H(Y)]}^2 ; \end{aligned}$$

589 where the second equality is by the definition of $H(y)$, the third by $H(y) < 1$, and the fourth by the
590 definition of L_2 norm.

This implies that

$$(B.28) = E_P \int_Y^Z R_f^{(1)}(y; \cdot; \cdot) \hat{(y)}^2 d[y] = O [H(Y)] \hat{[H(Y)]}^2 ;$$

591 where the first equality is by Eq. (B.30) and the second equality is by the above.
Also, Sum of (B.19,B.21,B.23,B.25) in (B.29) can be written as follows:

$$\begin{aligned} &\text{Sum of (B.19,B.21,B.23,B.25)} \\ &= \sum_{z \in \mathcal{Z}} O_P \kappa_z(W) \kappa_z(W) \kappa_z(W) \hat{(x; z; W)} [R_f(Y; \cdot; \cdot)] (x; z; W) [R_f(Y; \cdot; \cdot)] + \hat{(x; z; W)} \kappa_z(W) \kappa_z(W) \\ &= \sum_{z \in \mathcal{Z}} O_P \kappa_z(W) \kappa_z(W) \kappa_z(W) \hat{(x; z; W)} [H(Y)] (x; z; W) [H(Y)] + \hat{(x; z; W)} \kappa_z(W) \kappa_z(W) : \end{aligned}$$

592

For simplicity, we assume, for any z ,

$$\begin{aligned} O_P \kappa_x(z; W) \hat{(x; z; W)} \kappa_x(z; W) \kappa_x(z; W) \hat{(x; z; W)} &= \sum_{z \in \mathcal{Z}} O_P \kappa_x(z; W) \hat{(x; z; W)}^2 ; \text{ and} \\ O_P \kappa_x(z; W) \hat{(x; z; W)} \kappa_x(z; W) \kappa_x(z; W) [H(Y)] \hat{(x; z; W)} [H(Y)] \\ &= \sum_{z \in \mathcal{Z}} O_P \kappa_x(z; W) \hat{(x; z; W)} \kappa_x(z; W) [H(Y)] \hat{(x; z; W)} [H(Y)] : \end{aligned}$$

The other part of Eq. (B.29) is given as

$$\begin{aligned}
 E_P & \frac{X}{\Lambda_X} 1 [R_f(Y; \theta_0; \hat{\Lambda}) - \hat{\Lambda} [R_f(Y; \theta_0; \hat{\Lambda})]] \\
 & = O_P \left(\frac{X}{\Lambda_X} \right) [R_f(Y; \theta_0; \hat{\Lambda}) - \hat{\Lambda} [R_f(Y; \theta_0; \hat{\Lambda})]] \\
 & = O_P \left(\frac{X}{\Lambda_X} \right) \left[\frac{Y^X [R_f(Y; \theta_0; \hat{\Lambda})] - \hat{\Lambda}^{Y^X} [R_f(Y; \theta_0; \hat{\Lambda})]}{X} + \frac{\hat{\Lambda}^{Y^X} [R_f(Y; \theta_0; \hat{\Lambda})]}{X} - \frac{\hat{\Lambda}^{Y^X} [R_f(Y; \theta_0; \hat{\Lambda})]}{\hat{\Lambda}_X} \right] \\
 & = O_P \left(\frac{X}{\Lambda_X} \right) \left[\frac{Y^X [R_f(Y; \theta_0; \hat{\Lambda})] - \hat{\Lambda}^{Y^X} [R_f(Y; \theta_0; \hat{\Lambda})]}{X} + \frac{1}{X} - \frac{1}{\hat{\Lambda}_X} \right] \\
 & = O_P \left(\frac{X}{\Lambda_X} \right) \left[\frac{Y^X [R_f(Y; \theta_0; \hat{\Lambda})] - \hat{\Lambda}^{Y^X} [R_f(Y; \theta_0; \hat{\Lambda})]}{X} + \frac{1}{X} - \frac{1}{\hat{\Lambda}_X} \right] \\
 & = O_P \left(\frac{X}{\Lambda_X} \right)^2 + O_P \left(\frac{X}{\Lambda_X} \right) \left[\frac{Y^X [R_f(Y; \theta_0; \hat{\Lambda})] - \hat{\Lambda}^{Y^X} [R_f(Y; \theta_0; \hat{\Lambda})]}{X} \right] \\
 & = O_P \left(\frac{X}{\Lambda_X} \right)^2 + O_P \left(\frac{X}{\Lambda_X} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \\
 & = \sum_z O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] + \sum_z \frac{X}{\Lambda_X(z; W)} \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \hat{\Lambda}(x; z; W) [H(Y)] ;
 \end{aligned}$$

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

Similarly we assume, for any z ,

$$\begin{aligned}
 O_P & \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \hat{\Lambda}(x; z; W) [H(Y)] \\
 & = \sum_{z \in \mathcal{Z}} O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right]^2 :
 \end{aligned}$$

We have

$$\begin{aligned}
 O_P & \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right]^2 \\
 & = O_P \left(\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right)^2 \\
 & = O_P \left(\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right)^2 + \frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \left(\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right) \\
 & = \sum_{z \in \mathcal{Z}} O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right]^2 + \sum_{z \in \mathcal{Z}} O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \\
 & + \sum_{z \in \mathcal{Z}} O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \hat{\Lambda}(x; z; W) [H(Y)] :
 \end{aligned}$$

Finally

$$\begin{aligned}
 \text{Eq. (B.14)} & = \sum_z O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \hat{\Lambda}(x; z; W) [H(Y)] + \sum_z \frac{X}{\Lambda_X(z; W)} \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right]^2 \\
 & + \sum_z O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \hat{\Lambda}(x; z; W) [H(Y)] \\
 & + \sum_z O_P \left(\frac{X}{\Lambda_X(z; W)} \right) \left[\frac{Y^X [H(Y)] - \hat{\Lambda}^{Y^X} [H(Y)]}{X} \right] \hat{\Lambda}(x; z; W) [H(Y)] : \quad (\text{B.31})
 \end{aligned}$$

Therefore, with Eq. (B.13), the following holds

$$b_0 = M^{-1} E_D [m(V; \theta_0; \hat{\Lambda})] + \text{Eq. (B.31)} + o_P(n^{-1/2});$$

594 where Eq. (B.31) $\in R_2^n$.

595

□

596 Corollary 4 (Restated Corol. 4) If nuisances $\hat{\Lambda}; \hat{g}$ converges at $n^{-1/4}$ rate, then the target
 597 estimator $\hat{\Lambda}$ converges to Λ_0 at $n^{-1/2}$ rate.

598 Proof. If all nuisances converge at $1/n^4$ rate, then the R_2^m term in Thm. 3 converges at $1/n^2$ rate.
 599 Also, $E_D [m(\theta_0; (\beta_0; \theta_0))]$ converges in distribution to $N(\theta_0; \text{var}(m(\theta_0; (\beta_0; \theta_0))))$ at $1/n$ -rate. So
 600 $\hat{\theta}$ converges to θ_0 at $1/n$ -rate by Thm. 3. \square

601 C Details of empirical applications

602 C.1 Data generating processes for synthetic datasets

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

$$U \sim N(0, 1)$$

$$f_W(U) = 2U - 1 + W; \text{ where } W \sim N(0, 1)$$

$$f_Z(W) = 1(0.25W + Z > 0); \text{ where } Z \sim N(0, 1)$$

$$f_X(W; Z; U) = 1(Z + 0.25W + 0.25U + X > 0.5) (1 - Z) + Z; \text{ where } X \sim N(0, 1):$$

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

$$f_Y(W; X; U) = 0.6501(W - (2X - 1) + 2U + 0.374):$$

For Fig. 2b,

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

For Fig. 2c,

$$f_Y(W; X; U) = 1.08541(W < 0)(2X - 1 + 0.1U) + 1(0 - W < 1)(2X - 1 + 0.1U) \\ + 1.0854 - 0.9163(1 - (W - 1))(3(2X - 1) + 0.2U + 0.3) - 0.122$$

For Fig. 2d,

$$f_Y(W; X; U) = 0.7865 - 1.0628 - 1(W < -1)(0.8(2X - 1) + 0.1U) + 1(1 - W < 0)(2(2X - 1) + 0.1U) \\ + 0.7865 - 1.0628 - 1(0 - W < 1)(2(2X - 1) + 0.2U) + 1(W > 1)(0.5(2X - 1) + 0.2U) + 0.0525 \\ + 1.0628 - 0.104$$

603 C.2 Application to 401(k) data

604 We use the 401(k) dataset that is initially introduced by Specially, we used the version of the
 605 data named 'The Woodridge Data Set' [originally entitled '401ksu.dta' in STATA format (available
 606 in <https://www.stata.com/texts/eacsap/>). In the dataset, we use netffa (net financial asset
 607 in \$1000) as Y, p401k (participation in 401(k), participation = 1) as Z, e401k (eligibility for 401(k),
 608 eligible = 1) as W, and W = f(W₁; W₂; W₃; W₄; W₅) = f(agesq, fszsq, male, marr, incsq),
 609 where agesq means the square of the age, fszsq the family size, male the gender (male = 1), marr
 610 the marital status (married = 1) and incsq the square of the income.

611 C.3 Density plots illustrating uncertainty

612 In this section, we present the density plots corresponding to Figs. (3,4) illustrating uncertainty of the
 613 results. The same data generating processes as used for Figs. (3,4) are leveraged.

Synthetic dataset. To represent the uncertainty, we generate 100 synthetic datasets $\{D_k\}_{k=1}^{100}$, each of which has $N = 50000$ samples (i.e. $|D_k| = 50000$), from the same data generating process used for the simulation for Fig. 3. After learning the density estimation with each dataset, we obtain a vector of density values $\{p_1^k; p_2^k; \dots; p_m^k\}$ at m equi-spaced points for each method ('Moment', 'MLTE', 'Kernel-smoothing', 'KLTE'). For the model-based approach (Moment, MLTE), s is set to 1000. For the kernel-based approach (Kernel-smoothing, KLTE), s is set to 25. For each {model, kernel}-based approach, we have estimates of the density in the form of a matrix $\{p_{ij}^k; p_m^k\}_{k=1}^{100}$. Let p_{avgj} denote the average of p_{ij}^k for $k=1$ to 100. Let g_j denote the standard deviation of p_{ij}^k for $k=1$ to 100. Then, we take

$$p_{avgj} = \frac{1}{100} \sum_{k=1}^{100} p_{ij}^k \\ p_{upperj} = p_{avgj} + g_j \\ p_{lowerj} = p_{avgj} - g_j$$

614 A set of density plots corresponding to Fig. 3 with uncertainty information in Fig. C.5. For each
 615 density estimate in Fig. C.5, the middle dark-colored dotted line shows \mathbf{p}_{avg} , and the above,below
 616 light-colored solid line shows $\mathbf{p}_{\text{upper}}$, $\mathbf{p}_{\text{lower}}$ respectively.

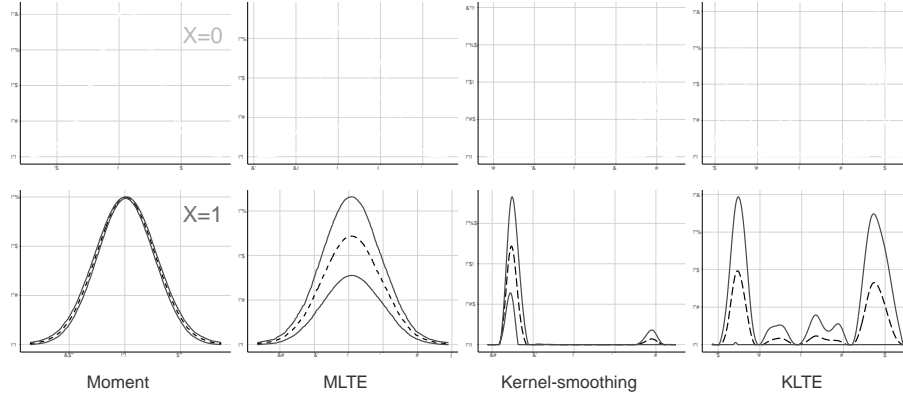


Figure C.5: LTE estimation with a synthetic dataset. The middle dark-colored dotted line denotes \mathbf{p}_{avg} , and the upper and lower light-colored solid lines represents $\mathbf{p}_{\text{upper}}$, $\mathbf{p}_{\text{lower}}$, respectively.

617 **Application to 401(k) data.** To represent the uncertainty, we randomly resample (with replacement)
 618 the dataset from the original dataset D , where the k th regenerated dataset is denoted D_k . We conducted
 619 this data regeneration process for 100 times and have $\hat{f}_{D_k} g_{k=1}^{100}$. After learning the density estimation
 620 with D_k , we obtain a vector of density values $(p_1^k, p_2^k, \dots, p_m^k)$ at m equi-spaced points for each
 621 method ('Moment', 'MLTE', 'Kernel-smoothing', 'KLTE'). For the model-based approach (Moment,
 622 MLTE), m is set to 1000. For the kernel-based approach (Kernel-smoothing, KLTE), m is set to 25.
 623 For each {model, kernel}-based approach, we have estimates of the density in the form of a matrix
 624 $\hat{f}(p_1^k, \dots, p_m^k) g_{k=1}^{100}$. Let $p_{\text{avg};i}$ denote the average of $\hat{f}_i p_i^k g_{k=1}^{100}$. Let σ_i denote the standard deviation of
 625 $\hat{f}_i p_i^k g_{k=1}^{100}$. Then, we take $\mathbf{p}_{\text{avg}} = \hat{f} p_{\text{avg};i} g_{i=1}^m$, $\mathbf{p}_{\text{upper}} = \hat{f} p_{\text{avg};i} + \sigma_i g_{i=1}^m$ and $\mathbf{p}_{\text{lower}} = \hat{f} p_{\text{avg};i} - \sigma_i g_{i=1}^m$.
 626 A set of density plots corresponding to Fig. 3 with uncertainty information in Fig. C.6. For each
 627 density estimate in Fig. C.6, the middle dark-colored dotted line shows \mathbf{p}_{avg} , and the above,below
 628 light-colored solid line shows $\mathbf{p}_{\text{upper}}$, $\mathbf{p}_{\text{lower}}$ respectively.

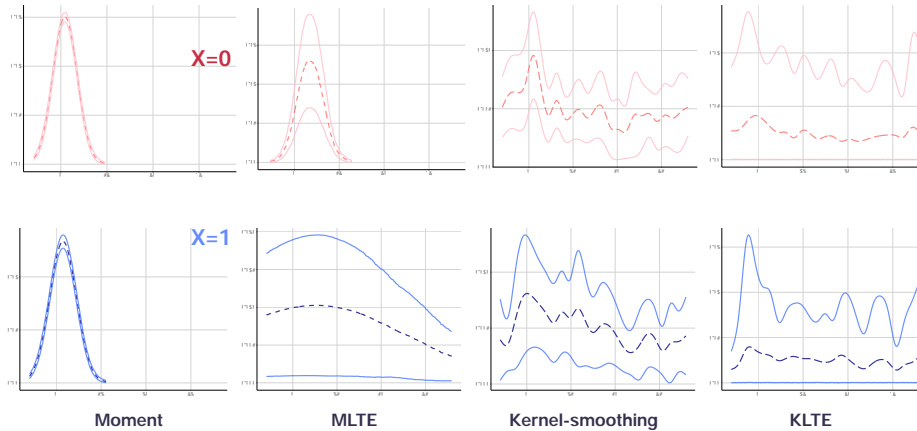


Figure C.6: LTE of 401(k) participation (X) on net financial asset (Y). The middle dark-colored dotted line denotes \mathbf{p}_{avg} , and the upper and lower light-colored solid lines represents $\mathbf{p}_{\text{upper}}$, $\mathbf{p}_{\text{lower}}$, respectively.