

ESTIMATING A NONPARAMETRIC TRIANGULAR MODEL WITH BINARY ENDOGENOUS REGRESSORS^{*†}

SUNG JAE JUN,[‡] JORIS PINKSE[§] AND HAIQING XU[¶]

Center for the Study of Auctions, Procurements and Competition Policy

Department of Economics

The Pennsylvania State University and the University of Texas

May 2, 2016

We consider identification and estimation in a nonparametric triangular system with a binary endogenous regressor and nonseparable errors. For identification we take a control function approach utilizing the Dynkin system idea developed in [Jun, Pinkse, and Xu \(2011, JPX11\)](#) and extended in [Kédagni and Mourifié \(2014, KM14\)](#). We articulate various tradeoffs, including continuity, monotonicity, and differentiability. For estimation, we use the idea of local instruments under smoothness assumptions, as in e.g. [Carneiro and Lee \(2009, CL09\)](#), [Heckman and Vytlačil \(1999\)](#) but we do not assume additive separability in latent variables. Our estimator uses nonparametric kernel regression techniques and its statistical properties are derived using the functional delta method. We establish that it is $n^{2/7}$ -consistent and has a limiting normal distribution. We apply the method to estimate the returns on a college education. Unlike existing work, notably [CL09](#) and [Carneiro, Heckman, and Vytlačil \(2011, CHV11\)](#), we find that returns on a college education are consistently positive. The returns curves we estimate are moreover inconsistent with the shape restrictions imposed in those papers.

Key Words: Triangular Models; Endogeneity; Nonparametric Estimation.

JEL Classification Codes: C14; C31; C35; C36.

^{*}This paper is based on research supported by NSF grant SES-0922127. We thank the Human Capital Foundation for their support of the Penn State economics graduate program. We thank Ismael Mourifié and seminar participants at the 2014 Canadian Econometrics Study Group and departmental seminars at Georgetown, Indiana, LSE, Michigan, Stanford, and UCL.

[†]This research was conducted with restricted access to Bureau of Labor Statistics (BLS) data. The views expressed here do not reflect the views of the BLS.

[‡](corresponding author) 303 Kern Graduate Building, University Park 16802, sjun@psu.edu

[§]joris@psu.edu

[¶]h.xu@austin.utexas.edu

1. INTRODUCTION

We consider identification and estimation in a nonparametric triangular system with a binary endogenous regressor and nonseparable errors. For identification we take a control function approach utilizing the Dynkin system idea developed in [Jun, Pinkse, and Xu \(2011, JPX11\)](#) and extended in [Kédagni and Mourifié \(2014, KM14\)](#). We articulate various tradeoffs, including continuity, monotonicity, and differentiability. For estimation, we use the idea of local instruments under smoothness assumptions, as in e.g. [Carneiro and Lee \(2009, CL09\)](#), [Heckman and Vytlačil \(1999\)](#) but we do not assume additive separability in latent variables. Our estimator uses nonparametric kernel regression techniques and its statistical properties are derived using the functional delta method. We establish that it is $n^{2/7}$ -consistent and has a limiting normal distribution. We apply the method to estimate the returns on a college education. Unlike existing work, notably [CL09](#) and [Carneiro, Heckman, and Vytlačil \(2011, CHV11\)](#), we find that returns on a college education are consistently positive. The returns curves we estimate are moreover inconsistent with the shape restrictions imposed in those papers.

There is a related—but-different approach in the literature, namely the instrumental variable (IV) quantile model of e.g. [Chernozhukov and Hansen \(2005, 2006\)](#). Unlike the control function approach, the IV quantile model does not treat discrete endogenous regressors differently. Therefore, binary endogenous regressors do not cause any extra issues in identification and estimation. However, the control function approach is more suitable when it is of interest to explicitly control for potential heterogeneity from the *first* stage equation.

Most papers on nonparametric triangular models with nonseparable errors focus on nonparametric identification rather than estimation, especially when there is a discrete endogenous regressor. [Hoderlein and Mammen \(2007\)](#) consider a nonseparable model with exogenous regressors but without monotonicity. [Chesher \(2003, CH03\)](#) and [Imbens and Newey \(2009\)](#) study nonseparable models with continuous endogenous regressors under monotonicity. Based on the identification result of e.g. [CH03](#), [Ma and Koenker \(2006\)](#) and [Jun \(2009\)](#) propose parametric and semiparametric estimation methods, respectively, but both require that endogenous regressors be continuous. [Chesher \(2005, CH05\)](#) establishes partial identification of the structural function at a given value in a triangular system with a discrete endogenous regressor, but [CH05](#) contains little discussion on estimation and inference. [JPX11](#) reconsider [CH05](#)'s result and provide tighter bounds under a weaker rank condition using an independence assumption on instruments, but also stop short of estimation and inference.

A related but different literature uses an instrumental variable (IV) quantile approach (e.g. [Chernozhukov and Hansen, 2005, 2006](#)). Unlike the control function approach, the IV quantile model treats discrete endogenous regressors in the same way it treats continuous endogenous regressors. Therefore, binary endogenous regressors do not cause additional issues in identification and estimation. However, we use the control function approach, because it allows us to control for potential heterogeneity arising from the *first* stage equation explicitly.

[JPX11](#) and [CH05](#) serve as our starting point. The two papers study the same model, albeit that [JPX11](#) use a global independence condition of instruments and errors to weaken [CH05](#)'s rank condition and to tighten identification bounds. The difference between the two approaches is most profound in the presence of a binary endogenous regressor since [CH05](#), unlike [JPX11](#), does not allow for a binary endogenous regressor. [JPX11](#) in fact establish conditions under which point identification obtains in the case of continuous instruments.

One of our objectives in this paper is to show how the Dynkin system idea can be used to obtain point identification in the triangular model with a binary endogenous regressor. Unlike the local instrument approach of [Heckman and Vytlacil \(1999, 2001\)](#) and [CL09](#), we show that it is the continuous variation, rather than the differentiability, of the propensity score that is essential for identification, albeit that the latter is useful for estimation. A more formal discussion can be found in section 2.

Once we have articulated conditions for identification we turn our attention to nonparametric estimation. We propose a kernel-based nonparametric estimator, which allows for the full flexibility of the triangular model with nonseparable errors. We then develop limit results for the proposed estimator. These results are derived in section 3.

Finally, in section 5 we implement our estimator using the same NLSY-based data set that is used in [CL09](#) and [CHV11](#) and an index model described in section 4. Our specification is different from theirs and so are our conclusions. [CL09](#) and [CHV11](#) find that returns to a college education can be negative whereas we find them to be consistently positive and substantial. Further, the returns curves that we estimate do not satisfy the additive separability of errors assumption imposed by [CL09](#) and [CHV11](#), albeit that the specification in [CL09](#) is not nested by ours. Nevertheless, the shape restrictions in the existing literature appear to be restrictive and should be studied in greater depth with a larger data set.

2. IDENTIFICATION

[JPX11](#) show that the identified bounds provided in [CH05](#) can be substantially tightened under a weaker rank condition when instruments are independent of the errors in a two equation triangular system. In an extreme case with independent and continuous instruments, the structural function evaluated at particular values of its arguments can even be point-identified. In this section we show that this general result is in fact closely related to existing results on the identification of treatment effects (e.g., [CL09](#) and [Heckman and Vytlacil \(1999, 2001\)](#)) and also to the results for continuous endogenous regressors of [CH03](#).

The approach taken in [JPX11](#) and [CH05](#) is general in that partial identification is discussed under the setup of a triangular system with discrete endogenous regressors. The source of the weaker rank condition and the tighter bounds of [JPX11](#) is the independence between instrumental variables and unobserved errors, which makes it possible to combine multiple values of the instrumental variables to obtain identified bounds. The idea is best explained when an endogenous regressor is binary, which is the case we focus on in the current paper.

Consider the model

$$\begin{cases} \mathbf{y} &= g(\mathbf{x}, \mathbf{u}), \\ \mathbf{x} &= \phi(\mathbf{z}, \mathbf{v}), \end{cases} \quad (1)$$

where \mathbf{x} is a binary regressor, $\mathbf{z} \in \mathcal{Z} \subset \mathbb{R}^{d_z}$ is a vector of observed ‘demographics,’ and \mathbf{y} is a scalar-valued outcome.¹ We omit covariates in the identification analysis but they will be introduced in the estimation section. We permit the errors to enter nonadditively, so $g(1, \mathbf{u}) - g(0, \mathbf{u})$ can vary with \mathbf{u} . Consequently there is a distinction between the difference between (conditional) quantiles of $g(1, \mathbf{u})$ and $g(0, \mathbf{u})$ and the (conditional) quantiles of the difference $g(1, \mathbf{u}) - g(0, \mathbf{u})$. We follow [Doksum \(1974\)](#) and many others and focus on the former.

Thus, for generic random variables \mathbf{a} and \mathbf{b} let $\mathbb{Q}_{\mathbf{a}|\mathbf{b}}(\tau|b)$ be the τ quantile of \mathbf{a} given $\mathbf{b} = b$. The parameter of interest is

$$\psi^* = \psi^*(x^*, \tau^*|v^*) = g\{x^*, \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*)\} = \mathbb{Q}_{g(x^*, \mathbf{u})|\mathbf{v}}(\tau^*|v^*), \quad (2)$$

¹Our analysis can be adapted to the potential outcome framework, where the two potential outcomes may have different errors. There is no loss of generality in assuming a scalar-valued error in a potential outcome variable by the Cantor–Schroeder–Bernstein theorem. But it becomes restrictive once monotonicity conditions are imposed, as we will do later. For a detailed discussion of issues in models with vector-valued errors, see [Kasy \(2011\)](#).

for given values of x^* , τ^* , v^* , where the last equality follows from assumptions **D** and **E** below.² The function ψ^* can be used to define causal parameters of interest. For instance, we will call

$$\text{QMTE}(\tau^*|v^*) = \psi^*(1, \tau^*|v^*) - \psi^*(0, \tau^*|v^*) \quad (3)$$

the *quantile marginal treatment effect*, which is the quantile version of the marginal treatment effect of e.g. Heckman and Vytlacil (2005). Integrating QMTE over τ^* yields the marginal treatment effect (MTE). Integrating the MTE over v^* with various weight functions is discussed in Heckman and Vytlacil (2005), and for one such choice results in the average treatment effect.

We make the following model assumptions, which are based on those in JPX11 and CH05.

Assumption A. u, v have (marginal) $U(0, 1]$ -distributions.

Assumption B. u, v are independent of z .

Assumption C. $\phi(z, v)$ is left-continuous and nondecreasing in v for all values of z .

Assumption D. $g(x^*, u)$ is nondecreasing in u on $(0, 1]$.

Assumption E. $g(x^*, u)$ is left-continuous in u at $\mathbb{Q}_{u|v}(\tau^*|v^*)$.

Assumption F. For all $\tau \in (0, 1]$, $\mathbb{Q}_{u|v}(\tau|v)$ is nondecreasing in v .

Assumption **A** is fairly standard in the literature and is essentially a normalization. Assumption **B** is strong, but indispensable here. The conditions on ϕ in assumption **C** are also common. Assumptions **A** to **C** imply that one can represent the relationship between x, z, v as $x = \mathbb{1}\{v > p(z)\}$, where $\mathbb{1}$ is the indicator function and $p(z) = \mathbb{P}(x = 0|z = z)$ is one minus the propensity score; see e.g. Vytlacil (2006). Assumptions **D** and **E** are needed for the last equality in (2), as noted before. Note that assumption **D** is weaker than strict monotonicity of g in u on $(0, 1]$. In particular, if for instance y represents earnings then assumptions **D** and **E** allow for the case where there is a mass point at the minimum wage and the minimum wage is below the desired quantile. Note also that $g(x^*, \cdot)$ is allowed to have a discontinuous jump at $\mathbb{Q}_{u|v}(\tau^*|v^*)$.

The positive dependence condition in assumption **F** is used in both JPX11 and CH05.³ We use assumption **F** to obtain identifiable bounds for ψ^* , but it is not needed to establish point identification of ψ^* if there are continuous instruments and $g(x^*, u)$ is continuous at $\mathbb{Q}_{u|v}(\tau^*|v^*)$.

²We define quantiles in the standard way, i.e. $\mathbb{Q}_{a|b}(\tau|b) = \inf\{a : \mathbb{P}(a \leq a|b = b) \geq \tau\}$.

³Negative dependence can be dealt with similarly. The essence of this assumption is the monotonicity of $\mathbb{Q}_{u|v}(\tau|\cdot)$.

Let V_L, V_U be arbitrary subsets (of positive measure) of $(0, v^*]$ and $(v^*, 1]$, respectively. Then assumptions **A**, **D** and **F** imply that

$$g\{x^*, \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V_L)\} \leq \psi^* \leq g\{x^*, \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V_U)\}, \quad (4)$$

which can be seen by inverting the inequality e.g.

$$\begin{aligned} \mathbb{P}\{g(x^*, \mathbf{u}) \leq y | \mathbf{v} \in V_U\} &= \frac{1}{\mathbb{P}(\mathbf{v} \in V_U)} \int_{V_U} \mathbb{P}\{g(x^*, \mathbf{u}) \leq y | \mathbf{v} = v\} dv \\ &\leq \frac{1}{\mathbb{P}(\mathbf{v} \in V_U)} \int_{V_U} \mathbb{P}\{g(x^*, \mathbf{u}) \leq y | \mathbf{v} = v^*\} dv = \mathbb{P}\{g(x^*, \mathbf{u}) \leq y | \mathbf{v} = v^*\}. \end{aligned}$$

The bounds in (4) are discussed in detail in **CH05** and **JPX11**. Since it is important for our discussion to understand when these bounds are identified, we briefly discuss **CH05** and **JPX11** focusing on the case $x^* = 0$.

Let $\mathcal{V}(0) = \{(0, p(z)] : z \in \mathfrak{Z}_z\}$ and $\mathcal{V}(1) = \{(p(z), 1] : z \in \mathfrak{Z}_z\}$. If $V_L \in \mathcal{V}(0)$ then assumptions **B**, **D** and **E** imply that

$$g\{0, \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V_L)\} = \mathbb{Q}_{g(0, \mathbf{u})|\mathbf{v}}(\tau^*|V_L) = \mathbb{Q}_{y|\mathbf{x}, z}(\tau^*|0, z),$$

which is identified, where $z \in \mathfrak{Z}_z$ is such that $V_L = (0, p(z)]$. However, exclusively relying on sets of the form $(0, p(z)]$ leads to a trivial upper bound of the identified set because there is no set of that form that lies in its entirety above v^* . Similarly, relying on $\mathcal{V}(1)$ leads to a trivial lower bound in the case of $x^* = 1$. **CH05** stops here and interprets this problem as a violation of his rank condition.

JPX11 go on to show that the bounds in (4) are identified when sets not belonging to $\mathcal{V}(0)$ are utilized. For instance, suppose that there exist z and \tilde{z} in \mathfrak{Z}_z such that $v^* \leq p(\tilde{z}) < p(z)$. Then $(p(\tilde{z}), p(z)] = (0, p(z)] - (0, p(\tilde{z})] \geq v^*$, i.e. all elements in the majorant side set are no less than v^* . Hence one can choose $V_U = (p(\tilde{z}), p(z)]$ in (4) to obtain an upper bound, namely the τ^* quantile of the conditional distribution given by

$$\begin{aligned} \mathbb{P}\{g(0, \mathbf{u}) \leq y | \mathbf{v} \in V_U\} \\ = \frac{1}{p(z) - p(\tilde{z})} \{ \mathbb{P}(y \leq y | \mathbf{x} = 0, z = z) p(z) - \mathbb{P}(y \leq y | \mathbf{x} = 0, z = \tilde{z}) p(\tilde{z}) \}. \quad (5) \end{aligned}$$

A Dynkin system $\mathcal{D}(x^*)$ generated by $\mathcal{V}(x^*)$ can be obtained by applying various set operations to $\mathcal{V}(x^*)$ and ensures that $g\{x^*, \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V)\}$ is identified whenever $V \in \mathcal{D}(x^*)$. Such a Dynkin

system can be used to identify the tightest bounds in (4). The following definition was first introduced in JPX11.

Definition 1 (Dynkin System, JPX11). A Dynkin system $\mathcal{D}(x^*)$ is defined by the collection \mathcal{D}_∞ in the following iterative scheme. Let $\mathcal{D}_0 = \mathcal{V}(x^*)$. Then for all $t \geq 0$, \mathcal{D}_{t+1} consists of all sets A^* such that at least one of the following three conditions is satisfied.

- (i) $A^* \in \mathcal{D}_t$,
- (ii) $\exists A_1, A_2 \in \mathcal{D}_t : A_1 \subset A_2, \mu(A_2 - A_1) > 0, A^* = A_2 - A_1$,
- (iii) $\exists A_1, A_2 \in \mathcal{D}_t : A_1 \cap A_2 = \emptyset, \mu(A_1 \cup A_2) > 0, A^* = A_1 \cup A_2$.

Since $\{\mathcal{D}_t : t = 0, 1, \dots\}$ is an increasing sequence of collections of sets, we have $\mathcal{D}(x^*) = \bigcup_{t=0}^{\infty} \mathcal{D}_t$. It can be shown that the conditional distribution function of $g(x^*, \mathbf{u})$ given $\mathbf{v} \in V$ is identified whenever $V \in \mathcal{D}(x^*)$.

Let $\mathcal{D}_L(x^*, v^*) = \{V \in \mathcal{D}(x^*) : V \leq v^*\}$ and let $\mathcal{D}_U(x^*, v^*)$ be similarly defined. Then JPX11 have shown⁴ that under assumptions A to F identified bounds for ψ^* are given by

$$\sup_{V \in \mathcal{D}_L(x^*, v^*)} g\{x^*, \mathbb{Q}_{\mathbf{u}|v}(\tau^*|V)\} \leq \psi^* \leq \inf_{V \in \mathcal{D}_U(x^*, v^*)} g\{x^*, \mathbb{Q}_{\mathbf{u}|v}(\tau^*|V)\}. \quad (6)$$

KM14 have shown that there are circumstances under which the above bounds can be tightened further.

We now discuss how additional continuity conditions can be used to obtain the point identification of ψ^* .

Assumption G. $\mathbb{Q}_{\mathbf{u}|v}(\tau^*|v)$ is left-continuous in v at v^* .

Assumption H. There exists a sequence $\{z_t\}$ in \mathcal{S}_z such that $p(z_t)$ is (strictly) increasing in t with supremum v^* .

Left-continuity in assumptions E and G can be replaced with right-continuity. Please note that left-continuity does not rule out the presence of discrete jumps in the function g and it hence allows for mass points in the distribution of \mathbf{y} . A sufficient condition for assumption H is that at least one element of the vector of instruments \mathbf{z} has continuous variation (given the other elements) and p is continuous in that element of \mathbf{z} .

We can now strengthen the result in (6).

⁴The conditions in JPX11 are slightly different from the ones here.

Theorem 1. (i) Suppose that assumptions *E* to *H* are satisfied. Then

$$\sup_{V_L \in \mathcal{D}_L(x^*, v^*)} g\{x^*, \mathbb{Q}_{\mathbf{u}|v}(\tau^* | V_L)\} = \psi^* \leq \inf_{V_U \in \mathcal{D}_U(x^*, v^*)} g\{x^*, \mathbb{Q}_{\mathbf{u}|v}(\tau^* | V_U)\}. \quad (7)$$

(ii) If moreover assumptions *A* to *D* are satisfied then ψ^* is point-identified.

Theorem 1 implies that the Dynkin system approach of JPX11 can be used to achieve point identification of ψ^* using continuous variation in the propensity score $p(\mathbf{z})$. This result can be compared to CH03 where point identification is achieved assuming strong monotonicity of $\phi(\mathbf{z}, v)$ in v , which implies that \mathbf{x} is continuously distributed. If left-continuity (assumptions *E* and *G*) is strengthened to continuity then the inequality in (7) becomes an equality such that the intersection bounds in (6) collapse to ψ^* .

Continuity of $g(x^*, u)$ at $u = \mathbb{Q}_{\mathbf{u}|v}(\tau^* | v^*)$ is more helpful for the identification of ψ^* than continuity of $\mathbb{Q}_{\mathbf{u}|v}(\tau^* | v)$ at v^* . Indeed, continuity of $\mathbb{Q}_{\mathbf{u}|v}(\tau^* | v)$ is of at most modest help to relax assumption *H* whereas continuity of $g(x^*, u)$ obviates the need for monotonicity assumptions on $\mathbb{Q}_{\mathbf{u}|v}(\tau | v)$ in v (assumption *F*) for the identification of ψ^* , as theorem 2 demonstrates. However, as we have pointed out before, assuming continuity rather than left-continuity of g will impose restrictions on the support of the outcome, which would be undesirable.

Assumption I. $g(x^*, u)$ is continuous at $\mathbb{Q}_{\mathbf{u}|v}(\tau^* | v^*)$.

Theorem 2. Suppose that assumptions *A* to *D* and *G* to *I* are satisfied. Then ψ^* is identified.

Please note that theorems 1 and 2 rely on monotonicity/continuity but not on smoothness (i.e. differentiability). We now discuss how these results are connected with the existing results of *identification via local instruments* in the treatment effects literature.

It is well-known in the treatment effects literature that differentiability can result in point identification of the distribution functions of counterfactual outcomes conditional on $p(\mathbf{z}) = v^*$.⁵ We now explain how this treatment effects literature result is related to the results in theorems 1 and 2. Let $\partial_z p(z^*)$ be the partial derivative of p with respect to z at z^* and $G^*(y | x^*, z) = \mathbb{P}(y \leq y | \mathbf{x} = x^*, z = z)$.

Assumption J. For any $y \in \mathbb{R}$, $G^*(y | x^*, z)$ is continuously differentiable in z at z^* .

Assumption K. For some z^* in the interior of \mathcal{Z} , (i) $p(z^*) = v^*$ and (ii) $\partial_z p(z^*) \neq 0$.⁶

⁵See e.g. CL09 and Heckman and Vytlačil (1999, 2001).

⁶Please note that $\partial_z p$ is vector-valued and we only require one of its elements to be nonzero.

It is useful to compare assumptions **E**, **G** and **H** with assumptions **J** and **K**. If the propensity score is differentiable then it follows from (8) that assumption **J** is equivalent to continuity of $F_{\mathbf{u}|\mathbf{v}}(u|v)$ in v at v^* and indeed to the differentiability of $F_{\mathbf{u}\mathbf{v}}(u, v)$ in v at v^* .

We now show that the smoothness conditions in assumptions **J** and **K** provide an alternative path to identification. Suppose that \mathbf{z} is scalar-valued. Since $\mathbf{x} = 0$ and $\mathbf{z} = z$ is equivalent to $\mathbf{v} \in (0, p(z)]$ and $\mathbf{z} = z$, we have that for any $y \in \mathbb{R}$,

$$G^*(y|x^*, z) = \mathbb{P}\{g(0, \mathbf{u}) \leq y | \mathbf{v} \in (0, p(z)]\} = \frac{1}{p(z)} \int_0^{p(z)} \mathbb{P}\{g(0, \mathbf{u}) \leq y | \mathbf{v} = v\} dv. \quad (8)$$

Differentiating both sides in (8) and evaluating at z^* yields

$$\mathbb{P}\{g(0, \mathbf{u}) \leq y | \mathbf{v} = v^*\} = G^*(y|x^*, z^*) + v^* \frac{\partial_z G^*(y|x^*, z^*)}{\partial_z p(z^*)}. \quad (9)$$

The right hand side in (9) is identified and ψ^* is defined as the smallest value of y for which the left hand side in (9) is equal to τ^* . An expression similar to (9) can be found in [CL09](#).

For vector-valued \mathbf{z} it is more natural to work with the propensity score.⁷ Thus, let $G(y|x, p) = \mathbb{P}\{\mathbf{y} \leq y | \mathbf{x} = x, p(\mathbf{z}) = p\}$ so that assumption **B** implies $G^*(y|x, z) = G\{y|x, p(z)\}$. Then, as was shown by [CL09](#), we have

$$\begin{aligned} \mathbb{P}\{g(0, \mathbf{u}) \leq y | \mathbf{v} = v^*\} &= G(y|0, v^*) + v^* \partial_p G(y|0, v^*), \\ \mathbb{P}\{g(1, \mathbf{u}) \leq y | \mathbf{v} = v^*\} &= G(y|1, v^*) - (1 - v^*) \partial_p G(y|1, v^*). \end{aligned} \quad (10)$$

Theorem 3. *If assumptions **A** to **D**, **J** and **K** are satisfied then ψ^* is identified.*

Theorems 1 to 3 articulate a trade-off between monotonicity, continuity, and smoothness assumptions. Continuity of $F_{\mathbf{u}|\mathbf{v}}(u|v)$ in v and differentiability of the propensity score are convenient for estimation but neither condition is necessary for identification.

Finally, we note that the Dynkin system idea in theorems 1 and 2 has applications far beyond the simple binary endogenous variable model of this paper: see e.g. [JPX11](#); [Jun, Pinkse, and Xu \(2012\)](#); [Jun, Pinkse, Xu, and Yildiz \(2010\)](#).

3. ESTIMATION

3.1. Assumptions. We now proceed to describe and motivate our estimation procedure, for which we will focus on the case $x^* = 0$. We add a subscript i to $y, \mathbf{x}, \mathbf{z}, \mathbf{u}, \mathbf{v}$ and assume that we have an

⁷Assumption **J** is sufficient but not necessary for the differentiability of $G(y|x, \cdot)$.

i.i.d. sample of size n . We allow for the presence of exogenous covariates $\mathbf{a}_i \in \mathbb{R}^{d_a}$ in the function g , i.e. we now consider

$$\mathbf{y}_i = g(\mathbf{x}_i, \mathbf{a}_i, \mathbf{u}_i). \quad (11)$$

The covariates \mathbf{a}_i are contained in the vector of instruments \mathbf{z}_i , which contains one or more additional elements \mathbf{q}_i and is assumed to be independent of $\mathbf{u}_i, \mathbf{v}_i$, as is formally assumed here. However, in our proofs, we explicitly allow for the possibility that \mathbf{a}_i and \mathbf{z}_i have the same dimension, because in the semiparametric version of our estimator introduced in section 4 we consider $\mathbf{a}_i^\top \theta_0$ and $\mathbf{z}_i^\top \gamma_0$ (with \mathbf{a}_i a subvector of \mathbf{z}_i) in lieu of \mathbf{a}_i and \mathbf{z}_i , respectively.

Assumption L. *Assumptions A to D and K are for some q^* satisfied with $\mathbf{z}_i = [\mathbf{q}_i^\top, \mathbf{a}_i^\top]^\top$, $z^* = [q^{*\top}, a^{*\top}]^\top$, and with $g(x^*, a^*, u)$ in lieu of $g(x^*, u)$.*

Let $\mathcal{F}_j : \mathbb{R}^{d_z} \rightarrow \mathbb{R}$ denote the class of functions which are j times continuously differentiable on \mathcal{Z} and $j + 2$ times boundedly differentiable with respect to z_1 . We replace assumption J with the stronger assumption M.

Assumption M. $G^*(y|0, z) \in \mathcal{F}_2$.

The addition of the covariates \mathbf{a}_i does not complicate the identification argument much. Indeed, one can simply condition on $\mathbf{a}_i = a^*$ in which case the entire argument of section 2 can be repeated with \mathbf{q}_i assuming the role of \mathbf{z}_i . For estimation we adopt the identification argument of theorem 3, because it is the most convenient. Assumption M is introduced to obtain the desired convergence rate.

Thus, let $\mathbf{b}_i = [\mathbf{p}_i, \mathbf{a}_i^\top]^\top$, $\mathbf{z}_i = [\mathbf{q}_i^\top, \mathbf{a}_i^\top]^\top$, $\mathbf{p}_i = p(\mathbf{z}_i)$, and (re)define

$$G(y|x, a, p) = \mathbb{P}(\mathbf{y}_i \leq y | \mathbf{x}_i = x, \mathbf{a}_i = a, \mathbf{p}_i = p). \quad (12)$$

We start by estimating

$$\psi^* = \psi^*(x^*, \tau^* | a^*, v^*) = \mathbb{Q}_{g(x^*, a^*, \mathbf{u}_i) | \mathbf{v}_i}(\tau^* | v^*) = g\{x^*, a^*, \mathbb{Q}_{\mathbf{u}_i | \mathbf{v}_i}(\tau^* | v^*)\}. \quad (13)$$

We propose estimating ψ^* by inverting the functions $H(\cdot | a^*, v^*)$ defined by

$$H(y|a, v) = \mathbb{P}\{g(0, a, \mathbf{u}_i) \leq y | \mathbf{v}_i = v\},$$

which under assumption [L](#) satisfies

$$H(y|a, v) = G(y|0, a, v) + v\partial_p G(y|0, a, v).$$

So whereas the estimator in [CL09](#) is semiparametric and the object of interest is the mean,⁸ our approach is nonparametric and we estimate quantiles which entails an additional inversion step which requires some empirical process theory. However, in section [4](#) we discuss the possibility of using single index restrictions allowing for the possibility of semiparametric estimation, albeit in a more structural fashion than [CL09](#).

Let $G(y|a, v) = G(y|0, a, v)$ and by assumptions [B](#) to [D](#) for $\mathbf{w}_i = w(\mathbf{z}_i)$ for some function w to be introduced later,

$$G(y|a^*, v^*) = \frac{\mathbb{E}\{\mathbb{1}(y_i \leq y)\mathbf{w}_i | \mathbf{x}_i = 0, \mathbf{a}_i = a^*, \mathbf{p}_i = v^*\}}{\mathbb{E}(\mathbf{w}_i | \mathbf{x}_i = 0, \mathbf{a}_i = a^*, \mathbf{p}_i = v^*)}. \quad (14)$$

Since the function p is estimable, so is the function G , and thence H . We propose estimating both G and $\partial_v G$ by nonparametric kernel (derivative) regression estimation and inverting the resulting estimator of $H(\cdot|a^*, v^*)$ to estimate ψ^* .

It is well-known that kernel regression estimation has problems in the tails of the distribution, or more precisely wherever the density of conditioning variables is close to zero. In the estimation we hence only use observations i for which \mathbf{z}_i belongs to some convex and compact set \mathcal{Z} on which the density f of \mathbf{z}_i is bounded away from zero and which is further constrained below. Not using all data does have efficiency implications, but the commonly used alternative of sample-size dependent trimming is practically cumbersome, technically messy, and any meaningful gains of such a procedure in empirical work are phantasmic. In what follows we will assume \mathbf{z}_i to be continuously distributed even though in empirical work discrete covariates and instruments are prevalent. Kernel estimation with discrete regressors can be accommodated (see e.g. [Delgado and Mora, 1995](#)) at the expense of longer proofs. However, because in practice the index version of the estimator proposed in section [4](#) will often be more attractive and since for the index version only one of the elements of \mathbf{a}_i and one of the elements in \mathbf{z}_i that are not in \mathbf{a}_i must be continuously distributed, we do not weaken the assumption here.

⁸or a quantile under an additional additive separability assumption.

The function w in (14) is chosen to be nonnegative on \mathcal{Z} and zero elsewhere. Let $\mathbf{I}_{xi} = \mathbb{1}(x_i = 0)\mathbf{w}_i$, $\mathbf{I}_i(y) = \mathbf{I}_{xi}\mathbb{1}(y_i \leq y)$, let f_{ap} be the joint density of $\mathbf{a}_i, \mathbf{p}_i$, and let

$$\begin{aligned} S_{0x} &= S_{0x}(a^*, v^*) = \mathbb{E}(\mathbf{I}_{xi} | \mathbf{a}_i = a^*, \mathbf{p}_i = v^*) f_{ap}(a^*, v^*), & S_{1x} &= \partial_v S_{0x}, \\ S_0(y) &= S_0(y; a^*, v^*) = \mathbb{E}\{\mathbf{I}_i(y) | \mathbf{a}_i = a^*, \mathbf{p}_i = v^*\} f_{ap}(a^*, v^*), & S_1(y) &= \partial_v S_0(y). \end{aligned}$$

Then, noting that $G(y|a^*, v^*) = S_0(y)/S_{0x}$, it follows that

$$H(y|a^*, v^*) = \frac{S_0(y)S_{0x} + v^*S_1(y)S_{0x} - v^*S_0(y)S_{1x}}{S_{0x}^2}.$$

We now develop our estimator. Let k be a kernel, K be a product kernel based on k whose dimension is determined by its argument, and let h_0, h_1, h_z be bandwidths. Define $\mathbf{K}_{zi}(z) = K\{(z - z_i)/h_z\}/h_z^{d_z}$, and $\hat{\mathbf{p}}_i = \hat{\mathbf{p}}(z_i)$, where

$$\hat{\mathbf{p}}(z) = \frac{\sum_{i=1}^n \mathbf{K}_{zi}(z) \mathbb{1}(x_i = 0)}{\sum_{i=1}^n \mathbf{K}_{zi}(z)}. \quad (15)$$

Let further $\mathbf{K}_{aij} = K\{(a^* - \mathbf{a}_i)/h_j\}/h_j^{d_a}$, $\mathbf{k}_{ij}^{(s)} = k^{(s)}\{(v^* - \mathbf{p}_i)/h_j\}/h_j^{s+1}$, $\hat{\mathbf{k}}_{ij}^{(s)} = k^{(s)}\{(v^* - \hat{\mathbf{p}}_i)/h_j\}/h_j^{s+1}$, and

$$\hat{\mathbf{S}}_s(y; p) = \frac{1}{n} \sum_{i=1}^n \mathbf{k}_{is}^{(s)} \mathbf{K}_{ais} \mathbf{I}_i(y), \quad \hat{\mathbf{S}}_s(y; \hat{p}) = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{k}}_{is}^{(s)} \mathbf{K}_{ais} \mathbf{I}_i(y), \quad (16)$$

The proposed estimator is then given by

$$\hat{\mathbf{H}}(y|a^*, v^*) = \frac{\hat{\mathbf{S}}_0(y) \hat{\mathbf{S}}_{0x} + v^* \hat{\mathbf{S}}_1(y) \hat{\mathbf{S}}_{0x} - v^* \hat{\mathbf{S}}_0(y) \hat{\mathbf{S}}_{1x}}{\hat{\mathbf{S}}_{0x}^2}, \quad (17)$$

where $\hat{\mathbf{S}}_0(y) = \hat{\mathbf{S}}_0(y; \hat{p})$, $\hat{\mathbf{S}}_{0x} = \hat{\mathbf{S}}_0(\infty)$, $\hat{\mathbf{S}}_1(y) = \hat{\mathbf{S}}_1(y; \hat{p})$, and $\hat{\mathbf{S}}_{1x} = \hat{\mathbf{S}}_1(\infty)$.

The bandwidths h_0, h_1 , and h_z vary with n according to

$$h_0 \sim n^{-\eta_0}, h_1 \sim n^{-\eta_1}, \text{ and } h_z \sim n^{-\eta_z} \quad (18)$$

for some $\eta_0, \eta_1, \eta_z > 0$ to be constrained in assumption S.

So there are a total of five different input parameters here: a kernel, the w -function, and three bandwidths. The number of bandwidths can be reduced to two by choosing $h_0 = h_z$, but our conditions require that h_0 and h_z converge to zero faster than h_1 .

We make the following assumptions.

Assumption N. $G(\cdot|a^*, v^*)$ and $\partial_p G(\cdot|a^*, v^*)$ are differentiable in y , and hence so is $H(\cdot|a^*, v^*)$.

Assumption N presumes a continuous outcome, which is not necessary for identification but is convenient for estimation. It is sufficient for the quantile of interest to be uniquely defined and is needed for the empirical process results that are used.

Assumption O. $\mathcal{Z} = \mathcal{Z}_1 \times \tilde{\mathcal{Z}}$ is a subset of the interior of the support \mathcal{S}_z of $\mathbf{z}_i = [z_{i1}, \tilde{\mathbf{z}}_i^\top]^\top$, for which $\tilde{\mathcal{Z}} \subset \mathbb{R}^{d_z-1}$, $\mathcal{Z}_1 = [\underline{z}_1, \bar{z}_1]$ for some $\underline{z}_1 < \bar{z}_1$ are compact and convex. On \mathcal{Z} the density $f \in \mathcal{F}_2$ of \mathbf{z}_i is bounded away from zero. Finally, \mathcal{Z} contains points of the form (q, a^*) and for any such points and any vector $\xi \in \mathbb{R}^{d_a}$ there exists an $\epsilon > 0$ such that $(q, a^* - \epsilon\xi)$ and $(q, a^* + \epsilon\xi)$ are also in \mathcal{Z} .

Assumption P. $p \in \mathcal{F}_2$ is (strictly) increasing in its first argument and $0 < \mathbb{P}\{p(\bar{z}_1, \tilde{\mathbf{z}}_i) \leq v^*\} < \mathbb{P}\{p(\underline{z}_1, \tilde{\mathbf{z}}_i) \leq v^*\} < 1$.

Assumptions M, O and P are typical for the kernel derivative estimation literature, albeit that we require the existence of one extra derivative in the first argument. There is nothing special about the first argument (other than that it is an element of \mathbf{q}_i rather than \mathbf{a}_i); one of the instruments used must satisfy this condition, but there is no need to know, or indeed specify, which one. The number of required derivatives with respect to z_1 can possibly be reduced by one at the expense of much more restrictive conditions on the bandwidth sequences (assumption S) and permitted dimensions d_a, d_z .

The remaining assumptions (assumptions Q to S below) pertain to the choice of input parameters and are hence of lesser importance as long as input parameters that satisfy the properties exist.

Assumption Q. $w \in \mathcal{F}_2$ is positive on the interior of \mathcal{Z} , zero everywhere else, and nowhere greater than one.

Assumption Q is there both to ensure that only observations i with $\mathbf{z}_i \in \mathcal{Z}$ are used (the need for which was explained earlier) and to allow us to use standard kernel bias expansions by removing discontinuities on the boundaries of \mathcal{Z} .

We now state our conditions for the kernel and bandwidth choices.

Assumption R. The kernel k is even, everywhere nonnegative, infinitely many times boundedly differentiable, and integrates to one. It further satisfies $\kappa_{s2} = \int \{k^{(s)}(t)\}^2 dt < \infty$ for $s = 0, 1$, and $\kappa_2 = \int k(t)t^2 dt < \infty$.

Conditions on the kernel similar to those in assumption **R** are standard in the kernel estimation literature. Since we get to choose k , assumption **R** is innocuous. It is possible to require a smaller number of derivatives at the expense of longer proofs and possibly stronger restrictions on the bandwidths than those found in assumption **S**.

Assumption S. *The constants η_0, η_1, η_z defined in (18) are such that for $\eta^* = \max(2d_z\eta_z - 1, 0)$, $2\eta_0 < \min(4\eta_z, 1 - \eta_z d_z)$, and*

$$\max\{\eta^*, 1 - 4\max(\eta_z, \eta_0), (d_a + 1)\eta_0, 2(d_a + 2)\eta_0 + \eta^* - 1\} < (d_a + 3)\eta_1 < \min\{(d_a + 3)\min(\eta_0, \eta_z), 1 - \eta^*\}.$$

The choice of bandwidths in assumption **S** results in the convergence rate

$$\rho_n = n^{\{1-(3+d_a)\eta_1\}/2}. \quad (19)$$

While assumption **S** allows for undersmoothing, the choice of $\eta_1 = 1/(7 + d_a)$ leads to the optimal rate of $\rho_n = n^{2/(7+d_a)}$ for kernel derivative estimators (using second order kernels). Faster convergence rates are feasible under additional smoothness conditions (more derivatives) using bias reduction techniques such as higher order kernels or local polynomial estimation. Such an extension is a well-trodden path, which adds no new theoretical insights, and its promised performance improvements are not often realized in samples of finite size.

To see that η_0, η_1, η_z exist for many (but certainly not all) combinations of d_a, d_z , we present table 1 which for $\eta_1 = 1/(7 + d_a)$ lists the values of 1,000 times the values of η_0, η_z which are the ‘points of gravity’ of the regions of η_0, η_z combinations for which assumption **S** is satisfied and which are in some sense hence farthest from violating assumption **S**. If there is no entry in the table for a particular $d_a \leq d_z - 1$ combination then that means that for $\eta_1 = 1/(7 + d_a)$ there are no values of η_0, η_z to satisfy assumption **S**. Of course, η_0, η_1, η_z only indicate a rate; the constant multiplying $n^{-\eta_0}$ for instance still needs to be chosen.

3.2. Limit results for our estimator of H . Before stating our formal results, we introduce some notation. Let $\lambda_b = \lim_{n \rightarrow \infty} (\rho_n h_1^2)$, $\lambda_v = \lim_{n \rightarrow \infty} (\rho_n^2 / n h_1^{3+d_a})$, $p_y(y|z) = \mathbb{P}(y_i \leq y | \mathbf{x}_i = 0, \mathbf{z}_i = z)$,

$$\Gamma(y, a^*, v^*) = v^* f_{ap}(a^*, v^*) \kappa_{12} \mathbb{E}\{p_y(y|\mathbf{z}_i) \mathbf{w}_i^2 | \mathbf{a}_i = a^*, \mathbf{p}_i = v^*\},$$

$d_a \downarrow$	1	2	3	4	d_z 5	6	7	8	9
0	226	209	190	166					
	142	142	142	142					
	345	216	174	154					
1	181	180	176	166	145				
	125	125	125	125	125				
	380	221	169	145	133				
2		145	144	142	136	128			
		111	111	111	111	111			
		223	168	140	125	115			
3			122	121	120	115			
			100	100	100	100			
			160	132	116	107			
4				107	106	104	100		
				90	90	90	90		
				126	110	100	93		
5					96	95	92		
					83	83	83		
					105	95	88		
6						86	85	83	
						76	76	76	
						90	83	78	
7							79	77	
							71	71	
							80	75	
8								72	70
								66	66
								72	68

TABLE 1. Suggested choices for 1,000 times η_0, η_1, η_z for various combinations of d_a, d_z .

and

$$\begin{aligned} \mathcal{C}(y, y^*) = \mathcal{C}(y, y^*; a^*, v^*) &= \Gamma\{\min(y, y^*), a^*, v^*\} - \Gamma(y, a^*, v^*)G(y^*|a^*, v^*) \\ &\quad - \Gamma(y^*, a^*, v^*)G(y|a^*, v^*) + G(y|a^*, v^*)G(y^*|a^*, v^*)\Gamma(\infty, a^*, v^*). \end{aligned} \quad (20)$$

Let further

$$\mathcal{C}(y, y^*) = \mathcal{C}(y, y^*; a^*, v^*) = \lambda_v \frac{v^{*2} \mathcal{C}(y, y^*; a^*, v^*)}{S_{0x}^2(a^*, v^*)},$$

and

$$\mathcal{B}(y) = \mathcal{B}(y; a^*, v^*) = \frac{\lambda_b v^* \kappa_2}{2S_{0x}} \text{tr} \left\{ \partial_v \partial_{bb^\top} S_0(y; a^*, v^*) - G(y|a^*, v^*) \partial_v \partial_{bb^\top} S_{0x}(a^*, v^*) \right\}.$$

Theorem 4. *Under assumptions L to S,*

$$\rho_n \{ \hat{\mathbf{H}}(\cdot|a^*, v^*) - H(\cdot|a^*, v^*) \} \xrightarrow{w} \mathbf{G},$$

on the space of bounded functions on $\mathcal{Y} = \{y : \exists u \in \mathcal{U} : g(0, a^*, u) = y\}$, $\mathcal{L}^\infty(\mathcal{Y})$, where \mathbf{G} is a Gaussian process with mean \mathcal{B} and covariance kernel \mathcal{C} .

Please note that table 1 implies that it is possible for the limit distribution not to be affected by the first step estimation of p — the ‘oracle property’ — even in some cases in which $d_z > d_a + 1$. This may appear to be at odds with other results in the voluminous literature on nonparametric generated regressors (Rilstone, 1996; Pinkse, 2001; Mammen, Rothe, and Schienle, 2012, inter multa alia) in which nonparametrically estimated regressors do affect the optimal convergence rate unless the estimated regressors are functions whose vector of arguments is of smaller dimension than the vector of arguments of the function of interest. However, here we are not evaluating \hat{p} at a fixed point, say z^* , to obtain our estimate $\hat{\mathbf{H}}(y|a^*, v^*)$. Instead we *only* use \hat{p}_i ’s which are averaged in some sense which reduces their contribution to the variance, thereby allowing us to use smaller h_z values to reduce the bias, also.

3.3. Limit results for our estimator of ψ^* . We finally turn to our estimator of ψ^* itself. We use the standard definition of quantile using the estimated conditional distribution function, i.e.

$$\hat{\psi}^* = \inf \{ \tilde{\psi} : \hat{\mathbf{H}}(\tilde{\psi}|a^*, v^*) \geq \tau^* \}.$$

The asymptotic behavior can then be inferred from theorem 4. Indeed, we have theorem 5.

Theorem 5. *Under assumptions L to S, $\rho_n(\hat{\psi}^* - \psi^*) \xrightarrow{d} N(\mathcal{B}_\psi, \mathcal{V}_\psi)$, where*

$$\mathcal{B}_\psi = -\frac{\mathcal{B}(\psi^*)}{H'(\psi^*|a^*, v^*)}, \quad \mathcal{V}_\psi = \frac{\mathcal{C}(\psi^*, \psi^*)}{\{H'(\psi^*|a^*, v^*)\}^2}.$$

3.4. Bias and variance estimation. The bias and variance in theorem 5 can be consistently estimated by standard methods. Since the bias can be removed by undersmoothing, the Jackknife, or

other methods, we focus on estimation of the variance below. Note that

$$H'(y|a^*, v^*) = \frac{S_0^{(1)}(y)S_{0x} + v^*S_1^{(1)}(y)S_{0x} - v^*S_0^{(1)}(y)S_{1x}}{S_{0x}^2},$$

where letting $f_y(\cdot|x, z)$ be the conditional density of y_i given $\mathbf{x}_i = x$ and $\mathbf{z}_i = z$,

$$\begin{cases} S_0^{(1)}(y) = \partial_y S_0(y) = \mathbb{E}\{\mathbf{I}_{xi} f_y(y|\mathbf{x}_i, \mathbf{z}_i) | \mathbf{p}_i = v^*\} f_{ap}(a^*, v^*), \\ S_1^{(1)}(y) = \partial_y S_1(y) = \partial_{yv} S_0(y). \end{cases}$$

For $s = 0, 1$, we can estimate $S_s^{(1)}(y)$ by

$$\hat{S}_s^{(1)}(y) = \frac{1}{n} \sum_{i=1}^n \mathbf{K}_{ais} \hat{\mathbf{k}}_{is}^{(s)} \mathbf{k}_{yi}(y) \mathbf{I}_{xi},$$

where $\mathbf{k}_{yi}(y) = k\{(y - y_i)/h_y\}/h_y$ with h_y a bandwidth. Also, $\mathcal{C}(y, y; a^*, v^*)$ can be estimated by

$$\hat{\mathcal{C}}(y, y; a^*, v^*) = \frac{h_1^{1+d_a}}{n} \sum_{i=1}^n \mathbf{K}_{ai1}^2 \hat{\mathbf{k}}_{i1}^2 \mathbb{1}(x_i = 0) \mathbf{w}_i^2 \left\{ \mathbb{1}(y_i \leq y) - \frac{\hat{S}_0(y)}{\hat{S}_{0x}} \right\}^2.$$

The final estimator of \mathcal{V}_ψ can be obtained by using $\hat{S}_s^{(1)}$ and $\hat{\mathcal{C}}$ evaluated at $y = y^* = \hat{\psi}^*$. The following theorem establishes the consistency of $\hat{S}_s^{(1)}(\hat{\psi}^*)$ and $\hat{\mathcal{C}}(\hat{\psi}^*, \hat{\psi}^*)$.

Theorem 6. *Suppose that assumptions L and N to S are satisfied with $h_y = o(1)$, $1 = o(\rho_n h_y)$, and $\sup_s |k(s)| < \infty$. Then, for $s = 0, 1$,*

$$\hat{S}_s^{(1)}(\hat{\psi}^*) \xrightarrow{P} S_s^{(1)}(\psi^*) \quad \text{and} \quad \hat{\mathcal{C}}(\hat{\psi}^*, \hat{\psi}^*) \xrightarrow{P} \mathcal{C}(\psi^*, \psi^*, v^*).$$

4. INDEX

In most applications the dimensions of the $\mathbf{a}_i, \mathbf{z}_i$ vectors are too large for estimates to be sufficiently precise. One solution to this problem is to impose semiparametric restrictions on the g and p functions or, said differently, to assume that $\mathbf{a}_i, \mathbf{z}_i$ enter as indices. As a leading example, we consider⁹

$$\begin{cases} y_i = g\{\mathbf{x}_i, \mathbf{a}_i^\top \theta_0, \mathbf{u}_i\}, \\ \mathbf{x}_i = \mathbb{1}\{\mathbf{v}_i > p(\mathbf{z}_i^\top \gamma_0)\}, \end{cases} \quad (21)$$

albeit that in our application we allow the value of θ_0 in $g(1, \cdot, \cdot)$ to be different from that in $g(0, \cdot, \cdot)$.

⁹Other parametric link functions or multiple indices can be accommodated but they complicate the identification conditions.

It follows from the copious work on index models that several normalizations are needed. First, $\mathbf{a}_i, \mathbf{z}_i$ should not include a constant term and even then the vectors θ_0, γ_0 are (at best) identified up to scale. Second, one should be able to move \mathbf{x}_i exogenously without changing \mathbf{a}_i , i.e. at least one of the γ -coefficients on the \mathbf{q}_i component of \mathbf{z}_i should be nonzero. Indeed, if one lets $\mathbf{z}_i^\top \gamma_0 = \mathbf{a}_i^\top \gamma_{0a} + \mathbf{q}_i^\top \gamma_{0q}$ then the conditions of sections 2 and 3 can be verified conditional on $\mathbf{a}_i^\top \theta_0 = \mathbf{a}^{*\top} \theta_0$ and taking \mathbf{z} in sections 2 and 3 to equal $\mathbf{q}_i^\top \gamma_{0q}$; doing this requires that $\gamma_{0q} \neq 0$.

From now on, we take identification of ψ^* and that of γ_0, θ_0 as given. We also take as given that \sqrt{n} -consistent estimators $\hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\theta}}$ of γ_0, θ_0 exist. We are not generally fans of high-level assumptions. However, the structure of (21) fits well into the index model estimation literature of which Powell, Stock, and Stoker (1989); Ichimura (1993); Klein and Spady (1993) are prominent examples. Indeed, $\mathbb{P}(x_i = 1 | z_i = z) = p(z^\top \gamma_0)$, which yields an estimate of γ_0 . Further, $\mathbb{E}(y | x = x, z = z)$ is an unknown function of $x, z^\top \gamma_0, \mathbf{a}^\top \theta_0$, which can be used to construct an estimate of θ_0 .

The main task for this section, then, is to establish that the estimation of the nuisance parameters γ_0, θ_0 does not affect the limit distribution of the estimator of ψ^* .

Let $\hat{\boldsymbol{\psi}}^*$ be defined as $\hat{\boldsymbol{\psi}}^*$, replacing \mathbf{a}_i with $\mathbf{a}_i^\top \hat{\boldsymbol{\gamma}}$ and \mathbf{z}_i with $\mathbf{z}_i^\top \hat{\boldsymbol{\theta}}$.

Theorem 7. *Suppose that $\hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\theta}}$ are \sqrt{n} -consistent estimates of γ_0, θ_0 , respectively. Then theorems 5 and 6 hold for $\hat{\boldsymbol{\psi}}^*$ when $\mathbf{a}_i, \mathbf{z}_i$ are replaced with $\mathbf{a}_i^\top \hat{\boldsymbol{\gamma}}$ and $\mathbf{z}_i^\top \hat{\boldsymbol{\theta}}$, respectively.*

5. APPLICATION

5.1. Data. We now apply our method to estimate the returns to college education using the NLSY 1979 data, which were used by CL09 and CHV11 inter alia. Indeed, we use the code provided by CHV11 to obtain exactly the same data set and therefore we compare our results with CHV11: the results of CL09 are not qualitatively different from those in CHV11. The sample consists of 1,747 white males. The data set contains a number of demographic and education-related variables as can be gleaned from our results tables, tables 2 and 3. Even though the sample is fairly small relative to the nonparametric specification of the model that we are estimating, we are finding several interesting results.

A detailed description of the data can be found in CHV11 and its supplementary material,¹⁰ which also contains summary statistics separated by education group (college versus noncollege).

¹⁰See https://www.aeaweb.org/aer/data/oct2011/20061111_app.pdf

5.2. **Methodology.** As noted in the introduction, our approach differs from that of CL09 and CHV11 in several respects. The most important difference is that in both those papers the error in the earnings equation enters additively whereas in our model it enters nonparametrically, i.e. nonadditively. Hence in CL09 and CHV11 there can be no interaction between the earnings equation error and regressors such that potential wage curves (as a function of v^*) for different covariate values are constrained to be vertical shifts of each other, whereas ours can vary freely. This shape restriction is also used in CHV11, which takes a more parametric (and therefore more restrictive) approach than CL09.

Like CL09 and CHV11 we distinguish between two groups: those with a college education (group 1) and those without (group 0). For the schooling equation we use the second half of (21). For the wage equation we allow the θ -coefficients in (21) to be different across education groups, i.e. we consider

$$\begin{cases} y_i = g_1(\mathbf{a}_i^\top \theta_{10}, \mathbf{u}_{i1}) \mathbb{1}\{\mathbf{x}_i = 1\} + g_0(\mathbf{a}_i^\top \theta_{00}, \mathbf{u}_{i0}) \mathbb{1}\{\mathbf{x}_i = 0\}, \\ \mathbf{x}_i = \mathbb{1}\{\mathbf{v}_i > p(\mathbf{z}_i^\top \gamma_0)\}, \end{cases}$$

where the outcome variable is a multi-year (1989–1993) average of log hourly wages deflated to 1983\$: our identification analysis focuses on each of the two potential outcomes and it naturally extends to the case where the errors in each of the potential outcome equations are different.

Since \mathbf{x}_i is *increasing* in \mathbf{v}_i for fixed \mathbf{z}_i and \mathbf{v}_i is assumed independent of \mathbf{z}_i , we interpret \mathbf{v}_i as an unobserved measure of an individual's inclination to attend college, i.e. a measure of such inclination that is not accounted for by \mathbf{z}_i .

We estimate γ_0 (normalized to have a unit length) in the schooling equation using Ichimura's semiparametric least squares estimator (SLSE), which uses the fact that

$$\gamma_0 = \underset{\|\gamma\|=1}{\operatorname{argmin}} \mathbb{E} \left\{ \mathbf{x}_i - \mathbb{E}(\mathbf{x}_i | \mathbf{z}_i^\top \gamma) \right\}^2.$$

Ichimura (1993) states conditions under which the SLSE is \sqrt{n} -consistent.

The specification of the θ -parameters is in essence a double index model (Ichimura and Lee, 1991) since

$$\begin{cases} \mathbb{E}(y_i | \mathbf{x}_i = 0, \mathbf{a}_i = a, \mathbf{z}_i = z) = \mathbb{E}\{g_0(0, a^\top \theta_{00}, \mathbf{u}_{i0}) | \mathbf{v}_i \leq p(\mathbf{z}_i^\top \gamma_0)\} = G_0(a^\top \theta_{00}, \mathbf{z}_i^\top \gamma_0), \\ \mathbb{E}(y_i | \mathbf{x}_i = 1, \mathbf{a}_i = a, \mathbf{z}_i = z) = \mathbb{E}\{g_1(1, a^\top \theta_{10}, \mathbf{u}_{i1}) | \mathbf{v}_i > p(\mathbf{z}_i^\top \gamma_0)\} = G_1(a^\top \theta_{10}, \mathbf{z}_i^\top \gamma_0), \end{cases}$$

albeit that in our case we already have \sqrt{n} -consistent estimates of γ_0 . We therefore use [Escanciano, Jacho-Chávez, and Lewbel \(2010, EJL10\)](#) instead of [Ichimura and Lee \(1991\)](#); \sqrt{n} -consistency obtains under conditions stated in [EJL10](#).

Since our analysis is semiparametric, we scale the exogenous variables (i.e. controls and instruments) by their standard deviations. All computations are done in Matlab using its global optimization toolbox with 100 initial values. For the bandwidth choice we follow [Härdle, Hall, and Ichimura \(1993\)](#), which entails optimizing over the coefficients for each bandwidth and then doing a grid search over the bandwidth. The search range is (0.01, 0.41) noting that $2n^{-1/5} \approx 0.4$. There is, as is not unusual, some sensitivity to the choice of bandwidth, but changing the bandwidth does not affect the qualitative conclusions of our study.

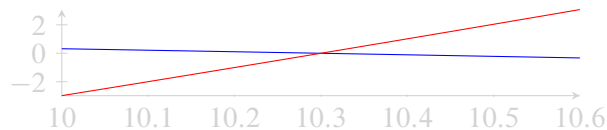


FIGURE 1. Effect of permanent local log earnings at age 17

5.3. Results pertaining to the schooling decision. We now turn to a discussion of our estimation results. Table 2 shows the estimates of the schooling equation coefficients. We consider two specifications: one includes proxies of intellectual ability (AFQT¹¹ and mother’s education) and the other does not. The full model is consistent with specifications used by others and is our main specification. We include the models without proxies mainly to demonstrate that (i) v^* should be interpreted as an inclination to attend college conditional on ability and that (ii) the AFQT scores and mother’s education do a great job proxying for ability.

The estimates can be interpreted as average partial effects up to scale normalization, because the schooling equation is a single index model. The coefficients in the model without proxies are normalized to have the same norm for the corresponding visible coefficients in the full model. The last column of table 2 contains the corresponding estimates of the first stage logit regression reported in the online appendix of [CHV11](#) normalized to have the same norm as the visible coefficients of the full model after correcting for our scale normalization of the covariates mentioned earlier. The most striking difference is in the effect of permanent local log earnings at 17 (PLLE17): figure 1 displays the index as a function of PLLE17 holding all other regressors constant, where the domain

¹¹The Armed Forces Qualification Test was administered to all subjects.

	Full model		No proxies		CHV11
	$\hat{\gamma}$	S.E.	$\hat{\gamma}$	S.E.	
Corrected AFQT	0.4630	(0.0421)**			0.2237
Corrected AFQT Squared	0.0411	(0.0634)			0.0118
Mother's Years of Schooling	-0.1852	(0.0532)**			-0.1897
Mother's Years of Schooling Squared	0.1883	(0.0680)**			0.0258
Number of Siblings	-0.4398	(0.0586)**	-0.4376	(0.0355)**	-0.3833
Number of Siblings Squared	0.0010	(0.0401)	-0.0235	(0.0302)	0.0016
Urban Residence at 14	-0.0130	(0.0330)	0.1243	(0.0596)*	0.0023
Permanent Local Log Earnings at 17	0.0559	(0.0239)*	0.0524	(0.0284)	-0.5112
Permanent Local Log Earnings at 17 Sq.	-0.0551	(0.0234)*	-0.0014	(0.0325)	0.5146
Permanent State Unempl. Rate at 17	0.0282	(0.0329)	-0.0415	(0.0323)	-0.0145
Permanent State Unempl. Rate at 17 Sq.	-0.0415	(0.0319)	0.0266	(0.0310)	-0.0139
Presence of a College at 14	0.0135	(0.0352)	0.0500	(0.0475)	0.0174
× AFQT	0.0046	(0.0495)			0.0059
× Mother's Education	0.0087	(0.0383)			-0.0180
× Number of Siblings	-0.0249	(0.0430)	-0.0212	(0.0396)	0.0019
Local Log Earnings at 17	-0.0618	(0.0431)	-0.0686	(0.0488)	-0.0305
× AFQT	-0.1882	(0.1100)			-0.1981
× Mother's Education	0.1665	(0.0293)**			0.1953
× Number of Siblings	0.5344	(0.0756)**	0.3375	(0.0779)**	0.3843
Local Unemployment Rate at 17	0.1097	(0.0428)*	-0.0214	(0.0372)	0.0377
× AFQT	-0.0099	(0.0499)			-0.0052
× Mother's Education	-0.0553	(0.0380)			-0.0398
× Number of Siblings	0.0292	(0.0447)	0.0072	(0.0355)	-0.0046
Tuition in 4 Year Public Colleges at 17	-0.1385	(0.0514)**	-0.0076	(0.0389)	-0.0279
× AFQT	0.1783	(0.0720)*			0.0081
× Mother's Education	0.1667	(0.0355)**			0.0297
× Number of Siblings	-0.2352	(0.0498)**	-0.1113	(0.0439)*	-0.0062
Bandwidth	0.03		0.04		

Notes: The birth year cohort dummy coefficients are omitted from the table. The numbers in parentheses are standard errors that are computed using the bootstrap with 500 replications. The coefficients in the full model are normalized to have norm one with the birth year cohort dummies included. The no proxies coefficients are normalized to have the same norm as the full model coefficients for the visible coefficients, e.g. the norm of the coefficients in the *no proxies* column is the same as that of the corresponding coefficients in the full model column. The CHV11 coefficients are taken from their online appendix, rescaled to be consistent with our choice to normalize covariates to have unit variance and then renormalized to have the same norm as the visible vector of full model coefficients. Significance at the 5% and 1% level is indicated with asterisks.

TABLE 2. Schooling equation coefficients.

is $\text{mean} \pm \text{two standard deviations}$ and where the functions are vertically centered at zero. Both are close to linear with one increasing and the other decreasing, albeit that the magnitudes of the CHV11 coefficients far exceeds ours. We have done further experiments (not tabulated) to determine the source of this difference, which suggest that the differences are mainly due to the nonparametric and

nonseparable second stage in our estimation method: if we replace our semiparametric first stage with a fully parametric first stage then the results are largely unchanged.

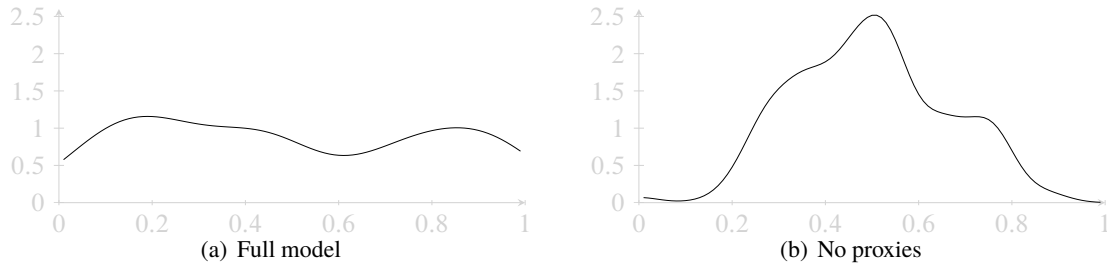


FIGURE 2. Density of $p(z_1^T \gamma_0)$

Figure 2 depicts the density estimates of $p(z_1^T \gamma_0)$. When the intelligence proxies are excluded, the density of the propensity score is small at extreme quantiles, i.e. at extreme levels of the objects' inclination to attend college. This is intuitive since there are no purely demographic variables that would explain the college education decision with (near) certainty; individuals with very low AFQT scores will not attend college with certainty but once AFQT scores and mother's education are omitted there are no covariates left to predict the outcome with an equal degree of certainty.

5.4. Results pertaining to returns to education. Regression coefficients for the second stage can be found in table 3. The last two columns contain the coefficients of the normal switching regressions reported in the online appendix of [CHV11](#), rescaled to make them comparable to ours. The coefficients in our full model, which as noted is our main specification, seem sensible. A few coefficients differ in sign from what is found in the [CHV11](#) regressions, especially for group 0. The most noticeable difference is in the effect of PLLE17, which is the same variable that had different coefficients in the schooling decision equation. But most coefficients, including those on PLLE17, are insignificant in [CHV11](#), it is unclear what one would expect the signs of the coefficients to be, and our model nests [CHV11](#) so the differences may be due to misspecification or indeed noise due to overspecification.

We now turn our attention to the potential wage curves depicted in figure 3.¹² Like in the existing literature potential wages for group 0 in the full model decrease as v^* increases. There are however some important differences: our potential wage curves for the two education groups do not cross —

¹²Note that averaging the covariates or v^* do not yield 'unconditional' quantiles of potential wages.

	Full model		No proxies		CHV11	
	no coll.	college	no coll.	college	no coll.	college
Experience	0.7998	0.6668	0.0131	0.6511	0.0575	0.5384
Experience Squared	-0.2504	-0.4453	-0.8404	-0.4946	0.0073	-0.4068
Log Local Earnings in 1991	0.2685	0.3839	-0.4730	0.4316	0.0283	0.2605
Local Unemployment in 1991	-0.1372	0.0748	0.0273	0.0766	-0.0014	0.0077
Corrected AFQT	-0.0321	0.1562			0.0143	0.2202
Corrected AFQT Squared	0.0169	0.1060			-0.0171	0.1202
Mother's Years of Schooling	-0.0433	-0.0990			-0.0129	-0.0429
Mother's Years of Schooling Sq.	0.0191	0.0336			0.0152	0.1508
Number of Siblings	0.0577	0.0024	-0.0648	-0.0131	0.0024	-0.0369
Number of Siblings Squared	0.0421	0.0306	0.0423	0.0121	-0.0036	0.0066
Urban Residence at 14	-0.0778	0.0595	0.0459	0.1747	0.0010	0.0452
Perm. Local Log Earnings at 17	-0.1781	-0.2023	-0.1534	-0.0462	0.7027	0.3201
Perm. Local Log Earnings at 17 Sq.	0.2842	0.2558	-0.0143	0.1173	-0.6936	-0.3144
Perm. State Unempl. Rate at 17	0.2230	0.1344	0.0744	0.0403	0.0210	0.2724
Perm. State Unempl. Rate at 17 Sq.	-0.1034	-0.1531	-0.0955	-0.1980	-0.0222	-0.3241
Bandwidth	0.21	0.39	0.40	0.29		

Notes: Birth year cohort dummy coefficients are not reported. The coefficients in the last two columns were scaled to have the same norm as the vector containing the coefficients to the same variables in the first two columns. The dependent variable is log hourly wages.

TABLE 3. Earnings equation coefficients

returns on a college education are consistently positive — and potential wages for group 1 decrease as v^* increases. There are several explanations for the decreasing potential wage curves.

One explanation is that the intelligence proxies adequately measure ability and that v^* should now be interpreted as one's inclination to attend college. Those with an unfulfilled desire to attend college may do poorly in the labor market, as do those who attend college despite their lack of ability (as measured by AFQT scores and mother's education level). Such explanations are consistent with the graphs for the model without the intelligence proxies in figure 3 which feature (mostly) increasing potential wage curves because absent proxies v^* is correlated with overall ability. We conclude that AFQT and mother's education do a good job proxying for ability and focus on the full model from hereon.

The return on a college education at different quantiles of education equation (v^*) and earnings equation (τ^*) unobservables are presented in figure 4, where the controls are fixed at their sample means. Figure 5 depicts the same information in greater detail at three values of τ^* and figure 6 at five values of τ^* but without confidence bands.

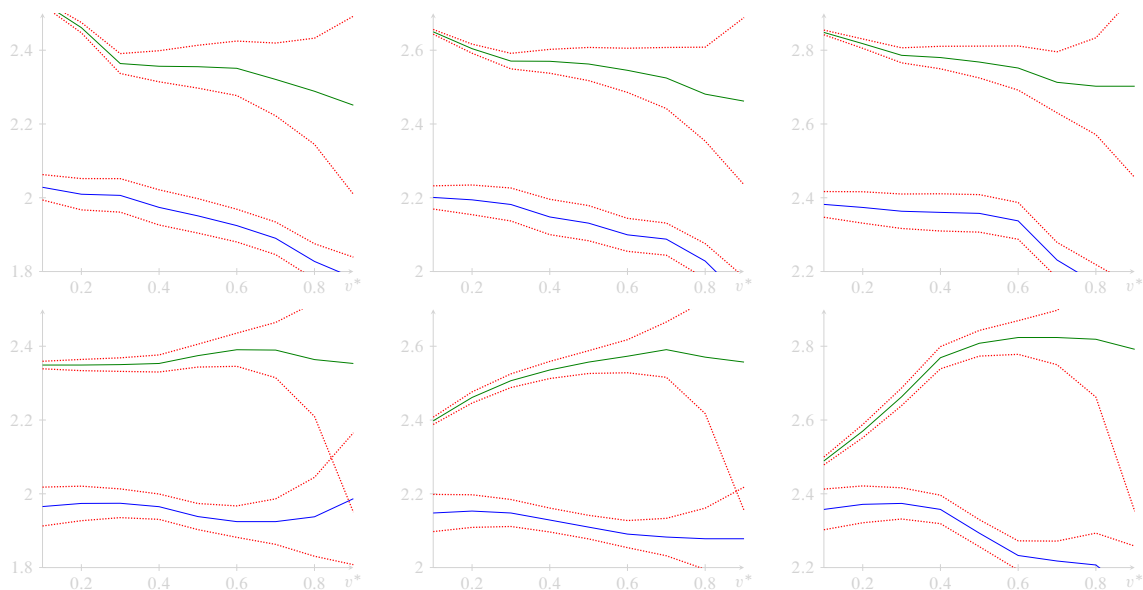


FIGURE 3. Potential wages as a function of v^* for $\tau^* = 0.3$ (left), $\tau^* = 0.5$ (middle), and $\tau^* = 0.7$ (right), with (top) and without (bottom) intelligence proxies. Covariates are evaluated at their means. Dotted lines show (pointwise) 95% confidence intervals. $2.2 \approx \$9.03$, $2.4 \approx \$11.02$, $2.6 \approx \$13.46$, $2.8 \approx \$16.44$.

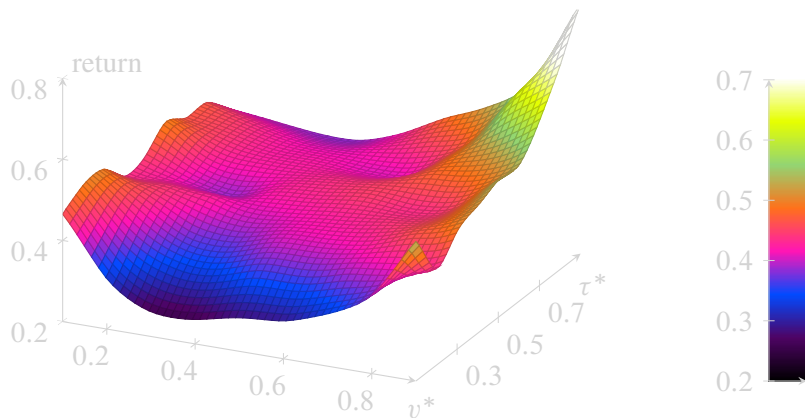


FIGURE 4. Return on a college education as a function of v^* , τ^* ; covariates evaluated at their means; 0.2 corresponds to a 22% premium, $0.4 \approx 49\%$, $0.6 \approx 82\%$, and $0.7 \approx 100\%$.

Returns to college vary substantially with both τ^* and v^* , albeit that figure 5 suggests that at the moderate τ^* quantiles depicted returns are similar except at high v^* values. More striking is the $\tau^* = 0.1$ graph in figure 6, which indicates that the returns on education are substantially less at the bottom end once inclination to attend college and demographics are controlled for.

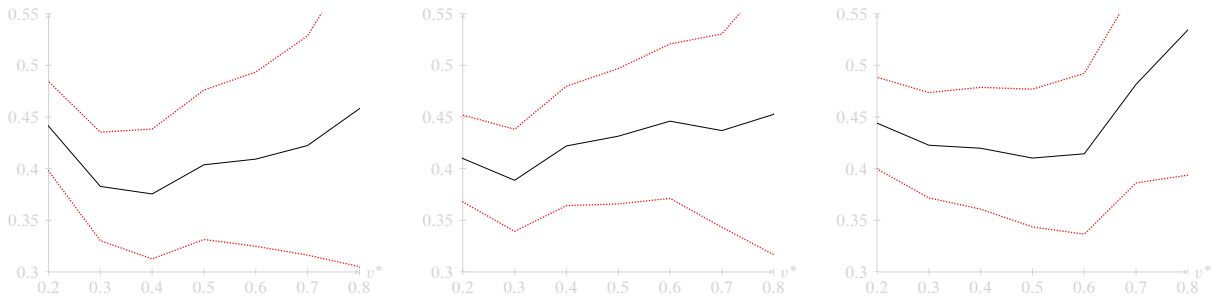


FIGURE 5. Returns to college as a function of v^* for $\tau^* = 0.3$ (left), $\tau^* = 0.5$ (middle), and $\tau^* = 0.7$ (right) for the full model. Covariates are evaluated at their means. Dotted lines show (pointwise) 95% confidence intervals

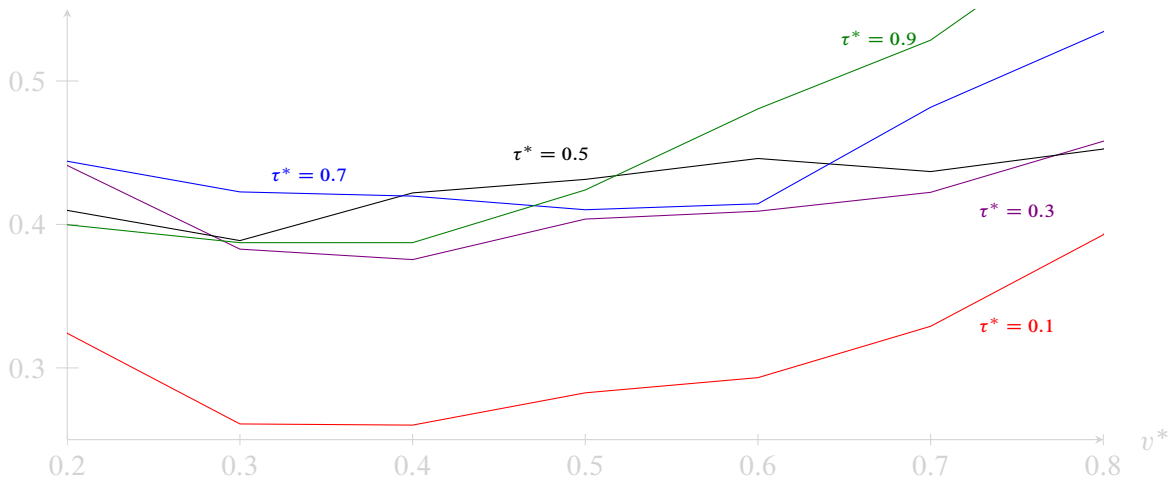


FIGURE 6. Returns for various τ^* values in the full model. Covariates are evaluated at their means.



FIGURE 7. Returns for $\tau^* = 0.3$ (left), $\tau^* = 0.5$ (middle), and $\tau^* = 0.7$ (right) in the full model at different covariate values. Covariates are evaluated as indicated.

All of the graphs discussed thus far have covariates evaluated at their means. Figure 7 introduces variation across covariate values and depicts returns as a function of v^* for three moderate values of τ^* where the covariates are measured at four different values: three quartiles and the mean.

CHV11 and CL09 impose shape invariance restrictions that are not imposed by us. The shape invariance restrictions in both of those papers imply that the returns curves (and in fact also the potential wage curves) are identical to each other up to a vertical shift if one changes covariate values. Such shape invariance conditions look inconsistent with figure 7.

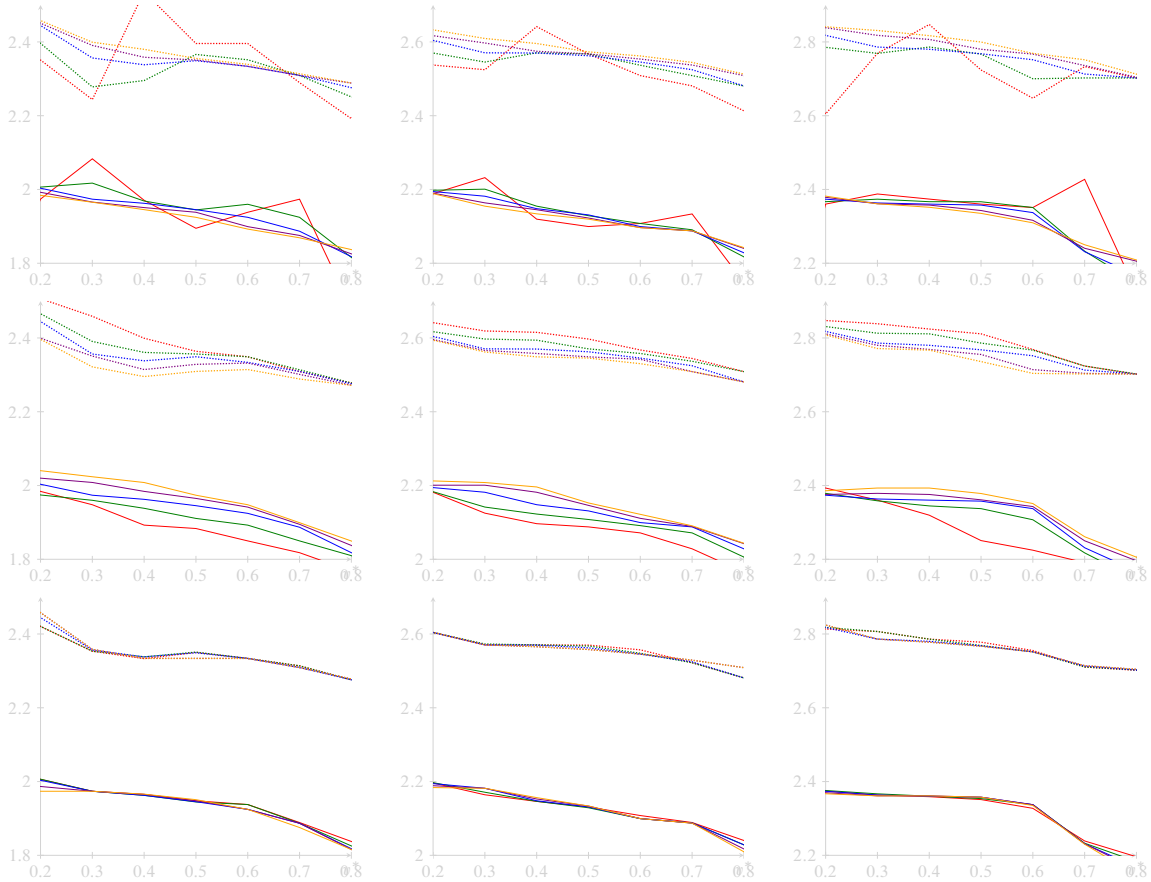


FIGURE 8. Potential wages for different values of τ^* (0.3, 0.5, 0.7 from left to right) and different bandwidth choices. At the top we vary $h_0 \in \{0.05, 0.10, 0.15, 0.20, 0.25\}$ with $h_1 = 0.35, h_z = 0.15$, in the middle we vary $h_1 \in \{0.25, 0.30, 0.35, 0.40, 0.45\}$ with $h_0 = 0.15, h_z = 0.15$, and at the bottom we vary $h_z \in \{0.05, 0.10, 0.15, 0.20, 0.25\}$ with $h_0 = 0.15, h_1 = 0.35$. Each colored curve denotes a different bandwidth ordered from small to large (red, green, blue, purple, and orange). Dotted curves are for the college educated case.

We conclude our discussion with an analysis of the sensitivity of our results to the choice of input parameters. The results are depicted in figure 8. The results are unusually robust to the choice of h_z , which we attribute to the fact that the semiparametric estimation procedure averages across \hat{p}_i 's.¹³

¹³The word 'averaging' here should be interpreted in the sense of inter alia Linton and Nielsen (1995), not an immediate sample mean. We average over functions in which \hat{p}_i enters as one of the arguments.

With severe undersmoothing the h_1 curves and especially the h_0 curves get the expected nonsmooth appearance which is exacerbated by the fact that we only computed estimates at 0.1 increments of the v^* -values.

REFERENCES

- ARCONES, M., AND E. GINÉ (1993): "Limit theorems for U-processes," *Annals of Probability*, 21(3), 1494–1542.
- CARNEIRO, P., J. J. HECKMAN, AND E. J. VYTLACIL (2011): "Estimating marginal returns to education," *American Economic Review*, 101, 2754–2781.
- CARNEIRO, P., AND S. LEE (2009): "Estimating distributions of potential outcomes using local instrumental variables with an application to changes in college enrollment and wage inequality," *Journal of Econometrics*, 149(2), 191–208.
- CHERNOZHUKOV, V., AND C. HANSEN (2005): "An IV model of quantile treatment effects," *Econometrica*, 73(1), 245–261.
- (2006): "Instrumental quantile regression inference for structural and treatment effect models," *Journal of Econometrics*, 132(2), 491–525.
- CHESHER, A. (2003): "Identification in nonseparable models," *Econometrica*, 71(5), 1405–1441.
- (2005): "Nonparametric identification under discrete variation," *Econometrica*, 73, 1525–1550.
- DE LA PEÑA, V., AND E. GINÉ (1999): *Decoupling: from dependence to independence*. Springer Verlag.
- DELGADO, M., AND J. MORA (1995): "Nonparametric and semiparametric estimation with discrete regressors," *Econometrica*, 63(6), 1477–1484.
- DOKSUM, K. (1974): "Empirical probability plots and statistical inference for nonlinear models in the two sample case," *Annals of Statistics*, 2(2), 266–277.
- ESCANCIANO, J. C., D. JACHO-CHÁVEZ, AND A. LEWBEL (2010): "Identification and estimation of semiparametric two step models," *Unpublished manuscript*.
- HÄRDLE, W., P. HALL, AND H. ICHIMURA (1993): "Optimal smoothing in single-index models," *Annals of Statistics*, 21(1), 157–178.
- HECKMAN, J., AND E. VYTLACIL (1999): "Local instrumental variables and latent variable models for identifying and bounding treatment effects," *Proceedings of the National Academy of Sciences*, 96(8), 4730.
- (2001): "Local instrumental variables," in *Nonlinear statistical modeling: Proceedings of the thirteenth international symposium in economic theory and econometrics: essays in honor of Takeshi Amemiya*, vol. 13, pp. 1–15. Cambridge University Press.

- (2005): “Structural equations, treatment effects, and econometric policy evaluation,” *Econometrica*, 73(3), 669–738.
- HODERLEIN, S., AND E. MAMMEN (2007): “Identification of marginal effects in nonseparable models without monotonicity,” *Econometrica*, 75(5), 1513–1518.
- ICHIMURA, H. (1993): “Semiparametric least squares (SLS) and weighted SLS estimation of single-index models,” *Journal of Econometrics*, 58(1), 71–120.
- ICHIMURA, H., AND L.-F. LEE (1991): “Semiparametric least squares estimation of multiple index models: single equation estimation,” in *Nonparametric and semiparametric methods in econometrics and statistics: Proceedings of the Fifth International Symposium in Economic Theory and Econometrics*. Cambridge, pp. 3–49.
- IMBENS, G. W., AND W. K. NEWEY (2009): “Identification and estimation of triangular simultaneous equations models without additivity,” *Econometrica*, 77(5), 1481–1512.
- JUN, S. (2009): “Local structural quantile effects in a model with a nonseparable control variable,” *Journal of Econometrics*, 151(1), 82–97.
- JUN, S. J., J. PINKSE, AND H. XU (2011): “Tighter bounds in triangular systems,” *Journal of Econometrics*, 161, 122–128.
- (2012): “Discrete endogenous variables in weakly separable models,” *Econometrics Journal*, 15, 288–303.
- JUN, S. J., J. PINKSE, H. XU, AND N. YILDIZ (2010): “Identification of average treatment effects in a triangular system of equations,” Discussion paper, Pennsylvania State University.
- KASY, M. (2011): “Identification in triangular systems using control functions,” *Econometric Theory*, 27, 663–671.
- KÉDAGNI, D., AND I. MOURIFIÉ (2014): “Tightening Bounds in Triangular Systems,” Discussion paper, University of Toronto.
- KLEIN, R., AND R. SPADY (1993): “An efficient semiparametric estimator for binary response models,” *Econometrica*, 61, 387–421.
- LEE, A. (1990): *U-statistics: theory and practice*. Marcel Dekker.
- LINTON, O., AND J. P. NIELSEN (1995): “A kernel method of estimating structured nonparametric regression based on marginal integration,” *Biometrika*, pp. 93–100.
- MA, L., AND R. KOENKER (2006): “Quantile regression methods for recursive structural equation models,” *Journal of Econometrics*, 134(2), 471–506.

- MAMMEN, E., C. ROTHE, AND M. SCHIENLE (2012): “Nonparametric regression with nonparametrically generated covariates,” *Annals of Statistics*, 40(2), 1132–1170.
- PAGAN, A., AND A. ULLAH (1999): *Nonparametric econometrics*. Cambridge University Press.
- PINKSE, J. (2001): “Nonparametric regression estimation using weak separability,” Discussion paper, UBC.
- POWELL, J. L., J. H. STOCK, AND T. M. STOKER (1989): “Semiparametric estimation of index coefficients,” *Econometrica: Journal of the Econometric Society*, pp. 1403–1430.
- RILSTONE, P. (1996): “Nonparametric estimation of models with generated regressors,” *International Economic Review*, pp. 299–313.
- VAN DER VAART, A. (2000): *Asymptotic statistics*. Cambridge University Press.
- VAN DER VAART, A., AND J. WELLNER (1996): *Weak convergence and empirical processes*. Springer.
- VYTLACIL, E. (2006): “A note on additive separability and latent index models of binary choice: representation results,” *Oxford Bulletin of Economics and Statistics*, 68(4), 515–518.

APPENDIX A. LEMMAS FOR IDENTIFICATION

Lemma A1. *Suppose that assumptions G and H are satisfied. Then,*

$$\lim_{t \rightarrow \infty} \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V_t) = \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*),$$

where $V_t = (p(z_t), p(z_{t+1}))$ and $\{p(z_t)\}$ is as in assumption H.

Proof. Choose $\epsilon > 0$. By assumption G there exists a $v_\epsilon < v^*$ such that for all $v \in (v_\epsilon, v^*]$,

$$\mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) - \epsilon < \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v) < \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) + \epsilon.$$

Recalling that $\mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*)$ is the smallest value of u for which $\mathbb{P}(\mathbf{u} \leq u|\mathbf{v} = v^*) \geq \tau^*$, it follows that for all $v \in (v_\epsilon, v^*]$,

$$\mathbb{P}\{\mathbf{u} \leq \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) - \epsilon|\mathbf{v} = v\} < \tau^* \leq \mathbb{P}\{\mathbf{u} \leq \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) + \epsilon|\mathbf{v} = v\}.$$

Hence, if one picks t large enough to ensure that $v_\epsilon < p(z_t) < p(z_{t+1}) < v^*$ then

$$\mathbb{P}\{\mathbf{u} \leq \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) + \epsilon|\mathbf{v} \in V_t\} = \frac{\int_{p(z_t)}^{p(z_{t+1})} \mathbb{P}\{\mathbf{u} \leq \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) + \epsilon|\mathbf{v} = v\} dv}{p(z_{t+1}) - p(z_t)} \geq \tau^*, \quad (24)$$

and similarly

$$\mathbb{P}\{\mathbf{u} \leq \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) - \epsilon|\mathbf{v} \in V_t\} < \tau^*. \quad (25)$$

Hence, it follows from (24) and (25) that

$$\mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) - \epsilon < \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V_t) \leq \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) + \epsilon.$$

Since $\epsilon > 0$ was arbitrarily chosen, the proof is done. \square

Lemma A2. *Suppose that assumptions F to H are satisfied. Then,*

$$\sup_{V \in \mathcal{D}_L(x^*, v^*)} \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V) = \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) \leq \inf_{V \in \mathcal{D}_U(x^*, v^*)} \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V). \quad (26)$$

Proof. By assumption F and definition of $\mathcal{D}_L(x^*, v^*)$ and $\mathcal{D}_U(x^*, v^*)$, we have

$$\sup_{V \in \mathcal{D}_L(x^*, v^*)} \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V) \leq \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) \leq \inf_{V \in \mathcal{D}_U(x^*, v^*)} \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V). \quad (27)$$

Now let $V_t = (p(z_t), p(z_{t+1}))$, where $\{p(z_t)\}$ be as in assumption G. Since $V_t \in \mathcal{D}_L(x^*, v^*)$, we have

$$\forall t, \quad \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V_t) \leq \sup_{V \in \mathcal{D}_L(x^*, v^*)} \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V).$$

Therefore, it follows from lemma A1 that

$$\mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|v^*) \leq \sup_{V \in \mathcal{D}_L(x^*, v^*)} \mathbb{Q}_{\mathbf{u}|\mathbf{v}}(\tau^*|V). \quad (28)$$

Combining (27) with (28) completes the proof. \square

APPENDIX B. TECHNICAL LEMMAS

Lemma B1. *Let $\Xi_n = \sum_{i=1}^n \xi_{ni}$, where $\{\xi_{ni}\}$ is an i.i.d. mean zero sequence of functions whose elements can depend on n . For any compact set Υ , suppose that $\tilde{\zeta}_n > 1$ is such that $\sup_{v \in \Upsilon} \|\partial_v \Xi_n(v)\| \leq \tilde{\zeta}_n$, let $\sigma_{n\xi}^2 = \sup_{v \in \Upsilon} \nabla \xi_{ni}(v)$ and let $\bar{\xi}_n$ be such that $\mathbb{P}\{\sup_{v \in \Upsilon} |\xi_{ni}(v)| > \bar{\xi}_n\} = 0$. If $\sigma_{n\xi}^2 < 1/n \log \tilde{\zeta}_n$ and $\bar{\xi}_n < 1/\log \tilde{\zeta}_n$ then $\sup_{v \in \Upsilon} |\Xi_n(v)| < 1$.*

Proof. Cover Υ using ζ_n balls $\Upsilon_1, \dots, \Upsilon_{\zeta_n}$ with centroids v_1, \dots, v_{ζ_n} , in such a way that for any n , $\max_{t=1, \dots, \zeta_n} \sup_{v \in \Upsilon_t} \|v - v_t\| \leq C/\zeta_n^{1/d_v}$ for some C independent of n . Then

$$\sup_{v \in \Upsilon} |\Xi_n(v)| \leq \max_{t=1, \dots, \zeta_n} \sup_{v \in \Upsilon_t} |\Xi_n(v) - \Xi_n(v_t)| + \max_{t=1, \dots, \zeta_n} |\Xi_n(v_t)|. \quad (29)$$

Choose any $\epsilon > 0$. For $\delta > 0$ to be chosen, let $\zeta_n = (\delta \tilde{\zeta}_n / \epsilon)^{d_v}$.

For RHS2 in (29) we have by the Bernstein inequality that

$$\mathbb{P}\left\{\max_{t=1,\dots,\zeta_n} |\Xi_n(\nu_t)| > \epsilon\right\} \leq \sum_{t=1}^{\zeta_n} \mathbb{P}\left\{|\Xi_n(\nu_t)| > \epsilon\right\} \leq 2\zeta_n \exp\left\{-\frac{\epsilon^2}{2(n\sigma_{n\xi}^2 + \bar{\xi}_n\epsilon)}\right\} < 1.$$

Finally, for RHS1 in (29) we have by the mean value theorem that

$$\begin{aligned} \mathbb{P}\left\{\max_{t=1,\dots,\zeta_n} \sup_{\nu \in \Upsilon_t} |\Xi_n(\nu) - \Xi_n(\nu_t)| > \epsilon\right\} &\leq \mathbb{P}\left\{\sup_{\nu \in \Upsilon} \|\partial_\nu \Xi_n(\nu)\| \max_{t=1,\dots,\zeta_n} \sup_{\nu \in \Upsilon_t} \|\nu - \nu_t\| > \epsilon\right\} \\ &\leq \mathbb{P}\left\{\sup_{\nu \in \Upsilon} \|\partial_\nu \Xi_n(\nu)\| > \frac{\epsilon \zeta_n^{1/d_\nu}}{C}\right\} \leq \mathbb{P}\left\{\sup_{\nu \in \Upsilon} \|\partial_\nu \Xi_n(\nu)\| > \frac{\delta \tilde{\zeta}_n}{C}\right\}. \end{aligned}$$

Let $n \rightarrow \infty$ followed by $\delta \rightarrow \infty$. □

Lemma B2. Let $\{\xi_{ni}^*\}$ be an i.i.d. sequence of mean zero functions defined on a compact set Υ for which for some $C < \infty$, $\sup_n [n^{-C} E\{\sup_{\nu \in \Upsilon} \|\partial_\nu \xi_{ni}^*(\nu)\|\}] < \infty$. Let further $\sigma_{n\xi^*}^2 = \sup_{\nu \in \Upsilon} \mathbb{V} \xi_{ni}^*(\nu)$ and $\sup_n \mathbb{P}(\sup_{\nu \in \Upsilon} |\xi_{ni}^*(\nu)| > \bar{\xi}_n^*) = 0$. Then for any $\zeta_n > \max(\sqrt{\sigma_{n\xi^*}^2 \log n/n}, \bar{\xi}_n^* \log n/n)$,

$$\sup_{\nu \in \Upsilon} \left| n^{-1} \sum_{i=1}^n \xi_{ni}^*(\nu) \right| < \zeta_n.$$

Proof. In lemma B1 take $\xi_{ni} = \xi_{ni}^*/n\zeta_n$. □

Lemma B3. Let $\{\xi_i\}$ be an i.i.d. sequence, let ξ_i include \mathbf{y}_i as an element, and let $\{\hat{\mathbf{A}}_i\}$ be such that $\hat{\mathbf{A}}_i = A_n(\xi_i, \xi_1, \dots, \xi_{i-1}, \xi_{i+1}, \dots, \xi_n)$ for arbitrary function A_n . If $\mathbb{P}(\|\hat{\mathbf{A}}_1\| > \epsilon) < 1/n$ and $\sup_{\mathbf{y} \in \mathcal{Y}} \mathbb{P}\left\{\left\|\sum_{i=1}^{n-1} \hat{\mathbf{A}}_i \mathbf{I}_i(\mathbf{y})\right\| > \epsilon\right\} < 1/n$, then $\sup_{\mathbf{y} \in \mathcal{Y}} \|\sum_i \hat{\mathbf{A}}_i \mathbf{I}_i(\mathbf{y})\| < 1$.

Proof.

$$\begin{aligned} \mathbb{P}\left\{\sup_{\mathbf{y} \in \mathcal{Y}} \left\|\sum_{i=1}^n \hat{\mathbf{A}}_i \mathbf{I}_i(\mathbf{y})\right\| > 2\epsilon\right\} &= \mathbb{P}\left\{\max_{t=1,\dots,n} \left\|\sum_{i=1}^n \hat{\mathbf{A}}_i \mathbf{I}_i(\mathbf{y}_t)\right\| > 2\epsilon\right\} \leq \sum_{t=1}^n \mathbb{P}\left\{\left\|\sum_{i=1}^n \hat{\mathbf{A}}_i \mathbf{I}_i(\mathbf{y}_t)\right\| > 2\epsilon\right\} \\ &\leq n \sup_{\mathbf{y} \in \mathcal{Y}} \mathbb{P}\left\{\left\|\sum_{i=1}^{n-1} \hat{\mathbf{A}}_i \mathbf{I}_i(\mathbf{y})\right\| > \epsilon\right\} + n \mathbb{P}\{\|\hat{\mathbf{A}}_1\| > \epsilon\} < 1. \quad \square \end{aligned}$$

APPENDIX C. V-STATISTICS

Let $\mathfrak{N}_n = \{1, \dots, n\}$, $\Upsilon_{n\ell} = \mathfrak{N}_n^\ell$, and let $\Upsilon_{n\ell j}$ be the set of vectors in $\Upsilon_{n\ell}$ with exactly j distinct elements. Let further for any $\iota \in \Upsilon_{n\ell}$, $\xi_\iota = (\xi_{\iota_1}, \dots, \xi_{\iota_\ell})^\top$.

Lemma C1. For $V_{n\ell} = \sum_{\iota \in \Upsilon_{n\ell}} m(\xi_\iota)$ and $U_n^{(\ell,j)} = \sum_{\iota \in \Upsilon_{n\ell j}} m(\xi_\iota)$, we have $V_{n\ell} = \sum_{j=1}^{\ell} U_n^{(\ell,j)}$, where $U_n^{(\ell,j)}$ is a U -statistic of order j whose kernel $m^{(\ell,j)}$ consists of a sum of $\sum_{t=1}^j (-1)^{j-t} t^{\ell-1} / \{(j-t)!(t-1)!\}$ elements.¹⁴

Proof. See Lee (1990, theorem 1 on p.183). \square

Lemma C2. For any symmetric j -th order U -statistic kernel $m^{(j)}$, let $U_{nj} = \sum_{\iota \in \Upsilon_{nj}} m^{(j)}(\xi_\iota)$. Let further for $0 \leq t \leq j$ and any a_1, \dots, a_t , $m_t^{(j)}(a) = \mathbb{E}m^{(j)}(a_1, \dots, a_t, \xi_1, \dots, \xi_{j-t})$, $U_{njt} = \sum_{\iota \in \Upsilon_{njt}} m_t^{(j)}(\xi_\iota)$, and U_{njt}^c the corresponding canonical U -statistic (De la Peña and Giné, 1999). Then if $\mu^{(j)} = \mathbb{E}m^{(j)}(\xi_1, \dots, \xi_j)$,

$$U_{nj} = \frac{n!}{(n-j)!} \mu^{(j)} + \sum_{t=1}^j \binom{j}{t} \frac{(n-t)!}{(n-j)!} U_{njt}^c.$$

Proof. This is essentially the Hoeffding decomposition (Lee, 1990, theorem 1 on p.26) combined with a rearrangement of terms.¹⁵ \square

Lemma C3. For U_{njt}^c defined in lemma C2, we have

$$\mathbb{P}(|U_{njt}^c| > \epsilon) \leq C_t \exp\left[-\frac{\epsilon^{2/t}/C_t}{n\sigma_{jt}^{2/t} + \beta_{jt}^{2/(t+1)} n^{(t-1)/(t+1)} \epsilon^{2/\{t(t+1)\}}}\right],$$

where C_t is a constant which only depends on t , $\beta_{jt} = \sup m_t^{(j)}(\cdot)$, $\sigma_{jt}^2 = \mathbb{V}m_t^{(j)}(\xi_1, \dots, \xi_t)$.

Proof. Follows from Arcones and Giné (1993, proposition 2.3(c)). \square

Lemma C4. For an ℓ -th order V -statistic $V_{n\ell}$ as defined in lemma C1 with symmetric kernel m , let for $1 \leq t \leq j \leq \ell$, $m^{(\ell,j)}$ be defined as in lemma C1, $\mu^{(\ell,j)} = \mathbb{E}m^{(\ell,j)}(\xi_1, \dots, \xi_j)$, $m_t^{(\ell,j)}(a) = \mathbb{E}m^{(\ell,j)}(a_1, \dots, a_t, \xi_1, \dots, \xi_{j-t})$, $\beta_t^{(\ell,j)} = \sup m_t^{(\ell,j)}$, and $\sigma_t^{(\ell,j)} = \sqrt{\mathbb{V}m_t^{(\ell,j)}(\xi_1, \dots, \xi_t)}$.

Then $\mathbb{P}(V_{n\ell} > \epsilon_n)$ decreases faster than any polynomial of n , where

$$\epsilon_n = \max_{1 \leq t \leq j \leq \ell} \left[(\log n)^{t+1} \max \left\{ n^{t/2} \sigma_t^{(\ell,j)}, n^{(t-1)/2} \beta_t^{(\ell,j)}, n^j \mu^{(\ell,j)} \right\} \right].$$

¹⁴These are Stirling numbers of the second kind.

¹⁵The representation is slightly different here from the one in Lee (1990) because the U -statistic kernel incorporates a number of permutations in his case.

Proof. In lemma C1 the V–statistic is separated into a number (independent of n) of U–statistics. Each of these U–statistics is further separated into a number (again independent of n) of canonical U–statistics in lemma C2 plus a mean. Finally, apply lemma C3 to each element individually.¹⁶ \square

APPENDIX D. \mathcal{Z}

Let π be such that $\pi\{p(z_1, \tilde{z}), \tilde{z}\} = z_1$ for all $(z_1, \tilde{z}) \in \mathcal{Z}$. The function π is well–defined by assumption P.

Lemma D1. *For all p and any $c > 0$, $f_p(p|\tilde{z})$ is four times boundedly differentiable with respect to p , uniformly in \tilde{z} for which $\tilde{f}(\tilde{z}) \geq c$.*

Proof. Note that $F_p(p|\tilde{z}) = \mathbb{P}(\mathbf{p}_i \leq p|\tilde{\mathbf{z}}_i = \tilde{z}) = \mathbb{P}\{z_{i1} \leq \pi(p, \tilde{z})|\tilde{\mathbf{z}}_i = \tilde{z}\}$, such that $f_p(p|\tilde{z}) = \partial_p \pi(p, \tilde{z}) f\{\pi(p, \tilde{z}), \tilde{z}\} / \tilde{f}(\tilde{z})$. The stated result then follows from assumption O. \square

Lemma D2. *For all p for which $\exists z \in \mathcal{Z} : p(z) = p$ and all t times boundedly differentiable functions v for which $v(z) = 0$ for all $z \notin \mathcal{Z}$, $\mathbb{E}\{v(\mathbf{z}_i)|\mathbf{p}_i = p\} f_p(p)$ is $\min(t, 3)$ times boundedly differentiable in p .*

Proof. Let π be as in lemma D1. Then for any $z \in \mathcal{Z}$ and $p = p(z)$,

$$\mathbb{E}\{v(\mathbf{z}_i)|\mathbf{p}_i = p\} f_p(p) = \int v\{\pi(p, \tilde{z}), \tilde{z}\} f\{\pi(p, \tilde{z}), \tilde{z}\} \partial_p \pi(p, \tilde{z}) d\tilde{z}. \quad \square$$

APPENDIX E. KERNELS

Lemma E1. *Let $\{(\xi_i, \mathbf{z}_i)\}$ be i.i.d., and suppose that $\mu(z) f(z)$ with $\mu(z) = \mathbb{E}(\xi_i|\mathbf{z}_i = z)$ has two bounded derivatives. Then*

$$\sup_{z \in \mathcal{Z}} \left| \mathbb{E}\{\mathbf{K}_{z_i}(z) \xi_i\} - \mu(z) f(z) \right| \leq h_z^2.$$

Proof. This follows from a standard kernel bias expansion. \square

Lemma E2 can be found, often in slightly different form, in many other sources, including Pagan and Ullah (1999).

¹⁶Because the number of canonical U–statistics has an upper bound independent of n , looking at each individual term separately is sufficient.

Lemma E2. Let $\{(\xi_i, z_i)\}$ be i.i.d., ξ_i uniformly bounded, and $\sigma_\xi^2(z) = \mathbb{V}(\xi_i | z_i = z)$ is continuous on \mathcal{Z} . Then

$$\sup_{z \in \mathcal{Z}} \left| \frac{1}{n} \sum_{i=1}^n [\mathbf{K}_{z_i}(z) \xi_i - \mathbb{E}\{\mathbf{K}_{z_i}(z) \xi_i\}] \right| < \frac{\log n}{\sqrt{nh_z^{d_z}}}.$$

Proof. Follows directly from lemma B2. \square

Let

$$\alpha_n = \log n / \sqrt{nh_z^{d_z}} + h_z^2. \quad (30)$$

Let further $\hat{r} = \hat{p} \hat{f}$, where \hat{f} is the kernel density estimator of f using bandwidth h_z and kernel K .

Lemma E3. (i) $\sup_{z \in \mathcal{Z}} |\hat{f}(z) - f(z)| \leq \alpha_n$, (ii) $\sup_{z \in \mathcal{Z}} |\hat{r}(z) - r(z)| \leq \alpha_n$, (iii) $\sup_{z \in \mathcal{Z}} |\hat{p}(z) - p(z)| \leq \alpha_n$,

Proof. The first two results follows by combining lemmas E1 and E2 and the third one from the first two by noting that for any $\tilde{z} \in \mathcal{Z}$,

$$\frac{r(\tilde{z}) - \sup_{z \in \mathcal{Z}} |\hat{r}(z) - r(z)|}{f(\tilde{z}) + \sup_{z \in \mathcal{Z}} |\hat{f}(z) - f(z)|} \leq \hat{p}(\tilde{z}) \leq \frac{r(\tilde{z}) + \sup_{z \in \mathcal{Z}} |\hat{r}(z) - r(z)|}{f(\tilde{z}) - \sup_{z \in \mathcal{Z}} |\hat{f}(z) - f(z)|}. \quad \square$$

Lemma E4. For some $\epsilon > 0$, $\lim_{n \rightarrow \infty} \mathbb{P}\{\inf_{z \in \mathcal{Z}} \hat{f}(z) < \epsilon\} = 0$.

Proof. Note that

$$\mathbb{P}\left\{\inf_{z \in \mathcal{Z}} \hat{f}(z) < \epsilon\right\} \leq \mathbb{P}\left\{\inf_{z \in \mathcal{Z}} f(z) < 2\epsilon\right\} + \mathbb{P}\left\{\sup_{z \in \mathcal{Z}} |\hat{f}(z) - f(z)| > \epsilon\right\}.$$

Apply lemma E3. \square

APPENDIX F. EXPANSIONS

Recalling (16), let $\mathcal{E}(y; p) = \mathbb{E}\{\mathbf{I}_i(y) | \mathbf{a}_i = a^*, \mathbf{p}_i = p\}$, $\mathcal{E}(y) = \mathcal{E}(y; v^*)$, and

$$\bar{\mathcal{S}}_s(y; p) = \frac{1}{n} \sum_{i=1}^n k_{is}^{(s)} \mathbf{K}_{ais} \mathcal{E}(y), \quad \bar{\mathcal{S}}_s(y; \hat{p}) = \frac{1}{n} \sum_{i=1}^n \hat{k}_{is}^{(s)} \mathbf{K}_{ais} \mathcal{E}(y). \quad (31)$$

Let further $\mathbf{K}_{zij} = \mathbf{K}_{zi}(z_j)$.

Lemmas F1 to F3 serve as inputs into establishing two results, namely

$$\sup_y \left| \hat{\mathcal{S}}_s(y; \hat{p}) - \bar{\mathcal{S}}_s(y; \hat{p}) - \hat{\mathcal{S}}_s(y; p) + \bar{\mathcal{S}}_s(y; p) \right| < 1/\rho_n \quad (32)$$

for $s = 0, 1$ and

$$\sup_y |\bar{\mathbf{S}}_0(y; \hat{p}) - \bar{\mathbf{S}}_0(y; p)| < 1/\rho_n, \quad (33)$$

i.e. lemmas F6 and F7. Each of these expression is expanded using the mean value theorem to some order J to apply lemmas F1 to F3. For instance, by assumption R, the RHS of (32) is bounded above by

$$\begin{aligned} & \sum_{j=1}^{J-1} \frac{1}{j!} \sup_{y \in \mathcal{Y}} \left| \frac{1}{n} \sum_{i=1}^n k_{i_s}^{(s+j)} \mathbf{K}_{ais}(\mathbf{p}_i - \hat{\mathbf{p}}_i)^j \{ \mathbf{I}_{i_0}(y) - \mathcal{E}(y) \} \right| \\ & + \sup_{y \in \mathcal{Y}} \frac{1}{n h_s^{s+1+J} J!} \sum_{i=1}^n \left| \mathbf{K}_{ais}(\mathbf{p}_i - \hat{\mathbf{p}}_i)^J \{ \mathbf{I}_{i_0}(y) - \mathcal{E}(y) \} \right|, \quad (34) \end{aligned}$$

where the first term is covered by lemmas F1 and F2 and the second term is dealt with in lemma F3.

Lemma F1. Let $\zeta_1, \dots, \zeta_j, \zeta^* \in \mathcal{F}_2$, and $\mathcal{T}_j \subset \mathbb{R}^j$ consist of vectors whose elements are all either equal to one or zero and let $\mathbf{u}_{\ell_i}^*$ be such that $\mathbb{E}(\mathbf{u}_{\ell_i}^* | \mathbf{z}_i) = 0$ a.s. and $\mathbb{E}(\mathbf{u}_{\ell_i}^{*2} | \mathbf{z}_i = z)$ is continuous on \mathcal{Z} . Then for $s = 0, 1$, $j = 1, 2, \dots$, and all $t \in \mathcal{T}_j$,

$$\begin{aligned} & \sup_{y \in \mathcal{Y}} \left| \frac{1}{n^{j+1}} \sum_{i_0=1}^n \dots \sum_{i_j=1}^n k_{i_0 s}^{(s+j)} \mathbf{K}_{ai_0 s} \zeta_{i_0}^* \{ \mathbf{I}_{i_0}(y) - \mathcal{E}(y) \} \prod_{\ell=1}^j \mathbf{K}_{z_{i_0 \ell}} \mathbf{u}_{\ell_i}^{*t_\ell} (\zeta_{\ell_i} - \zeta_{\ell_{i_0}})^{1-t_\ell} \right| \\ & < 1/\rho_n, \quad (35) \end{aligned}$$

where $\zeta_{\ell_i} = \zeta_\ell(\mathbf{z}_i)$ and similarly for other ζ symbols.

Proof. As will become apparent in lemma F6, for every j the LHS in (35) corresponds to the j -th term in a Taylor expansion of $\hat{\mathbf{S}}_s(y; \hat{p}) - \bar{\mathbf{S}}_s(y; \hat{p})$ around $\hat{\mathbf{S}}_s(y; p) - \bar{\mathbf{S}}_s(y; p)$; see (34). Because by lemma E3 $\hat{\mathbf{p}}_i - \mathbf{p}_i$ converges faster (uniformly in i) than the extra $1/h_s$ incurred for each additional derivative, the convergence rate is slowest for $s = j = 1$, so we establish convergence at the promised rate for that case; all other cases can be verified similarly, albeit sometimes more painfully.

Thus, we use lemma C4 to obtain a rate for

$$\sup_{y \in \mathcal{Y}} \left| \frac{1}{n^2} \sum_{i_0=1}^n \sum_{i_1=1}^n k''_{i_0 1} \mathbf{K}_{ai_0 1} \zeta_{i_0}^* \{ \mathbf{I}_{i_0}(y) - \mathcal{E}(y) \} \mathbf{K}_{z_{i_0 i_1}} \mathbf{u}_{1 i_1}^{*t_1} (\zeta_{1 i_1} - \zeta_{1 i_0})^{1-t_1} \right|. \quad (36)$$

Let ξ_i contain all random variables pertaining to observation i . Noting that (36) is a V statistic and that lemma C4 is based on a decomposition of the V statistic into a sum of U statistics, we have

for the m -symbols of lemma C4 and for some $\bar{\xi}, \tilde{\xi} \in \mathcal{F}_0$,

$$m^{(2,2)}(\xi_{i_0}, \xi_{i_1}) = \frac{1}{2n^2} \left[k''_{i_0 1} \mathbf{K}_{ai_0 1} \xi_{i_0}^* \{I_{i_0}(y) - \mathcal{E}(y)\} \mathbf{K}_{z i_0 i_1} \mathbf{u}_{1 i_1}^{* t_1} (\xi_{1 i_1} - \xi_{1 i_0})^{1-t_1} \right. \\ \left. + k''_{i_1 1} \mathbf{K}_{ai_1 1} \xi_{1 i_1}^* \{I_{i_1}(y) - \mathcal{E}(y)\} \mathbf{K}_{z i_0 i_1} \mathbf{u}_{1 i_0}^{* t_1} (\xi_{1 i_0} - \xi_{1 i_1})^{1-t_1} \right], \\ m_1^{(2,2)}(\xi_i) = \mathbb{E} \{ m^{(2,2)}(\xi_i, \xi_{i_1}) | \xi_i \} \quad (i_1 \neq i)$$

$$\simeq \begin{cases} \frac{h_z^2}{2n^2} \left[k''_{i_1} \mathbf{K}_{ai_1} \xi_i^* \{I_i(y) - \mathcal{E}(y)\} \bar{\xi}_i, \right. \\ \left. + k''_{i_1} \mathbf{K}_{ai_1} \{ \mathcal{E}(y; \mathbf{p}_i, \mathbf{a}_i) - \mathcal{E}(y; v^*, a^*) \} \tilde{\xi}_i \right], & t_1 = 0, \\ \frac{h_z^2}{2n^2} k''_{i_1} \mathbf{K}_{ai_1} \xi_i^* \{I_i(y) - \mathcal{E}(y)\} \bar{\xi}_i \mathbf{u}_{i_1}^*, & t_1 = 1, \end{cases} \\ m_1^{(2,1)}(\xi_i) = \begin{cases} 0, & t_1 = 0, \\ \frac{1}{n^2} k''_{i_1} \mathbf{K}_{ai_1} \xi_i^* \{I_i(y) - \mathcal{E}(y)\} \mathbf{K}_{z i i} \mathbf{u}_{i_1}^*, & t_1 = 1. \end{cases}$$

$$n^2 \mu^{(2,2)} \leq h_z^2, \quad (37)$$

$$n \mu^{(2,1)} \leq \frac{1}{n h_z^{d_z}}, \quad (38)$$

$$\sqrt{n} \beta_2^{(2,2)} \leq \frac{1}{n^{3/2} h_z^{d_z} h_1^{3+d_a}}, \quad (39)$$

$$\beta_1^{(2,2)} \leq \frac{h_z^2}{n^2 h_1^{3+d_a}}, \quad (40)$$

$$\beta_1^{(2,1)} \leq \frac{1}{n^2 h_z^{d_z} h_1^{3+d_a}}, \quad (41)$$

$$n \sigma_2^{(2,2)} \leq \frac{1}{n h_1^{3+d_a}}, \quad (42)$$

$$\sqrt{n} \sigma_1^{(2,2)} \leq \frac{h_z^2}{n^{3/2} h_1^{(5+d_a)/2}}, \quad (43)$$

$$\sqrt{n} \sigma_1^{(2,1)} \leq \frac{1}{n^{3/2} h_1^{(5+d_a)/2} h_z^{d_z}}. \quad (44)$$

Sufficient conditions for (37) to (44) to converge at a rate faster than $1/\rho_n$ are respectively

$$(3 + d_a) \eta_1 > 1 - 4 \eta_z, \quad (45)$$

$$(3 + d_a) \eta_1 > 2 d_z \eta_z - 1, \quad (46)$$

$$(3 + d_a) \eta_1 < 2 - 2 \eta_z d_z, \quad (47)$$

$$(3 + d_a)\eta_1 < 3 + 4\eta_z, \quad (48)$$

$$(3 + d_a)\eta_1 < 3 - 2\eta_z d_z, \quad (49)$$

$$(3 + d_a)\eta_1 < 1, \quad (50)$$

$$\eta_1 < 1 + 2\eta_z, \quad (51)$$

$$\eta_1 < 1 - \eta_z d_z. \quad (52)$$

Conditions (48) to (52) are implied by (47) and/or standard kernel estimation conditions needed for consistency of the estimator of H without nuisance parameters. Thus, only (45) to (47) are potentially relevant and the lemma statement holds if

$$\max(1 - 4\eta_z, 2d_z\eta_z - 1, 0) < (3 + d_a)\eta_1 < \min(2 - 2\eta_z d_z, 1),$$

which is satisfied by assumption S. \square

Lemma F2. *Using essentially the same notation and conditions as in lemma F1, for $j = 1, 2, \dots$,*

$$\left| \frac{1}{n^{j+1}} \sum_{i_0=1}^n \cdots \sum_{i_j=1}^n \mathbf{k}_{i_0}^{(j)} \mathbf{K}_{ai_0} \prod_{\ell=1}^j \mathbf{K}_{zi_0 i_\ell} \mathbf{u}_{i_\ell}^{*t_\ell} (\xi_{i_\ell} - \xi_{i_0})^{1-t_\ell} \right| < 1/\rho_n. \quad (53)$$

Proof. This lemma is used in lemma F7 to deal with the expansion of $\bar{\mathbf{S}}_0(y; \hat{p})$ around $\bar{\mathbf{S}}_0(y; p)$.

Using the same strategy and rationale for focusing on the case $s = 1$ as in lemma F1, we have

$$m^{(2,2)}(\xi_{i_0}, \xi_{i_1}) = \frac{1}{2n^2} \left[\mathbf{k}'_{i_0} \mathbf{K}_{ai_0} \mathbf{K}_{zi_0 i_1} \mathbf{u}_{i_1}^{*t} (\xi_{i_1} - \xi_{i_0})^{1-t} + \mathbf{k}'_{i_1} \mathbf{K}_{ai_1} \mathbf{K}_{zi_0 i_1} \mathbf{u}_{i_0}^{*t} (\xi_{i_0} - \xi_{i_1})^{1-t} \right],$$

$$m_1^{(2,2)}(\xi_i) \simeq \begin{cases} \frac{h_z^2}{n^2} \mathbf{k}'_{i_0} \mathbf{K}_{ai_0} \bar{\xi}_i, & t = 0, \\ \frac{1}{2n^2} \mathbf{k}'_{i_0} \mathbf{K}_{ai_0} \bar{\xi}_i \mathbf{u}_{i_1}^*, & t = 1. \end{cases}$$

$$m_1^{(2,1)}(\xi_i) = \begin{cases} 0, & t = 0, \\ \frac{1}{n^2} \mathbf{k}'_{i_0} \mathbf{K}_{ai_0} \mathbf{K}_{zii} \mathbf{u}_{i_1}^*, & t = 1, \end{cases}$$

$$n^2 \mu^{(2,2)} \leq h_z^2, \quad (54)$$

$$n \mu^{(2,1)} \leq \frac{1}{n h_z^{d_z}}, \quad (55)$$

$$\sqrt{n} \beta_2^{(2,2)} \leq \frac{1}{n^{3/2} h_0^{d_a+2} h_z^{d_z}}, \quad (56)$$

$$\beta_1^{(2,2)} \leq \frac{1}{n^2 h_0^{d_a+2}}, \quad (57)$$

$$\beta_1^{(2,1)} \leq \frac{1}{n^2 h_0^{d_a+2} h_z^{d_z}}, \quad (58)$$

$$n\sigma_2^{(2,2)} \leq \frac{1}{n h_0^{d_a+2}}, \quad (59)$$

$$\sqrt{n}\sigma_1^{(2,2)} \leq \frac{1}{n^{3/2} h_0^{(d_a+3)/2}}, \quad (60)$$

$$\sqrt{n}\sigma_1^{(2,1)} \leq \frac{1}{n^{3/2} h_0^{(d_a+3)/2}}. \quad (61)$$

Some of the numbered equations above were already shown to be $< 1/\rho_n$ in lemma F1. The remaining conditions are implied by

$$(3 + d_a)\eta_1 > 2\eta_0(d_a + 2) + 2\eta_z d_z - 2, \quad (62)$$

$$(3 + d_a)\eta_1 > 2\eta_0(d_a + 2) - 1, \quad (63)$$

$$(3 + d_a)\eta_1 > (d_a + 3)\eta_0 - 2, \quad (64)$$

which follow from assumption S. \square

Lemma F3. *Using the same notation as in lemma F1, for some $0 < J < \infty$ and $s = 0, 1$,*

$$\sup_{y \in \mathcal{Y}} \frac{1}{n h_s^{s+1+J}} \sum_{i=1}^n |\mathbf{K}_{ais} \xi_i^* \{I_i(y) - \mathcal{E}(y)\}| \sup_{z \in \mathcal{Z}} |\hat{p}(z) - p(z)|^J < 1/\rho_n.$$

Proof. By lemma E3, the stated result is implied by α_n/h_s decreasing to zero at a polynomial rate since for any polynomial π_n^* then $\pi_n^*(\alpha_n/h_s)^J < 1$. The requirement that α_n/h_s decrease to zero at a polynomial rate is guaranteed by $\eta_s < \min\{2\eta_z, (1 - \eta_z d_z)/2\}$, which was assumed in assumption S. \square

Lemma F4.

$$\sup_{z \in \mathcal{Z}} \left| \hat{p}(z) - p(z) - \frac{\hat{r}(z) - r(z)}{f(z)} \sum_{j=0}^{J-2} \left(\frac{f(z) - \hat{f}(z)}{f(z)} \right)^j - p(z) \sum_{j=1}^{J-1} \left(\frac{f(z) - \hat{f}(z)}{f(z)} \right)^j \right| \leq \alpha_n^J.$$

Proof. From the recursion of $f/\hat{f} = 1 + (f - \hat{f})/\hat{f}$, we have $f/\hat{f} = \sum_{j=0}^{J-1} \{(f - \hat{f})/\hat{f}\}^j + \{(f - \hat{f})/\hat{f}\}^J f/\hat{f}$. Therefore, the LHS in the lemma statement is bounded by

$$\sup_{z \in \mathcal{Z}} \left| \left(\frac{\hat{r}(z) - r(z)}{f(z)} \right) \left(\frac{f(z) - \hat{f}(z)}{f(z)} \right)^{J-1} \frac{f(z)}{\hat{f}(z)} + \left(\frac{f(z) - \hat{f}(z)}{f(z)} \right)^J \frac{f(z)}{\hat{f}(z)} \right|.$$

Apply lemmas E3 and E4. \square

Let $\bar{f}(z) = \mathbb{E} \hat{f}(z)$, $\bar{r}(z) = \mathbb{E} \hat{r}(z)$, $\hat{f}^* = \hat{f}/f$, $\hat{r}^* = \hat{r}/f$, $\bar{r}^* = \bar{r}/f$, $\bar{f}^* = \bar{r}/f$.

Lemma F5. For given j, J , let Λ_j be the collection of vectors ℓ of dimension four containing non-negative integers satisfying $\ell_1 + \ell_2 \leq j$ and $\ell_1 + \ell_2 + \ell_3 + \ell_4 < J$. Then for all sufficiently large n and any $1 \leq j < J$ and some constants $C_{j\ell}$ independent of n, z ,

$$\sup_{z \in \mathcal{Z}} \left| \{\hat{p}(z) - p(z)\}^j - \sum_{\ell \in \Lambda_j} C_{j\ell} \{\hat{r}^*(z) - \bar{r}^*(z)\}^{\ell_1} \{\bar{r}^*(z) - p(z)\}^{\ell_2} \{\hat{f}^*(z) - \bar{f}^*(z)\}^{\ell_3} \{\bar{f}^*(z) - 1\}^{\ell_4} \right| \leq \alpha_n^J.$$

Proof. Follows directly from lemma F4 combined with the multinomial theorem. \square

Lemma F6. For $s = 0, 1$, $\sup_{y \in \mathcal{Y}} |\hat{\mathbf{S}}_s(y; \hat{p}) - \bar{\mathbf{S}}_s(y; \hat{p}) - \hat{\mathbf{S}}_s(y; p) + \bar{\mathbf{S}}_s(y; p)| < 1/\rho_n$.

Proof. Let J be sufficiently large as in lemma F3. Then, expand the LHS of the lemma statement to order J using the mean value theorem to obtain an upper bound of (34). The second term in (34) is covered by lemma F3 and the first term in (34) is dealt with in lemma F1, using lemmas F4 and F5. \square

Lemma F7. $\sup_{y \in \mathcal{Y}} |\bar{\mathbf{S}}_0(y; \hat{p}) - \bar{\mathbf{S}}_0(y; p)| < 1/\rho_n$.

Proof. The proof is entirely analogous to that of lemma F6, albeit using lemma F2 instead of lemma F1, and is hence omitted. \square

Below we will write $S_s(y; p)$ for $S_s(y) = S_s(y; a^*, v^*)$ for $s = 0, 1$.

Lemma F8.

$$\sup_{y \in \mathcal{Y}} |\hat{\mathbf{S}}_0(y; p) - S_0(y; p)| < 1/\rho_n. \quad (65)$$

Proof. By standard kernel estimation theory, the squared LHS in (65) is $\leq h_0^4 + 1/n h_0^{1+da} < 1/\rho_n^2$ by assumption S. \square

Lemma F9.

$$\sup_{y \in \mathcal{Y}} |\hat{\mathbf{S}}_0(y; \hat{p}) - S_0(y; p)| < 1/\rho_n. \quad (66)$$

Proof. The LHS in (66) is bounded above by the sum of

$$\sup_{y \in \mathcal{Y}} |\hat{\mathcal{S}}_0(y; \hat{p}) - \bar{\mathcal{S}}_0(y; \hat{p}) - \hat{\mathcal{S}}_0(y; p) + \bar{\mathcal{S}}_0(y; p)|, \quad (67)$$

$$\sup_{y \in \mathcal{Y}} |\bar{\mathcal{S}}_0(y; \hat{p}) - \bar{\mathcal{S}}_0(y; p)|. \quad (68)$$

$$\sup_{y \in \mathcal{Y}} |\hat{\mathcal{S}}_0(y; p) - \mathcal{S}_0(y; p)|, \quad (69)$$

Apply lemmas F6 to F8. □

Lemma F10. For all $y \in \mathcal{Y}$, $\{\bar{\mathcal{S}}_1(y; \hat{p}) - \bar{\mathcal{S}}_1(y; p)\} \mathcal{S}_0(\infty; p) = \{\bar{\mathcal{S}}_1(\infty; \hat{p}) - \bar{\mathcal{S}}_1(\infty; p)\} \mathcal{S}_0(y; p)$.

Proof. Trivial. □

Lemma F11. Let $\hat{H}(y; \hat{p}) = \hat{H}(y|a^*, v^*)$ and $H(y; p) = H(y|a^*, v^*)$. Then

$$\sup_{y \in \mathcal{Y}} \left| \hat{H}(y; \hat{p}) - H(y; p) - v^* \frac{\{\hat{\mathcal{S}}_1(y; p) - \mathcal{S}_1(y; p)\} \mathcal{S}_0(\infty; p) - \{\hat{\mathcal{S}}_1(\infty; p) - \mathcal{S}_1(\infty; p)\} \mathcal{S}_0(y; p)}{\mathcal{S}_0^2(\infty; p)} \right| < 1/\rho_n. \quad (70)$$

Proof. For the remainder of this lemma, let \simeq_ρ mean that the left and right hand sides differ by a term $< 1/\rho_n$, uniformly in y . By lemma F9,

$$\hat{H}(y; \hat{p}) \simeq_\rho \frac{\mathcal{S}_0(y; p) \mathcal{S}_0(\infty; p) + v^* \{\hat{\mathcal{S}}_1(y; \hat{p}) \mathcal{S}_0(\infty; p) - \hat{\mathcal{S}}_1(\infty; \hat{p}) \mathcal{S}_0(y; p)\}}{\mathcal{S}_0^2(\infty; p)}.$$

Since by lemma F6 $\hat{\mathcal{S}}_1(y; \hat{p}) \simeq_\rho \hat{\mathcal{S}}_1(y; p) + \bar{\mathcal{S}}_1(y; \hat{p}) - \bar{\mathcal{S}}_1(y; p)$, it follows from lemma F10 that

$$\hat{H}(y; \hat{p}) \simeq_\rho \frac{\mathcal{S}_0(y; p) \mathcal{S}_0(\infty; p) + v^* \{\hat{\mathcal{S}}_1(y; p) \mathcal{S}_0(\infty; p) - \hat{\mathcal{S}}_1(\infty; p) \mathcal{S}_0(y; p)\}}{\mathcal{S}_0^2(\infty; p)}. \quad (71)$$

Claim (70) then follows by subtracting and adding $\mathcal{S}_1(y; p)$ and $\mathcal{S}_1(\infty; p)$ in the numerator of (71). □

APPENDIX G. WEAK CONVERGENCE

Let $\hat{\mathcal{S}}_s(y) = \hat{\mathcal{S}}_s(y; p)$. We first show the weak convergence of $\hat{\mathcal{C}}_{ns}^*(\cdot) = \sqrt{nh_s^{2s+1+d_a}} \{\hat{\mathcal{S}}_s(\cdot) - \mathbb{E} \hat{\mathcal{S}}_s(\cdot)\}$ in $\mathcal{L}^\infty(\mathcal{I})$, where \mathcal{I} is an arbitrary compact subset of \mathbb{R} . Let $\omega_{nsc}(x, y, z, p) = w(z) \mathbb{1}(x = 0) \mathbb{1}(y \leq c) K \{(a^* - a)/h_s\} k^{(s)} \{(v^* - p)/h_s\} / \sqrt{h_s^{1+d_a}}$ and consider

$$\mathcal{F}_{ns} = \mathcal{F}_{ns}(\mathcal{I}) = \{(x, y, z, p) \mapsto \omega_{nsc}(x, y, z, p) : c \in \mathcal{I}\}.$$

Define \mathcal{E}_{ns} by $\mathcal{E}_{ns}(x, y, z, p) = |K\{(a^* - a)/h_s\}k^{(s)}\{(v^* - p)/h_s\}|/\sqrt{h_s^{1+d_a}}$ so that it is an envelope function of \mathcal{F}_{ns} . Below we will write $\mathcal{E}_{ns}(a, p)$ for $\mathcal{E}_{ns}(x, y, z, p)$ given that $\mathcal{E}_{ns}(x, y, z, p)$ depends only on a, p .

Lemma G1. For $s = 0, 1$, $\mathbb{E}\mathcal{E}_{ns}^2(\mathbf{a}_i, \mathbf{p}_i) \leq 1$. Also, for any $\epsilon > 0$, $\mathbb{E}[\mathcal{E}_{ns}^2(\mathbf{a}_i, \mathbf{p}_i)\mathbb{1}\{\mathcal{E}_{ns}(\mathbf{a}_i, \mathbf{p}_i) > \epsilon\sqrt{n}\}] < 1$.

Proof. The first statement follows from a change of variables and assumption R. The second statement follows from $\mathbb{1}\{\mathcal{E}_{ns}(a, p) > \epsilon\sqrt{n}\} \leq \mathbb{1}\{\sup_{t_1, t_2} |K(t_1)k^{(s)}(t_2)| > \epsilon\sqrt{nh_s^{1+d_a}}\} = 0$ for sufficiently large n by assumption S. \square

Lemma G2. For any $\delta_n < 1$ and $s = 0, 1$,

$$\sup_{|c-c^*| \leq \delta_n} \mathbb{E} \left[\mathbf{I}_{xi} \{ \mathbb{1}(y_i \leq c) - \mathbb{1}(y_i \leq c^*) \} K \left(\frac{a^* - a_i}{h_s} \right) k^{(s)} \left(\frac{v^* - p_i}{h_s} \right) \right]^2 / h_s^{1+d_a} < 1.$$

Proof. The LHS of the lemma statement is bounded by twice of

$$\begin{aligned} \sup_{|c-c^*| \leq \delta_n} \mathbb{E} \left[\mathbb{1}\{\min(c, c^*) < y_i \leq \max(c, c^*)\} w(z_i) K \left(\frac{a^* - a_i}{h_s} \right) k^{(s)} \left(\frac{v^* - p_i}{h_s} \right) \right]^2 / h_s^{1+d_a} \\ \leq C \delta_n \sup_{y, z, p} f_{yzp}(y, z, p) < 1, \end{aligned}$$

where C is a constant and f_{yzp} is the density of y_i, z_i, p_i . \square

Lemma G3. For $s = 0, 1$, \mathcal{F}_{ns} is a Vapnik–Cervonenkis (VC) class with VC index uniformly bounded in n .

Proof. Let $\mathcal{J} = \{y \mapsto \mathbb{1}(y \leq c) : c \in \mathcal{I}\}$ and let $\chi_{ns}(x, z, p) = \mathbb{1}(x = 0)w(z)K\{(a^* - a)/h_s\}k^{(s)}\{(v^* - p)/h_s\}/\sqrt{h_s^{1+d_a}}$. Then, by van der Vaart and Wellner (1996, lemma 2.6.18), the VC index of $\mathcal{F}_{ns} = \chi_{ns} \cdot \mathcal{J} = \{\chi_{ns}\bar{J} : \bar{J} \in \mathcal{J}\}$ is bounded by the VC index of \mathcal{J} times 2 minus 1. Therefore, the VC index of \mathcal{F}_{ns} is bounded and independent of n , because \mathcal{J} is a VC class that does not depend on n . \square

Lemma G4. For $s = 0, 1$, $\hat{\mathbf{G}}_{ns}^* \xrightarrow{w} \mathbf{G}_s^*$ in $\mathcal{L}^\infty(\mathbb{R})$, where \mathbf{G}_s^* is a mean-zero Gaussian process.

Proof. Convergence of finite marginals easily follows by a central limit theorem. Now, for $\ell = 1, 2$, let $\mathcal{F}_{ns, \delta}^\ell$ be a set defined by

$$\left\{ (x, y, z, p) \mapsto \{ \omega_{nsc}(x, y, z, p) - \omega_{nsc^*}(x, y, z, p) \}^\ell : |c - c^*| < \delta, \omega_{nsc}, \omega_{nsc^*} \in \mathcal{F}_{ns}(\mathcal{I}) \right\}$$

Since $\{\mathbb{1}(z \in \mathcal{Z}, x = 0, y \leq c) - \mathbb{1}(z \in \mathcal{Z}, x = 0, y \leq c^*)\}^\ell$ is left- or right-continuous for every c, c^* and since \mathcal{I} is separable, $\mathcal{F}_{ns,\delta}^\ell$ contains a countable subclass $\mathcal{G}_{ns,\delta}^\ell$ such that for every $\chi \in \mathcal{F}_{ns,\delta}^\ell$ there exists a sequence $\{\chi_j\} \subset \mathcal{G}_{ns,\delta}^\ell$ with $\chi_j(x, y, z, p) \rightarrow \chi(x, y, z, p)$. Therefore, by the same reasoning as [van der Vaart and Wellner \(1996, example 2.3.4\)](#), $\mathcal{F}_{ns,\delta}^\ell$ for $\ell = 1, 2$ is a measurable class for every $\delta > 0$. Therefore, it follows from lemmas [G1](#) to [G3](#) and [van der Vaart and Wellner \(1996, theorem 2.11.22\)](#) that $\hat{\mathbf{G}}_{ns}^* \xrightarrow{w} \mathbf{G}_s^*$ in $\mathcal{L}^\infty(\mathcal{I})$. Since \mathcal{I} is an arbitrary compact set in \mathbb{R} , we know by [van der Vaart and Wellner \(1996, theorem 1.6.1\)](#) that $\hat{\mathbf{G}}_{ns}^* \xrightarrow{w} \mathbf{G}_s^*$ in $\mathcal{L}^\infty(\mathcal{I}_1, \mathcal{I}_2, \dots)$, where $\{\mathcal{I}_j\}$ is an increasing sequence of compact sets such that $\cup_j \mathcal{I}_j = \mathbb{R}$. Finally note that for all n , $\hat{\mathbf{G}}_{ns}^*, \mathbf{G}_s^* \in \mathcal{L}^\infty(\mathbb{R}) \subset \mathcal{L}^\infty(\mathcal{I}_1, \mathcal{I}_2, \dots)$. \square

Lemma G5. $\hat{\mathbf{G}}_{n1}^*(\cdot) - \hat{\mathbf{G}}_{n1}^*(\infty)G(\cdot|a^*, v^*) \xrightarrow{w} \mathbf{G}^*(\cdot)$ in $\mathcal{L}^\infty(\mathbb{R})$, where \mathbf{G}^* is a mean-zero Gaussian process with the covariance kernel given by \mathcal{C} in [\(20\)](#).

Proof. By a central limit theorem, $\hat{\mathbf{G}}_{n1}^*(\infty)G(y|a^*, v^*) \xrightarrow{d} \Psi G(y|a^*, v^*)$ for a mean-zero normal random variable Ψ . Since $G(\cdot|a^*, v^*)$ is uniformly continuous in \mathcal{I} , where \mathcal{I} is an arbitrary compact set in \mathbb{R} , we have $\hat{\mathbf{G}}_{n1}^*(\infty)G(\cdot|a^*, v^*) \xrightarrow{w} \Psi G(\cdot|a^*, v^*)$ in $\mathcal{L}^\infty(\mathcal{I})$. Therefore, by [van der Vaart and Wellner \(1996, theorem 1.6.1\)](#), we have $\hat{\mathbf{G}}_{n1}^*(\infty)G(\cdot|a^*, v^*) \xrightarrow{w} \Psi G(\cdot|a^*, v^*)$ in $\mathcal{L}^\infty(\mathcal{I}_1, \mathcal{I}_2, \dots)$, where $\{\mathcal{I}_j\}$ is an increasing sequence of compact sets such that $\cup_j \mathcal{I}_j = \mathbb{R}$. Now note that for all n , $\hat{\mathbf{G}}_{n1}^*(\infty)G(\cdot|a^*, v^*)$ and $\Psi G(\cdot|a^*, v^*)$ are in $\mathcal{L}^\infty(\mathbb{R}) \subset \mathcal{L}^\infty(\mathcal{I}_1, \mathcal{I}_2, \dots)$ and the lemma statement follows from the continuous mapping theorem. \square

Lemma G6. For $s = 0, 1$,

$$\sup_{y \in \mathcal{Y}} \left| \mathbb{E} \left\{ \mathbf{K}_{ais} \mathbf{k}_{is}^{(s)} \mathbf{I}_i(y) \right\} - \partial_v^s S_0(y) - \frac{h_s^2 \kappa_2}{2} \text{tr} \left\{ \partial_v^s \partial_{bb^\top} S_0(y; a^*, v^*) \right\} \right| < \frac{1}{\rho n}.$$

Proof. This is nothing but a standard kernel bias expansion after noting that $\mathbb{E}\{\mathbf{I}_i(y)|\mathbf{a}_i = a, \mathbf{p}_i = p\} = S_0(y; a, p)/f_{ap}(a, p)$. \square

APPENDIX H. SEMIPARAMETRIC ESTIMATION

Let $\hat{\mathbf{p}}(z) = \hat{\mathbf{r}}_L(z)/\hat{\mathbf{f}}_L(z)$ and redefine $\hat{\mathbf{p}}(z) = \hat{\mathbf{r}}_L(z)/\hat{\mathbf{f}}_L(z)$ (compare with [\(15\)](#)), where

$$\begin{cases} \hat{\mathbf{f}}_L(z) = n^{-1} \sum_{i=1}^n \hat{\mathbf{K}}_{zLi}, & \hat{\mathbf{r}}_L(z) = n^{-1} \sum_{i=1}^n \hat{\mathbf{K}}_{zLi} \mathbb{1}(x_i = 0), \\ \hat{\mathbf{f}}_L(z) = n^{-1} \sum_{i=1}^n \mathbf{K}_{zLi}, & \hat{\mathbf{r}}_L(z) = n^{-1} \sum_{i=1}^n \mathbf{K}_{zLi} \mathbb{1}(x_i = 0), \end{cases}$$

with $\hat{\mathbf{K}}_{zLi} = \hat{\mathbf{K}}_{zLi}(z) = K\{(z - z_i)^\top \hat{\boldsymbol{\gamma}}/h_z\}/h_z$ and $\mathbf{K}_{zLi} = K\{(z - z_i)^\top \boldsymbol{\gamma}_0/h_z\}/h_z$. $\hat{\mathbf{K}}_{aLi}$ and \mathbf{K}_{aLi} are similarly defined.

Lemma H1. *Let*

$$\tilde{\mu}(z) = \frac{1}{\gamma_{01}^2 f_L(z)} \int \partial_{z_1} p\left(\frac{z^\top \boldsymbol{\gamma}_0 - \tilde{t}^\top \tilde{\boldsymbol{\gamma}}_0}{\gamma_{01}}, \tilde{t}\right) f\left(\frac{z^\top \boldsymbol{\gamma}_0 - \tilde{t}^\top \tilde{\boldsymbol{\gamma}}_0}{\gamma_{01}}, \tilde{t}\right) \left(\frac{z^\top \boldsymbol{\gamma}_0 - \tilde{t}^\top \tilde{\boldsymbol{\gamma}}_0}{\gamma_{01}}, \tilde{t}\right) d\tilde{t}$$

Then

$$\sup_{z \in \mathcal{Z}} \left| \hat{\boldsymbol{\rho}}(z) - \hat{\boldsymbol{\rho}}(z) - \tilde{\mu}^\top(z)(\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_0) \right| \leq \frac{1}{\pi_n^* \sqrt{n}},$$

for some π_n^* increasing as fractional power of n .

Proof. Note that both $\hat{\mathbf{f}}_L(z) - \hat{\mathbf{f}}_L(z)$ and $\hat{\mathbf{r}}_L(z) - \hat{\mathbf{r}}_L(z)$ can be expanded as

$$\begin{aligned} & \sum_{j=1}^{J-1} \frac{1}{n j!} \sum_{i=1}^n \mathbf{K}_{zLi}^{(j)} \boldsymbol{\xi}_i \{(z - z_i)^\top (\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_0)\}^j \\ & + \frac{1}{n J! h_z^{J+1}} \sum_{i=1}^n K_{zLi}^{(J)}(\cdot) \boldsymbol{\xi}_i \{(z - z_i)^\top (\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_0)\}^J, \quad (72) \end{aligned}$$

for bounded $\boldsymbol{\xi}_i$.

The J -th order term in (72) is of order $n^{-J/2} h_z^{-J-1}$ (uniformly in $z \in \mathcal{Z}$) which, for sufficiently large J , is $< 1/\sqrt{n}$. For $1 \leq j \leq J-1$, note that for any $1 \leq \ell_1, \dots, \ell_j \leq d_z$,

$$\sup_{z \in \mathcal{Z}} \left| \mathbb{E} \left\{ \mathbf{K}_{zLi}^{(j)} \boldsymbol{\xi}_i \prod_{t=1}^j (z_{\ell_t} - z_{i\ell_t}) \right\} \right| \leq \frac{1}{h_z^{\max(0, j-2)}}, \quad \sup_{z \in \mathcal{Z}} \left| \mathbb{V} \left\{ \mathbf{K}_{zLi}^{(j)} \boldsymbol{\xi}_i \prod_{t=1}^j (z_{\ell_t} - z_{i\ell_t}) \right\} \right| \leq \frac{1}{h_z^{2j+1}}.$$

Consequently, analogous to lemma E2, the j -th term in (72) is for $j \geq 2$ of order no greater than

$$\frac{1}{n^{j/2}} \max \left(\frac{1}{h_z^{\max(0, j-2)}}, \frac{\log n}{\sqrt{n} h_z^{2j+1}} \right) < \frac{1}{\sqrt{n}},$$

uniformly in $z \in \mathcal{Z}$.

Finally for $j = 1$ note that by standard kernel theory for any function $\mu \in \mathcal{F}_1$,

$$\sup_{z \in \mathcal{Z}} \left| \frac{1}{h_z^2} \int k' \left(\frac{(z-t)^\top \boldsymbol{\gamma}_0}{h_z} \right) \mu(t) dt - \frac{1}{\gamma_{01}^2} \int \partial_{z_1} \mu \left(\frac{z^\top \boldsymbol{\gamma}_0 - \tilde{t}^\top \tilde{\boldsymbol{\gamma}}_0}{\gamma_{01}}, \tilde{t} \right) d\tilde{t} \right| \leq h_z^2.$$

Hence it follows that as in lemma F4

$$\sup_{z \in \mathcal{Z}} \left| \hat{\mathbf{p}}(z) - \hat{\mathbf{p}}(z) - \tilde{\boldsymbol{\mu}}^\top(z)(\hat{\boldsymbol{y}} - \boldsymbol{\gamma}_0) \right| \simeq$$

$$\sup_{z \in \mathcal{Z}} \left| \frac{\{\hat{\mathbf{r}}_L(z) - \hat{\mathbf{r}}_L(z)\} - p(z)\{\hat{\mathbf{f}}_L(z) - \hat{\mathbf{f}}_L(z)\}}{f_L(z)} - \tilde{\boldsymbol{\mu}}^\top(z)(\hat{\boldsymbol{y}} - \boldsymbol{\gamma}_0) \right| \leq \frac{1}{\pi_n^* \sqrt{n}},$$

for some π_n^* increasing as a fractional power of n . \square

Lemma H2. Let $\boldsymbol{\xi}_i(y)$ be of the form $\zeta^*(\mathbf{z}_i)\mathbf{I}_i(y) + \zeta^{**}(\mathbf{z}_i)\boldsymbol{\mathcal{E}}(y)$ with ζ^*, ζ^{**} depending only on \mathbf{z}_i and be such that $\mathbb{E}\{\boldsymbol{\xi}_i(y)|\mathbf{z}_i = z\} \in \mathcal{F}_2$. Then for $s = 0, 1$,

$$\sup_{y \in \mathcal{Y}} \left| \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{k}}_{is}^{(s)} \hat{\mathbf{K}}_{aLi} - \hat{\mathbf{k}}_{is}^{(s)} \mathbf{K}_{aLi}) \boldsymbol{\xi}_i(y) \right| < 1/\rho_n.$$

Proof. We have to deal both with the presence of $\hat{\mathbf{p}}_i$ in lieu of $\hat{\mathbf{p}}_i$ and with $\hat{\boldsymbol{\theta}}^\top \mathbf{a}_i$ in lieu of $\theta_0^\top \mathbf{a}_i$.

Since the former is more difficult than the latter, we shall establish below that

$$\sup_{y \in \mathcal{Y}} \left| \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{k}}_{is}^{(s)} - \hat{\mathbf{k}}_{is}^{(s)}) \mathbf{K}_{aLi} \boldsymbol{\xi}_i(y) \right| < 1/\rho_n, \quad (73)$$

where the remaining results can be established similarly but more simply. We again use lemma B3 repeatedly. By the mean value theorem, the LHS average in (73) can be expanded as

$$\sum_{j=1}^{J-1} \frac{1}{n j!} \sum_{i=1}^n \hat{\mathbf{k}}_{is}^{(s+j)} \mathbf{K}_{aLi} \boldsymbol{\xi}_i(y) (\hat{\mathbf{p}}_i - \hat{\mathbf{p}}_i)^j$$

$$+ \frac{1}{n J! h_s^{s+1+J}} \sum_{i=1}^n k^{(s+J)}(\cdot) \mathbf{K}_{aLi} \boldsymbol{\xi}_i(y) (\hat{\mathbf{p}}_i - \hat{\mathbf{p}}_i)^J. \quad (74)$$

The last term in (74) is of order $h_s^{-s-1-J} n^{-J/2}$ (uniformly in y) which, for sufficiently large J , is $< 1/\rho_n$. Further, for $1 \leq j < J$, we expand $\hat{\mathbf{k}}_{is}^{(s+j)}$ around $(v^* - \mathbf{p}_i)/h_s$ to obtain

$$\frac{1}{n} \sum_{i=1}^n \hat{\mathbf{k}}_{is}^{(s+j)} \mathbf{K}_{aLi} \boldsymbol{\xi}_i(y) (\hat{\mathbf{p}}_i - \hat{\mathbf{p}}_i)^j = \sum_{j^*=0}^{J^*-1} \frac{1}{n j^*!} \sum_{i=1}^n \mathbf{k}_{is}^{(s+j+j^*)} \mathbf{K}_{aLi} \boldsymbol{\xi}_i(y) (\hat{\mathbf{p}}_i - \hat{\mathbf{p}}_i)^j (\hat{\mathbf{p}}_i - \mathbf{p}_i)^{j^*}$$

$$+ \frac{1}{n J^*! h_s^{s+1+j+J^*}} \sum_{i=1}^n k^{(s+j+J^*)}(\cdot) \mathbf{K}_{ai} \boldsymbol{\xi}_i(y) (\hat{\mathbf{p}}_i - \hat{\mathbf{p}}_i)^j (\hat{\mathbf{p}}_i - \mathbf{p}_i)^{J^*} \quad (75)$$

The last term in (75) is $< 1/\rho_n$ (uniformly in y) for sufficiently large J^* because α_n/h_s vanishes as a (fractional) power of n . For $j + j^* > 1$ the terms in the RHS sum in (75) are of order $h_s^{-s-j-j^*} n^{-j/2} \alpha_n^{j^*} < h_s^{-s-j} n^{-j/2} \leq 1/\rho_n$.

For $j = 0$ the expansion in (75) is redundant, which leaves the case $j = 1, j^* = 0$. Thus, we must analyze

$$\frac{1}{n} \sum_{i=1}^n k_{is}^{(s+1)} \mathbf{K}_{aLi} \xi_i(y) (\hat{p}_i - \hat{p}_i),$$

which by lemma H1 and standard kernel estimation arguments equals

$$\frac{1}{n^{3/2}} \sum_{i=1}^n k_{is}^{(s+1)} \mathbf{K}_{aLi} \xi_i(y) \tilde{\mu}(z_i) \quad (76)$$

plus a term of order $1/\rho_n \pi_n^* \prec 1/\rho_n$ (uniformly in y). Finally, (76) is of order $1/\sqrt{n}$, uniformly in y . \square

APPENDIX I. PROOFS OF THEOREMS

Proof of Theorem 1. Part (i) follows from lemma A2 and assumption E. For (ii), please recall that the LHS in (6) was shown to be identified in JPX11. \square

Proof of Theorem 2. It follows from lemma A1 and assumption I that $g\{x^*, \mathbb{Q}_{\mathbf{u}|v}(\tau^*|V_t)\} \rightarrow \psi^*$. Identification of $g\{x^*, \mathbb{Q}_{\mathbf{u}|v}(\tau^*|V_t)\}$ follows from the fact that $V_t \in \mathcal{D}(x^*)$. \square

Proof of Theorem 3. It follows from (10) and the monotonicity of g . \square

Proof of Theorem 4. By lemma F11,

$$\begin{aligned} & \rho_n \{ \hat{\mathbf{H}}(y|a^*, v^*) - H(y|a^*, v^*) \} \\ & \simeq \rho_n v^* \frac{ \{ \hat{\mathbf{S}}_1(y; p) - S_1(y; p) \} - \{ \hat{\mathbf{S}}_1(\infty; p) - S_1(\infty; p) \} G(y; a^*, v^*) }{ S_0(\infty; p) } \\ & = \frac{v^*}{S_{0x}} \left[\{ \hat{\mathbf{G}}_{n1}^*(y) - \hat{\mathbf{G}}_{n1}^*(\infty) G(y|a^*, v^*) \} \right. \\ & \quad \left. + \{ \mathbb{E} \hat{\mathbf{S}}_1(y; p) - S_1(y; p) \} - \{ \mathbb{E} \hat{\mathbf{S}}_1(\infty; p) - S_1(\infty; p) \} G(y|a^*, v^*) \right]. \end{aligned}$$

The stated result then follows from lemmas G5 and G6. \square

Proof of Theorem 5. Let \mathcal{D} be a collection of CADLAG functions and define a mapping $T : \mathcal{D} \times \mathcal{U} \rightarrow \mathbb{R}$ such that for $F^* \in \mathcal{D}$ and $\tau \in \mathcal{U}$, $T(F^*, \tau) = \inf\{\tilde{\psi} \in \mathcal{Y} : F^*(\tilde{\psi}) \geq \tau\}$. We then have $\hat{\psi}^* = T\{\hat{\mathbf{H}}(\cdot|a^*, v^*), \tau^*\}$ and $\psi^* = T\{H(\cdot|a^*, v^*), \tau^*\}$. We now use the functional delta-method; see e.g. Van der Vaart (2000, theorem 20.8). In particular, by Van der Vaart (2000, lemma 21.3), $T(\cdot, \tau^*)$ is Hadamard-differentiable at $H(\cdot|a^*, v^*)$ tangentially to the set of functions $F \in \mathcal{D}$

that are continuous at ψ^* with derivative $T_H(F, \tau^*) = -F(\psi^*)/H'(\psi^*|a^*, v^*)$. Therefore, by the functional delta–method and theorem 4, we have

$$\rho_n [T\{\hat{H}(\cdot|a^*, v^*), \tau^*\} - T\{H(\cdot|a^*, v^*), \tau^*\}] \xrightarrow{d} T_H(\mathbb{G}, \tau^*) = -\frac{\mathbb{G}(\psi^*)}{H'(\psi^*|a^*, v^*)}. \quad \square$$

Proof of Theorem 6. First consider $\hat{\mathcal{C}}$. We have

$$\begin{aligned} \hat{\mathcal{C}}(y, y) &= \frac{1}{nh_1^{1+d_a}} \sum_{i=1}^n K^2\left(\frac{a^* - a_i}{h_1}\right) k^2\left(\frac{v^* - \hat{p}_i}{h_1}\right) \mathbf{I}_{x_i} \mathbf{I}_i(y) \\ &\quad - \frac{2}{nh_1^{1+d_a}} \sum_{i=1}^n K^2\left(\frac{a^* - a_i}{h_1}\right) k^2\left(\frac{v^* - \hat{p}_i}{h_1}\right) \mathbf{I}_{x_i} \mathbf{I}_i(y) \frac{\hat{S}_{0(y)}}{\hat{S}_{0x}} \\ &\quad + \frac{1}{nh_1^{1+d_a}} \sum_{i=1}^n K^2\left(\frac{a^* - a_i}{h_1}\right) k^2\left(\frac{v^* - \hat{p}_i}{h_1}\right) \mathbb{1}(x_i = 0) w_i^2 \frac{\hat{S}_{0(y)}^2}{\hat{S}_{0x}^2}. \end{aligned} \quad (77)$$

By lemma F9 and as in the proof of theorem 4, we have $\sup_{y \in \mathcal{Y}} |\hat{S}_{0(y)}/\hat{S}_{0x} - S_0(y)/S_{0x}| < 1$ and

$$\begin{aligned} \sup_{y \in \mathcal{Y}} \left| \frac{1}{nh_1^{1+d_a}} \sum_{i=1}^n K^2\left(\frac{a^* - a_i}{h_1}\right) k^2\left(\frac{v^* - \hat{p}_i}{h_1}\right) \mathbf{I}_{x_i} \mathbf{I}_i(y) \right. \\ \left. - \frac{1}{nh_1^{1+d_a}} \sum_{i=1}^n K^2\left(\frac{a^* - a_i}{h_1}\right) k^2\left(\frac{v^* - p_i}{h_1}\right) \mathbf{I}_{x_i} \mathbf{I}_i(y) \right| < 1. \end{aligned}$$

Therefore, $\hat{p}_i, \hat{S}_{0(y)}, \hat{S}_{0x}$ in (77) can be replaced with $p_i, S_0(y), S_{0x}$ without changing the (uniform) probability limit of $\hat{\mathcal{C}}$. Then, standard kernel estimation theory the uniform consistency of $\hat{\mathcal{C}}$.

For $\hat{S}_s^{(1)}$, let $\bar{S}_s^{(1)}$ be defined as $\hat{S}_s^{(1)}$ with \hat{p}_i replaced with p_i . Noting that $\sup_{y \in \mathcal{Y}} |k\{(y - y_i)/h_y\}|/h_y \leq C/h_y$ for some C , a slight modification of lemma F9 shows that

$$\sup_{y \in \mathcal{Y}} |\hat{S}_s^{(1)}(y; a^*, v^*) - \bar{S}_s^{(1)}(y; a^*, v^*)| < \frac{1}{\rho_n h_y} < 1,$$

using $\rho_n h_y > 1$. Then, standard kernel estimation theory shows the uniform consistency of $\bar{S}_s^{(1)}(\cdot; a^*, v^*)$. \square

Proof of Theorem 7. Redefine $\hat{H}(\cdot|a^*, v^*)$ in (17) as using $\mathbf{a}_i^\top \gamma_0$ and $\mathbf{z}_i^\top \theta_0$ in lieu of \mathbf{a}_i and \mathbf{z}_i . Let $\hat{\hat{H}}(\cdot|a^*, v^*)$ be identically defined but using $\hat{\boldsymbol{\gamma}}$ and $\hat{\boldsymbol{\theta}}$ instead of γ_0 and θ_0 . It suffices to show that

$$\sup_{y \in \mathcal{Y}} |\hat{\hat{H}}(y|a^*, v^*) - \hat{H}(y|a^*, v^*)| = o_p(1/\rho_n),$$

which follows from lemma H2. □