# Inference on Welfare and Value Functionals under Optimal Treatment Assignment[*]

Xiaohong Chen[†], Zhenxiao Chen[‡], and Wayne Yuan Gao[§]

October 30, 2025

## Abstract

We provide theoretical results for the estimation and inference of a class of welfare and value functionals of the nonparametric conditional average treatment effect (CATE) function under optimal treatment assignment, i.e., treatment is assigned to an observed type if and only if its CATE is nonnegative. For the optimal welfare functional defined as the average value of CATE on the subpopulation with nonnegative CATE, we establish the $\sqrt{n}$ asymptotic normality of the semiparametric plug-in estimators and provide an analytical asymptotic variance formula. For more general value functionals, we show that the plug-in estimators are typically asymptotically normal at the 1-dimensional nonparametric estimation rate, and we provide a consistent variance estimator based on the sieve Riesz representer, as well as a proposed computational procedure for numerical integration on submanifolds. The key reason underlying the different convergence rates for the welfare functional versus the general value functional lies in that, on the boundary subpopulation for whom CATE is zero, the integrand vanishes for the welfare functional but does not for general value functionals. We demonstrate in Monte Carlo simulations the good finite-sample performance of our estimation and inference procedures, and conduct an empirical application of our methods on the effectiveness of job training programs on earnings using the JTPA data set.

# 1 Introduction

In this paper, we study the estimation and inference on welfare and value functionals of a given treatment under optimal ("first-best") treatment assignment.

Let $D_i \in \{0,1\}$ denote a certain binary treatment for subject $i$, $(Y_i(0), Y_i(1))$ denote the potential outcomes of interest, and $Y_i := Y_i(D_i) \in \mathbb{R}$ denote the observed outcome. Let $X_i \in \mathbb{R}^d$ denote subject $i$'s observable characteristics of $i$, which is distributed with density $f_0$. We suppose that researchers have access to a random sample of training data $\{(D_i, Y_i, X_i)\}_{i=1}^n$.

Under the standard conditional unconfoundedness assumption $(Y_i(0), Y_i(1)) \perp D_i | X_i$ and the overlap condition $p_0(x) := \mathbb{E}[D_i| X_i = x] \in (0,1)$, the conditional average treatment effect (CATE) defined by

$$\text{CATE}(x) := \mathbb{E}[Y_i(1) - Y_i(0)| X_i = x]$$

is identified from data by

$$\text{CATE}(x) \equiv h_0(x) := \mu_0(x,1) - \mu_0(x,0), \tag{1}$$

where $h_0 : \mathbb{R}^d \mapsto \mathbb{R}$, and

$$\mu_0(x,d) := \mathbb{E}[Y_i| X_i = x, D_i = d] \tag{2}$$

is the nonparametric regression function of the outcome $Y_i$ on $X_i$ and $D_i$. We maintain the conditional unconfoundedness assumption and the overlap condition, and will thereafter simply refer to $h_0$ as the CATE function. We further assume that $h_0$ belongs to a Holder class of functions with smoothness $s > 1$.

We consider a standard scenario where policymakers can assign treatments based on co-variates, and focus on the following two core types of welfare and value parameters. The first type is the maximized welfare of the target population under optimal treatment assignment:

$$W(h_0) := \int [h_0(x)]_+ f(x) \, dx \tag{3}$$

where $f$ is the marginal density of $x$ in the target population and $[t]_+ := \max(t,0)$ is the rectified linear unit (ReLU) function. $W(h_0)$ averages the CATE over the population under the "first-best" treatment assignment rule: "treat type $x$ if and only if $\text{CATE}(x) \geq 0$," and is thus often referred to as the welfare under optimal treatment assignment. We will thereafter

refer to $W(h_0)$ as the welfare functional.

Alternatively, one may also be interested in evaluating the average of a value other than CATE over the population under optimal treatment assignment:

$$V(h_0) := \int \mathbb{1}\{h_0(x) \geq 0\} v_0(x) f(x) dx \qquad (4)$$

where $v_0 : \mathbb{R}^d \mapsto \mathbb{R}$ is a user-defined function that may be a utility function, a cost function, or any economically meaningful function of the observed covariate $x$. For example, setting $v_0(x) = a'x$ endows $V(h_0)$ with the interpretation as certain aggregate characteristics of the treated population under optimal treatment assignment, and setting $v_0 \equiv 1$ implies that $V(h_0) \equiv \mathbb{P}_f(h_0(X_i) \geq 0)$ becomes the share of the target population to be treated (i.e. with nonnegative CATE).

In the formulation of $W(h_0)$ and $V(h_0)$ above, we take the density $f(x)$ to be known. This is in itself relevant in settings where the covariate density of the target population is configured or known/estimated from other sources than the training sample used to estimated CATE. For example, CATE may be estimated from a smaller pilot program, while the policymakers are contemplating to implement the policy on a statewide or nation-wide basis with a much larger population. That said, in this paper we also consider an important case where $f$ is unknown and set to be the covariate density in the underlying population of the training sample. In this case, $f$ does not need to be estimated, as the integral with respect to $f$ can be naturally approximated via sample average in the training sample. This corresponds more closely to the "empirical welfare" as considered in Kitagawa and Tetenov (2018). We provide results for this setting as well.

In this paper, we establish inference results for the welfare and value functionals $W(h_0)$ and $V(h_0)$ with nonparametric estimated CATE $h_0$. The results can be summarized informally as follows.

For the welfare functional $W(h_0)$, we establish semiparametric plug-in estimators of the welfare functional are asymptotically normal at the parametric $\sqrt{n}$ rate and derive closed-form asymptotic variance formulas, along with consistent asymptotic variance estimators. The key insight of the $\sqrt{n}$ rate is geometric: by the definition of the welfare functional, the integrand $h_0(x)$ vanishes on the boundary of the integration $\{x : h_0(x) = 0\}$, neutralizing the non-smoothness of the indicator function $\mathbb{1}\{h_0(x) \geq 0\}$.

In contrast, for the value functional $V(h_0)$ with general weight $v_0$ that does not vanish on the boundary $\{x : h_0(x) = 0\}$, we show that the rate of convergence is slower than $\sqrt{n}$, and is instead given by the 1-dimensional nonparametric regression rate $n^{-\frac{s}{2s+1}}$ under appropriate conditions. We establish asymptotic normality of the semiparametric plug-in estimator under

this irregular convergence rate, and provide a consistent variance estimator based on the sieve Riesz representer. In particular, the consistent variance estimator features a Hausdorff integral on the boundary submanifold $\{x : h_0(x) = 0\}$, for which we provide a numerical integration and differentiation procedure for the computation of submanifold integrals.

We conduct an array of Monte Carlo experiments to document the good finite-sample accuracy of our theoretical inferential results. We show that the proposed standard error estimators perform well in finite sample, and the corresponding confidence intervals based on the asymptotic normality result and the standard error estimators have coverage probabilities close to their nominal levels. These findings hold not only for the $\sqrt{n}$-estimable welfare functional, but also for the value functional, which is estimated at slower-than-$\sqrt{n}$ rate with standard errors computed through numerical integration and differentiation.

We also apply our results to empirical data from the Job Training Partnership Act (JTPA) data set. Following Kitagawa and Tetenov (2018), we take 30-month post-program earning as the outcome variable and consider two covariates: pre-program earning and education. We then provide empirical estimates and confidence intervals for two parameters: the welfare under first-best treatment assignment, which is $\sqrt{n}$ estimable, and the share of population to be treated under first-best treatment assignment, which is not $\sqrt{n}$-estimable. As in Kitagawa and Tetenov (2018), we also consider two different scenarios: one with the cost of the treatment incorporated, and one without. These parameters have also been estimated in Kitagawa and Tetenov (2018) under the label of "nonparametric plug-in rule" using kernel first-stages, but Kitagawa and Tetenov (2018) only provide point estimates with no confidence intervals for them. We use sieve (B-spline) first-stage nonparametric estimators and find similar results to those in Kitagawa and Tetenov (2018), and further provide informative confidence intervals for both the welfare and the share parameters.

**Closely Related Literature**

This work is a companion paper to Chen and Gao (2025), and directly applies the general theoretical results there to handle differentiation with respect to region of integration and semiparametric estimation of integrals over submanifolds. Specifically, this work focuses on the important context of treatment assignment problems, and deal with two features that are not discussed in Chen and Gao (2025). First, CATE is defined as the difference of two nonparametric regression functions between the treated and untreated subpopulation (or the difference of two point evaluations of a nonparametric function with treatment status defined as an argument too), and is often estimated as the difference of two first-stage nonparametric estimators. Second, in treatment assignment problems, researchers may face two different scenario in terms of the distribution of the covariates: sometimes, this dis-

tribution can be treated as known and may be different from the covariate distribution in the experimental/observational population from which CATE is estimated; in other times, one may want to treat the covariate distribution as unknown and the same as the experimental/observational population, and uses sample averages to automatically incorporate the covariate distribution. In this paper, we takes into account these special structures of the problem, and establish the inference results by providing lower-level sufficient conditions to the general theory in Chen and Gao (2025).[1]

Our paper makes new contributions to the literature on estimation and inference on a general value functional of a policy under optimal treatment assignment. To the best of our knowleadge, all the existing work on limiting distributions of functionals of a policy under optimal treatment assignments considered $\sqrt{n}$-normality only. Our paper is the first to establish slower-than-root-n limiting distribution and inference results for irregular value functionals of a policy under optimal treatment assignment. Previously, Bhattacharya and Dupas (2012) establishes $\sqrt{n}$-normality of the optimal welfare value under budget constraint, using Nadaraya-Waston kernel estimator for first-stage estimation of CATE. The important work of Kitagawa and Tetenov (2018) focuses on a related but slightly different topic: empirical welfare maximization within a constrained class of policy rules. That said, Kitagawa and Tetenov (2018) also considers and reports empirical estimates on the nonparametric plug-in rule, which can be interpreted as a plug-in estimator of the first-best welfare, though there were no theoretical results or confidence intervals for this estimate. The recent papers by Park (2025) and Whitehouse, Austern and Syrgkanis (2025) use soft-max functions to smooth the max function, and apply the debiased machine learning approach to establish asymptotic normality and provide inferential results, and Whitehouse, Austern and Syrgkanis (2025) establishes $\sqrt{n}$ asymptotic normality of their soft-max welfare functional. Feng, Hong and Nekipelov (2025) analyzes binary treatment assignment under constraints and the asymptotic property of the welfare of the policy under optimal cutoff choice, and establishes root-$n$ asymptotic normality of the functionals. Our paper complements these existing work in several ways. First, we clarify that the $\sqrt{n}$-estimability of the welfare functional is due to the fact that the integrand (CATE) by construction vanishes on the boundary of the optimally treated population $\{x : \mathrm{CATE}\,(x) = 0\}$, which is specific to the welfare functional but generally not satisfied for other types of value functionals. Second, we demonstrate that the

---

[1]Two concurrent papers by Cattaneo, Titiunik and Yu (2025a,b) also feature submanifold integrals, but focuses on 1-dimensional cases that arise from the specific context of boundary discontinuity designs. Our current paper also differs significantly from Cattaneo, Titiunik and Yu (2025a,b), who focus on boundary discontinuity designs and boundary treatment effects, which are very different objects from the welfare and value functionals considered here under first-best treatment assignments. Consequently, the submanifolds (boundaries) in their settings are given or known based on the locations or a distance function, while in our current project the boundary submanifold is defined by the unknown and estimated CATE function.

$\sqrt{n}$-estimability of the welfare functional can be attained without the use of smoothing/soft-max functions. Third and most importantly, our results extend well beyond the welfare functional, and cover generic value functionals that may be slower than $\sqrt{n}$-estimable.

The rest of the paper is organized as follows. Section 2 lays out the main model setup, and provides a conceptual explanation of why the welfare functional $W(h_0)$ can be $\sqrt{n}$-estimable while the value functional $V(h_0)$ is not in general. Section 3 then establishes the inference results for the welfare functional $W(h_0)$, while Section 4 provides corresponding results for the value function $V(h_0)$. We report numerical results from Monte Carlo simulations in Section 5, and conduct an empirical illustration in Section 6. Proofs of theoretical results in the main text are available in Appendix A.

# 2  Model Setup and Functional Derivatives

We first introduce the standard treatment effect model, the general welfare functional of interest and the maintained assumptions in Subsection 2.1

## 2.1  The Model and the Parameters of Interest

We first state the maintained assumption in the paper. Let $\hat{h}(x) := \hat{\mu}(x, 1) - \hat{\mu}(x, 0)$ be a nonparametric estimator of the CATE function, in which $\hat{\mu}(x, d)$ is a first-stage estimator of the nonparametric regression model

$$Y_i = \mu_0(X_i, D_i) + \epsilon_i, \quad \mathbb{E}[\epsilon_i | X_i, D_i] = 0, \quad \mathbb{E}\left[\epsilon_i^2 \Big| X_i, D_i\right] < \infty. \tag{5}$$

We impose the following basic assumptions in this paper.

**Assumption 1** (Model)**.** *The training data and the model satisfy:*

*(a) Training sample: the training data $\{(Y_i, D_i, X_i)\}_{i=1}^n$ is a random sample drawn from $(Y, D, X) \in \mathbb{R} \times \{0, 1\} \times \mathcal{X}$ satisfying Model (18), where $\mathcal{X}$ is a bounded rectangular set (say $[0, 1]^d$) in $\mathbb{R}^d$, and $X_i$ has its true unknown marginal density $f_0$ supported on $\mathcal{X}$.*

*(b) Overlap: $0 < p_0(x) := \mathbb{E}[D_i | X_i = x] < 1$.*

*(c) Smoothness of Regression Function: For $d \in \{0, 1\}$, $\mu_0(\cdot, d) \in \Lambda^s(\mathcal{X})$ with $s > 1$.*

We first provide a heuristic overview of our theoretical analysis, and explain the key intuition why the welfare functional $W(h_0)$ may be $\sqrt{n}$-estimable (i.e. $W$ is a regular functional) while $V(h_0)$ is generally not (i.e., $V$ is a irregular functional).

We first introduce a more general value functional $\Phi(h_0)$ that nests $W(h_0)$ and $V(h_0)$ as special cases:

$$\Phi(h_0) := \int \mathbb{1}\{h_0(x) \geq 0\} \phi(h_0(x), x) f(x) \, dx, \tag{6}$$

where $\phi : \mathbb{R} \times \mathbb{R}^d \to \mathbb{R}$ is a known measurable mapping. We note that

- $\Phi(h_0) = W(h_0)$ when $\phi(h_0(x), x) = h_0(x)$;

- $\Phi(h_0) = V(h_0)$ when $\phi(h_0(x), x) = v_0(x)$ with $\partial_1 \phi(h_0(x), x) = 0$.

**Assumption 2** (Functional)**.** *The functional $\Phi$ satisfies*

*(a) $\phi : \mathbb{R} \times \mathbb{R}^d \to \mathbb{R}$ is continuously differentiable with respect to its first argument.*

*(b) Target Density: The target density $f$ of $X$ is absolutely continuous w.r.t. $f_0$ with uniformly bounded Radon-Nikodym derivative $\lambda := f/f_0$.*

*(c) Regular Level Set: $h_0(\cdot) := \mu_0(\cdot, 1) - \mu_0(\cdot, 0)$ satisfies $\|\nabla_x h_0(x)\| \geq \underline{c} > 0$ on the level set $\{x \in \mathcal{X} : h_0(x) = 0\}$.*

## 2.2 Functional Derivatives via Generalized Leibniz rule

By the standard semiparametric theory on the estimation of functionals of nonparametric regression functions, the asymptotic property of the semiparametric plug-in estimator $\Phi(\hat{h})$ can be analyzed via the functional derivative of $\Phi(h)$ w.r.t. $h_0$ in the direction of $h - h_0$, i.e., writing $h_t := h_0 + t(h - h_0)$,

$$D_h \Phi(h_0)[h - h_0] := \frac{d}{dt} \Phi(h_t) \Big|_{t=0} = \frac{d}{dt} \int \mathbb{1}\{h_t(x) \geq 0\} \phi(h_t(x), x) f(x) \, dx \Big|_{t=0}. \tag{7}$$

$$D_h W(h_0)[h - h_0] := \frac{d}{dt} W(h_t) \Big|_{t=0} = \frac{d}{dt} \int [h_t(x)]_+ f(x) \, dx \Big|_{t=0} \tag{8}$$

$$D_h V(h_0)[h - h_0] := \frac{d}{dt} V(h_t) \Big|_{t=0} = \frac{d}{dt} \int \mathbb{1}\{h_t(x) \geq 0\} v_0(x) f(x) \, dx \Big|_{t=0} \tag{9}$$

The presence of the ReLU/max function $[t]_+$ in $D_h W(h_0)$ and the indicator function $\mathbb{1}\{t \geq 0\}$ induces a point of nonsmoothness at $t = 0$, where $[t]_+$ is nondifferentiable and $\mathbb{1}\{t \geq 0\}$ is discontinuous. This complicates the calculation of the functional derivatives in (8) and (9), though to different degrees.

We now provide an overview of the key difference between the welfare and value functions from the perspective of the generalized Leibniz rule, which has the following generic form regarding the total time derivative of integrals with changing integrand and changing region

of integration:

$$\frac{d}{dt} \int_{\Omega_t} G_t(x)\, dx. \tag{10}$$

The generalized Leibniz rule,[2] states that, under mild regularity conditions,

$$\frac{d}{dt} \int_{\Omega_t} G_t(x)\, dx = \underbrace{\int_{\Omega_t} \frac{\partial}{\partial t} G_t(x)\, dx}_{\textbf{(I)}} + \underbrace{\int_{\partial \Omega_t} \langle \mathbf{n}_t(x), \mathbf{v}_t(x) \rangle G_t(x)\, dS_{\partial \Omega_t}(x)}_{\textbf{(II)}} \tag{11}$$

where term (I) captures the effect of the change in the integrand $G_t(x)$ with the region of integration $\Omega_t$ held fixed, while term (II) captures the effect of the change in the region of integration $\Omega_t$ with the integrand $G_t(x)$ held fixed. The somewhat "nonstandard" term (II) warrants some more explanations: $\partial \Omega_t$ denotes the boundary of $\Omega_t$, $\mathbf{n}_t(x)$ is the outward-pointing unit normal vector, $\mathbf{v}_t(x)$ is the velocity vector associated with the time movement in the $\partial \Omega_t$, and $S_{\partial \Omega_t}(x)$ denotes the surface measure on the boundary $\partial \Omega_t$. Note that, when $x$ is one-dimensional and $\Omega_t = [a_t, b_t]$, (11) specializes to the standard Leibniz rule:

$$\frac{d}{dt} \int_{a_t}^{b_t} G_t(x)\, dx = \int_{a_t}^{b_t} \frac{\partial}{\partial t} G_t(x)\, dx + G_t(b_t) \frac{d}{dt} b_t - G_t(a_t) \frac{d}{dt} a_t.$$

We observe that $D_h \Phi(h_0)$, $D_h W(h_0)$ and $D_h V(h_0)$ are all of the form (10), with the same parametrized region of integration

$$\Omega_t := \left\{ x \in \mathbb{R}^d : h_t(x) \geq 0 \right\}, \quad \Omega_0 := \left\{ x \in \mathbb{R}^d : h_0(x) \geq 0 \right\}$$

and $\partial \Omega_0 = \left\{ x \in \mathbb{R}^d : h_0(x) = 0 \right\}$ under mild regularity conditions on $h_0$.

Applying the generalized Leibniz rule (11) to the functional $\Phi(h)$, we obtain:

$$D_h \Phi(h_0)[h - h_0] = \underbrace{\int_{\Omega_0} \partial_1 \phi(h_0(x), x)(h(x) - h_0(x)) f(x)\, dx}_{\textbf{(I)}}$$
$$+ \underbrace{\int_{\partial \Omega_0} \langle \mathbf{n}_0(x), \mathbf{v}_0(x) \rangle \phi(h_0(x), x) f(x)\, dS_{\partial \Omega_0}(x)}_{\textbf{(II)}} \tag{12}$$

where the first term (I) is a full-dimensional Lebesgue integral (in $\mathbb{R}^d$), while the second term (II) is a lower-dimensional boundary integral. In particular, when $\phi(h_0(x), x) f(x)$ does not vanish on $\partial \Omega_0 = \left\{ x \in \mathbb{R}^d : h_0(x) = 0 \right\}$, the second term (II) cannot be ignored, despite $\partial \Omega_0$ has Lebesgue measure zero in $\mathbb{R}^d$. In fact, under Assumption 2(c) (see, e.g., Chen and Gao (2025)), the second term (II) of (12) can be expressed as

$$\int_{\partial \Omega_0} \langle \mathbf{n}_0(x), \mathbf{v}_0(x) \rangle \phi(h_0(x), x) f(x)\, dS_{\partial \Omega_0}(x) = \int_{\partial \Omega_0} \frac{h(x) - h_0(x)}{\|\nabla_x h_0(x)\|} \phi(h_0(x), x) f(x)\, d\mathcal{H}^{d-1}(x),$$
$$\tag{13}$$

---

[2]See, for example, Theorem 4.2 of Delfour and Zolésio (2001).

8

where $\mathcal{H}^{d-1}$ denotes the $(d-1)$ dimensional Hausdorff measure (see Chen and Gao (2025)), which coincides with the $(d-1)$ dimensional Lebesgue measure in $\mathbb{R}^{d-1}$. Thus the boundary integral term (II) of (12) is a lower-dimensional integral functional that only extracts information about $h_0$ on a Lebesgue measure-0 set (in $\mathbb{R}^d$).

For the welfare function $\Phi(h_0) = W(h_0)$, plugging $\phi(h_0(x), x) = h_0(x)$ into (12) yields

$$D_h W(h_0)[h - h_0] = \int \mathbb{1}\{h_0(x) \geq 0\}[h(x) - h_0(x)]f(x)\,dx, \qquad (14)$$

where the term (II) vanishes since $\phi(h_0(x), x) = h_0(x) = 0$ on the boundary $\partial\Omega_0$, and the term (I) is a full-dimensional Lebesgue integral functional of $h - h_0$.

For the value function $\Phi(h_0) = V(h_0)$, plugging $\phi(h_0(x), x) = v_0(x)$ with $\partial_1 \phi(h_0(x), x) = 0$ into (12) and (13) yields

$$D_h V(h_0)[h - h_0] = \int_{\{x \in \mathbb{R}^d : h_0(x) = 0\}} \frac{(h(x) - h_0(x))}{\|\nabla_x h_0(x)\|} v_0(x)f(x)\,d\mathcal{H}^{d-1}(x). \qquad (15)$$

which is not 0 as long as $v_0(x)f(x)$ does not vanish on the boundary $\partial\Omega_0 = \{x \in \mathbb{R}^d : h_0(x) = 0\}$. For example, setting $v_0(x) \equiv 1$ yields $V(h_0) = P_f(h_0(X_i) \geq 0)$, the share of population with nonnegative CATE, and the boundary integral term (II) does not vanish. In fact, as long as $v_0(x)f(x) \neq 0$ on the boundary $\partial\Omega_0 = \{x \in \mathbb{R}^d : h_0(x) = 0\}$, $D_h V(h_0)$ is a non-zero $(d-1)$ dimensional integral functional that only extracts information about $h_0$ on a Lebesgue measure-0 set (in $\mathbb{R}^d$), akin to a point evaluation of a nonparametric estimation.

## 2.3   Key Difference between the Welfare and Value Functionals

Let $L^2(f)$ denote the Hilbert space of square integrable (against $f$) functions with the inner product $\langle g, h \rangle_{2,f} := \int g(x) h(x) f(x)\,dx$. For the CATE function $h_0 \in L^2(f)$, it is well-known that a linear functional $L[h - h_0]$ is bounded (or equivalently, continuous) if and only if

$$\sup_{\nu \neq 0, \nu \in L^2(f)} \frac{|L[\nu(\cdot)]|^2}{E_f[|\nu(X)|^2]} < \infty$$

which is a necessary and sufficient condition for the existence of a Riesz representer $\nu^* \in L^2(f)$ such that

$$L[\nu] = \langle \nu^*, \nu \rangle_{2,f} \quad \text{for all } \nu \in L^2(f)$$

This in turn is a necessary condition for any plug-in estimator of the linear functional $L[\hat{h} - h_0] = \langle \nu^*, \hat{h} - h_0 \rangle_{2,f}$ to converge to zero at a root-$n$ rate.

For the welfare functional, its linear directional derivative functional $L[\nu] = D_h W(h_0)[\nu]$ given in (14), we immediately see that $\nu^*(x) := \mathbb{1}\{h_0(x) \geq 0\}$ is the Riesz representer of the linear functional $D_h W(h_0)[\nu]$ in the Hilbert space $L^2(f)$, and that this Riesz representer

9

has bounded norm

$$\|\nu^*\|^2 := \int \mathbb{1}^2 \{h_0(x) \geq 0\} f(x) dx \leq 1,$$

and thus, by well-known results in, say, Chen and Liao (2014), Chen, Liao and Sun (2014) and Chen and Pouzo (2015), the linear functional $D_h W(h_0)[\nu]$ is a regular (i.e., $\sqrt{n}$-estimable) functional under appropriate conditions.

In contrast, for the general value functional, the linear functional corresponding to its directional derivative $D_h V(h_0)[\nu]$ given in (15) does not have a well-defined Riesz representer in the Hilbert space $L^2(f)$. It is well-known that, according to Lemma 3.3 of Chen and Pouzo (2015), Consequently, the functional $V$ becomes an irregular functional that cannot be estimated at $\sqrt{n}$ rate.

The above provides an intuitive explanation of why the welfare functional $W(h_0)$ is very special relative to general types of value functionals $V(h_0)$ or $\Phi(h_0)$, and clarifies why $W(h_0)$ could be $\sqrt{n}$-estimable while others generally cannot. In subsequent sections, we provide formal conditions and theorems that establish the $\sqrt{n}$-normality of plug-in estimators of the welfare functional $W(h_0)$, as well as the slower-than-$\sqrt{n}$ asymptotic normality for the value functional $V(h_0)$.

**Remark 1** (An Alternative View). *We also briefly discuss alternative view of the determinant of $\sqrt{n}$-estimability of the welfare fucntional $W(h_0)$, based on the Lipchitz continuity of the ReLU/max function $[t]_+$. We note that, under mild conditions ensuring that the level set $\{x : h_0(x) = 0\}$ has Lebesgue measure $0$, we may interchange the order of differentiation and integral based on the almost sure differentiability of $[h_t(x)]_+$ and the dominant convergence theorem:*

$$
\begin{aligned}
D_h W(h_0)[h - h_0] &= \frac{d}{dt} \int [h_t(x)]_+ f(x) dx \bigg|_{t=0} \\
&= \int \frac{d}{dt} [h_t(x)]_+ \bigg|_{t=0} f(x) dx \qquad (16) \\
&= \int \mathbb{1}\{h_0(x) \geq 0\} (h(x) - h_0(x)) f(x) dx,
\end{aligned}
$$

*which yields the same formula as in (14). However, for general value functional $V(h_0)$, the indicator function $\mathbb{1}\{t \geq 0\}$ is no longer Lipchitz, and there is no analog of (16): the functional derivative $D_h V(h_0)$ needs to be derived using the generalized Leibniz rule as described above (or its many variants or generalized forms in differential geometry and geometric measure theory).*

# 3 Estimation and Inference of the Welfare Functional

In this section, we focus on the welfare functional defined in (3), which can also be equivalently written as a functional of $\mu_0$ as follows:

$$W\left(h_0\right) := \int \left[h_0\left(x\right)\right]_+ f\left(x\right) dx$$
$$\equiv \overline{W}\left(\mu_0\right) := \int \left[\mu_0\left(x,1\right) - \mu_0\left(x,0\right)\right]_+ f\left(x\right) dx \tag{17}$$

We write out the two equivalent definitions of the functionals based on $h_0$ and $\mu_0$, since each representation has its own merit. The representations $W\left(h_0\right)$ based on the CATE function $h_0$ is clearer in terms of its interpretation: the condition $h_0\left(x\right) \geq 0$ is a direct optimal treatment assignment, and this representation has been adopted in previous work such as Kitagawa and Tetenov (2018). On the other hand, the representations $\overline{W}\left(\mu_0\right)$ based on $\mu_0$ is clearer in terms of the underlying nonparametric regression function, which is notationally easier to work with in our subsequent semiparametric asymptotic analysis.

We consider two different setups for the estimation of $W\left(h_0\right)$, depending on how we treat the marginal density $f\left(x\right)$.

In the first setup, we treat $f$ as known and the functional $W\left(\cdot\right)$ as a known transformation of $h_0$. This is relevant in cases where the covariate density $f\left(x\right)$ of the target population is either configured or known/estimated from other sources than the training sample used to estimated CATE. In the formulation of $W\left(h_0\right)$ and $V\left(h_0\right)$ above, we take the density $f\left(x\right)$ to be known. For example, CATE may be estimated from a smaller pilot program (with sample size $n$), while the policymakers are contemplating to implement the policy on a statewide or nation-wide basis with a much larger target population, whose covariate density may either be known at the population level or estimated from an alternative data source with much larger sample size $N >> n$ (so that the sampling uncertainty in the estimation of $f$ becomes negligible relative to that in the estimation of $h_0$).

In the second setup, we take $f$ to be unknown and set it to $f_0$, the covariate density in the underlying population of the training sample. In this case, $f$ does not need to be explicitly estimated, but the integral in $W\left(\cdot\right)$ with respect to $f$ can be naturally approximated via sample average in the training sample. This corresponds more closely to the "empirical welfare" as considered in Kitagawa and Tetenov (2018). We also provide results for this setting as well.

## 3.1 Welfare Functional Under Known Density $f$

We start with the first setup, where the covariate density $f$ is taken to be known and $W(\cdot)$ is treated as a deterministic known integral of $h_0$ w.r.t. $f$. In this case, we can define a simple nonparametric plug-in estimators of $W(h_0)$ as

$$W\left(\hat{h}\right) \equiv \overline{W}\left(\hat{\mu}\right) = \int \left[\hat{h}(x)\right]_+ f(x)\, dx.$$

We first state some key assumptions. Let $\hat{h}(x) := \hat{\mu}(x,1) - \hat{\mu}(x,0)$ be a nonparametric estimator of the CATE function, in which $\hat{\mu}(x,d)$ is a first-stage estimator of the nonparametric regression model

$$Y_i = \mu_0(X_i, D_i) + \epsilon_i, \quad \mathbb{E}[\epsilon_i | X_i, D_i] = 0, \quad \mathbb{E}\left[\epsilon_i^2 \Big| X_i, D_i\right] < \infty. \tag{18}$$

**Assumption 3.** *Fist-Stage Convergence:* $\|\hat{\mu} - \mu_0\|_\infty = o_p\left(n^{-1/4}\right)$.

**Theorem 1** ($\sqrt{n}$-asymptotic normality for the welfare functional). *Under Assumptions 1 and 2(b)(c), we have:*

$$\nu^*(x,d) := \mathbb{1}\{h_0(x) \geq 0\}\, \lambda(x) \left(\frac{d}{p_0(x)} - \frac{1-d}{1-p_0(x)}\right)$$

*is the Riesz representer for the linear functional $D_\mu \overline{W}(\mu_0)[\cdot]$.*

*If furthermore Assumption 3 holds, we have:*

$$\sqrt{n}\left(\overline{W}(\hat{\mu}) - \overline{W}(\mu_0)\right) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \nu^*(X_i, D_i)\, \epsilon_i + o_p(1).$$

*Then:*

$$\sqrt{n}\left(W\left(\hat{h}\right) - W(h_0)\right) \equiv \sqrt{n}\left(\overline{W}(\hat{\mu}) - \overline{W}(\mu_0)\right) \xrightarrow{d} \mathcal{N}\left(0, \sigma_W^2\right),$$

*with*

$$\sigma_W^2 := \mathbb{E}\left[\left(\nu^*(X_i, D_i)\, \epsilon_i\right)^2\right] = \mathbb{E}\left[\frac{\mathbb{1}\{h_0(X_i) \geq 0\}\, \lambda^2(X_i)\, \sigma_\epsilon^2(X_i)}{p_0(X_i)(1 - p_0(X_i))}\right],$$

*where* $\sigma_\epsilon^2(x) := \mathbb{E}\left[\epsilon_i^2 \big| X_i = x\right]$.

Theorem 1 suggests the following natural estimator for the asymptotic variance $\sigma_W^2$:

$$\hat{\sigma}_W^2 := \frac{1}{N} \sum_i \frac{\mathbb{1}\left\{\hat{h}(X_i) \geq 0\right\} \lambda^2(X_i)\, \hat{u}_i^2}{\hat{p}(X_i)(1 - \hat{p}(X_i))} \tag{19}$$

where $\hat{u}_i := Y_i - \hat{h}(X_i)$. This requires nonparametric estimation of propensity score function $p(x)$, as well as the knowledge (or nonparametric estimation) of the density $f_0$ or the Radon-Nikodym derivative $\lambda(x)$.

Alternatively and more preferably, we can use the sieve-based asymptotic variance estimator, which does not require estimation or knowledge of $p(x)$ and $\lambda(x)$. To do so, we first

12

clarify some subtlety in the definition of the sieve in the first-stage nonparametric estimation of CATE as $\hat{\mu}(x, 1) - \hat{\mu}(x, 0)$, where $\hat{\mu}$ is a linear sieve estimator of $\mu_0(x, d)$ under random design on $(X_i, D_i)$ with $D_i$ being binary.

In practice, $\hat{\mu}(x, 1)$ is estimated with $K_1$ linear series in the treated subsample, while $\hat{\mu}(x, 0)$ is estimated separately with potentially different $K_0$ linear series in the untreated subsample. However, since we treat $D_i$ as a random variable, we cannot directly treat $\hat{\mu}(x, 1)$ and $\hat{\mu}(x, 0)$ as two completely separate nonparametric estimators with exogenously given sample sizes. Instead, we treat $\hat{\mu}$ as the least square estimator of

$$Y_i = D_i \psi^{(K_1)}(X_i)' \beta_1 + (1 - D_i) \psi^{(K_0)}(X_i)' \beta_0 + u_i \tag{20}$$

where $\psi^{(K_1)}(x) = (\psi_1(x), \ldots, \psi_{K_1}(x))'$ and $\psi^{(K_0)}(x) = (\psi_{K_1+1}(x), \ldots, \psi_{K_1+K_0}(x))'$ denote the B-spline basis functions used to estimate $\mu_0(x, 1)$ and $\mu_0(x, 0)$, respectively, with sieve dimensions $K_1$ and $K_0$.

Define the vector-valued function $\overline{\psi}(x)$ such that its $k$th component equals $d\psi_k(x)$ for $k \in \{1, \ldots, K_1\}$ and $(1-d)\psi_k(x)$ for $k \in \{K_1+1, \ldots, K_1+K_0\}$. Following Chen and Christensen (2018), specifically equations (6) and (7), the variance (or standard error) estimator in our setting takes the form

$$\hat{\sigma}_W^2 := D_\mu \overline{W}(\hat{\mu}) \left[\overline{\psi}\right]' \hat{\Omega} D_\mu \overline{W}(\hat{\mu}) \left[\overline{\psi}\right], \tag{21}$$

where $\hat{\Omega}$ denotes the estimated asymptotic covariance matrix of the OLS estimators in (20) given by

$$\hat{\Omega} := \left(\Psi^{(2K)} \Psi^{(2K)'}\right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n u_i^2 \Psi^{(2K)} \Psi^{(2K)'}\right) \left(\Psi^{(2K)} \Psi^{(2K)'}\right)^{-1}$$

and the estimated directional derivative vector $D_\mu \overline{W}(\hat{\mu}) \left[\overline{\psi}\right]$ is given by

$$D_\mu \overline{W}(\hat{\mu}) \left[\overline{\psi}\right] = \begin{pmatrix} \int_{\{\hat{\mu}(x,1)-\hat{\mu}(x,0)\geq 0\}} \psi^{(K_1)}(x) f(x) dx \\ - \int_{\{\hat{\mu}(x,1)-\hat{\mu}(x,0)\geq 0\}} \psi^{(K_0)}(x) f(x) dx \end{pmatrix},$$

where the minus sign in front of the sieve terms in $\psi^{K_0}(x)$ is due to the presence of the minus sign in front of $\mu(x, 0)$ in $\text{CATE}(x) = \mu(x, 1) - \mu(x, 0)$.

We use Sobol points to numerically compute the integral above. See the simulation section for more details.

## 3.2 Welfare Functional Under Unknown Density $f = f_0$

In certain cases, such as in Kitagawa and Tetenov (2018), one might be interested in the welfare functional under the original distribution of covariates $F_0$, which may not be known or controlled. Given the random sample of $(Y_i, D_i, X_i)_{i=1}^n$ used to estimate the CATE $h_0$, it

is natural to use the sample mean $\frac{1}{n}\sum_{i=1}^{n}[\cdot]$ as an estimator of the population expectation $\int[\cdot]dF_0$, in which case a natural plug-in estimator of $W(h_0) \equiv \overline{W}(\mu_0)$ is given by

$$\hat{W}\left(\hat{h}\right) \equiv \hat{\overline{W}}\left(\hat{\mu}\right) := \frac{1}{n}\sum_{i=1}^{n}\left[\hat{\mu}\left(X_i,1\right) - \hat{\mu}\left(X_i,0\right)\right]_+ \equiv \frac{1}{n}\sum_{i=1}^{n}\left[\hat{h}\left(X_i\right)\right]_+.$$

Clearly, the approximation of the integral introduces an additional source of randomness, but the result established in the last subsection continues to be useful in this case.

**Theorem 2.** *Under Assumptions 1, 2(b)(c) and 3,*

$$\sqrt{n}\left(\hat{W}\left(\hat{h}\right) - W\left(h_0\right)\right) \equiv \sqrt{n}\left(\hat{\overline{W}}\left(\hat{\mu}\right) - \overline{W}\left(\mu_0\right)\right)$$

$$= \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\left(\left[h_0\left(X_i\right)\right]_+ - W\left(h_0\right) + \nu^*\left(X_i,D_i\right)\epsilon_i\right) + o_p\left(1\right) \xrightarrow{d} \mathcal{N}\left(0,\overline{\sigma}_W^2\right)$$

*where* $\overline{\sigma}_W^2 := \mathbb{E}\left[\left(\left[h_0\left(X_i\right)\right]_+ - W(h_0) + \nu^*\left(X_i,D_i\right)\epsilon_i\right)^2\right] = Var\left(\left[h_0\left(X_i\right)\right]_+\right) + \sigma_W^2.$

The standard errors can be computed similarly, based on straightforward adaptations of the analytical formula (19) or the sieve-based formula (21) in Section 3.1.

# 4    Estimation and Inference of the Value Functional

## 4.1    Value Functional Under Known Density $f$

We now turn to the general value functional given by

$$V(h_0) := \int \mathbb{1}\left\{h_0\left(x\right) \geq 0\right\} v_0\left(x\right) f\left(x\right) dx$$

$$\equiv \overline{V}\left(\mu_0\right) := \int \mathbb{1}\left\{\mu_0\left(x,1\right) - \mu_0\left(x,0\right) \geq 0\right\} v_0\left(x\right) f\left(x\right) dx.$$

The simple plug-in estimators are defined as

$$V\left(\hat{h}\right) := \int \mathbb{1}\left\{\hat{h}\left(x\right) \geq 0\right\} v_0\left(x\right) f\left(x\right) dx$$

$$\equiv \overline{V}\left(\hat{\mu}\right) := \int \mathbb{1}\left\{\hat{\mu}\left(x,1\right) - \hat{\mu}\left(x,0\right) \geq 0\right\} v_0\left(x\right) f\left(x\right) dx.$$

Under Assumption 1 and 2(b)(c) the functional derivative of $\overline{V}(\mu_0)[\nu]$ is given by

$$D_\mu\overline{V}\left(\mu_0\right)\left[\nu\right] := \int_{\left\{x\in\mathbb{R}^d:h_0(x)=0\right\}} \frac{\left(\nu\left(x,1\right) - \nu\left(x,0\right)\right)}{\left\|\nabla_x h_0\left(x\right)\right\|} v_0\left(x\right) f\left(x\right) d\mathcal{H}^{d-1}\left(x\right).$$

As explained in Section 2, $V(h_0) \equiv \overline{V}(\mu_0)$ is generally not $\sqrt{n}$-estimable, in particular, the $D_\mu\overline{V}(\mu_0)[\nu]$ are not bounded (or continuous) linear functionals on $L^2(f)$, which means that they do not have a Riesz representer on the whole Hilbert space $L^2(f)$. Nevertheless, sieve Riesz representer is well-defined (see Chen and Gao (2025)), which will be the key ingredient

14

for the sieve variance term for the properties of $\overline{V}(\hat{\mu}) - \overline{V}(\mu_0)$. In particular we need to study the asymptotic property of the plug-in estimator of the submanifold integral of form (15), which has been studied in Chen and Gao (2025) using linear series (Bspline) first stage: in particular, Section 4.3 of Chen and Gao (2025) analyzes the integral on upper contour set of the form $V(h_0)$ specifically. We use the result in Chen and Gao (2025) without repeating it here, but focus on the adaptation required for the standard error computation.

**Assumption 4.** *Suppose that:*

*(a)* $\|\nabla_x^2 h_0(x)\| \leq M < \infty$.

*(b)* $\|\hat{\mu} - \mu_0\|_\infty \|\nabla(\hat{\mu} - \mu_0)\|_\infty = o_p\left(\sqrt{\frac{1}{n}K_n^{\frac{1}{d}}}\right)$.

**Theorem 3.** *Suppose that Assumptions 1 and 2(b)(c) hold. Let $\hat{\mu}$ be a linear sieve estimator of $\mu_0$ and suppose that Assumptions 6, 8 and 11 in Chen and Gao (2025) hold along with Assumption 4 above. Then:*

$$\frac{\sqrt{n}\left(\overline{V}(\hat{\mu}) - \overline{V}(\mu_0)\right)}{\sigma_{V,n}} \xrightarrow{d} \mathcal{N}(0,1), \quad \text{with } \sigma_{V,n}^2 \asymp K_n^{\frac{1}{d}}$$

The standard error estimates can be computed based on the linear sieve first stage in a way similar to that described in Section 3.1, with the following adaptions. Again, we use the formula

$$\hat{\sigma}_V^2 := \hat{D}_\mu \overline{V}(\hat{\mu})\left[\overline{\psi}\right]' \hat{\Omega} \hat{D}_\mu \overline{V}(\hat{\mu})\left[\overline{\psi}\right], \tag{22}$$

where the pathwise derivative $D_\mu \overline{V}(\hat{\mu})\left[\overline{\psi}\right]$ given by

$$D_\mu \overline{V}(\hat{\mu})\left[\overline{\psi}\right] = \begin{pmatrix} \int_{\{x \in \mathcal{X}: \hat{h}(x)=0\}} \frac{\psi^{(K_1)}(x)}{\|\nabla_x h_0(x)\|} v_0(x)f(x)\, d\mathcal{H}^{d-1}(x) \\ -\int_{\{x \in \mathcal{X}: \hat{h}(x)=0\}} \frac{\psi^{(K_0)}(x)}{\|\nabla_x h_0(x)\|} v_0(x)f(x)\, d\mathcal{H}^{d-1}(x) \end{pmatrix},$$

can be approximated via

$$\hat{D}_\mu \overline{V}(\hat{\mu})\left[\overline{\psi}\right] = \begin{pmatrix} \frac{1}{2\epsilon}\int_{\{x \in \mathcal{X}: -\epsilon < \hat{h}(x) < \epsilon\}} \psi^{(K_1)}(x)v_0(x)f(x)\, dx \\ -\frac{1}{2\epsilon}\int_{\{x \in \mathcal{X}: -\epsilon < \hat{h}(x) < \epsilon\}} \psi^{(K_0)}(x)v_0(x)f(x)\, dx \end{pmatrix}, \tag{23}$$

based on the mathematical result[3] that

$$\lim_{\epsilon \searrow 0} \frac{1}{2\epsilon}\int_{\{x \in \mathcal{X}: \ -\epsilon < h(x) < \epsilon\}} \omega(x)\, dx = \int_{\{x \in \mathcal{X}: \ h(x)=0\}} \frac{\omega(x)}{\|\nabla_x h(x)\|} d\mathcal{H}^{d-1}(x).$$

Again, we use Sobol points for numerical integration. See the simulation section for details, as well as robustness checks with respect the choice of $\epsilon$ in the numerical differentiation step.

---

[3]See Theorem 3.13.(iii) of Evans and Gariepy (2015).

## 4.2  Value Functional Under Unknown Density $f = f_0$

We now consider the case where $F = F_0$ and population expectation $\mathbb{E}[\cdot]$ is estimated by the sample average $\frac{1}{n}\sum_{i=1}^{n}[\cdot]$ , and seek to characterize the asymptotic behavior of the natural plug-in estimator of $V(h_0) \equiv \overline{V}(\mu_0)$ is given by

$$\hat{V}\left(\hat{h}\right) := \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}\left\{\hat{h}(X_i) \geq 0\right\}v_0(X_i)$$

$$\equiv \hat{\overline{V}}(\hat{\mu}) := \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}\left\{\hat{\mu}(X_i,1) - \hat{\mu}(X_i,0) \geq 0\right\}v_0(X_i).$$

It turns out that the additional error in the approximation of $V\left(\hat{h}\right)$ by $\hat{V}\left(\hat{h}\right)$ is asymptotically negligible relative to $V\left(\hat{h}\right) - V(h_0)$, which converges at a slower-than-$\sqrt{n}$ rate.

**Theorem 4.** *The asymptotic distribution of $\hat{V}\left(\hat{h}\right)$ coincides with that $V(\hat{h})$ in Theorem 3.*

The standard error can be computed using formula (22), with $\hat{D}_\mu\overline{V}(\hat{\mu})\left[\overline{\psi}\right]$ given by the following adapted estimator,

$$\hat{D}_\mu\overline{V}(\hat{\mu})\left[\overline{\psi}\right] = \begin{pmatrix} \frac{1}{2\epsilon}\int_{\{x\in\mathcal{X}:-\epsilon<\hat{h}(x)<\epsilon\}}\psi^{(K_1)}(x)v_0(x)\hat{f}(x)\,dx \\ -\frac{1}{2\epsilon}\int_{\{x\in\mathcal{X}:-\epsilon<\hat{h}(x)<\epsilon\}}\psi^{(K_0)}(x)v_0(x)\hat{f}(x)\,dx \end{pmatrix}, \tag{24}$$

where $\hat{f}(x)$ is a nonparametric density estimator of $f_0(x)$. We then again use Sobol points to numerically evaluate the integral in (24).

**Remark 2.** *Even though integrals of the form $\int w(x)f_0(x)dx$ can be estimated using the sample average $\frac{1}{N}\sum_i w(X_i)$ without the need of a nonparametric density estimator, we choose instead to estimate it using $\int w(x)\hat{f}(x)dx$ using the nonparametric density estimator $\hat{f}(x)$ together with numerical integration based on Sobol points in (24), because we need to compute the integral in (24) on a small constrained domain $\{-\epsilon < \hat{h}(x) < \epsilon\}$ for numerical differentiation. As a result, given the empirical data $(X_i)_{i=1}^{N}$, for small choices of $\epsilon$, there might be very few, or even no, data points that fall into the small band, making the sample average $\frac{1}{N}\sum_i w(X_i)\mathbb{1}\{-\epsilon < \hat{h}(X_i) < \epsilon\}$ very discrete (and even identically zero) for small values of $\epsilon$. The use of a nonparametric density estimator $\hat{f}(x)$, along with numerical integration, effectively smooths out such discreteness and results in much better numerical approximation of the derivative.*

# 5 Simulations

## 5.1 Results for Theorem 1

We first report the finite-sample performance of the semiparametric estimator, its associated standard error estimator, and the resulting confidence interval, based on the theoretical results in Theorem 1, which concern the welfare functional under a known target distribution. The model specifications used in the simulations are summarized in Table 1. For each specification, random samples of size $n$ are drawn with covariates $X_i \sim F_0$, and treatment status is assigned according to the propensity score function $p_0(X_i)$. Outcomes are then generated as $Y_i = \mu_0(X_i, D_i) + \epsilon_i$, with $\epsilon_i \sim N(0, 1)$.

Table 1: Theorem 1: Model Specifications

| Model | $F_0$ | $F$ | $\mu_0(x, d)$ | $p_0(x)$ | $\sigma^2$ |
|-------|-------|-----|---------------|----------|------------|
| M1 | $U[-0.2, 1.2]$ | $U[0, 1]$ | $5\sin(2\pi x)\cos(2\pi x)$ $+d(-0.4 + 2x^2)$ | $\frac{1}{1+e^{-(1-2x)}}$ | 1 |
| M2 | $U[-0.2, 1.2]$ | $U[0, 1]$ | $0.5|x| + d(0.5 - x^2)$ | $\frac{1}{1+e^{-(-0.5+x)}}$ | 1 |
| M3 | $U[-0.2, 1.2]$ | $U[0, 1]$ | $x^2 + d(1 - x)$ | $\frac{1}{1+e^{-(0.5-x)}}$ | 1 |
| M4 | $U[-0.2, 1.2]^2$ | $U[0, 1]^2$ | $(1 - x_1^2 - x_2^2)(4 + \sin x_1 x_2 + \cos x_2)$ $+d(0.5x_1 - 0.4x_2)$ | $\frac{1}{1+e^{-(x_1-x_2)}}$ | 1 |
| M5 | $U[-0.2, 1.2]^2$ | $U[0, 1]^2$ | $(1 - x_1 x_2)(3 + \sin(\pi x_1)\cos(\pi x_2))$ $+d(0.3x_1 - 0.3x_2)$ | $\frac{1}{1+e^{-(x_1-x_2)}}$ | 1 |
| M6 | $U[-0.2, 1.2]^2$ | $U[0, 1]^2$ | $\log(1 + x_1 + x_2) + d(x_1 - 0.7x_2)$ | $\frac{1}{1+e^{-(1.5x_1-0.5x_2)}}$ | 1 |
| M7 | $U[-0.2, 1.2]^2$ | $U[0, 1]^2$ | $(x_1^2 + x_2^2)e^{-(x_1+x_2)}$ $+d(0.5 - x_2)$ | $\frac{1}{1+e^{-(-0.5+x_1+2x_2)}}$ | 1 |

The semiparametric estimator of the welfare functional is constructed in two steps. In the first step, $\mu_0(x, 1)$ and $\mu_0(x, 0)$ are estimated separately for the treated group ($D_i = 1$) and the control group ($D_i = 0$) using B-spline regressions. In the second step, the welfare functional is approximated by numerically integrating over $M = 5{,}000$ Sobol points $\{X_j^{Sobol}\}_{j=1}^{M}$ drawn from the target distribution $F$:

$$\hat{\overline{W}}(\hat{\mu}) = \frac{1}{M}\sum_{j=1}^{M}\left[\hat{\mu}(X_j^{Sobol}, 1) - \hat{\mu}(X_j^{Sobol}, 0)\right]_+,$$

where $\hat{\mu}(x, 1)$ and $\hat{\mu}(x, 0)$ denote the first-stage nonparametric estimators.

Let $\hat{h}(x) = \hat{\mu}(x, 1) - \hat{\mu}(x, 0)$ denote the estimated CATE function, and let $\hat{p}(x)$ be a B-spline sieve estimator of the propensity score. The 95% confidence interval for the welfare

functional takes the usual form,

$$\left[ \hat{\overline{W}}(\hat{\mu}) - 1.96\frac{\hat{\sigma}_W}{\sqrt{n}}, \quad \hat{\overline{W}}(\hat{\mu}) + 1.96\frac{\hat{\sigma}_W}{\sqrt{n}} \right],$$

with

$$\hat{\sigma}_W^2 = \frac{1}{N} \sum_i \frac{\mathbb{1}\left\{ \hat{h}\left(X_i\right) \geq 0 \right\} \lambda^2\left(X_i\right)\left(Y_i - \hat{h}(X_i)\right)^2}{\hat{p}\left(X_i\right)\left(1 - \hat{p}\left(X_i\right)\right)}.$$

To evaluate performance, we simulate each model 2,000 times at sample sizes $n = 1500, 3000,$ and $6000$. Table 2 reports the true welfare functional ($W_{\text{true}}$), the average bias of the estimator $\hat{\overline{W}}(\hat{\mu})$ (Bias), its sampling standard deviation (SD), the average estimated standard error (SE), the standard deviation of SE across iterations (SD(SE)), and the empirical coverage probability of the associated 95% confidence interval (Coverage). The results[4] indicate that the nominal coverage rate is attained in nearly all cases, even with relatively small samples, and improves further as sample size increases. The bias of the estimator also becomes negligible relative to its sampling variability, and the proposed SE estimator closely tracks the sampling standard deviation with high precision.

**Remark 3.** *To improve computational efficiency, we predetermine the sieve dimensions for estimating $\mu_0(x,0)$ and $\mu_0(x,1)$. For each model specification, we first generate a dataset $(Y_i, X_i, D_i)_{i=1}^{n=6000}$ and apply an adaptation the sieve dimension selection procedure of Chen, Christensen and Kankanala (2025) separately to the treated and control groups.[5] As demonstrated in their paper, this approach ensures that the resulting estimators of $\mu_0(x,1)$ and $\mu_0(x,0)$ converge at the fastest possible (i.e., minimax) rates in the sup-norm. The selected sieve dimensions are then used in the simulation designs with $n = 1500, 3000, 6000$. In principle, one could implement data-driven dimension selection within each simulation iteration, but doing so would substantially increase computational cost.*

**Remark 4.** *The construction of the 95% confidence interval requires estimating the propensity score function $p_0(x)$. To this end, we again use the B-spline sieve estimator, regressing treatment status on the covariates. The sieve dimension is similarly predetermined using an adaptation of Chen, Christensen and Kankanala (2025) predetermined based on the full dataset. A potential concern is that the fitted propensity score $\hat{p}(x)$ may take values outside the unit interval for some observations, which is likely to indicate a violation of the strict overlap assumption. Accordingly, when computing the asymptotic standard deviation $\hat{\sigma}_W$, we trim observations with estimated propensity scores lying outside $[0,1]$.*

---

[4]Additional simulations based on GAM are presented in the Appendix B.1.

[5]The adaptation selects a larger sieve dimension than that selected by the CCK procedure (and the npiv R package) to achieve undersmoothing.

Table 2: Theorem 1: Simulation Results

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|-------|------|---------|--------|--------|--------|--------|----------|
| M1    | 1500 | 0.3857  | 0.0161 | 0.0470 | 0.0480 | 0.0033 | 0.9415   |
|       | 3000 | 0.3857  | 0.0080 | 0.0338 | 0.0337 | 0.0018 | 0.9480   |
|       | 6000 | 0.3857  | 0.0037 | 0.0229 | 0.0237 | 0.0010 | 0.9535   |
| M2    | 1500 | 0.2358  | 0.0030 | 0.0484 | 0.0518 | 0.0031 | 0.9615   |
|       | 3000 | 0.2358  | 0.0028 | 0.0352 | 0.0366 | 0.0015 | 0.9555   |
|       | 6000 | 0.2358  | 0.0009 | 0.0243 | 0.0259 | 0.0007 | 0.9630   |
| M3    | 1500 | 0.5001  | 0.0046 | 0.0579 | 0.0605 | 0.0025 | 0.9620   |
|       | 3000 | 0.5001  | 0.0025 | 0.0407 | 0.0430 | 0.0013 | 0.9595   |
|       | 6000 | 0.5001  | 0.0000 | 0.0293 | 0.0305 | 0.0006 | 0.9625   |
| M4    | 1500 | 0.1033  | 0.0185 | 0.0478 | 0.0561 | 0.0074 | 0.9735   |
|       | 3000 | 0.1033  | 0.0088 | 0.0348 | 0.0400 | 0.0040 | 0.9745   |
|       | 6000 | 0.1033  | 0.0043 | 0.0248 | 0.0284 | 0.0021 | 0.9695   |
| M5    | 1500 | 0.0499  | 0.0282 | 0.0418 | 0.0521 | 0.0093 | 0.9700   |
|       | 3000 | 0.0499  | 0.0158 | 0.0303 | 0.0369 | 0.0056 | 0.9720   |
|       | 6000 | 0.0499  | 0.0094 | 0.0227 | 0.0262 | 0.0032 | 0.9660   |
| M6    | 1500 | 0.2315  | 0.0268 | 0.0559 | 0.0622 | 0.0048 | 0.9550   |
|       | 3000 | 0.2315  | 0.0123 | 0.0394 | 0.0444 | 0.0025 | 0.9655   |
|       | 6000 | 0.2315  | 0.0050 | 0.0288 | 0.0315 | 0.0014 | 0.9695   |
| M7    | 1500 | 0.1250  | 0.0462 | 0.0505 | 0.0595 | 0.0073 | 0.9370   |
|       | 3000 | 0.1250  | 0.0218 | 0.0346 | 0.0406 | 0.0043 | 0.9610   |
|       | 6000 | 0.1250  | 0.0101 | 0.0245 | 0.0278 | 0.0022 | 0.9680   |

*Notes:* (1) $W_{\text{true}}$ denotes the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\widehat{\overline{W}}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

### 5.1.1 Sieve Variance Estimation

One potential drawback of the plug-in approach to estimating the asymptotic variance of the welfare functional is that it requires nonparametric estimation of the propensity score function $p_0(x)$, which can introduce additional sampling noise and numerical instability. As an alternative, one may employ the sieve variance estimator as in Chen and Liao (2014) or Chen and Liao (2015), whose formula depends only on the pathwise derivative of the welfare functional and a linear regression of the outcome variable on the sieve basis. When the target distribution $F$ is known, the pathwise derivative of the welfare functional can be numerically approximated using Sobol points drawn from the target population combined with importance sampling. When the target distribution is unknown, the pathwise derivative can instead be evaluated directly at the observed data and approximated by a sample average. Simulation results and the empirical application based on the sieve variance estimator are shown in the Appendix B.5.

## 5.2 Results for Theorem 2

We now investigate the finite-sample performance of our estimation and inference procedure for the welfare functional in the case where the target distribution $F$ coincides with the population distribution $F_0$, though both remain unknown. The relevant model specifications are presented in Table 3. In contrast to the designs considered in Section 5.1, these specifications explicitly impose that $F$ and $F_0$ share identical supports.

Table 3: Theorem 2: Model Specifications

| Model | $F_0$ | $F$ | $\mu_0(x,d)$ | $p_0(x)$ | $\sigma^2$ |
|-------|-------|-----|--------------|----------|------------|
| M8 | $U[0,1]$ | $U[0,1]$ | $5\sin(2\pi x)\cos(2\pi x)$ $+d(-0.4+2x^2)$ | $\frac{1}{1+e^{-(1-2x)}}$ | 1 |
| M9 | $U[0,1]$ | $U[0,1]$ | $0.5|x|+d(0.5-x^2)$ | $\frac{1}{1+e^{-(-0.5+x)}}$ | 1 |
| M10 | $U[0,1]$ | $U[0,1]$ | $x^2+d(1-x)$ | $\frac{1}{1+e^{-(0.5-x)}}$ | 1 |
| M11 | $U[0,1]^2$ | $U[0,1]^2$ | $(1-x_1^2-x_2^2)(4+\sin x_1 x_2+\cos x_2)$ $+d(0.5x_1-0.4x_2)$ | $\frac{1}{1+e^{-(x_1-x_2)}}$ | 1 |
| M12 | $U[0,1]^2$ | $U[0,1]^2$ | $(1-x_1 x_2)(3+\sin(\pi x_1)\cos(\pi x_2))$ $+d(0.3x_1-0.3x_2)$ | $\frac{1}{1+e^{-(x_1-x_2)}}$ | 1 |
| M13 | $U[0,1]^2$ | $U[0,1]^2$ | $\log(1+x_1+x_2)+d(x_1-0.7x_2)$ | $\frac{1}{1+e^{-(1.5x_1-0.5x_2)}}$ | 1 |
| M14 | $U[0,1]^2$ | $U[0,1]^2$ | $(x_1^2+x_2^2)e^{-(x_1+x_2)}$ $+d(0.5-x_2)$ | $\frac{1}{1+e^{-(-0.5+x_1+2x_2)}}$ | 1 |

The plug-in estimator $\widehat{\overline{W}}(\hat{\mu})$ proposed here differs from that in Section 5.1 only in the second step: instead of using Sobol points to approximate the integral, we take the sample

average of $\left[ \hat{\mu}(X_i, 1) - \hat{\mu}(X_i, 0) \right]_+$ over the observed data:

$$\hat{\overline{W}}(\hat{\mu}) \;=\; \frac{1}{n} \sum_{i=1}^{n} \left[ \hat{\mu}(X_i, 1) - \hat{\mu}(X_i, 0) \right]_+.$$

Relative to the known-$F$ case, the asymptotic variance of $\hat{\overline{W}}(\hat{h})$ includes an additional component, $\mathrm{Var}([h_0(X_i)]_+)$. We estimate the asymptotic standard deviation $\hat{\hat{\sigma}}_W$ by substituting the B-spline sieve estimates for the nuisance functions and replacing the population mean with its sample analog. As before, we restrict attention to observations with estimated propensity scores in $[0, 1]$. A 95% confidence interval is then given by

$$\left[ \hat{\overline{W}}(\hat{\mu}) - 1.96 \frac{\hat{\hat{\sigma}}_W}{\sqrt{n}}, \;\; \hat{\overline{W}}(\hat{\mu}) + 1.96 \frac{\hat{\hat{\sigma}}_W}{\sqrt{n}} \right].$$

The simulation results[6] based on 2,000 iterations with sample sizes $n = 1500, 3000$, and 6000 are reported in Table 4. Overall, the coverage of the proposed confidence intervals converges to the nominal 95% level as sample size increases, while the decreasing bias and standard deviation indicate a reduction in mean squared error.

**Remark 5.** *Extrapolation bias may arise when B-spline fits are evaluated outside the support of their training samples-for instance, when $\hat{\mu}(x, 1)$ is evaluated using control group data or $\hat{\mu}(x, 0)$ using treated group data. To mitigate this issue, we trim observations that fall outside the common support of treated and control groups, estimate the nuisance functions on the trimmed sample, and then compute the welfare functional. Since only a small fraction of observations are removed, this adjustment has a negligible effect on the results.*

## 5.3 Results for Theorem 3

To assess the finite-sample properties of our estimation inference procedure for the value functional $\overline{V}(\hat{\mu})$ under a known target distribution $F$, we analyze the model described in Table 5:

Table 5: Theorem 3: Model Specification

| Model | $F_0$ | $F$ | $\mu_0(x_1, x_2, d)$ | $\nu_0(x_1, x_2)$ | $p_0(x_1, x_2)$ | $\sigma^2$ |
|---|---|---|---|---|---|---|
| M15 | U($[-2, 2]^2$) | U($[-1.5, 1.5]^2$) | $d \cdot (1 - x_1^2 - x_2^2) \cdot (4 + \sin(x_1)x_2 + \cos(x_2))$ | 1 | $\frac{1}{1 + e^{-(x_1 - x_2)}}$ | 1 |

Our parameter of interest is a scaled value functional under known $F$:

$$V(h_0) \;=\; 3^2 \int \mathbb{1}\left\{ \left( 1 - x_1^2 - x_2^2 \right) \left( 4 + \sin(x_1)x_2 + \cos(x_2) \right) \;\geq\; 0 \right\} dF(x_1, x_2),$$

---

[6]Additional simulations based on GAM are presented in the Appendix B.2.

Table 4: Theorem 2: Simulation Results

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|---|---|---|---|---|---|---|---|
| M8 | 1500 | 0.3857 | 0.0068 | 0.0414 | 0.0423 | 0.0025 | 0.9515 |
| | 3000 | 0.3857 | 0.0029 | 0.0296 | 0.0297 | 0.0013 | 0.9550 |
| | 6000 | 0.3857 | 0.0006 | 0.0206 | 0.0210 | 0.0007 | 0.9575 |
| M9 | 1500 | 0.2358 | 0.0042 | 0.0425 | 0.0440 | 0.0025 | 0.9605 |
| | 3000 | 0.2358 | 0.0029 | 0.0304 | 0.0311 | 0.0012 | 0.9590 |
| | 6000 | 0.2358 | 0.0012 | 0.0209 | 0.0221 | 0.0006 | 0.9555 |
| M10 | 1500 | 0.5001 | 0.0067 | 0.0511 | 0.0515 | 0.0018 | 0.9495 |
| | 3000 | 0.5001 | 0.0037 | 0.0357 | 0.0366 | 0.0009 | 0.9600 |
| | 6000 | 0.5001 | 0.0013 | 0.0253 | 0.0260 | 0.0005 | 0.9600 |
| M11 | 1500 | 0.1033 | 0.0307 | 0.0365 | 0.0402 | 0.0038 | 0.9245 |
| | 3000 | 0.1033 | 0.0156 | 0.0264 | 0.0286 | 0.0021 | 0.9395 |
| | 6000 | 0.1033 | 0.0084 | 0.0192 | 0.0203 | 0.0012 | 0.9470 |
| M12 | 1500 | 0.0499 | 0.0418 | 0.0316 | 0.0373 | 0.0045 | 0.8790 |
| | 3000 | 0.0499 | 0.0232 | 0.0229 | 0.0263 | 0.0029 | 0.9135 |
| | 6000 | 0.0499 | 0.0126 | 0.0168 | 0.0186 | 0.0017 | 0.9355 |
| M13 | 1500 | 0.2315 | 0.0177 | 0.0432 | 0.0459 | 0.0032 | 0.9565 |
| | 3000 | 0.2315 | 0.0088 | 0.0309 | 0.0323 | 0.0014 | 0.9510 |
| | 6000 | 0.2315 | 0.0041 | 0.0221 | 0.0229 | 0.0007 | 0.9525 |
| M14 | 1500 | 0.1250 | 0.0251 | 0.0382 | 0.0421 | 0.0059 | 0.9460 |
| | 3000 | 0.1250 | 0.0117 | 0.0258 | 0.0287 | 0.0030 | 0.9585 |
| | 6000 | 0.1250 | 0.0054 | 0.0182 | 0.0198 | 0.0013 | 0.9595 |

*Notes:* (1) $W_{\text{true}}$ is the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\widehat{\overline{W}}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

where the scaling factor $3^2$ is chosen so that the integral evaluates to $\pi$. Equivalently, this setup can be interpreted as a Monte Carlo experiment for estimating the area of the unit circle by uniformly throwing darts over a $3 \times 3$ square.

The plug-in estimator $\hat{\overline{V}}(\hat{\mu})$ for $\overline{V}(\mu_0)$ is constructed in two steps. In the first step, we estimate $\mu_0(x,1) = \mu_0(x_1, x_2, 1)$ and $\mu_0(x,0) = \mu_0(x_1, x_2, 0)$ separately using B-spline sieve estimators, fitted on the treated and control groups, respectively, with sieve dimensions pre-determined as described in subsection 5.1. In the second step, $\mathbb{1}\{(\hat{\mu}_0(x,1) - \hat{\mu}_0(x,0)) \geq 0\}$ is numerically integrated using 5000 Sobol points drawn from $F$. The resulting semiparametric two-step estimator is

$$\hat{\overline{V}}(\hat{\mu}) \;=\; \frac{1}{M} \sum_{j=1}^{M} \mathbb{1}\{\hat{\mu}(X_j^{Sobol}, 1) - \hat{\mu}(X_j^{Sobol}, 0)\}.$$

The 95% confidence interval is then constructed as

$$\left[ \hat{\overline{V}}(\hat{\mu}) - 1.96 \frac{\hat{\sigma}_V}{\sqrt{n}}, \;\; \hat{\overline{V}}(\hat{\mu}) + 1.96 \frac{\hat{\sigma}_V}{\sqrt{n}} \right],$$

where the asymptotic variance estimate $\hat{\sigma}_V^2$ is given by

$$\hat{\sigma}_V^2 = \hat{D}_\mu \overline{V}(\hat{\mu}) \left[\overline{\psi}\right]' \hat{\Omega} \hat{D}_\mu \overline{V}(\hat{\mu}) \left[\overline{\psi}\right],$$

with $\hat{\Omega}$ being the estimated asymptotic covariance matrix for the OLS estimators in the linear regression model 20, and

$$\hat{D}_\mu \overline{V}(\hat{\mu}) \left[\overline{\psi}\right] = 3^2 \begin{pmatrix} \frac{1}{2\epsilon} \int_{\{x \in [-1.5,1.5]^2 : -\epsilon < \hat{h}(x) < \epsilon\}} \psi^{(K_1)}(x) \frac{1}{3^2} \, dx \\ -\frac{1}{2\epsilon} \int_{\{x \in [-1.5,1.5]^2 : -\epsilon < \hat{h}(x) < \epsilon\}} \psi^{(K_0)}(x) \frac{1}{3^2} \, dx \end{pmatrix}.$$

We set the tuning parameter $\epsilon = 0.005$ to mitigate bias in $\hat{\overline{V}}(\hat{\mu})$ and approximate $\hat{D}_\mu \overline{V}(\hat{\mu})[\nu]$ using the sample average over $M$ Sobol draws from $F$. Because draws from $F$ are unlikely to fall within the set $\{x \in [-1.5, 1.5]^2 : -\epsilon < \hat{h}(x) < \epsilon\}$ when $\epsilon$ is small, we use $M = 1,000,000$ Sobol points to ensure accuracy. Simulation results[7] based on 2,000 iterations with sample sizes $n = 1500, 3000,$ and $6000$ are reported in Table 6. The results show that the coverage rate reaches the nominal level with relatively modest sample sizes, even though the value functional is not $\sqrt{n}$-estimable. Also, as sample size increases, both bias and standard error decrease, and the plug-in estimator for the standard error provides a close estimate of the sampling standard deviation.

---

[7]Additional simulations under alternative model specifications, as well as robustness checks for different values of $\epsilon$, are presented in Appendix B.3.

Table 6: Theorem 3: Simulation Results

| Model | $n$ | $V_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|-------|------|-------------------|--------|--------|--------|--------|----------|
| M15 | 1500 | 3.1416 | 0.0076 | 0.0710 | 0.0711 | 0.0092 | 0.9420 |
| | 3000 | 3.1416 | 0.0080 | 0.0486 | 0.0499 | 0.0029 | 0.9475 |
| | 6000 | 3.1416 | 0.0062 | 0.0337 | 0.0353 | 0.0016 | 0.9490 |

*Notes:* (1) $V_{\text{true}}$ is the true value functional, equal to $\pi$. (2) Bias is the average deviation from $V_{\text{true}}$. (3) SD is the sampling standard deviation of the estimator. (4) SE is the average estimated asymptotic standard error across iterations. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% confidence interval.

# 6  Empirical Application

We revisit the empirical application analyzed in Kitagawa and Tetenov (2018, KT18) using Job Training Parternship Act (JTPA) dataset. The JTPA study randomized whether applicants were eligible to receive a mix of training, job-search assistance, and other services provided under the program for a period of 18 months. A detailed description of the study and an assessment of average program effects for five major subgroups of the target population are provided in Bloom, Orr, Bell, Cave, Doolittle, Lin and Bos (1997).

We evaluate welfare using two outcome measures, following the approach in KT18. The first is total earnings over the 30 months following treatment assignment. The second adjusts this measure by subtracting \$774 for individuals assigned to treatment, thereby incorporating program costs. These outcomes are considered from an intention-to-treat perspective, meaning we focus on eligibility assignment rather than treatment effects among compliers. The available covariates include applicants' pre-program earnings, years of education, and treatment status. Our objective is to estimate and conduct inference on the first-best welfare functional and the optimal fraction of the population that should receive treatment.

As the first step in the estimation and inference procedure, we trim observations outside the common support of the treated and control groups to enforce the overlap assumption and to avoid extrapolation when applying sieve estimators. For either outcome measure, we then estimate $\mu_0(x, 1)$ nonparametrically using B-spline sieves fitted on the treated sample, and $\mu_0(x, 0)$ analogously on the control sample. The sieve dimensions are selected in a data-driven manner following Chen, Christensen and Kankanala (2025). Combining these estimates on the common support yields the CATE estimate $\hat{h}(x) = \hat{\mu}(x, 1) - \hat{\mu}(x, 0)$. Taking sample averages of $[\hat{h}(x)]_+$ and $\mathbb{1}\{\hat{h}(x) > 0\}$ over the trimmed dataset produces the estimators $\overline{\hat{W}}(\hat{\mu})$ and $\overline{\hat{V}}(\hat{\mu})$, respectively. Confidence intervals for $\overline{W}(\mu_0)$ and $\overline{V}(\mu_0)$ are then constructed according to their asymptotic theories: $\left(\overline{\hat{W}}(\hat{\mu}) \pm 1.96\, \hat{\bar{\sigma}}_W / \sqrt{N}\right)$, and $\left(\overline{\hat{V}}(\hat{\mu}) \pm 1.96\, \hat{\sigma}_V / \sqrt{N}\right)$,

where $N$ denotes the sample size of the JTPA dataset.

While computing sieve estimate of $\hat{\bar{\sigma}}_W$ is straightforward, computing sieve estimate of $\hat{\sigma}_V$ involves several additional steps. Specifically, it requires a density estimate $\hat{f}(x)$ for the covariates as discussed in Remark 2. To this end, we use a Gaussian kernel density estimator, selecting the bandwidth matrix via the smoothed cross-validation method (`Hscv()` in the `ks` package) and applying a scaling factor of $s = 3$ to ensure adequate smoothness in the presence of discrete values for years of education. Furthermore, we need to specify a small hyperparameter $\epsilon$ to provide a close approximate for the pathwise derivative of the value functional in $\hat{\sigma}_V$. We set $\epsilon$ equal to a fraction $\iota = 0.01$ of the standard deviation of $\hat{h}(x)$ over the trimmed dataset. Robustness checks on the tuning parameters $s$ and $\iota$ are reported in Appendix C.1, and alternative density estimation approaches are examined in Appendix C.2.

Table 7 presents our estimation and inference results alongside the corresponding findings from KT18. Specifically, we report our estimated welfare gain and the share of individuals to be treated, along with their confidence intervals, based on the trimmed dataset. For comparability, we also present results obtained using the untrimmed dataset-on which KT18 conducted their analysis-with the same choice of tuning parameters. The nonparametric plug-in estimates from KT18, which target the same welfare and share parameters as in our analysis, serve as an empirical benchmark. In addition, the linear rule estimates with their associated confidence intervals from KT18 are reported to assess how conservative our nonparametric inference procedure is relative to their parametric counterparts.

Across both the trimmed and untrimmed datasets, our estimates of the welfare gain and the optimal treatment share are broadly consistent with the nonparametric plug-in rule estimates reported in KT18, despite methodological difference in the first stage: we employ sieve estimators for the nuisance functions, whereas they use a Nadaraya-Watson estimator with an Epanechnikov kernel. A key contribution of our analysis is the provision of confidence intervals for both parameters of interest. By contrast, KT18 report confidence intervals only for the welfare gain under parametric rules, but not for the optimal treatment share. Although our confidence interval for the welfare gain appears wide, its length is comparable to that under the linear rule in KT18, underscoring that our inference procedure remains sharp even while relying on fully nonparametric methods.

Table 7: Estimated Welfare Gains and Share of Population to be Treated Under Nonparametric Plug-in Rule

(a) 30-Month Post-Program Earnings, No Treatment Cost

| Method | Share Treated | Est. Welfare Gain |
|---|---|---|
| Ours (With Trimming) | 0.89<br>(0.73, 1.05) | $1,519<br>($691, $2347) |
| Ours (Without Trimming) | 0.92<br>(0.76, 1.08) | $1,459<br>($840, $2078) |
| KT18 (Nonparametric Plug-in) | 0.91<br>NA | $1,693<br>NA |
| KT18 (Linear) | 0.96<br>NA | $1,180<br>($464, $1,896) |

(b) 30-Month Post-Program Earnings, $774 Cost per Treatment

| Method | Share Treated | Est. Welfare Gain |
|---|---|---|
| Ours (With Trimming) | 0.80<br>(0.53, 1.07) | $858<br>($152, $1564) |
| Ours (Without Trimming) | 0.85<br>(0.65, 1.05) | $768<br>($190, $1347) |
| KT18 (Nonparametric Plug-in) | 0.78<br>NA | $996<br>NA |
| KT18 (Linear) | 0.69<br>NA | 404<br>($-313,$1,121) |

*Note:* We present point estimates and confidence intervals based on both the trimmed and untrimmed datasets, along with the linear rule estimates (with its confidence intervals) and the nonparametric plug-in rule estimates from Kitagawa and Tetenov (2018).

# References

BHATTACHARYA, D. and DUPAS, P. (2012). Inferring welfare maximizing treatment assignment under budget constraints. *Journal of Econometrics*, **167** (1), 168–196.

BLOOM, H. S., ORR, L. L., BELL, S. H., CAVE, G., DOOLITTLE, F., LIN, W. and BOS, J. M. (1997). The benefits and costs of jtpa title ii-a programs: Key findings from the national job training partnership act study. *Journal of Human Resources*, **32** (3), 549–576.

CATTANEO, M. D., TITIUNIK, R. and YU, R. R. (2025a). Estimation and inference in boundary discontinuity designs: Distance-based methods. *arXiv preprint arXiv:2505.05670*.

—, — and — (2025b). Estimation and inference in boundary discontinuity designs: Location-based methods. *arXiv preprint arXiv:2505.05670*.

CHEN, X., CHRISTENSEN, T. and KANKANALA, S. (2025). Adaptive estimation and uniform confidence bands for nonparametric structural functions and elasticities. *The Review of Economic Studies*, **92** (1), 162–196.

— and CHRISTENSEN, T. M. (2018). Optimal sup-norm rates and uniform inference on nonlinear functionals of nonparametric iv regression. *Quantitative Economics*, **9** (1), 39–84.

— and GAO, W. Y. (2025). Semiparametric learning of integral functionals on submanifolds. *arXiv preprint arXiv:2507.12673*.

— and LIAO, Z. (2014). Sieve m inference on irregular parameters. *Journal of Econometrics*, **182** (1), 70–86.

— and — (2015). Sieve semiparametric two-step gmm under weak dependence. *Journal of Econometrics*, **189** (1), 163–186.

—, — and SUN, Y. (2014). Sieve inference on possibly misspecified semi-nonparametric time series models. *Journal of Econometrics*, **178**, 639–658.

— and POUZO, D. (2015). Sieve wald and qlr inferences on semi/nonparametric conditional moment models. *Econometrica*, **83** (3), 1013–1079.

DELFOUR, M. C. and ZOLÉSIO, J.-P. (2001). *Shapes and geometries: metrics, analysis, differential calculus, and optimization*. SIAM.

EVANS, L. C. and GARIEPY, R. F. (2015). *Measure Theory and Fine Properties of Functions.* CRC Press.

FENG, K., HONG, H. and NEKIPELOV, D. (2025). Statistical inference of optimal allocations i: Regularities and their implications. *arXiv preprint arXiv:2403.18248.*

KITAGAWA, T. and TETENOV, A. (2018). Who should be treated? empirical welfare maximization methods for treatment choice. *Econometrica*, **86** (2), 591–616.

PARK, G. (2025). Debiased machine learning when nuisance parameters appear in indicator functions. *arXiv preprint arXiv:2403.15934.*

SCHELLHASE, C. and KAUERMANN, G. (2012). Density estimation and comparison with a penalized mixture approach. *Computational Statistics*, **27**, 757–777.

STONE, C. J., HANSEN, M. H., KOOPERBERG, C. and TRUONG, Y. K. (1997). Polynomial splines and their tensor products in extended linear modeling. *The Annals of Statistics*, **25** (4), 1371–1425.

WHITEHOUSE, J., AUSTERN, M. and SYRGKANIS, V. (2025). Inference on optimal policy values and other irregular functionals via smoothing. *arXiv preprint arXiv:2507.11780.*

# A   Main Proofs

*Proof of Theorem 1.* For any $\nu$ s.t. $\int \nu^2(x) f(x) < \infty$, write $h_t := h_0 + t\nu$ and consider the functional derivative

$$D_h \overline{W}(h_0)[\nu] = \left. \frac{d}{dt}\overline{W}(h_t) \right|_{t=0}$$
$$= \left. \frac{d}{dt} \int [h_t(x)]_+ f(x)\, dx \right|_{t=0}.$$

We first show the interchangeability of the differentiation and integration holds:

$$\left. \frac{d}{dt} \int [h_t(x)]_+ f(x)\, dx \right|_{t=0} = \int \left. \frac{\partial}{\partial t}[h_t(x)]_+ \right|_{t=0} f(x)\, dx.$$

Notice that $[h_t(x)]_+$ is Lipchitz in $t$, we have

$$\left| [h_t(x)]_+ - [h_s(x)]_+ \right| \le |\nu(x)| \, |t - s|$$

with

$$\int |\nu(x)| f(x)\, dx < \infty.$$

28

Furthermore, $[h_t(x)]_+$ is almost everywhere differentiable in $t$ with

$$\left.\frac{\partial}{\partial t}[h_t(x)]_+\right|_{t=0} = \mathbb{1}\{h_0(x) \geq 0\}\nu(x) \tag{25}$$

for any $x$ s.t. $h_0(x) \neq 0$. Since $\{x : h(x) = 0\}$ has Lebesgue measure 0, we have

$$P_f(h_0(X_i) = 0) = 0.$$

Hence, (25) holds a.s.-$f$ in $x$, and thus by the dominated convergence theorem, we have

$$\left.\frac{d}{dt}\int[h_t(x)]_+ f(x)\,dx\right|_{t=0} = \int\left.\frac{\partial}{\partial t}[h_t(x)]_+\right|_{t=0} f(x)\,dx.$$
$$= \int \mathbb{1}\{h_0(x) \geq 0\}\nu(x)f(x)dx \tag{26}$$

We now switch the notation from $h_0$ to $\mu_0$ for subsequent analysis, and consider the functional derivative of $\overline{W}(\mu_0)$ w.r.t. $\mu$ in the direction of $\nu$. Note that $\mu(x,d)$ and $\nu(x,d)$ are functions of both $x$ and $d$. Applying (26), we have

$$D_\mu\overline{W}(\mu_0)[\nu] := \left.\frac{d}{dt}\overline{W}(\mu_0 + t\nu)\right|_{t=0}$$

$$= \int \mathbb{1}\{\mu_0(x,1) - \mu_0(x,0) \geq 0\}(\nu(x,1) - \nu(x,0))f(x)\,dx$$

$$= \int \mathbb{1}\{h_0(x) \geq 0\}(\nu(x,1) - \nu(x,0))\lambda(x)f_0(x)\,dx$$

$$= \mathbb{E}[\mathbb{1}\{h_0(X_i) \geq 0\}\lambda(X_i)(\nu(X_i,1) - \nu(X_i,0))]$$

$$= \mathbb{E}\left[\mathbb{1}\{h_0(X_i) \geq 0\}\lambda(X_i)\left(\frac{D_i}{p_0(X_i)}\nu(X_i,D_i) - \frac{1-D_i}{1-p_0(X_i)}\nu(X_i,D_i)\right)\right]$$

$$= \mathbb{E}\left[\mathbb{1}\{h_0(X_i) \geq 0\}\lambda(X_i)\left(\frac{D_i}{p_0(X_i)} - \frac{1-D_i}{1-p_0(X_i)}\right)\nu(X_i,D_i)\right]$$

$$= \mathbb{E}[\nu^*(X_i,D_i)\nu(X_i,D_i)] \tag{27}$$

where $\lambda := f/f_0$, $\mathbb{E}[\cdot]$ is expectation taken with respect to the training data distribution, and

$$\nu^*(x,d) := \mathbb{1}\{h_0(x) \geq 0\}\lambda(x)\left(\frac{d}{p_0(x)} - \frac{1-d}{1-p_0(x)}\right) \tag{28}$$

is the Riesz representer for the linear functional $D_\mu\overline{W}(\mu_0)[\cdot]$. Since

$$\mathbb{E}\left[\left.\left(\frac{D_i}{p_0(X_i)} - \frac{1-D_i}{1-p_0(X_i)}\right)^2\right|X_i\right]$$

$$= \mathbb{E}\left[\left.\left(\frac{D_i - D_ip_0(X_i) - p_0(X_i) + D_ip_0(X_i)}{p_0(X_i)(1-p_0(X_i))}\right)^2\right|X_i\right]$$

$$= \frac{\mathbb{E}\left[(D_i - p_0(X_i))^2 \mid X_i\right]}{p_0^2(X_i)(1 - p_0(X_i))^2} = \frac{p_0(X_i)(1 - p_0(X_i))}{p_0^2(X_i)(1 - p_0(X_i))^2}$$

$$= \frac{1}{p_0(X_i)(1 - p_0(X_i))}$$

under Assumption 1(b), the Riesz representer $\nu^*$ has finite norm

$$\|\nu^*\|^2 = \mathbb{E}\left[\mathbb{1}\{h_0(X_i) \geq 0\} \lambda^2(X_i)\left(\frac{D_i}{p_0(X_i)} - \frac{1 - D_i}{1 - p_0(X_i)}\right)^2\right]$$

$$= \mathbb{E}\left[\frac{\mathbb{1}\{h_0(X_i) \geq 0\} \lambda^2(X_i)}{p_0(X_i)(1 - p_0(X_i))}\right] < \infty$$

showing that $D_\mu \overline{W}(\mu_0)$ is a regular linear functional.

Next, we control the remainder term from the linearization. Specifically, note that

$$D_h W(h)[\nu] = \int \mathbb{1}\{h(x) \geq 0\} \nu(x) f(x) dx$$

which is exactly of the form of the generic value functional. We then have

$$D_h^2 W(h_0)[\nu, u] = \int_{\{x \in \mathbb{R}^d : h_0(x) = 0\}} \frac{1}{\|\nabla_x h_0(x)\|} \nu(x) u(x) f(x) d\mathcal{H}^{d-1}(x)$$

$$\leq \frac{1}{\epsilon} \|\nu\|_\infty \|u\|_\infty.$$

Similar bound applies to $D_\mu^2 \overline{W}(h_0)[\nu, u]$ as well. Thus we obtain:

$$\left|\overline{W}(\hat{\mu}) - \overline{W}(\mu_0) - D_\mu \overline{W}(\mu_0)[\hat{\mu} - \mu_0]\right| \leq M \|\hat{\mu} - \mu_0\|_\infty^2.$$

Hence, we have

$$\sqrt{n}\left(\overline{W}(\hat{\mu}) - \overline{W}(\mu_0)\right) = \sqrt{n} D_\mu \overline{W}(\mu_0)[\hat{\mu} - \mu_0] + \sqrt{n} O_p\left(\|\hat{\mu} - \mu_0\|_\infty^2\right)$$

$$= \frac{1}{\sqrt{n}} \sum_{i=1}^n \nu^*(X_i, D_i) \epsilon_i + o_p(1) + \sqrt{n} o_p\left(n^{-1/2}\right)$$

$$\xrightarrow{d} \mathcal{N}\left(0, \sigma_W^2\right)$$

where, writing $\sigma_\epsilon^2(x) := \mathbb{E}\left[\epsilon_i^2 \mid X_i = x\right]$,

$$\sigma_W^2 := \mathrm{Var}\left(\nu^*(X_i, D_i) \epsilon_i\right) = \mathbb{E}\left[\frac{\mathbb{1}\{h_0(X_i) \geq 0\} \lambda^2(X_i) \sigma_\epsilon^2(X_i)}{p_0(X_i)(1 - p_0(X_i))}\right]$$

$$= \int \frac{\mathbb{1}\{h_0(x) \geq 0\} \lambda^2(x) \sigma_\epsilon^2(x)}{p_0(x)(1 - p_0(X_i))} f_0(x) dx$$

$$= \int \frac{\mathbb{1}\{h_0(x) \geq 0\} \lambda(x) \sigma_\epsilon^2(x)}{p_0(x)(1 - p_0(X_i))} f(x) dx$$

$\square$

*Proof of Theorem 2.* We apply the derivations in Theorem 1 with $f = f_0$ and consequently

$\lambda \equiv 1$.

Consider the following standard empircal process decomposition

$$\sqrt{n}\left(\mathbb{P}_n g_\mu - P g_{\mu_0}\right) = \sqrt{n} P\left(g_{\hat{\mu}} - g_{\mu_0}\right) + \mathbb{G}_n g_{\mu_0} + \mathbb{G}_n\left(g_{\hat{\mu}} - g_{\mu_0}\right), \tag{29}$$

where $g_\mu(x) := \left[\mu(x,1) - \mu(x,0)\right]_+$, $\mathbb{P}_n g := \frac{1}{n}\sum_{i=1}^n g$, $P g = \int g dF_0$ and $\mathbb{G}_n := \sqrt{n}\left(\mathbb{P}_n - P\right)$. Note that the term $\sqrt{n} P\left(g_{\hat{h}} - g_{h_0}\right)$ corresponds the analysis of population expectation (integral) with respect to the true distribution $F = F_0$ in the last subsection. There are two additional terms that appear due to the use of the sample average: a "stochastic equicontinuity" term $\mathbb{G}_n\left(g_{\hat{\mu}} - g_{\mu_0}\right)$ that will be shown to be asymptotically negligible under the permanance of Donsker property, as well as a "CLT"-term $\mathbb{G}_n g_{\mu_0}$ that adds to the asymptotic variance of the estimator.

Under Assumption 1(c), $\hat{\mu} - \mu_0$ is assumed to belong to a Donsker class of functions, and the Lipchitz transformation of $\mu_0$ through the ReLU function $[\cdot]_+$ preserves the Donsker property. Hence, $g_{\hat{\mu}} - g_{\mu_0}$ also belongs to a Donsker class and thus $\mathbb{G}_n\left(g_{\hat{\mu}} - g_{\mu_0}\right) = o_p(1)$.

Then, by (29) we have,

$$\sqrt{n}\left(\hat{W}\left(\hat{h}\right) - W(h_0)\right) \equiv \sqrt{n}\left(\hat{\overline{W}}(\hat{\mu}) - \overline{W}(\mu_0)\right)$$

$$= \frac{1}{\sqrt{n}}\sum_{i=1}^n\left(\left[h_0(X_i)\right]_+ - W(h_0) + \nu^*(X_i, D_i)\epsilon_i\right) + o_p(1) \xrightarrow{d} \mathcal{N}\left(0, \bar{\sigma}_W^2\right)$$

where

$$\bar{\sigma}_W^2 := \mathrm{Var}\left(\left[h_0(X_i)\right]_+ - W(h_0) + \nu^*(X_i, D_i)\epsilon_i\right)$$

$$= \mathrm{Var}\left(\left[h_0(X_i)\right]_+\right) + \mathrm{Var}\left(\nu^*(X_i, D_i)\epsilon_i\right) + 2\mathrm{Cov}\left(\left[h_0(X_i)\right]_+, \nu^*(X_i, D_i)\epsilon_i\right)$$

$$= \mathrm{Var}\left(\left[h_0(X_i)\right]_+\right) + \mathrm{Var}\left(\nu^*(X_i, D_i)\epsilon_i\right)$$

$$= \mathrm{Var}\left(\left[h_0(X_i)\right]_+\right) + \sigma_W^2.$$

$\square$

*Proof of Theorem 3.* Recall that

$$\theta_0 = V(h_0) := \int \mathbb{1}\left\{h_0(x) \geq 0\right\} v_0(x) f(x) dx, \quad \hat{\theta} = V\left(\hat{h}\right) = \overline{V}(\hat{\mu}).$$

Taking the functional derivative according to Chen & Gao (2025), we obtain the following submanifold integral with submanifold dimension $m = d - 1$,

$$D_h V(h_0)[\nu] := \int_{\left\{x \in \mathbb{R}^d : h_0(x) = 0\right\}} \frac{\nu(x)}{\|\nabla_x h_0(x)\|} v_0(x) f(x) d\mathcal{H}^{d-1}(x).$$

Equivalently, using the notation $\mu_0$, we have

$$D_\mu \overline{V}(\mu_0)[\nu] := \int_{\{x \in \mathbb{R}^d : h_0(x) = 0\}} \frac{(\nu(x,1) - \nu(x,0))}{\|\nabla_x h_0(x)\|} v_0(x) f(x) d\mathcal{H}^{d-1}(x).$$

Applying Theorem 4 and Proposition 2 in Chen and Gao (2025), we obtain that

$$\frac{\sqrt{n}\left(\hat{\theta} - \theta_0\right)}{\sigma_{V,n}} \xrightarrow{d} \mathcal{N}(0,1) \quad \text{with } \sigma_{V,n}^2 \asymp K_n^{\frac{1}{d}}.$$

$\square$

*Proof of Theorem 4.* We now work with the following rescaled empirical process decomposition

$$\sqrt{\frac{n}{\sigma_{V,n}^2}} \left(\mathbb{P}_n g_\mu - P g_{\mu_0}\right) = \sqrt{\frac{n}{\sigma_{V,n}^2}} P\left(g_{\hat{\mu}} - g_{\mu_0}\right) + \sigma_{V,n}^{-1} \mathbb{G}_n g_{\mu_0} + \sigma_{V,n}^{-1} \mathbb{G}_n \left(g_{\hat{\mu}} - g_{\mu_0}\right), \tag{30}$$

where $g_\mu(x) := \mathbb{1}\{\mu(x,1) - \mu(x,0) \geq 0\}$.

Note that the term $P\left(g_{\hat{\mu}} - g_{\mu_0}\right) \equiv \overline{V}(\mu) - \overline{V}(\mu_0)$ has been analyzed in Section 4.1, where we have shown that

$$\sqrt{\frac{n}{\sigma_{V,n}^2}} P\left(g_{\hat{\mu}} - g_{\mu_0}\right) \xrightarrow{d} \mathcal{N}(0,1).$$

Given the above, the term $\sigma_{V,n}^{-1} \mathbb{G}_n g_{\mu_0} = O_p\left(\sqrt{K_n^{-1/d}}\right) = o_p(1)$ becomes asymptotically negligible. Below we seek to show that, in fact, the last term is also asymptotically negligible:

$$\sigma_{V,n}^{-1} \mathbb{G}_n \left(g_{\hat{\mu}} - g_{\mu_0}\right) = o_p(1).$$

Note that $g_\mu(x)$ involves a discontinuous indicator function, which does not preserve the Donsker property for the Holder class in general.[8] We thus directly derive the Donsker property via a maximal inequality on the functional class

$$\mathcal{G}_{a_n} := \{g_\mu - g_{\mu_0} : \|\mu - \mu_0\|_\infty \leq a_n\}.$$

We first obtain an envelope function for $\mathcal{G}_{a_n}$ and its second moment:

$$\begin{aligned}
&|g_\mu(x) - g_{\mu_0}(x)| \\
&= |\mathbb{1}\{\mu(x,1) - \mu(x,0) \geq 0\} - \mathbb{1}\{\mu_0(x,1) - \mu_0(x,0) \geq 0\}| \\
&= \mathbb{1}\{\mu(x,1) - \mu(x,0) \geq 0 > \mu_0(x,1) - \mu_0(x,0)\} \\
&\quad + \mathbb{1}\{\mu_0(x,1) - \mu_0(x,0) \geq 0 > \mu(x,1) - \mu(x,0)\} \\
&\leq \mathbb{1}\{\mu_0(x,1) - \mu_0(x,0) + 2\|\mu - \mu_0\|_\infty \geq 0 > \mu_0(x,1) - \mu_0(x,0)\}
\end{aligned}$$

---

[8]In contrast, the indicator function transformation preserves the Donsker property for VC classes of functions, in which case the Donsker property would deliver $\mathbb{G}_n\left(g_{\hat{\mu}} - g_{\mu_0}\right) = o_p(1)$, which implies the weaker condition $K_n^{-1/d} \mathbb{G}_n\left(g_{\hat{\mu}} - g_{\mu_0}\right) = o_p(1)$ required in this paper.

$$+ \mathbb{1}\left\{ \mu_0\left(x,1\right) - \mu_0\left(x,0\right) \geq 0 > \mu_0\left(x,1\right) - \mu_0\left(x,0\right) - 2\left\|\mu - \mu_0\right\|_\infty \right\}$$

$$\leq \mathbb{1}\left\{ \mu_0\left(x,1\right) - \mu_0\left(x,0\right) + 2a_n \geq 0 > \mu_0\left(x,1\right) - \mu_0\left(x,0\right) \right\}$$

$$+ \mathbb{1}\left\{ \mu_0\left(x,1\right) - \mu_0\left(x,0\right) \geq 0 > \mu_0\left(x,1\right) - \mu_0\left(x,0\right) - 2a_n \right\}$$

$$= \mathbb{1}\left\{ \left| \mu_0\left(x,1\right) - \mu_0\left(x,0\right) \right| \leq 2a_n \right\}$$

$$=: G_{a_n}$$

with

$$PG_{a_n}^2 = \mathbb{P}\left( \left| \mu_0\left(X_i,1\right) - \mu_0\left(X_i,0\right) \right| \leq 2a_n \right)$$

$$= \mathbb{P}\left( \left| h_0\left(X\right) \right| \leq 2a_n \right)$$

$$\leq Ma_n.$$

Then, provided a finite uniform entropy integral $J_{\mathcal{G}}$,

$$P \sup_{\left\|\mu - \mu_0\right\|_\infty \leq a_n} \left| \mathbb{G}_n\left(g_\mu - g_{\mu_0}\right) \right| \leq J_{\mathcal{G}}\sqrt{PG_{a_n}^2} \leq M\sqrt{a_n}$$

and thus

$$\sigma_{V,n}^{-1}\mathbb{G}_n\left(g_\mu - g_{\mu_0}\right) = O_p\left( \sqrt{K_n^{-1/d}a_n} \right) = o_p\left(1\right).$$

Hence,

$$\sqrt{\frac{n}{\sigma_{V,n}^2}}\left( \hat{V}\left(\hat{h}\right) - V\left(h_0\right) \right) \equiv \sqrt{\frac{n}{\sigma_{V,n}^2}}\left( \hat{\overline{V}}\left(\hat{\mu}\right) - \overline{V}\left(\mu_0\right) \right)$$

$$= \sqrt{\frac{n}{\sigma_{V,n}^2}}\left( \overline{V}\left(\hat{\mu}\right) - \overline{V}\left(\mu_0\right) \right) + o_p\left(1\right)$$

$$\xrightarrow{d} \mathcal{N}\left(0,1\right).$$

$\square$

# B  Additional Simulation Results

## B.1  Theorem 1

In principle, a variety of nonparametric estimators can be employed in the first step for the nuisance parameters, provided their convergence rates are sufficiently fast. In this subsection, we assess the estimation and inference performance of the welfare functional estimator under a known target distribution $F$, using a generalized additive model (GAM) as the first-stage estimator. Specifically, for Models 1-3, we use B-splines as the smooth terms in GAM, while for Models 4-7, we adopt the default tensor product smooths with cubic regression splines.

The simulation results are presented in Table 8. In Models 6 and 7, the confidence

33

intervals exhibit overcoverage relative to the nominal 95% level. This likely reflects the slow convergence of the variance plug-in estimator toward the true asymptotic variance of the welfare estimator, potentially due to suboptimal choices of smooth terms in the GAM specification.

Table 8: Theorem 1 Simulation Results for Models 1–7 (GAM)

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|-------|------|--------|---------|--------|--------|--------|----------|
| M1 | 1500 | 0.3857 | 0.0062 | 0.0474 | 0.0481 | 0.0035 | 0.9500 |
|    | 3000 | 0.3857 | 0.0025 | 0.0341 | 0.0339 | 0.0018 | 0.9525 |
|    | 6000 | 0.3857 | -0.0009 | 0.0231 | 0.0239 | 0.0009 | 0.9545 |
| M2 | 1500 | 0.2358 | -0.0092 | 0.0445 | 0.0513 | 0.0027 | 0.9670 |
|    | 3000 | 0.2358 | -0.0057 | 0.0335 | 0.0365 | 0.0013 | 0.9610 |
|    | 6000 | 0.2358 | -0.0048 | 0.0239 | 0.0259 | 0.0006 | 0.9615 |
| M3 | 1500 | 0.5001 | 0.0037 | 0.0587 | 0.0604 | 0.0022 | 0.9580 |
|    | 3000 | 0.5001 | 0.0021 | 0.0409 | 0.0430 | 0.0011 | 0.9605 |
|    | 6000 | 0.5001 | -0.0002 | 0.0294 | 0.0305 | 0.0006 | 0.9600 |
| M4 | 1500 | 0.1033 | 0.0303 | 0.0465 | 0.0550 | 0.0065 | 0.9615 |
|    | 3000 | 0.1033 | 0.0193 | 0.0338 | 0.0394 | 0.0036 | 0.9635 |
|    | 6000 | 0.1033 | 0.0130 | 0.0243 | 0.0281 | 0.0021 | 0.9610 |
| M5 | 1500 | 0.0499 | 0.0316 | 0.0407 | 0.0527 | 0.0087 | 0.9685 |
|    | 3000 | 0.0499 | 0.0217 | 0.0294 | 0.0372 | 0.0050 | 0.9680 |
|    | 6000 | 0.0499 | 0.0138 | 0.0218 | 0.0261 | 0.0030 | 0.9635 |
| M6 | 1500 | 0.2315 | 0.0035 | 0.0450 | 0.0634 | 0.0039 | 0.9930 |
|    | 3000 | 0.2315 | 0.0004 | 0.0314 | 0.0448 | 0.0019 | 0.9965 |
|    | 6000 | 0.2315 | -0.0008 | 0.0234 | 0.0317 | 0.0010 | 0.9920 |
| M7 | 1500 | 0.1250 | 0.0040 | 0.0380 | 0.0545 | 0.0064 | 0.9950 |
|    | 3000 | 0.1250 | 0.0026 | 0.0272 | 0.0385 | 0.0030 | 0.9920 |
|    | 6000 | 0.1250 | 0.0016 | 0.0199 | 0.0272 | 0.0015 | 0.9915 |

*Notes:* (1) $W_{\text{true}}$ denotes the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\hat{\hat{W}}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

## B.2    Theorem 2

We next examine the estimation and inference performance of the welfare functional estimator under an unknown target distribution $F$, using GAM as the first-stage estimator.

The GAM specifications are identical to those described in the previous subsection. The simulation results, reported in Table 9, demonstrate that our proposed inference procedure exhibits strong finite-sample performance.

Table 9: Theorem 2 Simulation Results for Models 1–7 (GAM)

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|---|---|---|---|---|---|---|---|
| M8 | 1,500 | 0.3857 | 0.0077 | 0.0412 | 0.0421 | 0.0025 | 0.951 |
|    | 3,000 | 0.3857 | 0.0036 | 0.0296 | 0.0297 | 0.0013 | 0.952 |
|    | 6,000 | 0.3857 | 0.0008 | 0.0206 | 0.0209 | 0.0007 | 0.956 |
| M9 | 1,500 | 0.2358 | 0.0000 | 0.0414 | 0.0436 | 0.0021 | 0.965 |
|    | 3,000 | 0.2358 | 0.0008 | 0.0297 | 0.0310 | 0.0010 | 0.962 |
|    | 6,000 | 0.2358 | 0.0002 | 0.0206 | 0.0220 | 0.0005 | 0.960 |
| M10 | 1,500 | 0.5001 | 0.0059 | 0.0512 | 0.0514 | 0.0018 | 0.952 |
|     | 3,000 | 0.5001 | 0.0033 | 0.0357 | 0.0366 | 0.0010 | 0.957 |
|     | 6,000 | 0.5001 | 0.0010 | 0.0252 | 0.0260 | 0.0005 | 0.962 |
| M11 | 1,500 | 0.1033 | 0.0266 | 0.0362 | 0.0396 | 0.0036 | 0.938 |
|     | 3,000 | 0.1033 | 0.0154 | 0.0261 | 0.0284 | 0.0020 | 0.940 |
|     | 6,000 | 0.1033 | 0.0094 | 0.0190 | 0.0202 | 0.0011 | 0.942 |
| M12 | 1,500 | 0.0499 | 0.0277 | 0.0324 | 0.0370 | 0.0052 | 0.941 |
|     | 3,000 | 0.0499 | 0.0172 | 0.0232 | 0.0263 | 0.0029 | 0.946 |
|     | 6,000 | 0.0499 | 0.0106 | 0.0169 | 0.0187 | 0.0016 | 0.941 |
| M13 | 1,500 | 0.2315 | 0.0062 | 0.0431 | 0.0456 | 0.0024 | 0.964 |
|     | 3,000 | 0.2315 | 0.0027 | 0.0308 | 0.0323 | 0.0012 | 0.959 |
|     | 6,000 | 0.2315 | 0.0007 | 0.0218 | 0.0229 | 0.0006 | 0.957 |
| M14 | 1,500 | 0.1250 | 0.0056 | 0.0367 | 0.0394 | 0.0043 | 0.957 |
|     | 3,000 | 0.1250 | 0.0023 | 0.0249 | 0.0276 | 0.0019 | 0.968 |
|     | 6,000 | 0.1250 | 0.0007 | 0.0176 | 0.0195 | 0.0009 | 0.969 |

*Notes:* (1) $W_{\text{true}}$ denotes the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\widehat{\overline{W}}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

## B.3   Theorem 3

In addition to the model specification considered in the main text for the Theorem 3 simulations, we conducted two auxiliary simulation studies with slightly modified designs: one with the target population $F \sim \text{U}([-1.75, 1.75]^2)$ and the other with $F \sim \text{U}([-1.9, 1.9]^2)$, keeping all other settings fixed. The results are reported in Table 10. We find that the

coverage rate approaches 95% as the sample size increases in both cases. However, in the latter case the supports of the treated and control groups do not fully nest the target population $F$, potentially leading to extrapolation bias of the first-stage estimator and a slight undercoverage.

Table 10: Sensitivity Checks for Model 15

| Case | $n$ | $W_{\text{true}}$ | Bias | SD | SE | Coverage |
|------|-----|-------------------|------|-----|-----|----------|
| (a) $F \sim \text{U}[-1.75, 1.75]^2$ | 1500 | 3.1416 | 0.0222 | 0.1289 | 0.1198 | 0.926 |
| | 3000 | 3.1416 | 0.0077 | 0.0801 | 0.0796 | 0.942 |
| | 6000 | 3.1416 | 0.0089 | 0.0542 | 0.0558 | 0.952 |
| (b) $F \sim \text{U}[-1.9, 1.9]^2$ | 1500 | 3.1416 | 0.0118 | 0.0746 | 0.0676 | 0.923 |
| | 3000 | 3.1416 | 0.0047 | 0.0474 | 0.0453 | 0.938 |
| | 6000 | 3.1416 | 0.0037 | 0.0334 | 0.0317 | 0.936 |

Besides the sieve dimensions for estimating $\mu_0(x, 1)$ and $\mu_0(x, 0)$, $\epsilon$ is an additional tuning parameter that determines how well the level set is approximated. To assess the sensitivity of our method to this parameter, we consider the alternative values $\epsilon = 0.0025$ and $\epsilon = 0.0075$ while keeping anything else fixed. Table 11 shows that the simulation results are not materially affected by the choice of $\epsilon$.

Table 11: Sensitivity Checks for Model 15

| Case | $n$ | $W_{\text{true}}$ | Bias | SD | SE | Coverage |
|------|-----|-------------------|------|-----|-----|----------|
| (a) $\epsilon = 0.0075$ | 1500 | 3.1416 | 0.0120 | 0.0642 | 0.0635 | 0.9395 |
| | 3000 | 3.1416 | 0.0126 | 0.0442 | 0.0448 | 0.9445 |
| | 6000 | 3.1416 | 0.0102 | 0.0311 | 0.0316 | 0.9445 |
| (b) $\epsilon = 0.0025$ | 1500 | 3.1416 | 0.0120 | 0.0642 | 0.0638 | 0.9420 |
| | 3000 | 3.1416 | 0.0126 | 0.0442 | 0.0450 | 0.9445 |
| | 6000 | 3.1416 | 0.0102 | 0.0311 | 0.0317 | 0.9450 |

## B.4 Simulation Results for DML

An alternative to the semi-parametric two-step estimation of the welfare functional is to use the double-debiased machine learning approach. The DML estimator is implemented via a cross-fitting scheme. Specifically, the sample is partitioned into K folds. For each fold, we use the observations in the remaining $K - 1$ folds to obtain nonparametric estimates of the nuisance functions $\mu_0(x, 1)$ and $\mu_0(x, 0)$ (in our case, using GAM). These estimates are then applied to the held-out fold to take an average of the fitted indicator $\mathbb{1}\{\hat{\mu}_0(x, 1) - \hat{\mu}_0(x, 0) \geq$

0}. This process is repeated so that each fold serves once as the testing set, and the final estimator is obtained by averaging the resulting values across all folds. A confidence interval could then be constructed as in Theorem 2 with the debiased ML estimate in replacement of the semi-parametric two-step estimate.

We simulate 1,000 iterations for Models 1 to 3 with $K = 5$, and the results are summarized in table 12. The confidence intervals exhibit slight undercoverage in small samples, but this issue diminishes as the sample size increases. This pattern is consistent with the asymptotic nature of the interval's construction.

Table 12: Theorem DML Simulation Results

| Model | n | $W_0$ | Bias | SD | SE | Coverage |
|---|---|---|---|---|---|---|
| Model 1 | 1,500 | 0.3857 | -0.01271 | 0.04583 | 0.04215 | 0.906 |
| Model 1 | 3,000 | 0.3857 | -0.00776 | 0.03184 | 0.02969 | 0.921 |
| Model 1 | 6,000 | 0.3857 | -0.00392 | 0.02193 | 0.02092 | 0.934 |
| Model 2 | 1,500 | 0.2358 | -0.01051 | 0.04900 | 0.04372 | 0.916 |
| Model 2 | 3,000 | 0.2358 | -0.00475 | 0.03319 | 0.03097 | 0.934 |
| Model 2 | 6,000 | 0.2358 | -0.00195 | 0.02297 | 0.02200 | 0.929 |
| Model 3 | 1,500 | 0.5001 | -0.00568 | 0.05636 | 0.05134 | 0.913 |
| Model 3 | 3,000 | 0.5001 | -0.00261 | 0.03701 | 0.03654 | 0.953 |
| Model 3 | 6,000 | 0.5001 | -0.00174 | 0.02628 | 0.02597 | 0.954 |

## B.5 Sieve Variance Estimation

# C Additional Results for Empirical Application

## C.1 Sensitivity Analysis for Tuning Parameters

When conducting inference for the value functional in the empirical analysis of the JTPA data, two tuning parameters must be specified: $\iota$ and $s$. The parameter $\iota$ determines the size of the set $\{x \in \chi : -\epsilon < \hat{h}(x) < \epsilon\}$. The parameter $s$ controls the smoothness of the estimated density function over the trimmed dataset by scaling the bandwidths used in kernel density estimation.

To evaluate robustness, we perform sensitivity analyses with respect to both tuning parameters. The results, reported in Table 18, show that the estimated value functional $\hat{\bar{V}}(\hat{\mu})$ is largely insensitive to the choice of $\iota$, but more responsive to the choice of $s$. This shows the importance of choosing an appropriate kernel density estimator.

Table 13: Theorem 1 Simulation Results for Models 1–3 (Sieve Variance)

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|-------|------|--------|---------|--------|--------|--------|----------|
| M1 | 1500 | 0.3857 | 0.0152 | 0.0469 | 0.0467 | 0.0035 | 0.9370 |
|    | 3000 | 0.3857 | 0.0076 | 0.0337 | 0.0329 | 0.0019 | 0.9390 |
|    | 6000 | 0.3857 | 0.0036 | 0.0229 | 0.0232 | 0.0011 | 0.9510 |
| M2 | 1500 | 0.2358 | 0.0039 | 0.0481 | 0.0488 | 0.0029 | 0.9480 |
|    | 3000 | 0.2358 | 0.0042 | 0.0350 | 0.0345 | 0.0014 | 0.9425 |
|    | 6000 | 0.2358 | 0.0025 | 0.0241 | 0.0245 | 0.0007 | 0.9535 |
| M3 | 1500 | 0.5001 | 0.0038 | 0.0580 | 0.0581 | 0.0023 | 0.9560 |
|    | 3000 | 0.5001 | 0.0019 | 0.0407 | 0.0413 | 0.0012 | 0.9560 |
|    | 6000 | 0.5001 | -0.0002 | 0.0293 | 0.0294 | 0.0006 | 0.9525 |

*Notes:* (1) $W_{\text{true}}$ denotes the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\widehat{W}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

Table 14: Theorem 1 Simulation Results for Models 4–7 (Sieve Variance)

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|-------|------|--------|--------|--------|--------|--------|----------|
| M4 | 1500 | 0.1033 | 0.0185 | 0.0478 | 0.0482 | 0.0088 | 0.9400 |
|    | 3000 | 0.1033 | 0.0088 | 0.0348 | 0.0350 | 0.0048 | 0.9445 |
|    | 6000 | 0.1033 | 0.0043 | 0.0248 | 0.0251 | 0.0024 | 0.9490 |
| M5 | 1500 | 0.0499 | 0.0282 | 0.0418 | 0.0431 | 0.0110 | 0.9310 |
|    | 3000 | 0.0499 | 0.0158 | 0.0303 | 0.0311 | 0.0068 | 0.9335 |
|    | 6000 | 0.0499 | 0.0094 | 0.0227 | 0.0224 | 0.0038 | 0.9290 |
| M6 | 1500 | 0.2315 | 0.0268 | 0.0559 | 0.0552 | 0.0057 | 0.9260 |
|    | 3000 | 0.2315 | 0.0123 | 0.0394 | 0.0402 | 0.0030 | 0.9480 |
|    | 6000 | 0.2315 | 0.0050 | 0.0288 | 0.0288 | 0.0015 | 0.9490 |
| M7 | 1500 | 0.1250 | 0.0462 | 0.0505 | 0.0502 | 0.0087 | 0.8920 |
|    | 3000 | 0.1250 | 0.0218 | 0.0346 | 0.0350 | 0.0046 | 0.9270 |
|    | 6000 | 0.1250 | 0.0101 | 0.0245 | 0.0245 | 0.0023 | 0.9400 |

*Notes:* (1) $W_{\text{true}}$ denotes the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\widehat{W}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

Table 15: Theorem 2 Simulation Results for Models 8–10 (Sieve Variance)

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|-------|------|--------|--------|--------|--------|--------|----------|
| M8  | 1500 | 0.3857 | 0.0068 | 0.0414 | 0.0417 | 0.0026 | 0.9475 |
|     | 3000 | 0.3857 | 0.0029 | 0.0296 | 0.0294 | 0.0014 | 0.9505 |
|     | 6000 | 0.3857 | 0.0006 | 0.0206 | 0.0208 | 0.0007 | 0.9555 |
| M9  | 1500 | 0.2358 | 0.0042 | 0.0425 | 0.0431 | 0.0026 | 0.9555 |
|     | 3000 | 0.2358 | 0.0029 | 0.0304 | 0.0305 | 0.0012 | 0.9475 |
|     | 6000 | 0.2358 | 0.0012 | 0.0209 | 0.0216 | 0.0006 | 0.9510 |
| M10 | 1500 | 0.5001 | 0.0067 | 0.0511 | 0.0511 | 0.0020 | 0.9480 |
|     | 3000 | 0.5001 | 0.0037 | 0.0357 | 0.0364 | 0.0011 | 0.9590 |
|     | 6000 | 0.5001 | 0.0013 | 0.0253 | 0.0259 | 0.0006 | 0.9580 |

*Notes:* (1) $W_{\text{true}}$ denotes the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\widehat{\overline{W}}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

Table 16: Theorem 2 Simulation Results for Models 11–14 (Sieve Variance)

| Model | $n$ | $W_{\text{true}}$ | Bias | SD | SE | SD(SE) | Coverage |
|-------|------|--------|--------|--------|--------|--------|----------|
| M11 | 1500 | 0.1033 | 0.0307 | 0.0365 | 0.0379 | 0.0041 | 0.9065 |
|     | 3000 | 0.1033 | 0.0156 | 0.0264 | 0.0272 | 0.0023 | 0.9290 |
|     | 6000 | 0.1033 | 0.0084 | 0.0192 | 0.0195 | 0.0013 | 0.9395 |
| M12 | 1500 | 0.0499 | 0.0418 | 0.0316 | 0.0344 | 0.0050 | 0.8410 |
|     | 3000 | 0.0499 | 0.0232 | 0.0229 | 0.0245 | 0.0031 | 0.8930 |
|     | 6000 | 0.0499 | 0.0126 | 0.0168 | 0.0175 | 0.0018 | 0.9205 |
| M13 | 1500 | 0.2315 | 0.0177 | 0.0432 | 0.0438 | 0.0030 | 0.9420 |
|     | 3000 | 0.2315 | 0.0088 | 0.0309 | 0.0313 | 0.0016 | 0.9440 |
|     | 6000 | 0.2315 | 0.0041 | 0.0221 | 0.0222 | 0.0008 | 0.9495 |
| M14 | 1500 | 0.1250 | 0.0251 | 0.0382 | 0.0386 | 0.0050 | 0.9230 |
|     | 3000 | 0.1250 | 0.0117 | 0.0258 | 0.0268 | 0.0024 | 0.9440 |
|     | 6000 | 0.1250 | 0.0054 | 0.0182 | 0.0188 | 0.0012 | 0.9415 |

*Notes:* (1) $W_{\text{true}}$ denotes the true welfare functional, computed numerically using 5,000 Sobol points drawn from $F$. (2) Bias is the average deviation from $W_{\text{true}}$. (3) SD is the sampling standard deviation. (4) SE is the average estimated asymptotic standard error of $\widehat{\overline{W}}(\hat{\mu})$. (5) SD(SE) is the standard deviation of SE across iterations. (6) Coverage is the empirical coverage probability of the 95% CI.

Table 17: Estimated Welfare Gains and Share of Population to be Treated Under Nonparametric Plug-in Rule (Sieve Variance)

(a) 30-Month Post-Program Earnings, No Treatment Cost

| Method | Share Treated | Est. Welfare Gain |
|---|---|---|
| Ours | 0.89<br>(0.73 1.05) | $1,519<br>($691, $2347) |
| KT (2018) | 0.91<br>NA | $1,693<br>NA |

(b) 30-Month Post-Program Earnings, $774 Cost per Treatment

| Method | Share Treated | Est. Welfare Gain |
|---|---|---|
| Ours | 0.80<br>(0.53, 1.07) | $858<br>($84, $1631) |
| KT (2018) | 0.78<br>NA | $996<br>NA |

*Note:* Two-sided 95% confidence intervals in parentheses, constructed based on the asymptotic distributions of the corresponding estimators.

Table 18: Sensitivity of $\hat{\bar{V}}(\hat{\mu})$ to Tuning Parameters $\iota$ and $s$ (Cost = F vs. Cost = T)

(a) Sensitivity to $\iota$

| $\iota$ | $\hat{V}$ | SE | CI Low | CI High | Num |
|---|---|---|---|---|---|
| **Cost = F** | | | | | |
| 0.005 | 0.8908 | 0.0866 | 0.7210 | 1.0606 | 475 |
| 0.0075 | 0.8908 | 0.0862 | 0.7219 | 1.0597 | 713 |
| 0.0100 | 0.8908 | 0.0828 | 0.7286 | 1.0530 | 931 |
| 0.0125 | 0.8908 | 0.0840 | 0.7262 | 1.0554 | 1203 |
| 0.0150 | 0.8908 | 0.0857 | 0.7228 | 1.0588 | 1463 |
| **Cost = T** | | | | | |
| 0.005 | 0.7971 | 0.1343 | 0.5338 | 1.0603 | 689 |
| 0.0075 | 0.7971 | 0.1303 | 0.5417 | 1.0525 | 1044 |
| 0.0100 | 0.7971 | 0.1373 | 0.5279 | 1.0662 | 1411 |
| 0.0125 | 0.7971 | 0.1399 | 0.5229 | 1.0712 | 1784 |
| 0.0150 | 0.7971 | 0.1406 | 0.5214 | 1.0727 | 2152 |

(b) Sensitivity to $s$

| $s$ | $\hat{V}$ | SE | CI Low | CI High | Num |
|---|---|---|---|---|---|
| **Cost = F** | | | | | |
| 1 | 0.8908 | 0.0896 | 0.7151 | 1.0665 | 931 |
| 2 | 0.8908 | 0.0860 | 0.7224 | 1.0593 | 931 |
| 3 | 0.8908 | 0.0828 | 0.7286 | 1.0530 | 931 |
| 4 | 0.8908 | 0.0809 | 0.7322 | 1.0494 | 931 |
| 5 | 0.8908 | 0.0793 | 0.7354 | 1.0462 | 931 |
| **Cost = T** | | | | | |
| 1 | 0.7971 | 0.1103 | 0.5809 | 1.0132 | 1411 |
| 2 | 0.7971 | 0.1254 | 0.5513 | 1.0428 | 1411 |
| 3 | 0.7971 | 0.1373 | 0.5279 | 1.0662 | 1411 |
| 4 | 0.7971 | 0.1420 | 0.5187 | 1.0754 | 1411 |
| 5 | 0.7971 | 0.1554 | 0.4925 | 1.1016 | 1411 |

*Notes:* (a) Cost = T if the outcome variable is the 30-month earnings minus an additional \$774 and F otherwise. (b) SE = $\hat{\sigma}_V/\sqrt{n}$. (c) $\epsilon = \iota \times \text{SD}(\hat{h}(x))$ over **X_trim**. (d) CI Low = $\hat{\bar{V}}(\hat{\mu}) - 1.96 \times$ SE, CI High = $\hat{\bar{V}}(\hat{\mu}) + 1.96 \times$ SE. (e) Num is the number of Sobol points whose $\hat{h}$ evaluations fall into $\{x \in \textbf{X\_trim} : |\hat{h}(x)| < \epsilon\}$.

## C.2 Sensitivity Analysis for Density Estimation

We employed a naive Gaussian kernel density estimator for the covariate distribution in the trimmed dataset, treating years of education as a continuous variable. While years of education is theoretically continuous, in practice it takes on only discrete values in the dataset. This motivates considering alternative density estimators that explicitly treat years of education as categorical in order to assess the robustness of our inference. To this end, we partitioned the trimmed dataset by educational level and examined three cases: (1) a naive Gaussian kernel density estimator with bandwidths selected by `Hscv()` and scaled by a small factor to ensure smoothness; (2) a logspline density estimator (Stone, Hansen, Kooperberg and Truong (1997)); (3) a penalized B-spline density estimator (Schellhase and Kauermann (2012)).

The results are presented in Table 19. As shown, when costs are taken into account, the Gaussian kernel density estimator discussed in the main text yields relatively conservative confidence intervals compared to the three alternatives. In contrast, when costs are not considered, the results remain largely unchanged across all the estimators. However, we should emphasize that these three alternative density estimators are applied based on the assumption that years of education is categorical, but in theory it is treated as a continuous variable, which provides justification for using the Gaussian kernel estimator in the main text.

Table 19: Sensitivity of $\hat{\overline{V}}(\hat{\mu})$ to Alternative Density Estimators

| Estimator | Cost = T | | | | Cost = F | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | $\hat{V}$ | SE | CI Low | CI High | $\hat{V}$ | SE | CI Low | CI High |
| Logspline | 0.7971 | 0.0932 | 0.6144 | 0.9797 | 0.8908 | 0.1041 | 0.6868 | 1.0948 |
| Kernel | 0.7971 | 0.0704 | 0.6590 | 0.9351 | 0.8908 | 0.0980 | 0.6988 | 1.0829 |
| Penalized B-spline | 0.7971 | 0.0601 | 0.6793 | 0.9148 | 0.8908 | 0.0691 | 0.7555 | 1.0262 |

*Notes:* (a) Cost = T if the outcome variable is 30-month earnings minus an additional \$774, and Cost = F otherwise. (b) SE = $\hat{\sigma}_V/\sqrt{n}$. (c) CI Low = $\hat{\overline{V}}(\hat{\mu}) - 1.96 \times$ SE, CI High = $\hat{\overline{V}}(\hat{\mu}) + 1.96 \times$ SE.

## C.3 Extrapolating Sieve Estimates

We also examine the case where the dataset is not trimmed. In this setting, interpolation is required to evaluate the estimated functions $\hat{\mu}(x, 0)$ and $\hat{\mu}(x, 1)$ outside the supports of the control and treated groups, respectively. The formulas for the welfare and value functional estimators, as well as their asymptotic variances, remain unchanged. The only difference is

that the sample averages are now computed over the full covariate support rather than the trimmed dataset. The results for this specification are reported in Table 20.

Table 20: Estimated Welfare Gains and Share of Population to be Treated Under Nonparametric Plug-in Rule Without Trimming

(a) 30-Month Post-Program Earnings, No Treatment Cost

| Method | Share Treated | Est. Welfare Gain |
|---|---|---|
| Ours | 0.92 (0.80 1.03) | $1,459 ($825, $2093) |
| KT (2018) | 0.91 NA | $1,693 NA |

(b) 30-Month Post-Program Earnings, $774 Cost per Treatment

| Method | Share Treated | Est. Welfare Gain |
|---|---|---|
| Ours | 0.85 (0.71, 0.99) | $768 ($164, $1373) |
| KT (2018) | 0.78 NA | $996 NA |

*Note:* Two-sided 95% confidence intervals in parentheses, constructed based on the asymptotic distributions of the corresponding estimators.

## C.4   Estimating $\hat{\Omega}$ Using the Trimmed Dataset

We used the full sample to estimate the asymptotic variance-covariance matrix $\hat{\Omega}$. For comparison, Table 21 reports the results obtained when $\hat{\Omega}$ is instead estimated using only the trimmed dataset. We observe that the CIs expands by a small factor.

## C.5   The Sieve Score Bootstrapping Procedure for Critical Values

Following Chen and Christensen (2018), one could also use the sieve score bootstrap procedure to calculate the critical value used in the construction of the 95% CI for $\overline{V}(\mu_0)$. The algorithm starts with making iid draws $\{\omega_i\}_{i=1}^n$ from a distribution independent of the data **df** with mean zero, unit variance, and finite third moment (e.g. a standard normal

Table 21: Estimated Share of Population to be Treated Under Nonparametric Plug-in Rule With $\hat{\Omega}$ Estimated Using Trimmed Data

(a) 30-Month Post-Program Earnings, No Treatment Cost

| Method | Share Treated |
|---|---|
| Ours | 0.89 |
| | (0.72, 1.06) |
| KT (2018) | 0.91 |
| | NA |

(b) 30-Month Post-Program Earnings, $774 Cost per Treatment

| Method | Share Treated |
|---|---|
| Ours | 0.80 |
| | (0.44, 1.15) |
| KT (2018) | 0.78 |
| | NA |

distribution), and then calculates the the bootstrap sieve t-statistic using the formula

$$Z_n^* = \frac{DV(\hat{\mu})[\nu]'(B'B/n)}{\hat{\sigma}_V/\sqrt{n}} \begin{bmatrix} \frac{1}{\sqrt{n}} \sum_{i=1}^{n_1} \psi^{K_1}(x_i)\hat{u}_i\omega_i \\ \frac{1}{\sqrt{n}} \sum_{j=n_1+1}^{n} \psi^{K_0}(x_j)\hat{u}_j\omega_j \end{bmatrix}$$

with $B$ being the design matrix of regression $Y_i$ on $D_i\psi^{K_1}(X_i)$ and $(1-D_i)\psi^{K_0}(X_i)$. The 95% quantile of the bootstrapped $|Z_n^*|$ could be used as the critical value in the construction of the CI. Setting the number of bootstrap equals 1000, the resulting critical values is 1.865758, not significantly different from 1.96.