



**UCD GEARY INSTITUTE
DISCUSSION PAPER
SERIES**

**Identification of Treatment
Effects Using Control Functions
in Models with Continuous,
Endogenous Treatment and
Heterogeneous Effects**

15th December 2008

The views expressed here do not necessarily reflect those of the Geary Institute.
All errors and omissions remain those of the author.

Geary WP/32/2008

Identification of Treatment Effects Using Control Functions in Models with Continuous, Endogenous Treatment and Heterogeneous Effects.*

J.P. Florens[†], J. J. Heckman[‡], C. Meghir[§] and E. Vytlacil[¶]

This draft, June 13, 2008

Abstract

We use the control function approach to identify the average treatment effect and the effect of treatment on the treated in models with a continuous endogenous regressor whose impact is heterogeneous. We assume a stochastic polynomial restriction on the form of the heterogeneity but, unlike alternative nonparametric control function approaches, our approach does not require large support assumptions.

*We thank two anonymous referees and Whitney Newey for their comments. We also thank participants at the Berkeley-Stanford (March 2001) workshop on non-parametric models with endogenous regressors as well as participants at the UCL empirical microeconomics workshop for useful comments. J. Heckman would like to thank NIH R01-HD043411 and NSF SES-024158 for research support. C. Meghir would like to thank the Centre for Economics of Education and the ESRC for funding through the Centre for Fiscal policy at the IFS and his ESRC Professorial Fellowship. E. Vytlacil would like to thank the NSF for financial support (NSF SES-05-51089). The views expressed in this paper are those of the authors and not necessarily those of the funders listed here. All errors are our own.

[†]IDEI, Toulouse

[‡]University of Chicago and University College Dublin

[§]IFS and UCL, c.meghir@ucl.ac.uk

[¶]Yale University

1 Introduction

There is a large and growing theoretical and empirical literature on models where the impacts of discrete (usually binary) treatments are heterogeneous in the population.¹ The objective of this paper is to analyze non-parametric identification of treatment effect models with continuous treatments when the treatment intensity is not randomly assigned. This generally leads to models that are non-separable in the unobservables and produces heterogeneous treatment intensity effects. Imposing a stochastic polynomial assumption on the heterogeneous effects, we use a control function approach to obtain identification without large support assumptions. Our approach has applications in a wide variety of problems, including demand analysis where price elasticities may differ across individuals; labor supply, where wage effects may be heterogeneous; or production functions, where the technology may vary across firms.

Other recent papers on semiparametric and nonparametric models with nonseparable error terms and an endogenous, possibly continuous, covariate include papers using quantile instrumental variable methods such as Chernozhukov and Hansen (2005) and Chernozhukov, Imbens, and Newey (2007), and papers using a control variate technique such as Altonji and Matzkin (2005), Blundell and Powell (2004), Chesher (2003), and Imbens and Newey (2002, 2007). Chesher (2007) surveys this literature. The analysis of Imbens and Newey (2002, 2007) is perhaps the most relevant to our analysis, with the key distinction between our approach and their approach being a tradeoff between making a stochastic polynomial assumption on the outcome equation versus assuming large support. We discuss the differences between our approach and their approach further in Section 3.2.

2 The Model, Parameters of Interest and the Observables.

Let Y_d denote the potential outcome corresponding to level of treatment intensity d . When the treatments are discrete this notation represents the two possible outcomes for a particular individual in the treated and non-treated state. In this paper, there are a continuum of alternatives as the treatment intensity varies. Define $\varphi(d) = E(Y_d)$ and $U_d = Y_d - \varphi(d)$, so that, by

¹See, e.g., Roy (1951); Heckman and Robb (1985, 1986); Björklund and Moffitt (1987); Imbens and Angrist (1994); Heckman (1997); Heckman, Smith, and Clements (1997); Heckman and Honoré (1990); Card (1999, 2001); Heckman and Vytlačil (2001, 2005, 2007a,b), who discuss heterogeneous response models.

construction,

$$Y_d = \varphi(d) + U_d. \quad (1)$$

We restrict attention to the case where the stochastic process U_d takes the polynomial form

$$U_d = \sum_{j=0}^K d^j \varepsilon_j, \quad \text{with } E(\varepsilon_j) = 0, \quad j = 0, \dots, K, \quad (2)$$

where $K < \infty$ is known.²

Let D denote the realized treatment, so that the realized outcome Y is given by $Y = Y_D$. We do not explicitly denote observable regressors that directly affect Y_d . All of our analysis implicitly conditions on such regressors. We assume

(A-1) $\varphi(D)$ is K times differentiable in D (a.s.), and the support of D does not contain any isolated points (a.s.).

This allows for heterogeneity of a finite set of derivatives of Y_d . This specification can be seen as a nonparametric, higher order generalization of the random coefficient model analyzed by Heckman and Vytlacil (1998) and Wooldridge (1997, 2003, 2007). The normalization $E(\varepsilon_j) = 0$, $j = 0, \dots, K$, implies that $\frac{\partial^j}{\partial d^j} E(Y_j) = \frac{\partial^j}{\partial d^j} \varphi(d)$.³

Equations (1) and (2) can be restated as follows to emphasize that we analyze a nonseparable model:

$$Y = h(D, \epsilon) = \varphi(D) + \sum_{j=0}^K D^j h_j(\epsilon) \quad (3)$$

where ϵ need not be a scalar random variable. The notation of equation (3) can be mapped into the notation of equations (1) and (2) by setting $\varepsilon_j = h_j(\epsilon)$. Notice that do not assume that ϵ is a scalar random variable, and h need not be monotonic in ϵ .

One parameter of interest in this paper is the Average Treatment Effect,

$$\Delta^{ATE}(d) = \lim_{\Delta d \rightarrow 0} \frac{E(Y_{d+\Delta d} - Y_d)}{\Delta d} \equiv \frac{\partial}{\partial d} E(Y_d) = \frac{\partial}{\partial d} \varphi(d) \quad (4)$$

which is the average effect of a marginal increase in treatment if individuals were randomly assigned to base treatment level d . Note that the average treatment effect depends on the

²As discussed later, we can test for the order of the polynomial as long as a finite upper bound on K is known. The question of identification with K infinite is left for future work.

³To see that $E(\varepsilon_j) = 0$, $j = 0, \dots, K$, is only a normalization, note that $\varphi(d) + \sum_{j=0}^K d^j \varepsilon_j = \left[\varphi(d) + \sum_{j=0}^K d^j E(\varepsilon_j) \right] + \sum_{j=0}^K d^j (\varepsilon_j - E(\varepsilon_j)) = \tilde{\varphi}(d) + \sum_{j=0}^K d^j \tilde{\varepsilon}_j$.

base treatment level, and for any of the continuum of possible base treatment levels we have a different average treatment effect. The average treatment effect is the derivative of the average structural function of Blundell and Powell (2004).

We also consider the effect of treatment on the treated (TT), given by

$$\begin{aligned}\Delta^{TT}(d) &= \lim_{\Delta d \rightarrow 0} \frac{E(Y_{d+\Delta d} - Y_d | D = d)}{\Delta d} \\ &\equiv E\left(\frac{\partial}{\partial d_1} Y_{d_1} | D = d_2\right) \Bigg|_{d=d_1=d_2} = \frac{\partial}{\partial d} \varphi(d) + \sum_{j=1}^K j d^{j-1} E(\varepsilon_j | D = d)\end{aligned}$$

which is the average effect of treatment for those currently choosing treatment level d of an incremental increase in the treatment holding their unobservables fixed at baseline values. This parameter corresponds to the local average response parameter considered by Altonji and Matzkin (2001, 2005).

We denote the choice equation (the assignment mechanism to treatment intensity) as

$$D = g(Z, V) \tag{5}$$

where Z are observed covariates that enter the treatment choice equation but are excluded from the equation for Y_d and V is a scalar unobservable. We make the following assumption:

(A-2) *V is absolutely continuous with respect to Lebesgue measure; g is strictly monotonically increasing in V ; and $Z \perp\!\!\!\perp (V, \varepsilon_0, \dots, \varepsilon_K)$.*

As long as D is a continuous random variable (conditional on Z), we can always represent D as a function of Z and a continuous scalar error term, with the function increasing in the error term and the error term independent of Z . To see this, set $V = F_{D|Z}(D|Z)$ and $g(Z, V) = F_{D|Z}^{-1}(V|Z)$. Thus, $D = g(Z, V)$ where g is strictly increasing in the scalar V which is distributed unit uniform and independent of Z . However, the assumption that $g(Z, V)$ is monotonic in a scalar unobservable V with $Z \perp\!\!\!\perp (V, \varepsilon_0, \dots, \varepsilon_K)$ is restrictive. The construction $V = F_{D|Z}(D|Z)$ and $D = F_{D|Z}^{-1}(V|Z) = g(Z, V)$ does not guarantee $Z \perp\!\!\!\perp (V, \varepsilon_0, \dots, \varepsilon_K)$.

Given assignment mechanism (5) and assumption (A-2), without loss of generality we can impose the normalization that V is distributed unit uniform. Given these assumptions and the normalization of V , we can follow Imbens and Newey (2002, 2007) and recover V from $V = F_{D|Z}(D|Z)$ and the function g from $g(Z, V) = F_{D|Z}^{-1}(V|Z)$. Assignment mechanism (5) and assumption (A-2) will not be directly used to prove identification. However, we use it to clarify the primitives underlying our identification assumptions.

2.1 Education and Wages: A Simple Illustration

To illustrate the type of problem we analyze in this paper, consider a simple model of educational choice. Suppose that the agent receives wages Y_d at direct cost C_d if schooling choice d is made. We work with discounted annualized earnings flows. We write wages for schooling level d , Y_d , as

$$Y_d = \varphi_0 + (\varphi_1 + \varepsilon_1)d + \frac{1}{2}\varphi_2d^2 + \varepsilon_0$$

and the cost function for schooling as

$$C_d = C_0(Z) + (C_1(Z) + v_1)d + \frac{1}{2}C_2(Z)d^2 + v_0 \quad (6)$$

where ε_s and v_s ($s = 0, 1$) are, respectively, unobserved heterogeneity in the wage level and in the cost of schooling. These unobserved heterogeneity terms are the source of the identification problem considered in this paper. We impose the normalizations that $E(\varepsilon_s) = 0$, $E(v_s) = 0$, for $s = 0, 1$. We implicitly condition on variables such as human capital characteristics that affect both wages and the costs of schooling. The Z are factors that only affect the cost of schooling, such as the price of education.

Assume that agents choose their level of education to maximize wages minus costs. Let D denote the resulting optimal choice of education. D solves the first order condition

$$(\varphi_1 - C_1(Z)) + (\varphi_2 - C_2(Z))D + \varepsilon_1 - v_1 = 0.$$

Assuming that $\varphi_2 - C_2(Z) < 0$ for all Z , the second order condition for a maximum will be satisfied. This leads to an education choice equation (assignment to treatment intensity rule)

$$D = \frac{\varphi_1 - C_1(Z) + \varepsilon_1 - v_1}{C_2(Z) - \varphi_2}.$$

This choice equation is produced as a special case of the model given by equations (1), (2) and (5), with

$$\begin{aligned} \varphi(d) &= \varphi_0 + \varphi_1d + \frac{1}{2}\varphi_2d^2 \\ U_d &= \varepsilon_0 + \varepsilon_1d \\ g(z, v) &= \frac{\varphi_1 - C_1(z) + F_{\varepsilon_1 - v_1}^{-1}(v)}{C_2(z) - \varphi_2} \end{aligned}$$

where $V = F_{\varepsilon_1 - v_1}(\varepsilon_1 - v_1)$ with $F_{\varepsilon_1 - v_1}$ the cumulative distribution function of $\varepsilon_1 - v_1$. The goal is to identify the average return to education: $\Delta^{ATE}(d) = \varphi_1 + \varphi_2d$, or TT, which is $\Delta^{TT}(d) = (\varphi_1 + E(\varepsilon_1|D = d)) + \varphi_2d$.

In this example, the treatment intensity is given by equation (5) with g strictly increasing in a scalar error term $V = F_{\varepsilon_1 - v_1}(\varepsilon_1 - v_1)$. The structure of the treatment intensity mechanism is sensitive to alternative specifications. Consider the same example as before, except now the second derivative of Y_d is also stochastic: $Y_d = \varphi_0 + (\varphi_1 + \varepsilon_1)d + \frac{1}{2}(\varphi_2 + \varepsilon_2)d^2 + \varepsilon_0$. The choice equation becomes $D = \frac{\varphi_1 - C_1(Z) + \varepsilon_1 - v_1}{C_2(Z) - \varphi_2 - \varepsilon_2}$. In this case, the structural model makes D a function of $V = (\varepsilon_1 - v_1, \varepsilon_2)$, which satisfies $Z \perp\!\!\!\perp (V, \varepsilon_0, \varepsilon_1, \varepsilon_2)$ but V is not a scalar error. We can still define $\tilde{V} = F_{D|Z}(D|Z)$ and the function \tilde{g} by $\tilde{g}(Z, \tilde{V}) = F_{D|Z}^{-1}(\tilde{V}|Z)$. With this construction, D is strictly increasing in a scalar error term \tilde{V} that is independent of Z . However, Z is not independent of $(\tilde{V}, \varepsilon_0, \varepsilon_1, \varepsilon_2)$. To see why, note that $\Pr(\tilde{V} \leq v|Z, \varepsilon_0, \varepsilon_1, \varepsilon_2) = \Pr\left[v_1 : \frac{\varphi_1 - C_1(Z) + \varepsilon_1 - v_1}{C_2(Z) - \varphi_2 - \varepsilon_2} \leq F_{D|Z}^{-1}(v)|Z, \varepsilon_0, \varepsilon_1, \varepsilon_2\right] \neq \Pr(\tilde{V} \leq v|\varepsilon_0, \varepsilon_1, \varepsilon_2)$. This is a case where assumption (A-2) does not hold. The fragility of the specification of equation (5) where g is strictly increasing in a scalar error term arises in part because, under rational behavior, heterogeneity in response to treatment (heterogeneity in the Y_d model) generates heterogeneity in selection into treatment intensity. This heterogeneity is absent if agents do not know their own treatment effect heterogeneity, which can happen if agents are uncertain at the time they make participation decisions (see Abbring and Heckman, 2007).

3 Identification Analysis.

An instrumental variable estimator IV does not identify ATE in the case of binary treatment with heterogeneous impacts (Heckman, 1997; Heckman and Robb, 1986) unless one imposes covariance restrictions between the errors in the assignment rule and the errors in the structural model. Following Newey and Powell (2003) and Darolles, Florens, and Renault (2002), consider a nonparametric IV strategy based on the identifying assumption that $E(Y - \varphi(D)|Z) = 0$. Suppose $K = 0$, which is the special case of no treatment effect heterogeneity. In this case, $U_D = \varepsilon_0$ so that $Y_D = \varphi(D) + \varepsilon_0$. We obtain the standard additive-in-unobservables model considered in the cited papers. The identification condition is $E(\varepsilon_0|Z) = 0$. However, in the general case of treatment effect heterogeneity ($K > 0$), the IV identification restriction implies special covariance restrictions between the error terms. For example, suppose $K = 1$ and that $D = g(Z) + V$. Then $E(Y - \varphi(D)|Z) = 0$ requires $E(\varepsilon_0|Z) = 0$ and $E(\varepsilon_1 D|Z) = 0$, with the latter restriction generically equivalent to $E(\varepsilon_1|Z) = 0$ and $E(\varepsilon_1 V|Z) = 0$. In other words, in addition to the more standard type of condition that ε_0 be mean independent of

the instrument, we now have a new restriction in the heterogeneous case that the covariance between the heterogeneous effect and the unobservables in the choice equation conditional on the instrument does not depend on the instrument.⁴ Instead of following an instrumental variables approach, we explore identification through a control function.⁵ We assume the existence of a (known or identifiable) control function \tilde{V} that satisfies the following conditions:

(A-3) *Control Function Condition:* $E(\varepsilon_j | D, Z) = E(\varepsilon_j | \tilde{V}) = r_j(\tilde{V})$.⁶

and

(A-4) *Rank condition:* D and \tilde{V} are measurably separated, i.e., any function of D almost surely equal to a function of \tilde{V} must be almost surely equal to a constant.

A necessary condition for assumption (A-4) to hold is that the instruments Z affect D .⁷ We return later in this section to consider sufficient conditions on the underlying model that implies the existence of such a control variate \tilde{V} . Under these assumptions, ATE and TT are identified.

Theorem 1. *Assume equations (1) and (2) hold with finite $K \geq 1$. Under assumptions (A-3) (control function condition), (A-4) (rank condition), and the smoothness and support condition (A-1), ATE and TT are identified.*

Proof. See Appendix. □

The control function assumption gives the basis for an empirical determination of the relevant degree of the polynomial in (2). If the true model is defined by a polynomial of degree ℓ we have that for any $k > \ell$

$$\frac{\partial^k}{\partial d^k} E(Y | D = d, \tilde{V} = v) = \frac{\partial^k \varphi(d)}{\partial d^k}$$

⁴See Heckman and Vytlačil (1998) and Wooldridge (1997, 2003, 2007).

⁵See Newey, Powell, and Vella (1999) for a control function approach for the case of separable models ($K = 0$). See Heckman and Vytlačil (2007b) for a discussion of the distinction between control functions and control variables. Technically “control function” is a more general concept. We adopt the recent nomenclature even though it is inaccurate. See the Matzkin (2007) paper for additional discussion.

⁶Note that our normalization $E(\varepsilon_j) = 0$, $j = 0, \dots, K$, implies the normalization that $E(r_j(\tilde{V})) = 0$, $j = 0, \dots, K$,

⁷Measurable separability, which we maintain in this paper is just one way of achieving identification. Alternatively, one could restrict the space of functions $\varphi(D)$ not to contain $r_j(\tilde{V})$ functions; this in turn can be achieved for example by assuming that $\varphi(D)$ is linear in D and r_j is non-linear as in the Heckman (1979) selection model. See also Heckman and Robb (1985, 1986) who discuss this condition.

which does not depend on v and thus is only a function of d . This property can be verified by checking whether the following equality holds almost surely:

$$\frac{\partial^k E(Y|D, \tilde{V})}{\partial D^k} \stackrel{a.s.}{=} E \left[\frac{\partial^k}{\partial D^k} E(Y|D, \tilde{V}) \Big| D \right], \quad k > \ell.$$

3.1 Primitive Conditions Justifying the Control Function Assumption.

In the previous section a control function is assumed to exist and satisfy certain properties. The analysis in the previous section did not use assignment rule (5) or condition (A-2). In this section, we use assignment rule (5) and condition (A-2) along with the normalization that V is distributed unit uniform. Under these conditions, consider using $V = F_{D|Z}(D|Z)$ as the control function. This leads to the following corollary to Theorem 1:

Corollary 3.1. *Assume equations (1) and (2) hold with finite $K \geq 1$ and assume smoothness and support condition (A-1). If D is generated by assignment equation (5) and condition (A-2) holds, and if V and D are measurably separable, (A-4), then ATE and TT are identified.*

In order for the conditions of Theorem 1 to be satisfied it is sufficient to verify that under the conditions in the corollary the control function assumption (A-3) is satisfied. Given that D satisfies assignment equation (5) and condition (A-2), from Imbens and Newey (2002, 2007) we obtain that $V = F_{D|Z}(D|Z)$ is a control variate satisfying assumption (A-3).

Next consider measurable separability condition (A-4). Measurable separability is a relatively weak condition, as illustrated by the following theorem.

Theorem 2. *Assume that (D, V) has a density with respect to Lebesgue measure in \mathbb{R}^2 and denote its support by S and let S^0 be the interior of the support. Further, assume that i) any point in S^0 has a neighborhood such that the density is strictly positive within it and ii) any two points within S^0 can be connected by a continuous curve that lies strictly in S^0 . Then measurable separability between D and V (A-4) holds.*

Proof. See Appendix. □

Measurable separability is a type of rank condition. To see this, consider the following heuristic argument. Consider a case where the condition is violated at some point in the interior of the support of (D, V) , i.e. $h(D) = l(V)$. Hence $h(g(Z, V)) = l(V)$. Differentiating both

sides of this expression with respect to Z , we obtain $\frac{\partial h}{\partial g} \frac{\partial g}{\partial Z} = 0$. If measurable separability fails, $\frac{\partial h}{\partial g} \neq 0$ and hence $\frac{\partial g}{\partial Z} = 0$ which means that g does not vary with Z . Note that the conditions in Theorem 2 are not very restrictive. For example, the conditional support of D can depend on V and vice versa.

Assignment rule (5) and condition (A-2) do not imply measurable separability (A-4). To show this, we consider two examples where equation (5) and condition (A-2) hold but D and V are not measurably separable. In the first example, Z is a discrete random variable. In the second example, $g(z, v)$ is a discontinuous function of v .

First, suppose $Z = 0, 1$ and suppose that $D = g(z, v) = z + v$, with $V \text{ Unif}[0, 1]$. Then (A-4) fails, i.e., D and V are not measurably separable. To see this, let $m_1(t) = t$ and let $m_2(t) = \mathbf{1}[t \leq 1]t + \mathbf{1}[t > 1](t - 1)$. Then $m_1(V) = m_2(D)$, but m_1 and m_2 are not a.s. equal to a constant. Now consider a second example. Suppose that $D = g_1(z) + g_2(v)$, where $g_2(t) = \mathbf{1}[t \leq .5]t + \mathbf{1}[t > .5](1 + t)$. Let g_1^{max} and g_1^{min} denote the maximum and minimum of the support of the distribution of $g_1(Z)$, and suppose that $g_1^{max} - g_1^{min} < 1$. Then (A-4) fails, i.e., D and V are not measurably separable. To see this, let $m_1(t) = \mathbf{1}[t \leq .5]$, let $m_2(t) = \mathbf{1}[t \leq .5 + g_1^{max}]$, and note that $m_1(V) = m_2(D)$ but that m_1 and m_2 do not (a.s.) equal a constant.

Assignment rule (5), condition (A-2), and regularity conditions that require Z to contain a continuous element and that g be continuous in v are sufficient to imply that measurable separability (A-4) holds. We prove the following theorem.

Theorem 3. *Suppose that D is determined by equation (5). Suppose that $g(z, v)$ is a continuous function of v . Suppose that, for any fixed v , the support of the distribution of $g(Z, v)$ contains an open interval. Then, under assumption (A-2), D and V are measurably separated (A-4) holds.*

Proof. See Appendix. □

Note that, for any fixed v , for the support of the distribution of $g(Z, v)$ to contain an open interval requires that Z contains a continuous element. A sufficient condition for the support of the distribution of $g(Z, v)$ to contain an open interval is that (a) Z contains an element whose distribution conditional on the other elements of Z contains an open interval, and (b) g is a continuous monotonic function of that element. Thus, under the conditions of Theorem 3, V is identified by $V = F(D|Z)$ and both the control function condition (A-3) and the rank

condition (A-4) hold with $\tilde{V} = V$.

3.2 Alternative Identification Analyses.

A general analysis of identification using the control function assumption without the polynomial structure is related to the work of Heckman and Vytlacil (2001) on identifying the marginal treatment effect (MTE). That paper considers a binary treatment model, but their analysis may be extended to the continuous treatment case. Similar approaches for semiparametric models with a continuous treatment D that is strictly monotonic in V is pursued for the Average Structural Function (ASF) by Blundell and Powell (2004) and for the Local Average Response (LAR) by Altonji and Matzkin (2001, 2005). The derivative of the ASF corresponds to our ATE, and the LAR corresponds to treatment on the treated.

Most relevant to this note is the analysis of Imbens and Newey (2002, 2007). They invoked the same structure as we do on the first stage equation for the endogenous regressor as our assignment mechanism (5) and they also invoke assumption (A-2). The control variate is V , with V identified, and with a distribution that can be normalized to be unit uniform. By the same reasoning, we have $E(Y_d|D = d, V = v) \stackrel{as}{=} E(Y_d|V = v)$. Furthermore, they assume that the support of (D, V) is the product of the support of the two marginal distributions, i.e., they assume rectangular support. Their assumption implies that the conditional support of D given V does not depend on V (and vice versa). It is stronger than the measurable separability assumption we previously used to establish identification. From these assumptions it follows that

$$E(Y|D = d, V = v) = E(Y_d|D = d, V = v) = E(Y_d|V = v)$$

and

$$E(Y_d) = \int E(Y_d|V = v)dF(v).$$

Then $\varphi(d) = \int E(Y|D = d, V = v)dF(v)$ and is identified. Identification of $\varphi(d)$ in turn implies identification of $\Delta^{ATE} = \frac{\partial}{\partial d}\varphi(d)$. The rectangular support condition is needed to replace $E(Y_d|V = v)$ by $E(Y|D = d, V = v)$ for all v in the unconditional support of V in the previous integral. The rectangular support condition may not be satisfied and in general requires a large support assumption as illustrated by the following example. Suppose $D = g_1(Z) + V$. Let \mathcal{G}_1 denote the support (Supp) of the distribution of $g_1(Z)$. If Z and V are independent, then

$$\text{Supp}(V|D = d) = \text{Supp}(V|g_1(Z) + V = d) = \text{Supp}(V|V = d - g_1(Z)) = \{d - g : g \in \mathcal{G}_1\}$$

where the last equality uses the condition that $Z \perp\!\!\!\perp V$. $\{d - g : g \in \mathcal{G}_1\}$ does not depend on d if and only if $\mathcal{G}_1 = \mathfrak{R}$. For example, if $\mathcal{G}_1 = [a, b]$, then $\{d - g : g \in \mathcal{G}_1\} = \{d - g : g \in [a, b]\} = [d - b, d - a]$ which does not depend on d if and only if $a = -\infty$ and $b = \infty$, i.e., if and only if $\mathcal{G}_1 = \mathfrak{R}$.

Instead of imposing $E(Y_d|D = d, V = v) = E(Y_d|V = v)$, one could instead impose

$$\frac{\partial}{\partial d}E(Y|D = d, V = v) = E\left(\frac{\partial}{\partial d}Y_d|D = d, V = v\right) = E\left(\frac{\partial}{\partial d}Y_d|V = v\right). \quad (7)$$

$E\left(\frac{\partial}{\partial d}Y_d|V = v\right)$ is the marginal treatment effect of Heckman and Vytlacil (2001), adapted to the case of a continuous treatment. Instead of integrating $E(Y_d|V = v)$ to obtain $\varphi(d)$, one could instead integrate $E\left(\frac{\partial}{\partial d}Y_d|V = v\right)$ to obtain ATE or TT:

$$\begin{aligned} \int \frac{\partial}{\partial d}E(Y|D = d, V = v)dF(v) &= \int \frac{\partial}{\partial d}E(Y_d|V = v)dF(v) = \Delta^{ATE}(d), \\ \int \frac{\partial}{\partial d_1}E(Y_{d_1}|D = d_2, V = v)dF(v|D = d_2)\Big|_{d=d_1=d_2} &= E\left(\frac{\partial}{\partial d_1}Y_{d_1}|D = d_2\right)\Big|_{d=d_1=d_2} \\ &= \frac{\partial}{\partial d_1}E(Y_{d_1}|D = d_2)\Big|_{d=d_1=d_2} = \Delta^{TT}(d). \end{aligned}$$

This is the identification strategy followed in Heckman and Vytlacil (2001), adapted to the case where D is a continuous treatment. As discussed in Heckman and Vytlacil (2001), a rectangular support condition is required in order to integrate up MTE to obtain ATE. Note that one does not require the rectangular support condition to integrate up $\frac{\partial}{\partial d}E(Y|D = d, V = v)$ to obtain TT. For TT, one only needs to evaluate $\frac{\partial}{\partial d}E(Y|D = d, V = v)$ for v in the support of V conditional on $D = d$, not in the unconditional support of V .

While a rectangular support condition is not required to integrate MTE to recover TT, a support condition is required for equation (7) to hold. That equation requires that $E(Y|D = d, V = v)$ can be differentiated with respect to d while keeping v fixed. This property is closely related to measurable separability between D and V . Assume that there exists a (differentiable) function of D , $h(D)$ equal (a.s.) to a function of V , $m(V)$, which is not constant. Then we obtain

$$E(Y|D = d, V = v) \stackrel{as}{=} E(Y_d|V = v) + h(d) - m(V)$$

and

$$\frac{\partial}{\partial d}E(Y|D = d, V = v) \stackrel{as}{=} \frac{\partial}{\partial d}E(Y_d|V = v) + \frac{\partial}{\partial d}h(d)$$

which implies that equation (7) is violated. Thus, for TT, we still need measurable separability between D and V in order for equation (7) to hold.

There are trade-offs between the approach presented in this note versus an approach that identifies MTE/MTE-like objects and then integrates them to obtain the object of interest. The approach developed here requires a stochastic polynomial structure on U_D of equation (2) and higher order differentiability. These conditions are not required by Imbens and Newey (2002, 2007) or Heckman and Vytlačil (2001). The approach of this note does not require the large support assumption required to implement these alternative approaches. As shown by Theorem 2, measurable separability between D and V is a relatively mild restriction on the support of (D, V) . As shown by Theorem 3, measurable separability between D and V follows from assignment mechanism (5) and Assumption (A-2) combined with a relatively mild regularity condition.

4 Estimation

Under the control function assumption we have:

$$\begin{aligned} E(Y|D = d, Z = z) &= E(Y|D = d, V = v) \\ &= \varphi(d) + \sum_{j=0}^K d^j h_j(v) . \end{aligned}$$

The method we propose is an extension of Newey, Powell and Vella (1999) and may also be viewed as an extension of estimation of additive models in a nonparametric context. The estimation is carried out in two steps: first estimate the residual v_i from the nonparametric regression $D = E(D|Z) + V$; then estimate φ and the h_j 's.

Define the estimation criterion

$$\min_{\varphi, h_0, \dots, h_K} E \left[Y - \varphi - \tilde{D}'h(v) \right]^2 \quad (8)$$

where $\tilde{D} = [1, D, D^2, \dots, D^K]'$, and $h = [h_0, h_1, \dots, h_K]'$. The first order conditions for the minimisation are

$$\begin{cases} E(Y|D = d) = \varphi(d) + \tilde{d}'E(h(V)|D = d) \\ E(\tilde{D}Y|V = v) = E(\tilde{D}\varphi(D)|V = v) + E(\tilde{D}\tilde{D}'|V = v)h \end{cases} \quad (9)$$

where $\tilde{d} = [1, d, \dots, d^K]'$.

This linear system in φ and h can easily be solved if the conditional expectations are replaced by their estimators (by kernels for example). In that case it is easily seen that (9) generates a

linear system with respect to the $\varphi(d_i)$ and the $h_j(v_i)$ ($i = 1, \dots, n; j = 0, \dots, K$) and this system may be solved by usual methods of linear equations. The equations in (9) are then used to compute $\varphi(d)$ and $h_j(v)$ at any point of evaluation. If we only wish to focus attention on φ , the vector h may be eliminated from (9) and we obtain:

$$\begin{aligned} \varphi(d) - \tilde{d}'E \left[E(\tilde{D}\tilde{D}'^{-1}E(\varphi(D)|V = v)|D = d) \right] \\ = E(Y|D = d) - \tilde{d}'E \left[E(\tilde{D}\tilde{D}'^{-1}E(\tilde{D}Y|V = v)|D = d) \right] . \end{aligned}$$

This equation has the form $(I - T)\varphi = \psi$ where T is, under very general conditions, a compact operator and ψ may be estimated. It is a Fredholm equation of type II which may be analysed using the methods in Carrasco, Florens, and Renault (2007, section 7). The original system (9) is also a Fredholm equation of type II and both systems generate well posed inverse problems. The asymptotic theory developed in Carrasco, Florens, and Renault (2007) applies with the exception that the v_i are now estimated. A precise analysis of this approach and some applications will be developed in future work.

5 Conclusions

This paper considers the identification and estimation of models with a continuous endogenous regressor and non-separable errors when continuous instruments are available. We present an identification result using a control function technique. Our analysis imposes a stochastic, finite-order polynomial restriction on the outcome model but does not impose a large support assumption.

6 Appendix: Proofs of Theorems

Proof of Theorem 1

Suppose that there are two sets of parameters $(\varphi^1, r_K^1, \dots, r_0^1)$ and $(\varphi^2, r_K^2, \dots, r_0^2)$ such that

$$E(Y|D = d, \tilde{V} = v) = \varphi^i(d) + \sum_{k=0}^K d^k r_k^i(v), \quad i = 1, 2,$$

where the conditional expectation on the left-hand side takes this form as a result of the control

function assumption (A-3). Then

$$[\varphi^1(d) - \varphi^2(d)] + \sum_{k=0}^K d^k [r_k^1(v) - r_k^2(v)] = 0. \quad (10)$$

Given smoothness assumption (A-1), this implies

$$\frac{\partial^K}{\partial d^K} \varphi^1(d) - \frac{\partial^K}{\partial d^K} \varphi^2(d) + (K!)(r_K^1(v) - r_K^2(v)) = 0.$$

Measurable separability assumption (A-4) implies that if any function of d is equal to a function of v (a.s.) then this must be a constant (a.s.). Hence, $r_K^1(v) - r_K^2(v)$ is a constant a.s.. Hence,

$$r_K^1(v) - r_K^2(v) = E \left[r_K^1(\tilde{V}) - r_K^2(\tilde{V}) \right].$$

This expression equals zero given our normalization that $E(\varepsilon_K) = 0$. Hence,

$$r_K^1(v) - r_K^2(v) \stackrel{a.s.}{=} 0.$$

Considering the $(K-1)^{\text{st}}$ derivative of equation (10), we find that

$$\frac{\partial^{K-1}}{\partial d^{K-1}} \varphi^1(d) - \frac{\partial^{K-1}}{\partial d^{K-1}} \varphi^2(d) + (K!)d \left[r_K^1(v) - r_K^2(v) \right] + ((K-1)!) \left[r_{K-1}^1(v) - r_{K-1}^2(v) \right] = 0.$$

We have already shown that $r_K^1(v) = r_K^2(v)$, and thus

$$\frac{\partial^{(K-1)}}{\partial d^{(K-1)}} \varphi^1(d) - \frac{\partial^{(K-1)}}{\partial d^{(K-1)}} \frac{\partial}{\partial d} \varphi^2(d) + ((K-1)!) (r_{K-1}^1(v) - r_{K-1}^2(v)) = 0.$$

Using the logic of the previous analysis, we can show that $r_{K-1}^1(v) - r_{K-1}^2(v) \stackrel{a.s.}{=} 0$. Iterating this procedure for $k = K-2, \dots, 0$, it follows that $r_k^1(v) - r_k^2(v) \stackrel{a.s.}{=} 0$ for all $k = 0, \dots, K$. Again appealing to equation (10), it follows that $\varphi^1(d) - \varphi^2(d) \stackrel{a.s.}{=} 0$, and thus ATE is identified. Using the fact that $\varphi^1(d) - \varphi^2(d) \stackrel{a.s.}{=} 0$ and $r_k^1(v) - r_k^2(v) \stackrel{a.s.}{=} 0$ for all $k = 0, \dots, K$, we also have that $\frac{\partial}{\partial d} \varphi^1 + \sum_{k=1}^K k d^{k-1} E[r_k^1(v)|d] = \frac{\partial}{\partial d} \varphi^2 + \sum_{k=1}^K k d^{k-1} E[r_k^2(v)|d] = 0$, and thus TT is identified. ■

Proof of Theorem 2

Let (d, v) be a point of the interior of the support S^0 . Let N^d denote a neighborhood of d and N^v a neighborhood of v such that $N^d \times N^v$ is included in S^0 . The distribution of (D, V) restricted to $N^d \times N^v$ is equivalent to Lebesgue measure (i.e. has the same null sets). Then using Theorem 5.2.7 of Florens, Mouchart, and Rolin (1990) (D, V) restricted to $N^d \times N^v$ are measurably separated. This implies that if within that neighborhood $h(D) \stackrel{as}{=} l(V)$, then $h(D)$

and $l(V)$ are *a.s.* constants. We need to show that this is true everywhere in the interior of the support. Consider any two points (d, v) and (d', v') in S^0 . The theorem will be true if $h(d) = h(d')$. As S^0 satisfies the property (ii) in the theorem and is open by definition, there exists a finite number of overlapping open sets with non-empty overlaps, i.e. \exists a finite sequence of neighborhoods $N_j^d \times N_j^v$, $j = 1, \dots, J$ such that each $N_j^d \times N_j^v \subset S^0$ and $N_j^d \cap N_{j+1}^d \neq \emptyset$ and similarly for N_j^v . The first point (d, v) is in $N_1^d \times N_1^v$ and the second point (d', v') is in $N_J^d \times N_J^v$. Take $d_1 \in N_1^d$ and in the next overlapping neighborhood $d_2 \in N_2^d$. From the previous result (D, V) are measurably separated on $N_1^d \times N_1^v$ and on $N_2^d \times N_2^v$. Thus $h(d_i)$ $i = 1, 2$ is constant on each and thus constant on the union implying $h(d_1) = h(d_2)$. Iterating in this way along the sequence of neighborhoods until $N_J^d \times N_J^v$, it follows that $h(d) = h(d')$. Hence $h(D)$ is *a.s.* constant and, because $h(D) \stackrel{a.s.}{=} l(V)$, $l(v)$ is *a.s.* constant.

■

Proof of Theorem 3

Let \mathcal{Z} denote the support of the distribution of Z . Consider any two functions m_1 and m_2 such that $m_1(D) = m_2(V)$ a.s. For (a.e. F_V) fixed v_0 , using the assumption that Z and V are independent, it follows that $m_1(g(z, v_0)) = m_2(v_0)$ for a.e. z conditional on $V = v_0$ implies that m_1 is (a.s. F_Z) constant on $\{g(z, v_0) : z \in \mathcal{Z}\}$. Likewise, for a v_1 close to v_0 , we have m_1 is constant on $\{g(z, v_1) : z \in \mathcal{Z}\}$. Using the fact that $g(z, v)$ is continuous in v and that $\{g(z, v) : z \in \mathcal{Z}\}$ contains an open interval for any v , we can pick v_1 sufficiently close to v_0 so that $\{g(z, v_0) : z \in \mathcal{Z}\}$ and $\{g(z, v_1) : z \in \mathcal{Z}\}$ have a nonnegligible intersection, and we thus conclude that m_1 is constant on $\{g(z, v) : z \in \mathcal{Z}, v = v_0, v_1\}$. Proceeding in this fashion, we conclude that m_1 is (a.s.) constant on $\{g(z, v) : z \in \mathcal{Z}, v \in [0, 1]\}$, and thus that m_1 is a.s. equal to a constant. ■

References

- ABBRING, J. H., AND J. J. HECKMAN (2007): “Econometric Evaluation of Social Programs, Part III: Distributional Treatment Effects, Dynamic Treatment Effects, Dynamic Discrete Choice, and General Equilibrium Policy Evaluation,” in *Handbook of Econometrics*, ed. by J. Heckman, and E. Leamer, vol. 6B, pp. 5145–5303. Elsevier, Amsterdam.
- ALTONJI, J. G., AND R. L. MATZKIN (2001): “Panel Data Estimators for Nonseparable Models with Endogenous Regressors,” Technical Working Paper t0267, NBER.
- (2005): “Cross Section and Panel Data Estimators for Nonseparable Models with Endogenous Regressors,” *Econometrica*, 73(4), 1053–1102.
- BJÖRKLUND, A., AND R. MOFFITT (1987): “The Estimation of Wage Gains and Welfare Gains in Self-Selection,” *Review of Economics and Statistics*, 69(1), 42–49.
- BLUNDELL, R., AND J. POWELL (2004): “Endogeneity in Semiparametric Binary Response Models,” *Review of Economic Studies*, 71(3), 655–679.
- CARD, D. (1999): “The Causal Effect of Education on Earnings,” in *Handbook of Labor Economics*, ed. by O. Ashenfelter, and D. Card, vol. 5, pp. 1801–1863. North-Holland, New York.
- (2001): “Estimating the Return to Schooling: Progress on Some Persistent Econometric Problems,” *Econometrica*, 69(5), 1127–1160.
- CARRASCO, M., J.-P. FLORENS, AND E. RENAULT (2007): “Linear Inverse Problems in Structural Econometrics Estimation Based on Spectral Decomposition and Regularization,” in *Handbook of Econometrics*, ed. by J. J. Heckman, and E. Leamer, vol. 6B, pp. 5633–5751. Elsevier, Amsterdam.
- CHERNOZHUKOV, V., AND C. HANSEN (2005): “An IV Model of Quantile Treatment Effects,” *Econometrica*, 73(1), 245–261.
- CHERNOZHUKOV, V., G. W. IMBENS, AND W. K. NEWEY (2007): “Instrumental Variable Estimation of Nonseparable Models,” *Journal of Econometrics*, 139(1), 4–14.
- CHESHER, A. (2003): “Identification in Nonseparable Models,” *Econometrica*, 71(5), 1405–1441.

- (2007): “Identification of Non-Additive Structural Functions,” in *Advances in Economics and Econometrics: Theory and Applications, Ninth World Congress*, ed. by W. K. N. Richard Blundell, and T. Persson, vol. 3, chap. 1. Cambridge University Press, New York, Presented at the Econometric Society Ninth World Congress. 2005. London, England.
- DAROLLES, S., J.-P. FLORENS, AND E. RENAULT (2002): “Nonparametric Instrumental Regression,” Working Paper 05-2002, Centre interuniversitaire de recherche en économie quantitative, CIREQ.
- FLORENS, J.-P., M. MOUCHART, AND J. ROLIN (1990): *Elements of Bayesian Statistics*. M. Dekker, New York.
- HECKMAN, J. J. (1979): “Sample Selection Bias as a Specification Error,” *Econometrica*, 47(1), 153–162.
- (1997): “Instrumental Variables: A Study of Implicit Behavioral Assumptions Used in Making Program Evaluations,” *Journal of Human Resources*, 32(3), 441–462, Addendum published vol. 33 no. 1 (Winter 1998).
- HECKMAN, J. J., AND B. E. HONORÉ (1990): “The Empirical Content of the Roy Model,” *Econometrica*, 58(5), 1121–1149.
- HECKMAN, J. J., AND R. ROBB (1985): “Alternative Methods for Evaluating the Impact of Interventions,” in *Longitudinal Analysis of Labor Market Data*, ed. by J. Heckman, and B. Singer, vol. 10, pp. 156–245. Cambridge University Press, New York.
- (1986): “Alternative Methods for Solving the Problem of Selection Bias in Evaluating the Impact of Treatments on Outcomes,” in *Drawing Inferences from Self-Selected Samples*, ed. by H. Wainer, pp. 63–107. Springer-Verlag, New York, Reprinted in 2000, Mahwah, NJ: Lawrence Erlbaum Associates.
- HECKMAN, J. J., J. A. SMITH, AND N. CLEMENTS (1997): “Making the Most Out Of Programme Evaluations and Social Experiments: Accounting for Heterogeneity in Programme Impacts,” *Review of Economic Studies*, 64(221), 487–536.

- HECKMAN, J. J., AND E. J. VYTLACIL (1998): “Instrumental Variables Methods for the Correlated Random Coefficient Model: Estimating the Average Rate of Return to Schooling When the Return Is Correlated with Schooling,” *Journal of Human Resources*, 33(4), 974–987.
- (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*, ed. by C. Hsiao, K. Morimune, and J. L. Powell, pp. 1–46. Cambridge University Press, New York.
- (2005): “Structural Equations, Treatment Effects and Econometric Policy Evaluation,” *Econometrica*, 73(3), 669–738.
- (2007a): “Econometric Evaluation of Social Programs, Part I: Causal Models, Structural Models and Econometric Policy Evaluation,” in *Handbook of Econometrics*, ed. by J. Heckman, and E. Leamer, vol. 6B, pp. 4779–4874. Elsevier, Amsterdam.
- (2007b): “Econometric Evaluation of Social Programs, Part II: Using the Marginal Treatment Effect to Organize Alternative Economic Estimators to Evaluate Social Programs and to Forecast Their Effects in New Environments,” in *Handbook of Econometrics*, ed. by J. Heckman, and E. Leamer, vol. 6B, pp. 4875–5144. Elsevier, Amsterdam.
- IMBENS, G. W., AND J. D. ANGRIST (1994): “Identification and Estimation of Local Average Treatment Effects,” *Econometrica*, 62(2), 467–475.
- IMBENS, G. W., AND W. K. NEWY (2002): “Identification and Estimation of Triangular Simultaneous Equations Models Without Additivity,” Technical Working Paper 285, National Bureau of Economic Research.
- (2007): “Identification and Estimation of Triangular Simultaneous Equations Models Without Additivity,” Unpublished manuscript, Harvard University and MIT.
- MATZKIN, R. L. (2007): “Nonparametric Identification,” in *Handbook of Econometrics*, ed. by J. Heckman, and E. Leamer, vol. 6B. Elsevier, Amsterdam.
- NEWY, W. K., AND J. L. POWELL (2003): “Instrumental Variable Estimation of Nonparametric Models,” *Econometrica*, 71(5), 1565–1578.

- NEWKEY, W. K., J. L. POWELL, AND F. VELLA (1999): “Nonparametric Estimation of Triangular Simultaneous Equations Models,” *Econometrica*, 67(3), 565–603.
- ROY, A. (1951): “Some Thoughts on the Distribution of Earnings,” *Oxford Economic Papers*, 3(2), 135–146.
- WOOLDRIDGE, J. M. (1997): “On Two Stage Least Squares Estimation of the Average Treatment Effect in a Random Coefficient Model,” *Economics Letters*, 56(2), 129–133.
- (2003): “Further Results on Instrumental Variables Estimation of Average Treatment Effects in the Correlated Random Coefficient Model,” *Economics Letters*, 79(2), 185–191.
- (2007): “Instrumental Variables Estimation of the Average Treatment Effect in Correlated Random Coefficient Models,” in *Advances in Econometrics: Modeling and Evaluating Treatment Effects in Econometrics*, ed. by D. Millimet, J. Smith, and E. Vytlacil, vol. 21. Elsevier, Amsterdam, Forthcoming.