

Endogenous Semiparametric Binary Choice Models with Heteroscedasticity

Stefan Hoderlein*

Boston College

First Draft: February 15, 2006

This Draft: September 29, 2014

Abstract

In this paper we consider endogenous regressors in the binary choice model under a weak median exclusion restriction, but without further specification of the distribution of the unobserved random components. Our reduced form specification with heteroscedastic residuals covers various heterogeneous structural binary choice models. We employ a control function IV assumption to establish identification of a slope parameter β in the reduced form model by the mean ratio of derivatives of two functions of the instruments. We propose a direct estimator based on sample counterparts, and discuss the large sample behavior of this estimator. In particular, we show \sqrt{n} consistency and derive the asymptotic distribution. As a particularly relevant example of a structural model where no semiparametric estimator has of yet been analyzed, we consider the binary random utility model with endogenous regressors and heterogeneous parameters. Moreover, we propose tests for heteroscedasticity, overidentification and endogeneity. We analyze the small sample performance through a simulation study. An application of the model to discrete choice demand data concludes this paper.

Keywords: Nonparametric, Discrete Choice, Heteroscedasticity, Average Derivative, Median, Random Coefficients.

*Boston College, Department of Economics, 140 Commonwealth Ave, Chestnut Hill, MA 02467, USA, email: stefan.hoderlein@bc.edu. I have received helpful comments from Richard Blundell, Joel Horowitz, Simon Lee, Arthur Lewbel, Oliver Linton, Amil Petrin, Ed Vytlacil and from seminar participants in Brown, Boston College, Georgetown, Iowa, Maryland, Minnesota, NYU, Penn, Princeton, St. Gallen, Tilburg and ESEM Milano. I am also grateful to Martin Hackmann and Vitor Hadad for research assistance.

1 Introduction

The Model: The binary choice model constitutes a workhorse of modern Microeconometrics and has found a great many applications throughout applied economics. It is commonly treated in a latent variable formulation, i.e.

$$\begin{aligned} Y^* &= X'\beta + U \\ Y &= \mathbb{I}\{Y^* > 0\}, \end{aligned} \tag{1.1}$$

where Y^* is an unobserved continuously distributed random variable, in the classical choice literature often utility or differences in utility, X is a random K -vector of regressors, β is a K -vector of fixed coefficients, and $\mathbb{I}\{\cdot\}$ denotes the indicator of an event. Throughout much of the literature, and indeed in this paper, interest centers on the coefficient β which summarizes the effect of a set of regressors X on the dependent variable. If U is assumed independent of X , and U follows a certain parametric distribution then $\mathbb{E}[Y|X] = F_U(X'\beta)$, where F_U is the known parametric cdf of U , and estimation is straightforward via ML. Both assumptions are restrictive in many economic applications and have therefore come under critique. In particular, invoking these assumptions rules out that model (1.1) is the reduced form of individual behavior in a heterogeneous population, where parameters vary across the population in an unrestricted fashion, and it rules out endogeneity.

This paper aims at weakening these critical assumptions, while retaining an estimator with a simple and interpretable structure. In particular, it establishes interpretation and constructive identification of β under assumptions that are compatible with, e.g., a heterogeneous population characterized by an unrestricted distribution of random utility parameters. The weakening of assumptions is twofold: First, we do not want to place restrictive parametric or full independence assumptions on the distribution of the unobservables (or indeed any random variable in this model), and employ instead relatively weak median exclusion restrictions. Second, due to its paramount importance in applications we want to handle the case of endogenous regressors, e.g., we want to allow for X to be correlated with U . The estimator we propose has a simple, “direct” structure, akin to average derivative estimator (ADE). A characteristic feature of this class of estimators is that they use a control function instrumental variables approach for identification. The identification result is constructive in the sense that it can be employed to yield a \sqrt{n} consistent semiparametric estimator for β .

Main Identification Idea: Throughout this paper, we will be concerned with model (1.1). However, we view model (1.1) as a reduced form of a structural model in a heterogeneous population. As a consequence, we will also be concerned with providing an example for the interpretation of β when employing sensible independence restrictions in the heterogeneous structural model.

The independence restriction we are invoking in model (1.1) is a conditional median exclusion restriction. Specifically, we introduce a L random vector of instruments, denoted Z , and assume that they are related to X via

$$X = \vartheta(Z) + V, \tag{1.2}$$

where ϑ is a smooth, but unknown function of Z . For instance, in the special case where out of K continuously distributed regressors (X^1, \dots, X^K) only X^1 is endogenous and there is exactly one additional instrument denoted by Z^1 (so that $Z = (Z^1, X^2, \dots, X^K)'$), ϑ could be the mean regression $m_X(z) = \mathbb{E}[X|Z = z] = (\mathbb{E}[X^1|Z = z], X^2, \dots, X^K)'$, in which case V would be the mean regression residuals in the first equation. Alternatively, it could also be a vector containing the conditional α quantile of X^1 conditional on Z as first element, as would for instance arise in a triangular nonseparable model with one endogenous regressor¹. This classifies our way of dealing with unobserved heterogeneity as a triangular systems of equations, which is standard in nonseparable models (see, e.g., Matzkin (2005), Hoderlein (2005, 2011)), but rules out simultaneity in the structural model. Contrary to the recent control function literature, e.g., Imbens and Newey (2009), we do not employ a nonseparable model with monotonicity in a scalar unobservable, as it rules out random coefficients or other high dimensional unobservables in the first stage equation. Random coefficients are still covered by our assumptions, because we allow for heteroscedasticity of V , conditional on Z .

If we let the conditional median of U given $Z = z$ and $V = v$ be denoted by $k_{U|ZV}^{0.5}(z, v)$, then we may formalize our identifying assumption as

$$k_{U|ZV}^{0.5}(z, v) = g(v),$$

for all (z, v) in its support. What does this assumption mean in economic terms, and why is it a sensible assumption if we think of a heterogeneous population? In the fourth section, we show that this assumption is implied by a heterogeneous random utility model with endogeneity arising from omitted variables. In this case, the median exclusion restriction is implied if instruments are (jointly) independent of omitted variables and of V , but it holds also under weaker restrictions. What economic interpretation of β is implied by our assumptions? Taking the binary choice model with random utility parameters as an example, in the fourth section we establish that β has the interpretation of a local average structural derivative (see Hoderlein (2011) and Hoderlein and Mammen (2007, 2009)).

Given that we have devised a sensible identification restriction and defined an interesting structural parameter, the question that remains to be answered is how to actually identify and

¹While the identification analysis proceeds on this level of generality, for the large sample theory we specify ϑ to be the mean regression.

estimate this parameter. To answer this question, we introduce the following notation: Let $\bar{Y} = k_{Y|ZV}^{0.5}(Z, V) = \mathbb{I}\{\mathbb{P}[Y = 0|Z, V] < 0.5\}$ denote the conditional median of Y given Z and V , and assume that $\vartheta(z) = \mathbb{E}[X|Z = z]$. Then, under assumptions to be detailed below,

$$\beta = \mathbb{E} \left[[D_z \mathbb{E}[X|Z]']^{-1} \nabla_z \mathbb{E}[\bar{Y}|Z] B(Z) \right], \quad (1.3)$$

where ∇_z and D_z denote gradient and Jacobian, and $B(z)$ denotes a bounded weighting function to be defined below. Intuitively, the identification follows by a combination of arguments employed to identify average derivatives (see Powell, Stock and Stoker (1989), PSS, for short), and the chain rule, and is only up to scale. This identification principle is constructive, and yields in a straightforward fashion a sample counterparts estimator, see equation (3.3) below. Because of its direct structure, the estimator shares all advantages of direct estimators. In particular, the estimator is robust to misspecification and avoids computationally difficult optimization problems involving nonconvex objective functions. Moreover, the estimator is \sqrt{n} consistent, and has a standard limiting distribution.

Additional Contributions: Beyond constructing an estimator for a sensible parameter in a heterogeneous population, the flexibility of the model enables us to check the specification for several issues that have not been considered exhaustively, if at all, in the literature on this type of models. For instance, we propose powerful tests for endogeneity and heteroscedasticity. Another important issue we discuss is overidentification. As will turn out, in a general non-separable setup overidentification is markedly different from the issue in the linear framework. In addition to clarifying the concept, we propose a Hausman type test for overidentification. We develop a semiparametric notion of weakness of the instruments, and establish how our approach allows to mitigate the problem of weak instruments. Finally, we show that our approach allows to handle discrete and continuous endogenous regressors.

Literature: The binary choice model (1.1) with exogenous regressors has been analyzed extensively in the semiparametric literature, most often via single index models. Since this paper employs a direct estimator, our approach is related to contributions by PSS (1989), Hristache, Juditsky and Spokoiny (2001) and Chaudhuri, Doksum and Samarov (1997), to mention just a few. The main alternative are “optimization”, or M -, estimators for β , including semiparametric LS (Ichimura (1993)), semiparametric ML (Klein and Spady (1993)), and general M -estimators (Delecroix and Hristache (1997)). None of these of estimators can handle general forms of heteroscedasticity even in the exogenous setting, and to do so one has to employ maximum score type estimators, see Manski (1975). But these estimators have a slow convergence rate and a nonstandard limiting distribution, and only the estimator of Horowitz (1992) achieves almost \sqrt{n} convergence.

Related in particular is a line of work which estimates single index random coefficient mod-

els nonparametrically, in particular Ichimura and Thompson (1998), Hoderlein, Klemelä, and Mammen (2010) and Gautier and Kitamura (2013). All of these models focus on estimating the entire distribution of a vector of random coefficients nonparametrically. The present paper nests this class of models as special case - in this setup we focus on the vector of centrality parameters of the distribution of random coefficients. In addition, Ichimura and Thompson (1998) do not deal with endogeneity, while Hoderlein, Klemelä, and Mammen (2010) only discuss continuous outcome variables. The closest is hence the elegant paper of Gautier and Kitamura (2013), who, however, when dealing with endogeneity also assume a linear and homogenous structure in the IV equation.

In spite of the wealth of literature about model (1.1) in the exogenous case with fixed parameters, and the importance of the concepts of endogeneity and instruments throughout econometrics, the research on model (1.1) with endogenous regressors has been relatively limited. However, there are important contributions that deserve mentioning. For the parametric case, we refer to Blundell and Smith (1986) and Rivers and Vuong (1988). For the semiparametric case, Lewbel proposes the concept of special regressors, i.e. one of the regressors is required to have infinite support, which is essential for identification (Lewbel (1998)). Our approach is more closely related to the work of Blundell and Powell (2004, BP, for short). Like BP, we use a control function assumption to identify the model, but as already mentioned in a different fashion, as we allow for heteroscedasticity. This makes our approach also weakly related to other control function models in the semiparametric literature, most notably Newey, Powell and Vella (1998) and Das, Newey and Vella (2003). Related is also Matzkin (1992, 2005), who was the first to consider nonparametric identification in the exogenous binary choice model, and Bajari, Fox, Kim and Ryan (2008), who propose a nonparametric estimator of the distribution of random coefficients. Finally, our work is also related to Ai and Chen (2001), Vytlacil and Yildiz (2007), and in particular the “Local Instruments” approach of Heckman and Vytlacil (2005) and Florens, Heckman, Meghir and Vytlacil (2008) for analyzing treatment effects.

Organization of Paper: We start out by stating formally the assumptions required for identification of β , and provide a discussion. Moreover, we establish identification both in the heteroscedastic as well as the homoscedastic case (we require the latter among other things to test the random utility specification). This identification principle is constructive in the sense that it yields direct estimators through sample counterparts. We derive the median exclusion restriction formally, and establish the interpretation of β stated above. The asymptotic distribution of these estimators is in the focus of the third section. Specifically, we establish \sqrt{n} consistency to a standard limiting distribution². Beyond suggesting a \sqrt{n} consistent estima-

²This is in stark contrast to the exogenous binary choice model, where single index estimators only allow for very limited forms of heteroscedasticity (namely that the distribution of $U|X$ is only a function of the index

tor, the general identification principle is fruitful in the sense that it allows to construct tests for endogeneity, heteroscedasticity and overidentification, and this will be our concern in the fourth section. We then provide a structural example by considering the case of a linear random utility model with heterogeneous parameters. A simulation study underscores the importance of correcting for endogeneity and heterogeneity and will be discussed in the fifth section. In the sixth section, we apply our methods to a real world discrete choice demand application: We consider the decision to subscribe to cable TV, using data similar to Goolsbee and Petrin (2004). Finally, this paper ends with a conclusion and an outlook.

2 Identification of β via Median Restrictions on U

In this section we discuss the identification of β in the baseline scenario. We start out by introducing notation, stating and discussing the assumptions, and then sketching the main arguments in the proof of identification, before stating and discussing the main identification theorem. This theorem lends itself in a natural way to the construction of a sample counterparts estimator which will be discussed in the subsequent section.

Notation: Let the $K \times L$ matrix of derivatives of a K -vector valued Borel function $g(z)$ be denoted by $D_z g(z)$, and let $\nabla_z g(z)$ denote the gradient of a scalar valued function. Denote by $m_{A|B}(a, b) = \mathbb{E}[A|B = b]$ the conditional expectation of a random vector A given B , and let $k_{A|B}^\alpha(b)$ denote the conditional α -quantile of a random variable A given $B = b$, i.e. for $\alpha \in (0, 1)$, $k_{A|B}^\alpha(b)$ is defined by $\mathbb{P}(Y \leq k_{A|B}^\alpha(b)|B = b) = \alpha$. Moreover, let $f_A(a)$, $f_{AB}(a, b)$ and $f_{A|B}(a; b)$ be the marginal, joint and conditional Radon-Nikodym density with respect to some underlying measure μ , which may be the counting or the Lebesgue measure, (i.e., A may be discretely or continuously distributed). Define the nonparametric score $Q_b(a, b) = \nabla_b \log f_{A|B}(a; b)$. Let G^- denote the Moore-Penrose pseudo-inverse of a matrix G . Finally, let $c_k, k = 1, 2, \dots$ denote generic constants, and we suppress the arguments of the functions whenever convenient.

Assumptions: As already discussed in the introduction, the main idea is now that instead of running a regression using Y , we employ $\bar{Y} = k_{Y|ZV}^{0.5}(Z, V)$, i.e. the conditional median of Y given Z and V (which is the L_1 -projection of Y on $\mathcal{Z} \times \mathcal{V}$), and consider the mean regression, i.e., the L_2 -projection, of \bar{Y} on \mathcal{Z} . Consequently, we consider weighting functions defined on \mathcal{Z} only. In the following two subsections we first list and discuss all assumptions that specify the true population distribution and the DGP, and then establish the role they play in identifying β . We will focus henceforth on the case where $\vartheta(z) = m_{X|Z}(z)$, but the arguments hold very analogously for alternative $\vartheta(z)$. Readers less interested in the econometric details of this model

$X'\beta$), and only maximum score type estimators allow for heteroscedastic errors of general form (Manski (1975, 1985), Horowitz (1992)), but those do not achieve \sqrt{n} rate of convergence.

may skip these subsections, and proceed directly to the main result (theorem 1).

2.1 Assumptions

Assumption 1. *The data $(Y_i, X_i, Z_i), i = 1, \dots, n$ are independent and identically distributed such that $(Y_i, X_i, Z_i) \sim (Y, X, Z) \in \mathcal{Y} \times \mathcal{X} \times \mathcal{Z} \subset \mathbb{R}^{1+K+L}$. The joint distribution of (Y, X, Z) is absolutely continuous with respect to a σ -finite measure μ on $\mathcal{Y} \times \mathcal{X} \times \mathcal{Z}$ with Radon-Nikodym density $f_{YXZ}(y, x, z)$. The underlying measure μ can be written as $\mu = \mu_Y \times \mu_{XZ}$, where μ_{XZ} is the Lebesgue measure.*

Assumption 2. *The weighting function $B(z)$ is nonzero and bounded with compact support $\mathcal{B} \subset \mathcal{Z}$, where usually $\mathcal{Z} = \mathbb{R}^L$.*

Assumption 3. *$m_{X|Z}(z)$ is continuously differentiable in the components of z for all $z \in \text{Int}(\mathcal{B})$. $[D_z m_{X|Z}(z)]^-$ exists and every element is bounded from below for all $z \in \mathcal{B}$. $[D_z m_{X|Z}(z)]^-$ is square integrable on \mathcal{B} .*

Assumption 4. *$m_{Y|Z}(z)$ is continuously differentiable in the components of z for all $z \in \text{Int}(\mathcal{B})$. $D_z m_{Y|Z}(Z)$ is square integrable on \mathcal{B} . $g(z, v) = F_{U|V}(m_{X|Z}(z)' \beta + v' \beta; v) f_{V|Z}(v; z)$ is bounded in absolute value by a nonnegative integrable function $q(z)$, for all $z \in \mathcal{B}$.*

Assumption 5. *$\mathbb{E}[\bar{Y}|Z = z] = m_{\bar{Y}|Z}(z)$ is continuously differentiable in the components of z for all $z \in \text{Int}(\mathcal{B})$. $D_z m_{\bar{Y}|Z}(Z)$ is square integrable on \mathcal{B} . Moreover $0 < \mathbb{P}[\bar{Y} = 1|Z = z] < 1$ for all $z \in \mathcal{B}$.*

For the stochastic terms U and V , the following holds:

Assumption 6. *U and V are jointly continuously distributed.*

In addition, either of the following hold:

Assumption 7. *U is independent of Z given V .*

Assumption 8. *1. $k_{U|ZV}^{0.5}(Z, V) = g(V)$.*

For the second part, let $\tilde{V} = l(V) = -(g(V) + V' \beta)$.

2. \tilde{V} is independent of Z and absolutely continuously with respect to Lebesgue measure with Radon-Nikodym density $f_{\tilde{V}}$. $f_{\tilde{V}}(\varpi)$ is differentiable for all $\varpi \in \text{im}(l)$. Finally, $f_{\tilde{V}}(D_z m_{X|Z}(Z)' \beta)$ is absolutely integrable on \mathcal{B} .

Remark 2.1 - Discussion of Assumptions: Starting with assumption 1, we assume to possess continuously distributed instruments and regressors. Strictly speaking, we do not even require continuous instruments for identification, but only for the specific direct estimator we propose. Indeed, we conjecture that a direct estimator akin to Horowitz and Haerdle (1998) in the exogenous single index model may be devised, but this is beyond the scope of this paper. The *iid* assumption is inessential and may be relaxed to allow for some forms of stationary time series dependence. Note that unlike Blundell and Powell (2004), we do not require the support of $V|X$ to be invariant in X , which is why endogenous binary regressors are ruled out in their case. However, note that the assumption that $l(V)$ be independent of Z in assumption 8.2 effectively rules out binary endogenous regressors in our case, too.

Regarding the choice of weighting function B , due to assumption 2 we delete all observations outside a fixed multivariate interval I_z . As such, the weighting is unrestrictive and merely serves as a device to simplify already involved derivations below. It could be abandoned at the price of a vanishing trimming procedure which we do not pursue here because we want to focus on the innovative part and want to keep the exposition concise. In addition we require that $[D_z m_{X|Z}(z)]^-$ exists and is bounded on \mathcal{B} (cf. assumption 3), and hence we choose $B(z) = \mathbb{I}\{z \in I_z\} \mathbb{I}\{|\det [D_z m_{X|Z}(z) D_z m_{X|Z}(z)']| \geq b\}$, with $b > 0$. By choosing the weighting function and the region \mathcal{B} appropriately, we may ensure that the instruments are not weak in the sense that $\det [D_z m_{X|Z}(z) D_z m_{X|Z}(z)'] \geq b$ for some subset of \mathcal{Z} with positive measure. If we view the derivative in a linear regression of X on Z as an average derivative, it may be the case that instruments are on average not strongly related to endogenous regressors, but are quite informative for β in certain areas of \mathcal{Z} space. We consider it to be an advantage of our nonparametric approach that we can concentrate on those areas, and hence suggest that a similar weighting be performed in applications. However, in applications \mathcal{B} is usually not known, implying that a threshold b be chosen and $D_z m_{X|Z}(z)$ be pre-estimated³.

Particularly novel is assumption 8.1. Instead of the full independence of U and Z conditional on V assumed in assumption 7 (and implying the Blundell and Powell (2004) assumption $U \perp X|V$) this assumption (only) imposes a conditional location restriction. Hence it allows for all other quantiles of U than the median to depend on Z and V , and thus on X , in an arbitrary fashion, which as we have seen in the introduction is sensible when unobserved heterogeneity is modelled. Assumption 8.2 covers the case when \tilde{V} is independent of Z , in which case the function l need not be restricted at all, and we leave the case where this assumption does not

³The trimming becomes then dependent on estimated quantities. We skip the large sample theory of such an approach, because it adds little new insight and makes the analysis more involved. An interesting situation arises when the instruments are weak everywhere. We conjecture that we may derive a generalized inverse by some type of regularization, e.g. by constructing a matrix $[D_z m_{X|Z}(Z)]^*$ that is analogous to, say Ridge regression. However, we do leave the explicit behavior of such a model for future research.

hold for an extension (see section 4.2 below).

2.2 Essential Arguments in the Identification of β in the Heteroscedastic Case

To see how assumptions 1–8 help in identifying β , rewrite the model as follows

$$Y = \mathbb{I} \left\{ (m_{X|Z}(Z) + V)' \beta + U > 0 \right\}. \quad (2.1)$$

Note first that under assumption 8.1, the conditional median \bar{Y} becomes

$$\bar{Y} = \mathbb{I} \left\{ m_{X|Z}(Z)' \beta + k_{U|ZV}^{0.5}(Z, V) + V' \beta > 0 \right\} = \mathbb{I} \left\{ m_{X|Z}(Z)' \beta > l(V) \right\}, \quad (2.2)$$

as \mathbb{I} is a monotonic function. This very much resembles the standard model, but with \bar{Y} instead of Y . However, note two complications: first $\tilde{V} = l(V)$ may not be fully independent of Z , second, l is unknown. We now establish that β is nevertheless constructively identified in the setup where \tilde{V} is independent of Z , while we leave the more general case as an extension (see again section 4.2). We start by noting that due to assumption 8.2, we have

$$m_{\bar{Y}|Z}(z) = \mathbb{E} [\bar{Y}|Z = z] = \mathbb{P} \left\{ m_{X|Z}(z)' \beta > l(V) \right\} \quad (2.3)$$

by standard arguments. To focus now on the essential arguments, we consider only a compact set $\mathcal{B} \subset \mathcal{Z}$ and a nonzero and bounded weighting function $B(z)$ with support \mathcal{B} , see assumption 2. Since, for all $z \in \mathcal{B}$, $m_{\bar{Y}|Z}(z)$ and $m_{X|Z}(z)$ are continuously differentiable in all components of Z we obtain by the chain rule that

$$\nabla_z m_{\bar{Y}|Z}(Z) = f_{l(V)} \left(m_{X|Z}(Z)' \beta \right) D_z m_{X|Z}(Z)' \beta, \quad (2.4)$$

with probability one. This steps rules out that X contains a constant. Moreover, note that $f_{l(V)}$ is a scalar valued function. Next, we premultiply equation (2.4) by the generalized inverse $[D_z m_{X|Z}(Z)']^-$, which exists on \mathcal{B} due to assumption 3, and the weighting function $B(z)$ to obtain

$$[D_z m_{X|Z}(Z)']^- \nabla_z m_{\bar{Y}|Z}(Z) B(Z) = \beta f_{l(V)} \left(m_{X|Z}(Z)' \beta \right) B(Z), \quad (2.5)$$

or, upon taking expectations,

$$\beta c_1 = \mathbb{E} \left\{ [D_z m_{X|Z}(Z)']^- \nabla_z \mathbb{E} [\bar{Y}|Z] B(Z) \right\}, \quad (2.6)$$

where $c_1 = \mathbb{E} [f_{l(V)} \left(m_{X|Z}(Z)' \beta \right) B(Z)]$. From now on, we will tacitly suppress this constant, i.e., set it to unity, or, put differently, identification is only up to scale. This last step is warranted, because the elementwise square integrability of all functions on \mathcal{B} (assumption 5), together with Cauchy-Schwarz ensures that the expectations exist.

2.3 Main Identification Results

The following theorem summarizes the discussion in the previous section and in the appendix:

Theorem 1. (i) *Let the true model be as defined in 1.1 and 1.2, and suppose that assumptions 1–3, 5, 6, 8.1 and 8.2 hold. Assume further that $\mathbb{E} [f_{\bar{V}|Z}(m_{X|Z}(Z)' \beta; Z) B(Z)] = 1$. Then β is identified by the relationship*

$$\beta = \mathbb{E} \left[[D_z m_{X|Z}(Z)]^{-1} \nabla_z \mathbb{E} [\bar{Y}|Z] B(Z) \right]. \quad (2.7)$$

(iii) *If we strengthen the conditional median independence assumption 8 to the full independence assumption 7 and assume that assumption 4 holds, we obtain that in addition to (2.7), β is (up to scale) identified by*

$$\mathbb{E} \left[[D_z m_{X|Z}(Z)]^{-1} \nabla_z m_{Y|Z}(Z) B(Z) \right], \quad (2.8)$$

Remark 2.2 - Interpretation of Theorem 1: First, consider the scenario where V and Z are independent which gives rise to (2.7). β is identified by a weighted average ratio of derivatives, involving the derivatives of the function $m_{X|Z}$, and of the mean regression of \bar{Y} (i.e., the conditional median given Z and V), on Z alone. Note that the control residuals V do not appear in this equation, however, the model relies on correct specification of the conditional median restriction and on the correct specification of the first stage functions as $m_{X|Z}$. Allowing the first stage model to be a conditional mean or a conditional quantile enables the applied researcher to choose between various specifications of the IV equation in order to select the one with the best economic interpretation .

It is instructive to compare the heteroscedastic case with the case when $U \perp Z|V$, i.e., the full (conditional) independence assumption 7. This assumption obviously implies assumption 8, so that equation (2.7) remains valid. But we obtain in addition that β is identified up to scale by (2.8). Under full independence, we have thus (at least) two estimating equations: we could either use directly an L_2 -projection of Y on Z , or use a two projection strategy, where we use L_1 , respectively, L_2 -projections of Y on (Z, V) in the first stage, and then use a L_2 -projection in the second stage. As shown below, we are able to obtain a test for heteroscedasticity out of this comparison.

3 A Sample Counterpart Estimator: Asymptotic Distribution and Conditions for \sqrt{n} Consistency

3.1 The Case for Direct Estimation

As mentioned above, the identification principle does not necessarily imply that we have to use a direct estimator. Indeed, in the case where assumption 8.2 holds (i.e., $\tilde{V} \perp Z$) which is not

the focus of this paper, we could base an optimization estimator on equation (2.2), i.e.

$$\bar{Y} = \mathbb{I} \left\{ m_{X|Z}(Z)' \beta > \tilde{V} \right\}. \quad (3.1)$$

However, there are a number of reasons to use direct estimators here. Several have already been mentioned: First, direct estimators are natural because they build upon sample counterparts of the identification result. Consequently, the role the assumptions play and the mechanics of the estimator are well-understood and transparent, which makes them accessible to applied people. Moreover, several related issues (like overidentification) can be discussed straightforwardly. Second, they avoid the optimization of a highly nonlinear function, which both may not lead to global maxima (sometimes not even to well defined ones, if the semiparametric likelihood is flat), and may be computationally very expensive. Finally, and perhaps most importantly, in the case discussed in this paper where \tilde{V} is not fully independent of Z , it is not even clear how to construct an optimization estimator.

There are, however, also reasons that speak against the use of kernel based direct estimators. One of the theoretical arguments against them is that they require higher order smoothness assumptions, as will be obvious below. Note, however, that in the general setup with unrestricted (nonparametric) IV equation $X = m_{X|Z}(Z) + V$, there is a “diminished smoothness gap”. Any optimization estimator depends on an estimator \hat{V} of V as a regressor. In the general nonparametric setup, this is a function of a nonparametric estimator for $m_{X|Z}$. Using results in Newey (1994), it is straightforward to see that for a \sqrt{n} consistent estimator of β we require that $\mathbb{E} [\hat{m}_{X|Z}] - m_{X|Z} = o_p(n^{-1/2})$ for the “no bias condition” to hold. In the kernel case this is, however, only possible under smoothness assumptions which are very similar to the ones we require to hold for our direct estimator, in particular undersmoothing.

The second main drawback of direct estimators is the lack of efficiency compared to optimization estimators. Improving the efficiency, however, is possible through a so called one step efficient estimator, taking the direct estimator as a starting point. Alternatively, as in Newey and Stoker’s (1994) analysis of the weighted average derivative estimators, we can define optimal weights.

3.2 A Sample Counterpart Estimator for β

In this section, we discuss the behavior of a sample counterpart estimator to (2.7). In this analysis, we take V as given and do not treat the effect of pre-estimation of V . However, given previous work on nonparametric regression with generated regressors, this is not restrictive under appropriate smoothness assumption which we are invoking anyway, see e.g., Sperlich (2009).

The first impression from looking at

$$\beta = \mathbb{E} \left[[D_z m_{X|Z}(Z)]^{-1} \nabla_z \mathbb{E} [\bar{Y}|Z] B(Z) \right]. \quad (3.2)$$

is that due to the non-smoothness in \bar{Y} , no fast enough first step estimator can be devised for an average derivative type estimator to become root n estimable. However, this is not the case. To see how the estimator is constructed and understand why it is \sqrt{n} consistent, note first that since Y is binary,

$$\bar{Y} = k_Y^{0.5}|_{ZV}(Z, V) = \mathbb{I} \{ \mathbb{P} [Y = 0|Z, V] < 0.5 \},$$

and consequently $\mathbb{E} [\bar{Y}|Z] = \mathbb{P} [\mathbb{P} [Y = 0|Z, V] < 0.5|Z]$. This suggests estimating $\nabla_z \mathbb{E} [\bar{Y}|Z = z]$ via

$$\sum_j \nabla_z W_j(z) \mathbb{I} \{ \hat{P}_j < 0.5 \},$$

where $W_j(z)$ are appropriate Kernel weights, e.g., $[\sum_j \mathcal{K}_{hj}(z)]^{-1} \mathcal{K}_{hj}(z)$, $\mathcal{K}_{hj}(z) = h^{-L} \mathcal{K}((Z_j - z)/h)$ and $\mathcal{K}((Z_j - z)/h) = \prod_{l=1, \dots, L} K((Z_j^l - z^l)/h)$ is a standard L -variate product kernel with standard univariate kernel function K . Moreover, \hat{P}_j denotes an estimator of $P_j = p(Z_j, V_j) = \mathbb{P} [Y_j = 0|Z_j, V_j]$. The problem with this estimator is that the pre-estimator \hat{P}_j appears within the nondifferentiable indicator, resulting in a potentially very difficult pre-estimation analysis. To improve upon the tractability of the problem, we replace the indicator \mathbb{I} by a smooth indicator, i.e., $\mathbb{K}(\xi) = \int_{-\infty}^{\xi} K(t) dt$. Then, a straightforward sample counterpart estimator to $\beta = \mathbb{E} \left[[D_z m_{X|Z}(Z)]^{-1} \nabla_z \mathbb{E} [\bar{Y}|Z] B(Z) \right]$, may be defined as follows:

$$\hat{\beta}_H = n^{-1} \sum_i [D_z \hat{m}_{X|Z}(Z_i)]^{-1} \sum_{j \neq i} \nabla_z W_j(Z_i) \mathbb{K} \left\{ \left(\hat{P}_j - 0.5 \right) / h \right\} B(Z_j), \quad (3.3)$$

where the subscript H indicates ‘‘heterogeneity’’. As is shown formally in theorems 2 and 3 below, the main result of this section is that under appropriate assumptions

$$\sqrt{n} \left(\hat{\beta}_H - \beta \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_H),$$

where Σ_H is defined as $\Sigma_H = \mathbb{E} \left(\sum_{k=1}^3 \sigma_k \sigma_k' \right) + 2 \mathbb{E} \left(\sigma_2 \sigma_3' \right) - \beta \beta'$, and

$$\begin{aligned} \sigma_1 &= [D_z m_{X|Z}(Z_i)]^{-1} \nabla_z m_{\bar{Y}|Z}(Z_i) B(Z_i), \\ \sigma_2 &= [D_z m_{X|Z}(Z_i)]^{-1} f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) V_i' [D_z m_{X|Z}(Z_i)]^{-1} \nabla_z m_{\bar{Y}|Z}(Z_i) B(Z_i), \\ \sigma_3 &= [D_z m_{X|Z}(Z_i)]^{-1} f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) (\bar{Y}_i - m_{\bar{Y}|Z}(Z_i)) B(Z_i). \end{aligned} \quad (3.4)$$

To understand the large sample behavior of this estimator, rewrite $\hat{\beta}_H$ as $\hat{\beta}_H = T_{1n} + T_{2n} + T_{3n}$,

where

$$\begin{aligned}
T_{1n} &= n^{-1} \sum_i [D_z \hat{m}_{X|Z}(Z_i)']^{-1} \sum_{j \neq i} \nabla_z W_j(Z_i) \bar{Y}_j B(Z_j), \\
T_{2n} &= n^{-1} \sum_i [D_z \hat{m}_{X|Z}(Z_i)']^{-1} \sum_{j \neq i} \nabla_z W_j(Z_i) [\mathbb{K}\{(P_j - 0.5)/h\} - \mathbb{I}\{P_j < 0.5\}] B(Z_j), \quad (3.5) \\
T_{3n} &= n^{-1} \sum_i [D_z \hat{m}_{X|Z}(Z_i)']^{-1} \sum_{j \neq i} \nabla_z W_j(Z_i) \left[\mathbb{K}\{(\hat{P}_j - 0.5)/h\} - \mathbb{K}\{(P_j - 0.5)/h\} \right] B(Z_j).
\end{aligned}$$

In this decomposition, T_{1n} is the leading term. It will dominate the asymptotic distribution. Its large sample behavior can be established using theorem 2, which also covers the sample counterparts estimator in the case defined by assumption 7, which yields the estimator:

$$\hat{\beta}_1 = \frac{1}{n} \sum_i [D_z \hat{m}_{X|Z}(Z_i)']^{-1} \nabla_z \hat{m}_{Y|Z}(Z_i) B(Z_i). \quad (3.6)$$

Hence, we will first discuss the behavior of $\hat{\beta}_1$ and T_{1n} . We proceed then by providing assumptions under which T_{2n} and T_{3n} tend to zero faster than the leading term. Essentially, these conditions are higher order smoothness conditions on the conditional cdf $F_{P|Z}$ and on f_Z , as well as the corresponding restrictions on the kernel (i.e., to be of higher order), so that fast enough rates of convergence are obtained. In this paper, we will not discuss the pre-estimation of V_i , again to focus on the essentials. We simply note that this can be neglected because it produces additional higher order terms if one assumes enough smoothness in the regression determining V_i (which we henceforth tacitly assume). In the simulations, we show that the repeated averaging removes the variance part of the estimator almost completely in finite samples.

3.3 The Large Sample Behavior of $\hat{\beta}_1$

When discussing the estimation of β using any form of regression it is important to clarify the properties of details of the estimator (3.6). This concerns in particular the kernel and bandwidth. As mentioned above, we use a product kernel in all regressions. Therefore we formulate our assumptions for the one-dimensional kernel functions K . To simplify things further, instead of a bandwidth vector $\mathbf{h} \in \mathbb{R}^L$, we assume that we have only one single bandwidth for each regression, denoted h . We shall make use of the following notation: Define kernel constants

$$\mu_k = \int u^k K(u) du \quad \text{and} \quad \kappa_k^2 = \int u^k K(u)^2 du.$$

In principle, we also have two bandwidths to consider, one in estimating $m_{X|Z}$, and one in estimating $m_{Y|Z}$. However, since the estimation problems are symmetric (i.e., in particular both mean regressions share the same regressors and have thus the same dimensionality), we

assume for ease of exposition the same kernel and the same bandwidth, denoted by K and h , in both regressions. Our assumptions regarding kernel and bandwidth are standard for average derivative estimation (cf. PSS); however, they mean that for any $L \geq 2$ we require higher order kernels.

Assumption 9. *Let $r = (L+4)/2$ if L is even and $r = (L+3)/2$ if L is odd. All partial derivatives of $\mathbb{E}[X|Z = z]$, $\mathbb{E}[Y|Z = z]$ and $f_Z(z)$ of order $r+1$ exist for all $z \in \mathcal{B}$. Moreover, the expectations of $[D_z m_{X|Z}(Z)]^- BY_l(Z)$ and $[D_z m_{X|Z}(Z)]^- BX_l(Z) [D_z m_{X|Z}(Z)]^- \nabla_z m_{Y|Z}(Z)$ exist for all $l = 1, \dots, r$, where BY_l (resp., BX_l) contains sums of products of all partial derivatives of $m_{Y|Z}$ and f_Z (resp. $m_{X|Z}$ and f_Z) such that the combined order of derivatives of the product is at most $l+1$.*

Assumption 10. *The one-dimensional kernel is Lipschitz continuous, bounded, has compact support, is symmetric around 0 and of order r (i. e. $\mu_k = \int u^k K(u) du = 0$ for all $k < r$ and $\int u^r K(u) du < \infty$).*

Assumption 11. *As $n \rightarrow \infty$, $h \rightarrow 0$, $nh^{L+2} \rightarrow \infty$ and $nh^{2r} \rightarrow 0$.*

The following theorem establishes asymptotic normality of the appropriate sample counterpart estimators, first in the case defined by assumption 7, which as already mentioned repeatedly can also be used to establish the large sample behavior of the leading term of $\hat{\beta}_H$.

Theorem 2. *Let the true model be as defined in 1.1 and 1.2. Suppose assumptions 1–4, 6–7, and 9–11 hold. Assume further that $V \perp Z$. For scale normalization, assume $\mathbb{E}[f_{U|V}(X'\beta; V) B(Z)] = 1$. Then,*

$$\sqrt{n} \left(\hat{\beta}_1 - \beta \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_1)$$

where

$$\Sigma_1 = \mathbb{E} \left(\sum_{k=1}^3 \sigma_k \sigma_k' \right) + 2\mathbb{E}(\sigma_2 \sigma_3') - \beta \beta'$$

and

$$\begin{aligned} \sigma_1 &= [D_z m_{X|Z}(Z_i)]^- \nabla_z m_{Y|Z}(Z_i) B(Z_i) \\ \sigma_2 &= [D_z m_{X|Z}(Z_i)]^- f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) V_i' [D_z m_{X|Z}(Z_i)]^- \nabla_z m_{Y|Z}(Z_i) B(Z_i) \\ \sigma_3 &= [D_z m_{X|Z}(Z_i)]^- f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) (Y_i - m_{Y|Z}(Z_i)) B(Z_i) \end{aligned}$$

Remark 3.1 – Discussion of Theorem 2: This results shows a number of parallels to PSS in the case of exogenous ADEs. Similar to PSS, we obtain \sqrt{n} consistency of our estimator for β , and we may be able to eliminate the bias under similar assumptions on the rate of convergence, as detailed in assumptions 9 and 11. The variance in turn is a more complicated

expression, but shares similar features with the PSS result, in particular in the first two terms. This result is of interest in its own right, but also serves as building block when we discuss the conditions under which we may include first stage projections of Y like the median regression which is required to deal with heteroscedasticity.

Remark 3.2 – Estimating Σ_1 : Estimation of the variance components is straightforward by sample counterparts. For instance, an estimator for $\Sigma_1^{23} = \mathbb{E}(\sigma_2\sigma_3')$ is given by

$$\begin{aligned} \widehat{\Sigma}_1^{23} &= n^{-1} \sum \hat{f}_Z(Z_i)^{-2} [D_z \hat{m}_{X|Z}(Z_i)']^{-} \nabla_z \hat{f}_Z(Z_i) (Y_i - \hat{m}_{Y|Z}(Z_i)) \\ &\quad \times \left\{ [D_z \hat{m}_{X|Z}(Z_i)']^{-} \nabla_z \hat{f}_Z(Z_i) (X_i - \hat{m}_{X|Z}(Z_i))' [D_z \hat{m}_{X|Z}(Z_i)']^{-} \nabla_z \hat{m}_{Y|Z}(Z_i) \right\}' B(Z_i), \end{aligned}$$

where the hats denote appropriate sample counterpart estimators. Consistency of this estimator can essentially be shown by appealing to a law of large numbers, but this analysis is beyond the scope of this paper.

3.4 The Large Sample Behavior of $\hat{\beta}_H$

We now extend theorem 2 to the heteroscedastic case. To treat heteroscedasticity, as outlined above we suggest the estimator

$$\hat{\beta}_H = n^{-1} \sum_i [D_z \hat{m}_{X|Z}(Z_i)']^{-} \sum_{j \neq i} \nabla_z W_j(Z_i) \mathbb{K} \left\{ \left(\hat{P}_j - 0.5 \right) / h \right\} B(Z_j),$$

Recall the decomposition $\hat{\beta}_H = T_{1n} + T_{2n} + T_{3n}$ in (3.5). The first term T_{1n} can be handled along exactly the same lines as the estimator in theorem 2, using some minor modifications in assumptions. It remains to be shown that the terms T_{2n} and T_{3n} tend to zero faster. To this end, we have to be precise about details of the estimator $\hat{\beta}_H$. First, there are several bandwidths: There is a bandwidth associated with the $\mathbb{K}\{\cdot\}$ function, as well as smoothness parameters when estimating $P_j = p(Z_j, V_j)$. To distinguish between the different kernels and bandwidths, we call the derivative of $\mathbb{K}\{\cdot\}$ K_1 , a kernel with bandwidth h_1 and order r_1 , and the univariate elements of a product kernel employed in the estimation of p as K_2 , with bandwidth h_2 and order r_2 .

Assumption 12. K_1 and K_2 are continuous, bounded, compactly supported, and symmetric functions of order r_1, r_2 (i. e. $\int u^k K(u) du = 0$ for all $k < r$ and $\int u^r K(u) du < \infty$).

Assumption 13. Let $r = (L + 4)/2$ if L is even and $r = (L + 3)/2$ if L is odd. All partial derivatives of $F_{P|Z}$ and $f_Z(z)$ of order $r + 1$ exist for all $z \in \mathcal{B}$. Moreover, the expectations of $[D_z m_{X|Z}(Z)]^{-} BF_l(Z)$ and $[D_z m_{X|Z}(Z)]^{-} BF_l(Z) [D_z m_{X|Z}(Z)]^{-} \nabla_z m_{Y|Z}(Z)$ exist for all $l = 1, \dots, r$, where BF_l contains sums of products of all partial derivatives of $F_{P|Z}$ and f_Z such that the combined order of derivatives of the product is at most $l + 1$.

Assumption 14. f_{ZV} is bounded and has bounded first partial derivatives with respect to all components of z , for all $z \in \mathcal{B}$.

Assumption 15. As $n \rightarrow \infty$, $h_1, h_2 \rightarrow 0$, $nh_1, nh_2^{L+\dim(V)+2} \rightarrow \infty$ and $nh_1^{2r_1}, nh_2^{L+\dim(V)+4} \rightarrow 0$.

Note that we require higher order smoothness conditions on $F_{P|Z}$ and f_Z which - in connection with higher order kernels - ensure that the bias terms $\sqrt{n}T_{2n}$ and $\sqrt{n}T_{3n}$ are $o_p(1)$.

Theorem 3. Let the true model be as defined in 1.1 and 1.2, and suppose that assumptions 1–3, 5–6, 8.1–8.2a, 9–15 are true. Assume further $\mathbb{E}[f_{\bar{V}|Z}(m_{X|Z}(Z)' \beta; Z) B(Z)] = 1$ holds. Then, $\sqrt{n}T_{2n} = o_p(1)$, $\sqrt{n}T_{3n} = o_p(1)$, and $\sqrt{n}(\hat{\beta}_H - \beta) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_H)$, where Σ_H is defined in equation (3.4).

This theorem characterizes the large sample behavior of our estimator. Under the smoothness and higher order bias reduction assumptions, it essentially behaves like the independence case estimator $\hat{\beta}_1$, with Y replaced by \bar{Y} , i.e., with known conditional median.

4 Examples and Extensions

4.1 A Structural Example: Binary Demand Decisions in a Heterogeneous Population

The question that should be answered for any reduced form microeconomic model is how it can be derived from individual behavior in a heterogeneous population. To answer this identification question for the one defined through (1.1), we start out with a general nonseparable model of a heterogeneous population as in Hoderlein (2011) or Hoderlein and Mammen (2007, 2009). The most general version of (1.1) has the structural unobservables (e.g., preferences) influencing the latent variable in a nonseparable fashion, i.e. $Y^* = \phi(X, A)$, where $A \in \mathfrak{A}$ denotes the unobservables. Here \mathfrak{A} is a Borel space, e.g., the space of piecewise continuous utility functions. Note that A may include objects like preferences, but also other omitted determinants. In our example, we denote the former by A_1 , while the remainder of A is denoted by A_2 . In discrete choice demand analysis for instance, A_2 are often omitted characteristics of the product.

Specific Model: While we could proceed to discuss the model on this level of generality, in this paper we restrict ourselves to linear models on individual level, largely because linear models are the dominating class of models in economic applications, and leave the most general case for a companion paper. A linear heterogeneous population with omitted variables A_2 may

then be formalized through a random coefficient model, i.e.,

$$\begin{aligned} Y^* &= X'\beta(A_1) + A_2'\gamma(A_1) \\ Y &= \mathbb{I}\{Y^* > 0\}, \end{aligned} \tag{4.1}$$

where $\theta(A_1) = (\beta(A_1)', \gamma(A_1)')$ is a mapping from the space of unobservables (say, preferences) $\mathfrak{A}_1 \subseteq \mathfrak{A}$ into \mathbb{R}^K . Since we assume the random elements A_1 (in our example, preferences) to vary across the population, $\theta(A_1)$ varies across the population, too. This model admits a reduced form representation as (1.1). What are now plausible stochastic conditions that we would like to impose on the reduced form (1.1) to identify β , and how can they be derived from restrictions in the structural model (4.1), and in particular on the random coefficients $\theta(A_1)$?

We answer this question under the assumption that the endogeneity arises from potential correlation of X and A_2 only, and that A_1 , e.g., the unobservable preferences that determine the parameters, are independent from all economic variables in the system, i.e. $A_1 \perp (X, Z, A_2)$. In discrete choice demand analysis for instance, this endogenous regressor is the own price of the good, which is assumed to be correlated with its omitted unobserved characteristics contained in A_2 , see Goolsbee and Petrin (2004), and Berry and Haile (2009), unless there is only a single good (and the outside alternative). As is straightforward to show by standard arguments in this literature, the triangular structure we propose (with additive errors in the pricing equation) arises if a monopolistic firm assumes the market share of a product to arise from aggregation of a linear probability model. We think that this is plausible as an approximation of real world behavior, however, if one believes that firms assume individuals' random utility model to follow, say, a probit or logit model then this has the consequence that the pricing equation cannot have an additive error, see Berry and Haile (2009). While we still believe the additive error specification to be a reasonable approximation in this case, we would like to point out this potential limitation.

Returning to the identification of our model, an unnecessarily strong, but economically plausible identifying independence restriction is the independence of instruments Z from all unobservables in the system, i.e., $Z \perp (A, V)$. Continuing our discrete choice example, the discrete choice literature suggests to use the wholesale price, franchise fees, or other regional supply side characteristics of a market as instruments. It is plausible that these instruments are independent of individual preferences and omitted characteristics of the product. This assumption implies that $X \perp A_2|V$, and recall that our maintained hypothesis is that $A_1 \perp (X, V, A_2)$ which is implied by $A_1 \perp (X, Z, A_2)$, see also Petrin and Train (2006) for similar arguments, and discussions on when an additive IV equation arises.

Interpretation of β : The following result states that these independence conditions imply an exclusion restriction that defines a sensible centrality parameter of the distribution of random

coefficients in (4.1). For the result, we require the notation $U = X'(\beta(A_1) - \beta) + A_2'\gamma(A_1)$, and let $\mathbb{E}_{XV}[\cdot]$ denotes integration over the distribution of (X, V) .

Theorem 4. *Let the model defined by equations (1.1) and (1.2) be the reduced form of the structural model defined in equations (1.2) and (4.1). Suppose that $A_1 \perp (X, Z, A_2)$ and $Z \perp (A, V)$ hold. Assume further that the median exclusion restriction $k_{U|ZV}^{0.5}(Z, V) = g(V)$ holds, and that the conditional median of $Y^*|Z, V$ satisfies the regularity conditions in Hoderlein and Mammen (2007). Then we obtain that $k_{U|XV}^{0.5}(X, V) = g(V)$ and*

$$\beta = \mathbb{E}_{XV} \left[\mathbb{E} \left[\beta(A_1) | X, V, Y^* = k_{Y^*|XV}^{0.5}(X, V) \right] \right]. \quad (4.2)$$

Therefore, if we assume the median exclusion restriction $k_{U|ZV}^{0.5}(Z, V) = g(V)$, we obtain that the coefficient β has the interpretation of a local average structural derivative, i.e., $\beta = \mathbb{E} \left[\beta(A_1) | X = x, V = v, Y^* = k_{Y^*|XV}^{0.5}(x, v) \right]$, for all $(x, v) \in \text{supp}(X) \times \text{supp}(V)$, even without assuming any symmetry assumption on the distribution of the error. Due to the linear random coefficient structure with exogenous A_1 , this quantity is invariant to changes in x, v , and hence we may integrate over x, v , keeping the quantile of the unobservable latent variable fixed at the median, i.e., at the center of the conditional distribution. If we identify this center of the distribution with a type of individuals (the “median” person), then we may speak of β as an average structural effect for this type⁴. Another more statistical interpretation of (4.2) is that of a best approximation to the underlying heterogeneous coefficient $\beta(A_1)$, conditioning on all the information that we have to our disposal in the data⁵.

4.2 Heteroscedasticity in the IV Equation

In this subsection we investigate the case where the residuals in the IV equation are not fully independent, as would for instance arise in the case of heteroscedasticity. We therefore replace assumption 8.2 by the following condition:

Assumption 16. *Recall that $\tilde{V} = l(V) = -(g(V) + V'\beta)$. Assume that there is one endogenous regressor X^k , and l is a continuous piecewise invertible function. Moreover, $f_{V|Z}(v, z)$ and its partial derivatives wrt the components of z are bounded on \mathcal{B} from below and above, i.e. $c_1 > \sup_{(v,z) \in \text{supp}(V) \times \mathcal{B}} f_{V|Z}(v, z) \geq \inf_{(v,z) \in \text{supp}(V) \times \mathcal{B}} f_{V|Z}(v, z) = c_2 > 0$, and $c_3 > \sup_{(v,z) \in \text{supp}(V) \times \mathcal{B}} \left\| \nabla_z f_{V|Z}(v, z) \right\| \geq \inf_{(v,z) \in \text{supp}(V) \times \mathcal{B}} \left\| \nabla_z f_{V|Z}(v, z) \right\| = c_4 > 0$. Finally, let*

⁴“Average” here means more precisely to average over the distribution of X, V

⁵See Hoderlein and Mammen (2007) for a related discussion in the case of a continuous dependent variable. As already mentioned, this result could be generalized to models of the form $Y^* = \phi(X, A) = m(X) + U$, with $U = \phi(X, A) - m(X)$ and $k_{U|XV}^{0.5}(X, V) = l(V)$, but due to the lack of relevance for applications we desist from discussing this more general case here.

$Q_z(V, Z)$ be absolutely integrable on $\text{supp}(V) \times \mathcal{B}$, and let $\tau(z) = \mathbb{E} [\bar{Y} Q_z(V, Z) | Z = z]$ be square integrable on \mathcal{B} .

With these modified assumptions, we are now in the position to state the effect of a heterogeneous IV equation on the identification of β :

Theorem 5. (i) Let the true model be as defined in 1.1 and 1.2, and suppose that assumptions 1–3, 5–6, 8.1, and 16 hold. Then we obtain that β is identified up to scale by

$$\beta = \mathbb{E} \left[[D_z m_{X|Z}(Z)]' \left\{ \nabla_z \mathbb{E} [\bar{Y} | Z] - \mathbb{E} [\bar{Y} Q_z(V, Z) | Z] \right\} B(Z) \right], \quad (4.3)$$

where $Q_z(V, Z)$ denotes the nonparametric score $\nabla_z \log f_{V|Z}(V; Z)$.

The proof can be found in the appendix. Observe the large similarities with the previous identification result. Indeed, the expression is identical, if it were not for the “bias correction” $\mathbb{E} [\bar{Y} Q_z(V, Z) | Z]$, which accounts for heterogeneity (i.e., higher order dependence) in the IV equation, and measures the correlation between the fitted values \bar{Y} , and the score $Q_z(V, Z) = \nabla_z \log f_{V|Z}(V; Z) = \nabla_z f_{V|Z}(V; Z) / f_{V|Z}(V; Z)$.

When moving to estimation by sample counterparts, the latter terms is straightforwardly estimated.

$$\sum_j W_j(z) \mathbb{K} \left\{ \left(\hat{P}_j - 0.5 \right) / h \right\} \hat{Q}_z(V, Z),$$

where \hat{Q} is an obvious estimator for $\nabla_z f_{V|Z}(V; Z) / f_{V|Z}(V; Z) = \nabla_z f_{VZ}(V, Z) / f_{VZ}(V, Z) - \nabla_z f_Z(Z) / f_Z(Z)$, i.e., Kernel density estimator and derivative of the Kernel density estimator of the densities $f_Z(Z)$ and $f_{V,Z}(V, Z)$. By similar arguments as in the proof of the large sample results (i.e., theorems 2 and 3), under appropriate smoothness assumptions on P_j and $f_{VZ}(V, Z)$ the terms involving pre-estimated terms vanish, as does the term involving the difference $\mathbb{K} - \mathbb{I}$, so that when considering the variance, we can replace \mathbb{K} , \hat{P}_j and \hat{Q}_z by their population counterparts.

Again by similar arguments as above, this adds one more variance components and modifies the existing ones in an obvious way. More specifically, we have $\Sigma_{HH} = \mathbb{E} \left(\sum_{k=1}^4 \sigma_k \sigma_k' \right) + 2 \sum_{j>k>1} \mathbb{E} \left(\sigma_k \sigma_j' \right) - \beta \beta'$, with

$$\begin{aligned} \sigma_1 &= [D_z m_{X|Z}(Z_i)]' \left[\nabla_z m_{\bar{Y}|Z}(Z_i) - \mathbb{E} [\bar{Y} Q_z(V, Z) | Z] \right] B(Z_i) \\ \sigma_2 &= [D_z m_{X|Z}(Z_i)]' f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) V_i' [D_z m_{X|Z}(Z_i)]' \left[\nabla_z m_{\bar{Y}|Z}(Z_i) - \mathbb{E} [\bar{Y} Q_z(V, Z) | Z] \right] B(Z_i) \\ \sigma_3 &= [D_z m_{X|Z}(Z_i)]' f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) (\bar{Y}_i - m_{\bar{Y}|Z}(Z_i)) B(Z_i) \\ \sigma_4 &= [D_z m_{X|Z}(Z_i)]' \left(\mathbb{E} [\bar{Y}_i Q_z(V_i, Z_i) | Z_i] - \bar{Y}_i Q_z(V_i, Z_i) \right) B(Z_i) \end{aligned}$$

4.3 Specification Testing

4.3.1 Overidentification: Issue and Test

The first question that we can analyze within our framework is how to treat overidentification if we have more instruments than regressors. In the linear model, overidentification allows to delete instruments and recover β by various different estimators that always only use a subset of instruments. In the (X, V) projection of the Blundell and Powell (2003) approach, as already noted by the authors a similar feature is missing. In our setup it may be introduced, and the linear model result may be better understood. We discuss in the following the full independence case, but all arguments may be trivially extended to the heteroscedastic case random coefficients case.

If we return to the theorem 1 and the associated assumptions, we see that β would be identified by taking the derivatives w.r.t. any subset of instruments Z_1 such that $Z = (Z'_1, Z'_{-1})'$ and $[D_{z_1} m_{X|Z}(z) D_{z_1} m_{X|Z}(z)']$ would be nonsingular for all $z \in \mathcal{B}$. By similar arguments as in theorem 1, the following result holds:

$$\beta = \mathbb{E} \left[[D_{z_1} m_{X|Z}(Z)]^{-1} \nabla_{z_1} m_{Y|Z}(Z) B(Z) \right]. \quad (4.4)$$

Consequently, the question of overidentification is **not** about exclusion of instruments in the regression. Instead the question of overidentification is about exclusion of **derivatives** of instruments, while the instruments should always be included in the regressions. Indeed, one can show that otherwise a nonvanishing bias term of the form $\mathbb{E}[\mathbb{E}[Y|Z] Q_{z_1} | Z_1, V]$, where $Q_{z_1} = \nabla_{z_1} \log f_{Z_{-1}|Z_1, V}(Z_{-1}; Z_1, V)$, is obtained. Excluding instruments is only possible if they can be excluded from both equations (using, say, a standard omission-of-variables test).

An overidentification test is straightforwardly constructed in the spirit of Hausman (1978): Suppose there exist M such partition of $Z = (Z'_1, Z'_{-1})'$, which may be obtained by successively deleting one or more derivatives in constructing the estimator, such that β is identified for each one of them, then we can simply compare their distance using some metric. To this end, we determine the joint distribution of $\mathcal{B} = (\beta^{(1)'}, \beta^{(2)'}, \dots, \beta^{(M)'})'$. The test would then consider $H_0 : \beta^{(1)} = \beta^{(2)} = \dots = \beta^{(M)}$. As a corollary from the large sample theory of this paper, $R' \mathcal{B} = 0$, $\widehat{\mathcal{B}} \xrightarrow{d} \mathcal{N}(0, \Sigma_I)$, where Σ_I is a covariance matrix with typical element Σ_{jk} . This element is given by

$$\Sigma_{jk} = \mathbb{E} \left(\sum_{l=1}^3 \sigma_k^j \sigma_k^{k'} \right) + 2 \mathbb{E} (\sigma_2^j \sigma_3^{k'}) - \beta \beta'$$

where for $h = j, k$.

$$\begin{aligned} \sigma_1^h &= [D_{z_h} m_{X|Z}(Z_i)]^{-1} \nabla_{z_h} m_{Y|Z}(Z_i) B(Z_i) \\ \sigma_2^h &= [D_{z_h} m_{X|Z}(Z_i)]^{-1} f_Z(Z_i)^{-1} \nabla_{z_h} f_Z(Z_i) V_i' [D_{z_h} m_{X|Z}(Z_i)]^{-1} \nabla_{z_h} m_{Y|Z}(Z_i) B(Z_i) \\ \sigma_3^h &= [D_{z_h} m_{X|Z}(Z_i)]^{-1} f_Z(Z_i)^{-1} \nabla_{z_h} f_Z(Z_i) (Y_i - m_{Y|Z}(Z_i)) B(Z_i) \end{aligned}$$

Then,

$$\widehat{\Gamma}_{OvId} = \left(R' \widehat{\mathcal{B}} \right)' \left[R \widehat{\Sigma}_I R' \right]^{-1} \left(R' \widehat{\mathcal{B}} \right) \xrightarrow{\mathcal{D}} \chi_{M-1}^2,$$

by standard arguments.

4.3.2 Testing for Heterogeneity under the Assumption of Endogeneity

The principle of comparing different coefficients as a means for testing a hypothesis about our specification can be maintained more generally. If we assume to be in the scenario with endogenous regressors, we can test whether we have a heteroscedastic error or not. To illustrate the main idea, suppose that $V \perp Z$, and hence, in the case of heteroscedasticity we know that a sample counterpart to

$$\beta = \mathbb{E} \left[\left[D_z m_{X|Z}(Z)' \right]^{-1} \nabla_z m_{\bar{Y}|Z}(Z) B(Z) \right], \quad (4.5)$$

where $\bar{Y} = k_{Y|Z,V}^{0.5}(Z, V)$ produces a \sqrt{n} consistent, asymptotically normal estimator regardless of heteroscedasticity of U , while a sample counterpart estimator based on

$$\beta = \mathbb{E} \left[\left[D_z m_{X|Z}(Z)' \right]^{-1} \nabla_z m_{Y|Z}(Z) B(Z) \right],$$

will be inconsistent under heteroscedasticity. However, under H_0 of homoscedasticity, we have again that both estimators should vary only by sampling error. Using the notation $\mathbb{E} [\bar{Y} - Y|Z] = m_{\bar{Y}-Y|Z}(Z)$, a straightforward test statistic is therefore suggested by the following reformulation of H_0 :

$$0 = \mathbb{E} \left[\left[D_z m_{X|Z}(Z)' \right]^{-1} \nabla_z m_{\bar{Y}-Y|Z} B(Z) \right] = \delta.$$

The theory of the obvious sample counterpart $\hat{\delta} = n^{-1} \sum \left[D_z \widehat{m}_{X|Z}(Z_i)' \right]^{-1} \nabla_z \widehat{m}_{\bar{Y}-Y|Z}(Z_i) B(Z_i)$ is a corollary to theorem 3. Specifically, $\sqrt{n} \hat{\delta} \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_\delta)$, where Σ_δ is defined as in equation (3.4), safe for the fact that \bar{Y} is replaced by $\bar{Y} - Y$. A test statistic for heteroscedasticity is then simply a Wald test of whether δ is greater than zero, i.e.

$$\widehat{\Gamma}_{het} = \hat{\delta}' \widehat{\Sigma}_\delta^{-1} \hat{\delta} \xrightarrow{\mathcal{D}} \chi_K^2,$$

where $\widehat{\Sigma}_\delta$ is an estimator for Σ_δ , and this test statistic may be used to assess whether our model is truly heteroscedastic.

4.3.3 Testing for Endogeneity

Finally, consider analyzing whether regressors are endogenous. There are a variety of options. As in Hoderlein (2005, 2011) and Hoderlein and Mammen (2009), we may compare the regression $\mathbb{E}[Y|X]$ with the regression $\mathbb{E}[Y|X, V]$. Under the null of exogeneity, the two functions

should be the same, and hence we use a standard nonparametric omission of variables test, with the only added difficulty that V is now a generated regressor. This test would be consistent regardless of whether the single index specification on the regressors is correct or not, and would deliver nonparametric test statistics that have local power against Pitman alternatives converging at a certain rate. This procedure can be seen as a nonparametric generalization of Hausman's (1978) second test for the inclusion of control functions as test of exogeneity in a linear model.

However, if we believe the index specification to be correct, than there are other, and in some instances better, options. Note that, under the null of exogeneity, a sample counterpart estimators to the average derivative identification principle $\beta = \mathbb{E}[\nabla_x \mathbb{E}[Y|X] C(X)]$ (where $C(\cdot)$ is a again a bounded weighting function), and an estimator based on our identification principle (say, $\beta = \mathbb{E}\left[[D_z m_{X|Z}(Z)]^{-1} \nabla_z \mathbb{E}[Y|Z] B(Z)\right]$), should yield estimators that vary only be sample randomness, while under the alternative they should differ significantly.

Therefore, a similar test as the original test in Hausman (1978) may be performed. Let $\hat{\beta}_{Ex}$ denote a sample counterpart estimator to $\mathbb{E}[\nabla_x \mathbb{E}[Y|X] C(X)]$ like the PSS ADE, $\hat{\beta}_{End}$ any of the sample counterpart estimator to $\mathbb{E}\left[[D_z m_{X|Z}(Z)]^{-1} \nabla_z \mathbb{E}[Y|Z] B(Z)\right]$ defined below, $\hat{\mathcal{B}} = (\hat{\beta}'_{Ex}, \hat{\beta}'_{End})$, and $G = (I, -I)$. Next, rewrite $H_0 : G'\mathcal{B} = 0$, and use the fact that $\hat{\mathcal{B}} \xrightarrow{d} \mathcal{N}(0, \Sigma_E)$, where Σ_E is a variance covariance matrix that is straightforwardly derived from the theory below, in particular theorem 2 (the subscript E is meant to denote endogeneity). Then, a Hausman-type test statistic for H_0 , $\hat{\Gamma}_1 = (G'\hat{\mathcal{B}})' [G\hat{\Sigma}_E G']^{-1} (G'\hat{\mathcal{B}})$ behaves asymptotically as follows:

$$\hat{\Gamma}_1 = (G'\hat{\mathcal{B}})' [G\hat{\Sigma}_E G']^{-1} (G'\hat{\mathcal{B}}) \xrightarrow{d} \chi_K^2. \quad (4.6)$$

What would be the main advantage of such a specification test? It has more power against certain alternatives. Indeed, because of the parametric rate of all estimators, we may detect local alternatives in the parameter vector that converge to H_0 at root n rate. In this case, this test will be superior, provided the misspecification due to endogeneity affects the index.

5 Simulation

The finite sample performance of the estimators we propose is best analyzed by a Monte Carlo simulation study. In this section, we are chiefly concerned with analyzing the behavior of $\hat{\beta}_H$. The main scenario we consider involves an asymmetric error distribution, such that conditional mean and median differ. Moreover, we assume that V in the IV equation is fully independent of Z , in which case there is no correction term, and the estimator takes the convenient ratio-of-coefficients form as in (3.3).

To obtain an idea of the behavior of our estimator, we analyze the performance of our estimator at different data sizes. We find that our estimator performs well for moderate data sizes, and as theory predicts, we find that the mean square error reduces as the sample size increases, but we observe a small bias even in relatively large samples. However, we establish that our estimator is superior to parametric and semiparametric estimators that do not account for heterogeneity. As examples for estimators that do not account for heterogeneity we consider the parametric estimator of River and Vuong (1988) $\hat{\beta}_{RV}$, and the full independence estimator $\hat{\beta}_1$. Moreover, we show that even an infeasible oracle estimator that uses some prior knowledge not available to the econometrician shows slow convergence behavior in this setup.

We consider the case of one endogenous regressor, w.l.o.g. X_{1i} , and denote the set of regressors by $X_i = (X_{1i}, X_{2i}, \dots, X_{5i})'$, and the set of all instruments $Z_i = (Z_{1i}, X_{2i}, \dots, X_{5i})'$. For the purpose of concreteness, we specify the DGP as the following 5 - dimensional regression:

$$\begin{aligned} Y_i &= \mathbb{I}\{\beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \beta_5 X_{5i} + U_i > 0\}, \\ X_{1i} &= Z_{1i} + V_i, \quad i = 1, \dots, n, \end{aligned}$$

where $\beta = (1, 0.5, 0.5, 0.5, 0.5)'$ and the data (Y_i, X_i, Z_i, U_i) , $i = 1, \dots, n$, are *iid* draws from the following distribution: For the error U_i , we assume that there is an omitted determinant called W_i such that

$$\begin{aligned} (\log(W_i), Z_i)' &\sim \mathcal{N}(\mu, \Sigma), \\ \log(V_i)' &\sim \mathcal{N}(0, 1), \end{aligned}$$

where

$$\mu = 0, \quad \Sigma = \begin{bmatrix} 2 & 1.5 & 0 & 0 & 0 & 0 \\ 1.5 & 2 & 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 1 & 1 & 1 \\ 0 & 1 & 1 & 2 & 1 & 1 \\ 0 & 1 & 1 & 1 & 2 & 1 \\ 0 & 1 & 1 & 1 & 1 & 2 \end{bmatrix},$$

and V_i is independent of $(\log(W_i), Z_i)'$. Observe that the W_i are in particular correlated with Z_{1i} . Next, the error U_i is defined through:

$$U_i = W_i - k_{W|Z}^{0.5}(Z_i) + V_i,$$

so that $k_{U|ZV}^{0.5}(Z_i, V_i) = V_i$. Hence, as we require the error U_i obeys the conditional median exclusion restriction, but depends on Z_i . As baseline, our estimator (3.3) is defined as a local

quadratic polynomial estimator, with Epanechnikov kernels⁶. Moreover, the “smooth indicator” is defined as the integral of the Epanechnikov kernel over the positive areas. The conditional probability is also estimated using a local quadratic polynomial estimator, with Epanechnikov kernel. The independence estimator $\hat{\beta}_1$ is defined similarly, with the exception that no “smooth indicator” is required. The oracle estimator is obtained by using fitted values Y_{1i} instead of either the conditional median $k_{Y|ZV}^{0.5}(Z_i, V_i)$ (as is the case in the $\hat{\beta}_H$) or the Y_i (as is the case in $\hat{\beta}_1$). The fitted values Y_{1i} are obtained in the following way: We assume that the oracle has knowledge of the Y_i^* , and compute the conditional median $k_{Y^*|ZV}^{0.5}(Z_i, V_i)$, and set $Y_{1i} = \mathbb{I} \left\{ k_{Y^*|ZV}^{0.5}(Z_i, V_i) > 0 \right\}$. Bandwidths for all estimators are obtained by doing a grid search for the bandwidth that minimizes the MSE in 100 repetitions. Finally, $\hat{\beta}_{RV}$ is obtained by estimating a probit model with control function residuals as additional regressors.

The result of applying our methods can be found in figures 1 - 3 in the graphs in the appendix. For each $j = 1, \dots, J$, $J = 500$, a new sample (X_i, Z_i, W_i, U_i) of size n is drawn from the distribution specified above. To illustrate the behavior of the estimator, in fig. 1 we plot the density of the estimated four elements of β for $n = 2500$ as solid line. Note that the first coefficient is normalized to one. The vertical line in all of these plots is at the true value of 0.5. The closer this distribution is to a spike centered at this value, the better the performance of the estimator. We compare the distribution of $\hat{\beta}_H$ with that of the parametric estimator $\hat{\beta}_{RV}$ (first dotted line), the independence estimator $\hat{\beta}_1$ (second dotted line), and the oracle estimator $\hat{\beta}_O$ (solid line).

First, the most obvious feature of the result is the clear ordering in terms of the performance of the estimators, regardless of data size. Obviously, in terms of the bias, $\hat{\beta}_H$ is less biased than any of the two alternative feasible estimators $\hat{\beta}_1$ and $\hat{\beta}_{RV}$, and only the infeasible oracle estimator show less bias. In terms of variance, all three estimators nonparametric estimators are approximately similar. The fact that the variance is not significantly affected by the first step estimation of the conditional median arises because the median estimator still uses all observations. While the parametric estimator $\hat{\beta}_{RV}$ shows less variance than any of the semiparametric estimators, its much larger bias results in a greatly increased MSE compared to any semiparametric estimator, a result of the double misspecification of $\hat{\beta}_{RV}$: First, the estimator erroneously imposes a parametric structure, second it erroneously assumes full conditional independence. Indeed, going from $\hat{\beta}_{RV}$ over $\hat{\beta}_1$ to $\hat{\beta}_H$ may be seen as first removing the parametric assumption, and then correcting for heteroscedasticity. As we see, the first step reduces the bias at the expense of increasing the variance somewhat. Still, the MSE is almost

⁶This may not be of sufficiently high order to remove the bias completely. As a consequence, all our estimators exhibit some bias even at large samples. However, we believe that in applications a higher order bias reduction than the one we perform is impracticable, and hence we concentrate on this case.

halved. The second step reduces again the bias while only affecting the variance marginally, and the MSE is again reduced significantly.

These graphs were obtained using the true V_i . Of course, in any real world application V_i is not known, and has to be replaced by a nonparametric first stage estimator. To estimate V_i , we use the fact that $V_i = X_i - m_{X|Z}(Z_i)$, and replace $m_{X|Z}$ by a nonparametric leave-one-out estimator. The results do not change in any significant way; in fact, if anything it seems to induce a small upward bias, which generally compensates some of the downward bias in particular at small sample sizes, and for all estimators except the oracle estimator. For instance, the MSE for $\hat{\beta}_H$ for $n = 2500$ decreases from approximately 0.016 to 0.014. For larger bandwidths, however, the effect is small, and generally the contribution to the variance seems negligible, probably since the averaging that happens in the regression, as well as the final averaging, washes out this effect. We conclude that correcting for the fact that V_i is pre-estimated is not an issue of great importance.

Another potentially important issue worth investigating in this paper is the use of a smooth indicator. To understand whether this is a purely theoretical device, or whether it has some real world bearings, we have conducted a series of experiments where we have let the bandwidth on this kernel tend from 0.1 to exactly 0. Again, we do not find any significant effect; as long as you choose a small value of h , the precise magnitude does not matter too much, with one minor caveat: At very small bandwidth including 0, the variance increases slightly, so that the average MSE increases for $n = 2500$ from 0.016 to 0.017.

It is also instructive to look in more detail at how the estimators behave as n varies. The heteroscedasticity robust estimator $\hat{\beta}_H$ significantly outperforms $\hat{\beta}_1$ and $\hat{\beta}_{RV}$ at moderate sample sizes ($n = 2500$); for smaller sample sizes the advantage in particular over the semiparametric competitor $\hat{\beta}_1$ becomes less pronounced. As such we find the familiar results in other simulation studies on the binary choice case (e.g., Frölich (2005)), namely that in binary choice models semiparametric methods require a significant amount of data to outperform misspecified models. Once, however, we have a significant amount of data, the advantages become obvious, see fig.2 and 3, who show the behavior with $n = 7500$ and $n = 15000$ observations. The bias of $\hat{\beta}_H$ starts to vanish, clearly visible in the bottom right panel of fig. 3. The same result is also obtained from the tables, cf tab 1-4 below, which provide the specific numerical results. We obtain for $\hat{\beta}_H$:

Coefficient	2	3	4	5
$n = 2500$	0.017288	0.015739	0.016421	0.016928
$n = 7500$	0.010301	0.009498	0.008755	0.008411
$n = 15000$	0.009207	0.008299	0.007690	0.007135

Table 1: MSE of $\hat{\beta}_H$ at Different Data Sizes

The reduction of the MSE with increasing sample size is obvious. Note also that due to the largely symmetric setup, all four coefficients are equally affected. A more detailed analysis shows a reduction in both bias and variance, as is also evident from the graphs, see fig. 1-3. Note, however, that the reduction in bias is quite slow. It is instructive to compare the estimator with the other estimators. For $\hat{\beta}_I$, we obtain the following result:

Coefficient	2	3	4	5
$n = 2500$	0.020972	0.020581	0.020385	0.021524
$n = 7500$	0.017027	0.016416	0.015744	0.015766
$n = 15000$	0.016466	0.016306	0.015324	0.014847

Table 2: MSE of $\hat{\beta}_I$ at Different Data Sizes

This result is clearly worse than the heteroscedasticity robust estimator $\hat{\beta}_H$, with an increase in MSE of roughly 25 - 100 %. In contrast, as was to be expected, the (infeasible) oracle estimator $\hat{\beta}_O$ outperforms both estimators:

Coefficient	2	3	4	5
$n = 2500$	0.012412	0.009253	0.011056	0.009996
$n = 7500$	0.004054	0.004245	0.005386	0.005296
$n = 15000$	0.002875	0.002583	0.003155	0.002967

Table 3: MSE of $\hat{\beta}_O$ at Different Data Sizes

When decomposing the MSE, we find that the variance remains very comparable across all semiparametric estimators given the data size, while it is the bias that causes the differences. Naturally, the oracle estimator starts out relatively unbiased and remains so. In contrast, the independence estimator exhibits a nonvanishing bias component. The heteroscedastic estimator starts out with a bias that diminishes with increasing sample size. Note that the difference between $\hat{\beta}_H$ and $\hat{\beta}_O$ can be seen as a measure of the degree of information loss associated with the indicator function. Viewing the indicator as a filter, we conclude that the information loss is quite severe, and that significant data sizes are required to distinguish between different structures within the indicator.

Finally, consider the estimator of Rivers and Vuong (1988), which is close to what an educated applied person would do. Unfortunately, due to the significant bias, this estimator

performs worst, see the following table 4.

Coefficient	2	3	4	5
$n = 2500$	0.036005	0.037955	0.036398	0.036653
$n = 7500$	0.035706	0.034574	0.034498	0.035680
$n = 15000$	0.034166	0.034278	0.034125	0.034720

Table 4: MSE of $\hat{\beta}_{RV}$ at Different Data Sizes

Indeed, by closer inspection we find that the MSE of $\hat{\beta}_{RV}$ is almost entirely due to the squared bias, and the variance contribution is quite small. Thus, the performance of the estimator shows only very little improvement with increasing sample size. Since all estimators correct for endogeneity, we can summarize the finding by saying that in this scenario correcting for heterogeneity and adopting a semiparametric procedure produces significantly better results. We conclude that the interaction between the various sources of misspecification makes at least in this setup semiparametric estimators quite attractive. And with a sufficient amount of data, it is also evident that allowing for a heterogeneous error structure improves the result significantly, and leads to a performance which is not much worse than that of an infeasible oracle estimator. Our application below will make the importance of being less restrictive in this part of the model also for real world data apparent.

The last point we want to discuss in this Monte Carlo is the performance of asymptotic standard errors and bootstrap standard errors. While we recommend using the bootstrap to perform inference, it is instructive to see how the large sample variance Σ_H performs. To do so, we have estimated the sample counterpart of the variance equation (3.4) nonparametrically. What follows are the details of the estimation of each of its components.

First, the density function of instruments f_Z was evaluated by a simple Gaussian kernel density estimator, using a bandwidth $h_{f_Z} = 3$. Similarly, its gradient, ∇f_Z was estimated via a Gaussian product kernel of the form

$$\widehat{\nabla f_Z}(z) \equiv \begin{bmatrix} \frac{1}{nh_{\nabla f_Z}^{K+1}} \sum_i^n K'_1\left(\frac{Z_{ik}-z_k}{h_{\nabla f_Z}}\right) \\ \vdots \\ \frac{1}{nh_{\nabla f_Z}^{K+1}} \sum_i^n K'_K\left(\frac{Z_{ik}-z_k}{h_{\nabla f_Z}}\right) \end{bmatrix}$$

where $h_{\nabla f_Z} = 3.5$ and K'_k is defined as the product kernel with a derivative on the k^{th} position, i.e., $K'_k(u) \equiv K(u_1)K(u_2) \cdots \partial_u K(u_k) \cdots K(u_K)$

Next, the estimate of the conditional expectation $m_{\bar{Y}|Z}$ was obtained by a second-order local polynomial estimator, again using Gaussian kernels, with a bandwidth of $h_{m_{\bar{Y}|Z}} = 3$. The second order is sensible since we are interested in the vector of first-order derivatives (i.e., the gradient $\nabla m_{\bar{Y}|Z}$). Estimates of the conditional expectation $m_{X|Z}$ and its Jacobian $D_z m_{X|Z}$ are

	$Q_{.10}$	$Q_{.25}$	$Q_{.50}$	$Q_{.75}$	$Q_{.90}$
β_2	0.062	0.094	0.148	0.234	0.450
β_3	0.067	0.096	0.146	0.236	0.509
β_4	0.067	0.099	0.147	0.243	0.461
β_5	0.066	0.098	0.140	0.222	0.522

Table 5: Quantiles of the Distributions of the Asymptotic Standard Errors for the resp. Coefficients.

	$Q_{.10}$	$Q_{.25}$	$Q_{.50}$	$Q_{.75}$	$Q_{.90}$
β_2	0.196	0.215	0.238	0.269	0.307
β_3	0.197	0.214	0.239	0.268	0.303
β_4	0.200	0.216	0.235	0.262	0.306
β_5	0.198	0.212	0.236	0.264	0.304

Table 6: Quantiles of the Distributions of the Bootstrap Standard Errors for the resp. Coefficients.

determined in the exact same manner, again using the same bandwidth $h_{m_{X|Z}} = 3$. Finally, the weighting function $B(z) = \mathbb{I}\{z \in I_z\} \mathbb{I}\{\det |D_z m_{X|Z}(z) D_z m_{X|Z}(z)'| \geq b\}$ was modelled by simply dropping observations that were not in $I_z = [-4, +4]^K$ or whose pre-estimated determinant inner product was smaller than $b = 0.1$.⁷

We ran 450 simulations with $N = 2500$ observations each and computed the standard errors of the coefficients for each one of them. The quantiles of the distribution can be found on the next table.

For completeness, we also include the same quantiles for the distribution of the bootstrap standard errors on Table 6. To produce this table, we drew 100 times with replacement from each of the 450 simulations above, then computed the corresponding values of β_H .

We also show a graphical comparison of the finite sample distributions of these standard errors in graphs 4 and 5. While the order of magnitude seems correct (the squared standard error is roughly of the same size as the variance component of the MSE), we see that the asymptotic standard errors exhibit more very small values. We conjecture thus that the bootstrap is tighter centered around the true value, and since it is easier to compute and does not rely on dropping observations, we recommend its usage in practise. This completes our simulation exercise.

⁷Further documentation and MATLAB code is available on request

6 Application to Discrete Consumer Choice

6.1 Description of Data and Variables

As an example for an application of our method to a structural model in a heterogeneous population, we use data that is very similar to the one employed in Goolsbee and Petrin (2004) about the choice of television transmission mode, see Table A.1 for an overview of all variables. The data comes from two data sources. First, from December 2000 until January 2001 NFO Worldwide fielded a household survey on television choices sponsored by Forrester Research as part of their Technographics 2001 program⁸. These households were randomly drawn from the NFO mail panel that is designed to be nationally representative.

The households that were surveyed basically have the choice between four different ways to receive television programming: local antenna, direct broadcast satellite (DBS), as well as basic and expanded cable, which we group into cable versus non-cable (satellite dish/local antenna). Local antenna reception is free but only carries the local broadcast stations⁹. DBS systems are national companies that deliver many of the cable channels that usually priced uniformly across the whole country (in 2001 the two leading companies DirectTV and DISH Network (Echostar) charged \$30 and \$32 per month respectively). Hence, there is almost no price variation in the alternative. Compared to cable, DBS provides a greater variety of channels and more pay per view options but bares the potential for signal interference and charges a higher price. The fair amount of regional variation in cable prices permits us to estimate own price effects, while the cross price effects are constant, and hence neglectable.

Other than the choices people make, the survey also provides information on various socio-economic household characteristics e.g. household income, household composition, education of the head of household and if applicable of the respective partner. Dropping observations with missing values in their choices or doubtful values in several household characteristics and removing outliers (recall that we also have to compactify our support) reduces the sample to approximately 15.900 observations. Table A.2 in the appendix provides summary statistics for the sub sample including renter status and whether households live in single unit dwellings. Both characteristics are known to influence the ability to receive satellite.

We also make use of a second source of data, which provides us with information on cable prices and cable franchise characteristics each household faces (within a specific cable franchise area). The data come from Warren Publishing's 2002 Television and Cable Factbook, and

⁸NFO was the largest custom market-research firm in the United States until it became part of the TNS Group in 2004.

⁹Looking at households that have a TV allows to assume that local antenna forms the chosen alternative for those who neither declare to subscribe to cable nor to DBS.

provides detailed information on the cable industry, which is divided into geographically separated cable systems. From this data source, we use the channel capacity of the cable system, whether pay per view is available, the price of basic plus expanded basic service, the price for premium channels (here we use the price for HBO) and the number of over-the-air channels (this corresponds to the number of local channels carried by the cable system).

As source of endogeneity, we follow Goolsbee and Petrin (2004), and assume that prices are correlated with unobserved characteristics like advertisement. To deal with endogeneity, we use variation in city franchise tax/fee to instrument cable prices (recall that the own price might be correlated with unobserved cable characteristics e.g. advertising or quality). Table A.3 presents summary statistic for the respective variables. Technically, we can match both data sources using Warren’s ICA system identification number, which is based on zip code information. Hence, we can assign a specific household to the adequate local cable company¹⁰ even though these individuals might not subscribe to cable.

6.2 Empirical Results

The focus in our empirical analysis is on the own price effect, and how the result is altered by the introduction of our method. The effect of household covariates is not of interest, and we use these variables merely as controls. Since we are not interested in their effect, we employ principal components to reduce them to some three orthogonal approximately continuous variables, mainly because we require continuous covariates for nonparametric estimation. While this has some additional advantages, it is arguably ad hoc. However, we performed some robustness checks like alternating the components or adding parametric indices to the regressions, and the results do not change in an appreciable fashion (nor is the remaining variation statistically significant).

To show the performance of our estimator, it is instructive to start out with the standard

¹⁰Typically only one cable company receives the right to serve a region as a result of a franchise agreement with a local government even though the household might not subscribe to cable.

practise of estimating a linear probability model and using 2SLS. We obtain the following result:

	Estimate	Std. Error	t value	p value
Intercept	0.697908	0.008805	79.266	0
Own Price	0.228026	0.020040	11.379	0
Income	0.028096	0.002513	11.181	0
PrinComp 1	- 0.025945	0.008904	- 2.914	0.003575
PrinComp 2	0.014143	0.004033	3.507	0.000454
PrinComp 3	- 0.018363	0.002663	- 6.895	0

Table 7: Linear Probability Model - 2SLS

There are two things noteworthy: First, quite in contrast to Economic theory, the model predicts that higher own price is associated with higher demand. Second, the income effect is positive, but small in absolute size. Due to the large sample size of $n = 15.918$ all variables are highly significant, with p -values of virtually zero. This holds true even for the - in absolute size - small income effect. This finding remains stable across specifications, however, the own price effect becomes progressively more plausible as we move to less obviously misspecified specifications.

The following tables show the behavior of the full independence estimator $\hat{\beta}_I$. Specifically, it shows the point estimate, as well as the 2.5 and 97.5 quantile of the bootstrap distribution¹¹ instead of the asymptotic distribution which is cumbersome to estimate. In this procedure, a coefficient is statistically not significant from zero if the confidence interval contains zero.

	Estimate	BS 0.025 value	BS 0.975 value
Own Price	- 2.10788	- 5.94305	1.76096
Income	1	1	1
PrinComp 1	- 3.31908	- 4.32975	- 2.49948
PrinComp 2	1.70843	1.23556	2.35323
PrinComp 3	- 0.35490	- 1.15326	0.48962

Table 8: Coefficients of $\hat{\beta}_I$ (Relative to Income) with Bootstrap Confidence Intervals

As we see from the results, this is the case for the own price effect, which is in absolute value only twice as strong as the income effect. Compared to the income effect, the estimate

¹¹We have performed $n = 200$ bootstrap repetitions with replacement from the same data. Since the choice of bandwidth is not clear (we conjecture that a second order expansion type of analysis can be performed), we have settled for a slightly smaller bandwidth in the bootstrap replications, because this is a common devise to mitigate small sample bias in the construction of pointwise confidence bands in nonparametric regression.

points in the opposite direction. And if we look at the non normalized results we also obtain that the income effect is positive (and actually of as small an order of magnitude as in the linear probability model), while the price effect is negative as it should be, but as mentioned insignificant. The first two principal components are significant, however not the third, and have generally the same sign and relative order of magnitude as in the linear probability model.

The heteroscedasticity robust estimator $\hat{\beta}_H$ produces the most sensible results:

	Estimate	BS 0.025 value	BS 0.975 value
Own Price	- 8.02943	- 12.91400	- 2.90706
Income	1	1	1
PrinComp 1	- 0.12521	- 1.04786	0.52044
PrinComp 2	1.38809	0.93442	2.01614
PrinComp 3	- 0.88721	- 2.02577	- 0.08665

Table 9: Coefficients of $\hat{\beta}_H$ (Relative to Income) with Bootstrap Confidence Intervals

Here we see that the own price effect is significantly negative. At first glance, the results appears to be slightly different from Goolsbee and Petrin (2004), who find a relatively low own price elasticity. However, as the income effect is rather weak (it is again of the same order of magnitude as in the linear probability model in the non normalized version. but recall that identification is only up to scale), this is not necessarily a contradiction. With respect to the application, we conclude that the likelihood that the average person in this population chooses cable reacts only modestly to an increase in income, which given the small fraction of total expenditures seems plausible (and is perhaps very different if one were to consider the demand for cars). However, given that price of cable is a significant variable in the marketing of this good, the average consumer seems to react more strongly to price incentives, and as theory predicts, a price increase reduces the probability of buying cable.

Finally, we use a specification test as outlined in section 4.3.2. To determine the critical values we use, however, the bootstrap. More specifically, we follow Hoderlein and Winter (2009) and use a hybrid bootstrap to obtain the (one sided) 0.05 confidence bound. The bootstrap p -value of the test statistics is .09, which means that we do not reject the null hypothesis formally that there is homogeneity at the 5% level, however, we would do so at the 10% level.

In summary, regarding the performance of various different estimators, we conclude that avoiding the misspecification associated with the linear probability model, as well as allowing for heterogeneous preferences (compared to the full independence estimator $\hat{\beta}_I$) substantially alters the result, and provides us with more plausible estimates for the (centrality) parameter of interest.

7 Summary and Outlook

The notion that we do not observe important determinants of individual behavior even in data sets with large cross section variation becomes more and more influential across microeconometrics. Indeed, it is widely believed now that unobserved tastes and preferences account for much more of the variation than observable characteristics. Hence, it is imperative to devise models that account for heterogeneity on individual level, in particular if the unobserved determinants and omitted variables are believed to be correlated with observables.

In these heterogeneous models, most often interest centers on average effects. In this paper, we analyze the binary choice model with random utility parameters under a median exclusion restriction that defines such a (local) average effect. It is moreover established that this effect coincides with the parameter β in the reduced form binary choice model with heteroscedastic errors under a median exclusion restriction. We show how to nonparametrically identify this parameter β , and we propose a \sqrt{n} consistent, asymptotically normal sample counterparts estimator. Moreover, based on our theory, we propose tests for overidentification, endogeneity as well as heterogeneity. Therefore we can provide means to check the specification, in addition to provide the first estimator for this parameter in this class of models.

In a Monte Carlo study we show that the performance of our estimator is superior to one that does not exploit the heterogeneity structure of the model. In an application, we show that our estimator uses significantly weaker assumptions than those employed in the literature, and through its use we may be able to reveal new and interesting features. How to extend this type of semiparametric approach from binary choice data to multinomial choice data and more complicated settings including simultaneity remains an interesting direction for future research. Our conjecture is that a similar estimation principle may be applicable to a large class of models.

References

- [1] Ai, Ch. 1997. A Semiparametric Maximum Likelihood Estimator, *Econometrica* **65**, 933–964.
- [2] Bajari, P., Fox, J., Kim K. and S. Ryan, 2012. The random coefficients logit model is identified, *Journal of Econometrics*, 166(2), 204-212.
- [3] Berry, S. and P. Haile, 2008. Nonparametric Identification of Multinomial Choice Demand Models with Heterogeneous Consumers, Working Paper, Cowles Foundation.

- [4] Blundell, R. and J. Powell. 2004. Endogeneity in semiparametric binary response models, *Review of Economic Studies* **71**, 655–679.
- [5] Blundell, R., and R. Smith. 1986. An Exogeneity Test for a Simultaneous Tobit Model, *Econometrica* **54**, 679–685.
- [6] Chaudhuri, P., K. Doksum and A. Samarov. 1997. On average derivative quantile regression, *Ann. Statist.* **25**, 715–744.
- [7] Das, M., Newey, W. and F. Vella. 2003. Nonparametric estimation of sample selection models, *Review of Economic Studies* **70**, 33–58.
- [8] Delecroix, M., and M. Hristache. 1999. M-estimateurs semi-paramétriques dans les modèles à direction révélatrice unique, *Bull. Belg. Maths. Soc.* **6**, 161–185.
- [9] Fan, J. and I. Gijbels. 1996. *Local Polynomial Modelling and Its Applications*, Chapman and Hall, London.
- [10] Florens, J.P., Heckman J, Meghir, C., and E. Vytlacil, 2008. Identification of Treatment Effects Using Control Functions in Models with Continuous, Endogenous Treatment and Heterogeneous Effects, *Econometrica*,
- [11] Gautier, E. and Y. Kitamura, 2013, Nonparametric Estimation in Random Coefficients Binary Choice Models, *Econometrica*, **81**, 581-607.
- [12] Goolsbee, A. and A. Petrin, 2004. The Consumer Gains from Direct Broadcast Satellites and the Competition with Cable TV, *Econometrica*, **72**, 351-381.
- [13] Heckman, J. and E. Vytlacil, 2005. Structural Equations, Treatment Effects, and Economic Policy Evaluation, *Econometrica*, **73**, 669-738..
- [14] Hoderlein, S., 2011. How many Consumers are Rational?, *Journal of Econometrics*, **164**, 294-309.
- [15] Hoderlein, S. , and E. Mammen, 2007. Identification of marginal effects in nonseparable models without monotonicity, *Econometrica* **75**, 1513–1518.
- [16] Hoderlein, S.,and J. Winter, 2010. Structural Measurement Errors in Nonseparable Models, *Journal of Econometrics*, **157**,
- [17] Hoderlein, S., J. Klemelä, and E. Mammen 2010: Analyzing the Random Coefficient Model Nonparametrically, *Econometric Theory*, **26**, 804–837

- [18] Horowitz, J. L. 1992, A smoothed maximum score estimator for the binary response model, *Econometrica* **60**, 505–531.
- [19] Horowitz, J., and W., Härdle. 1996. Direct Semiparametric Estimation of Single Index Models with Discrete Covariates, *JASA* **91**, 1632–1640.
- [20] Horowitz, J., 1997, Bootstrap Methods in Econometrics: Theory and Numerical Performance, In *Kreps, D.M. and Wallis, K.F. (eds), Advances in Economics and Econometrics: Theory and Applications*, pp. 188-222, Cambridge University Press.
- [21] Hristache, M., A. Juditsky, J. Polzehl, and V. Spokoiny. 2001. Structure Adaptive Approach for Dimension Reduction, *Ann. Statist.* **29**, 1537–1566.
- [22] Ichimura, H. 1993, Semiparametric Least Squares (SLS) and Weighted SLS Estimation of Single-Index Models, *Journal of Econometrics* **58**, 71–120.
- [23] Ichimura, H., and T. S. Thompson 1998: Maximum Likelihood Estimation of a Binary Choice Model With Random Coefficients of Unknown Distribution, *Journal of Econometrics*, 86, 269–295.
- [24] Klein, R.W. and R.S. Spady. 1993, An Efficient Semiparametric Estimator of the Binary Response Model, *Econometrica* **61**, 387–422.
- [25] Lewbel, A. 1998. Semiparametric Latent Variable Model Estimation with Endogenous or Mismeasured Regressors, *Econometrica*, **66**, 105-122.
- [26] Manski, C. F. 1975, Maximum score estimation of the stochastic utility model of choice, *Journal of Econometrics* **3**, 205–228.
- [27] Matzkin, R. 1992. Nonparametric and Distribution-Free Estimation of the Binary Threshold Crossing and the Binary Choice Models, *Econometrica*, **62**, 239-270.
- [28] Matzkin, R. 2005. Heterogeneous Choice, for *Advances in Economics and Econometrics*, edited by Richard Blundell, Whitney Newey, and Torsten Persson, Cambridge University Press; presented at the Invited Symposium on Modeling Heterogeneity, World Congress of the Econometric Society, London, U.K.
- [29] Newey, W.K., Powell, J.L. and Vella, F. 1999. Nonparametric Estimation of Triangular Simultaneous Equations Models, *Econometrica* **67**, 565–603.
- [30] Petrin, A. and K. Train. Control Function Corrections for Omitted Attributes in Differentiated Product Markets. 2006 Working paper.

- [31] Powell, J. L., J. H. Stock, and T. M. Stoker. 1989. Semiparametric estimation of index coefficients, *Econometrica* **57**, 1403–1430.
- [32] Rivers, D., and Q Vuong. 1988. Limited Information Estimators and Exogeneity Tests for Simultaneous Probit Models, *Journal of Econometrics*, 39, 347 - 364.
- [33] Stoker, T. M. 1986. Consistent Estimation of Scaled Coefficients, *Econometrica* **54**, 1461–1481.
- [34] Vytlacil, E. and N. Yildiz, 2007. "Dummy Endogenous Variables in Weakly Separable Models," *Econometrica*, **75**, 757-779.

8 Appendix 1: Technical Proofs

The Proof of Theorem 2

The Structure of the Proof

Neglect for a moment the weighting function A . Rewrite (3.6) as $\frac{1}{n} \sum_i \hat{G}_i^- \hat{B}_i$, where $G_i = D_z m_{X|Z}(Z_i)'$, $B_i = \nabla_z m_{Y|Z}(Z_i)$, $\hat{G}_i = D_z \hat{m}_{X|Z}(Z_i)'$ and $\hat{B}_i = \nabla_z \hat{m}_{Y|Z}(Z_i)$. Tedious, but straightforward manipulations lead to

$$\begin{aligned}
 \hat{G}_i^- \hat{B}_i &= G_i^- B_i + G_i^- \left[(G_i - \hat{G}_i) G_i^- B_i + (\hat{B}_i - B_i) \right] \\
 &\quad + G_i^- (G_i - \hat{G}_i) \left(G_i^- - \hat{G}_i^- \right) B_i \\
 &\quad + G_i^- (G_i - \hat{G}_i) G_i^- (\hat{B}_i - B_i) \\
 &\quad + G_i^- (G_i - \hat{G}_i) \left(G_i^- - \hat{G}_i^- \right) (\hat{B}_i - B_i).
 \end{aligned} \tag{8.1}$$

Now, in (8.1) the first two terms on the right hand side will provide us with the asymptotic distribution, while the terms from three to five will prove asymptotically negligible. In **Step 1**, we treat the behavior of the first two summands first in the case where $m_{X|Z}$ is a mean regression. Specifically, we show in **Step 1a** that

$$\tau_{1n} = n^{-1} \sum_i G_i^- B_i + G_i^- \left[(G_i - \hat{G}_i) G_i^- B_i + (\hat{B}_i - B_i) \right] = S_{1n} + S_{2n}$$

i.e., the sum can be decomposed into two terms, the first of which provides us with the asymptotic distribution, while the second one produces the bias. In **Step 1b**, we establish that the large sample theory of S_{1n} may be handled using projection arguments coming from U -statistic theory, while in **Step 1c** we show that the bias term S_{2n} will vanish under appropriate conditions on the bandwidths, as in PSS. Finally, in Step 1d we derive the asymptotic distribution.

In **Step 2**, we discuss the behavior of the higher order terms in (8.1), i.e., the behavior of terms three to five. In **Step 3** we establish under which conditions generated dependent variables do not matter for the asymptotic distribution of the estimator.

Step 1: The General Proof

Step 1a: Consider

$$\tau_{1n} = n^{-1} \sum_i \left\{ G_i^- B_i + G_i^- (G_i - \hat{G}_i) G_i^- B_i + G_i^- (\hat{B}_i - B_i) \right\} \quad (8.2)$$

$$= n^{-1} \sum_i G_i^- B_i + n^{-1} \sum_i G_i^- (G_i - \hat{G}_i) G_i^- B_i + n^{-1} \sum_i G_i^- (\hat{B}_i - B_i). \quad (8.3)$$

Since the first term has a trivial structure, and the second and third terms are similar, we start by considering the second term on the right hand side of (8.2) first. In the case where \hat{G}_i is a nonparametric Nadaraya Watson derivative estimator, it rewrites as

$$\left(\sum_{j \neq i} \mathcal{K}_{hj}(Z_i) \right)^{-1} \left[\sum_{j \neq i} \nabla_z \mathcal{K}_{hj}(Z_i) X_j' - \left(\sum_{j \neq i} \mathcal{K}_{hj}(Z_i) \right)^{-1} \sum_{j \neq i} \nabla_z \mathcal{K}_{hj}(Z_i) \sum_{j \neq i} \mathcal{K}_{hj}(Z_i) X_j' \right]. \quad (8.4)$$

Hence, $G_i - \hat{G}_i$ has a representation as

$$D_z m_{X|Z}(Z_i) - \left(\sum_{j \neq i} \mathcal{K}_{hj}(Z_i) \right)^{-1} \left[\sum_{j \neq i} [\nabla_z \mathcal{K}_{hj}(Z_i) - W_n(Z_i) \mathcal{K}_{hj}(Z_i)] [V_j' + m_{X|Z}(Z_j)'] \right]$$

where

$$W_n(Z_i) = \left(\sum_{s \neq i} \mathcal{K}_{hs}(Z_i) \right)^{-1} \sum_{j \neq i} \nabla_z \mathcal{K}_{hs}(Z_i).$$

Separate this expressions into the two parts, where

$$\begin{aligned} -P_{1i} &= \left(\sum_{j \neq i} \mathcal{K}_{hj}(Z_i) \right)^{-1} \left[\sum_{j \neq i} [\nabla_z \mathcal{K}_{hj}(Z_i) - W_n(Z_i) \mathcal{K}_{hj}(Z_i)] V_j' \right] \\ &= (n-1)^{-1} \sum_{j \neq i} \mathcal{W}_{jn}(Z_i) V_j' \end{aligned} \quad (8.5)$$

where $\mathcal{W}_{jn}(Z_i) = \left((n-1)^{-1} \sum_{j \neq i} \mathcal{K}_{hj}(Z_i) \right)^{-1} [\nabla_z \mathcal{K}_{hj}(Z_i) - W_n(Z_i) \mathcal{K}_{hj}(Z_i)]$ and

$$\begin{aligned} P_{2i} &= D_z m_{X|Z}(Z_i)' - \left(\sum_{j \neq i} \mathcal{K}_{hj}(Z_i) \right)^{-1} \left[\sum_{j \neq i} [\nabla_z \mathcal{K}_{hj}(Z_i) - W_n(Z_i) \mathcal{K}_{hj}(Z_i)] m_{X|Z}(Z_j)' \right] \\ &= D_z m_{X|Z}(Z_i)' - (n-1)^{-1} \sum_{j \neq i} \mathcal{W}_{jn}(Z_i) m_{X|Z}(Z_j)' \end{aligned} \quad (8.6)$$

Note that $\mathcal{W}_{jn}(Z_i) = -\mathcal{W}_{in}(Z_j)$ by the symmetry of the kernel. The first part, (8.5), will contribute to the variance of the estimators, whereas the second will produce the leading bias term for which we shall give conditions under which it vanishes. Rewriting

$$\begin{aligned} n^{-1} \sum_i G_i^- (G_i - \hat{G}_i) G_i^- B_i &= -(n(n-1))^{-1} \sum_i \sum_{j \neq i} G_i^- \mathcal{W}_{jn}(Z_i) V_j' G_i^- B_i \\ &\quad + (n(n-1))^{-1} \sum_i G_i^- P_{2i} G_i^- B_i \\ &= S_{1n}^2 + S_{2n}^2, \end{aligned}$$

where the superscript 2 denotes the second term in the expression (8.2). A similar decomposition may be performed on $n^{-1} \sum_i G_i^- (\hat{B}_i - B_i) = (n(n-1))^{-1} \sum_i \sum_{j \neq i} G_i^- \mathcal{W}_{jn}(Z_i) Q_j + (n(n-1))^{-1} \sum_i G_i^- P_{4i} = S_{1n}^3 + S_{2n}^3$, where $Q_i = Y_i - m_{Y|Z}(Z_i)$ and P_{4i} denotes again bias terms in the regression of Y on Z . In total, we obtain that

$$\tau_{1n} = n^{-1} \sum_i G_i^- B_i + S_{1n}^2 + S_{1n}^3 + S_{2n}^2 + S_{2n}^3 = S_{1n} + S_{2n}, \quad (8.7)$$

where $S_{1n} = n^{-1} \sum_i G_i^- B_i + S_{1n}^2 + S_{1n}^3$ collects all terms that affect the asymptotic distribution, while $S_{2n} = S_{2n}^2 + S_{2n}^3$ are all bias terms that vanish under appropriate conditions.

Step 1b: To analyze all terms that affect the distribution and are contained in S_{1n} , consider S_{1n}^2 first. Manipulating this expression produces

$$\begin{aligned} U_n &= (n(n-1))^{-1} \sum_i \sum_{j > i} \{G_i^- \mathcal{W}_{jn}(Z_i) V_j' G_i^- B_i - G_j^- \mathcal{W}_{jn}(Z_i) V_i' G_j^- B_j\} \\ &= (n(n-1))^{-1} \sum_i \sum_{j > i} p_n(S_i, S_j), \end{aligned}$$

where $S_i = (Y_i, X_i', Z_i)'$, with p_n symmetric, and we made use of $\mathcal{W}_{jn}(Z_i) = -\mathcal{W}_{in}(Z_j)$. To apply Lemma 3.1 of PSS which yields $\sqrt{n} (\hat{U}_n - U_n) = o_p(1)$, where

$$\hat{U}_n = \theta + n^{-1} \sum_i \mathbb{E} [p_n(S_i, S_j) | S_i], \quad (8.8)$$

we require that $\mathbb{E} (\|p_n(S_i, S_j)\|^2) = o(n)$. Following similar and straightforward, but more tedious arguments as in PSS, this is the case provided $nh^{L+2} \rightarrow \infty$. To analyze (8.8), note first that $\theta = \mathbb{E} [p_n(S_i, S_j)] = 0$, and consider first p_n^* which equals p_n save that in $\mathcal{W}_{jn}(Z_i)$, $(n-1)^{-1} \sum_{s \neq i} \mathcal{K}_{hs}(Z_i)$ and $(n-1)^{-1} \sum_{s \neq i} \nabla_z \mathcal{K}_{hs}(Z_i)$ are replaced with their probability limits,

$f_Z(Z_i)$. and $\nabla_z f_Z(Z_i)$. Then,

$$\begin{aligned}
& \mathbb{E} [p_n^*(S_i, S_j) | S_i = s_i] \\
&= \int h^{-(L+1)} (D_z m_{X|Z}(z_i)^- f_Z(z_i)^{-1} (\nabla_z \mathcal{K}((z_i - z)/h) - \nabla_z f_Z(z_i) f_Z(z_i)^{-1} h \mathcal{K}((z_i - z)/h)) \\
&\quad \times v'_i D_z m_{X|Z}(z_i)^- \nabla_z m_{Y|Z}(z_i) f_Z(z) dz \\
&= D_z m_{X|Z}(z_i)^- f_Z(z_i)^{-1} \int h^{-1} \nabla_z \mathcal{K}(\psi) f_Z(z_i + \psi h) d\psi v'_i D_z m_{X|Z}(z_i)^- \nabla_z m_{Y|Z}(z_i) \\
&\quad - D_z m_{X|Z}(z_i)^- f_Z(z_i)^{-2} \nabla_z f_Z(z_i) \int \mathcal{K}(\psi) f_Z(z_i + \psi h) d\psi v'_i D_z m_{X|Z}(z_i)^- \nabla_z m_{Y|Z}(z_i) \\
&= -D_z m_{X|Z}(z_i)^- f_Z(z_i)^{-1} \nabla_z f_Z(z_i) v'_i D_z m_{X|Z}(z_i)^- \nabla_z m_{Y|Z}(z_i) + \eta_{2i} \\
&= -g_i^- f_Z(z_i)^{-1} \nabla_z f_Z(z_i) v'_i g_i^- b_i + \eta_{2i},
\end{aligned} \tag{8.9}$$

where η_{2i} denotes higher order terms, for which, by standard arguments $n^{-1/2} \sum_i \eta_{2i} = o_p(1)$ (Here we use g_i^-, b_i to denote G_i^-, B_i at a fixed position z_i . We will now that we may replace p_n by p_n^* at the expense of a higher order term that vanishes as well (under boundedness assumptions on the densities), i.e.,

$$n^{-1/2} \sum_i \mathbb{E} [p_n(S_i, S_j) - p_n^*(S_i, S_j) | S_i] = o_p(1).$$

To see this, consider a typical expression in $\mathbb{E} [p_n(S_i, S_j) - p_n^*(S_i, S_j) | S_i]$. Using the right hand side of the third equality in equation (8.9),

$$\begin{aligned}
\rho_{ni} &= D_z m_{X|Z}(z_i)^- \left\{ \hat{f}_Z(z_i)^{-1} - f_Z(z_i)^{-1} \right\} \nabla_z f_Z(z_i) v'_i D_z m_{X|Z}(z_i)^- \nabla_z m_{Y|Z}(z_i) \\
&= \left\{ f_Z(z_i) - \hat{f}_Z(z_i) \right\} \hat{f}_Z(z_i)^{-1} D_z m_{X|Z}(z_i)^- f_Z(z_i)^{-1} \nabla_z f_Z(z_i) v'_i D_z m_{X|Z}(z_i)^- \nabla_z m_{Y|Z}(z_i),
\end{aligned}$$

where $\hat{f}_Z(z_i) = (n-1)^{-1} \sum_{s \neq i} \mathcal{K}_{hs}(Z_i)$. Next, write

$$n^{-1/2} \sum_i \rho_{ni} = n^{1/2} \int \frac{f_Z(z) - \hat{f}_Z(z)}{\hat{f}_Z(z)} \chi(z, v) \hat{F}_{ZV}(dz, dv), \tag{8.10}$$

where $\chi(z, v) = D_z m_{X|Z}(z)^- f_Z(z)^{-1} \nabla_z f_Z(z) v' D_z m_{X|Z}(z)^- \nabla_z m_{Y|Z}(z)$, and \hat{F}_{ZV} denotes the empirical cdf. Considering the denominator in (8.10), observe that

$$\frac{1}{\left| f_Z(z) + \hat{f}_Z(z) - f_Z(z) \right|} \leq \frac{1}{\left| f_Z(z) \right| - \left| \hat{f}_Z(z) - f_Z(z) \right|} \leq \frac{2}{b}, \tag{8.11}$$

since $f_Z(z) \geq b$ by the assumption that Z is continuously distributed RV on \mathcal{B} , with density bounded away from zero. Moreover, $\left| \hat{f}_Z(z) - f_Z(z) \right| \leq b/2$ with probability going to one, as $\hat{f}_Z(z)$ is consistent by assumptions on kernels and bandwidths. Hence, $n^{-1/2} \sum_i \rho_{ni}$ is, in absolute value, bounded by

$$c \sup_{z \in \mathcal{B}} \left| f_Z(z) - \hat{f}_Z(z) \right| n^{-1/2} \sum_i |\chi(Z_i, V_i)|, \tag{8.12}$$

But since $n^{-1/2} \sum_i |\chi(Z_i, V_i)|$ converges by a standard CLT for *iid* random variables to a normal limit (provided the second moment are finite which we tacitly assume), and $\sup_{z \in \mathcal{B}} |f_Z(z) - \hat{f}_Z(z)| = O_p\left(h^{2r} + (nh^L)^{-1/2} \log n\right) = o_p(1)$ under general conditions, it follows that $n^{-1/2} \sum_i \rho_{ni} = o_p(1)$. Similar arguments can be applied to any other term appearing in $\mathbb{E}[p_n(S_i, S_j) - p_n^*(S_i, S_j) | S_i]$, implying that the difference vanishes.

Repeating the same arguments as from the start of Step 1b, we can show that $G_i^-(\hat{B}_i - B_i) = G_i^- f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) Q_i + \eta_{3n}$, where $Q_i = Y_i - m_{Y|Z}(Z_i)$. Returning to (8.7)

$$\begin{aligned} S_{1n} &= n^{-1} \sum_i \left\{ G_i^- B_i + G_i^- \left[(G_i - \hat{G}_i) G_i^- B_i + (\hat{B}_i - B_i) \right] \right\} \\ &= n^{-1} \sum_i \left\{ G_i^- B_i - G_i^- f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) V_i' G_i^- B_i - G_i^- f_Z(Z_i)^{-1} \nabla_z f_Z(Z_i) Q_i \right\} + n^{-1} \sum_i T_{3i}, \end{aligned}$$

where T_{3i} denotes all higher order terms that. Note that $\sqrt{n} [n^{-1} \sum_i T_{3i}] = o_p(1)$, by arguments above..

Step 1c: To analyze all terms that affect the distribution and are contained in S_{2n} , consider S_{2n}^2 first. More specifically,

$$\begin{aligned} S_{2n}^2 &= n^{-1/2} \sum_i G_i^- \left\{ D_z m_{X|Z}(Z_i)' - (n-1)^{-1} \sum_{j \neq i} \mathcal{W}_{jn}(Z_i) m_{X|Z}(Z_j)' \right\} G_i^- B_i \\ &= \sqrt{n} \int \int [D_z m_{X|Z}(\zeta)']^- \\ &\quad \times \left\{ D_z m_{X|Z}(\zeta)' - \left[\frac{h^{-L-1} \nabla_z \mathcal{K}((z-\zeta)/h)}{\hat{f}_Z(\zeta)} - \frac{\nabla_z \hat{f}_Z(\zeta) h^{-L} \mathcal{K}((z-\zeta)/h)}{(\hat{f}_Z(\zeta))^2} \right] m_{X|Z}(z)' \right\} \\ &\quad \times [D_z m_{X|Z}(\zeta)']^- \nabla_z m_{Y|Z}(\zeta) \hat{F}_Z(dz) \hat{F}_Z(d\zeta). \end{aligned}$$

Next, let $S_{2n}^2 = A^1 + \omega_n$ where A^1 equals S_{2n}^2 with the exception that we replace \hat{F}_Z by F_Z , and we replace $\hat{f}_Z(\zeta)$ by $f_Z(\zeta)$. Hence we get a remainder term that contains expressions of the form $\hat{F}_Z - F_Z$ and $\hat{f}_Z(\zeta) - f_Z(\zeta)$. In the case of the replacement of \hat{F}_Z by F_Z , we can appeal to Glivenko-Cantelli together with the fact that \mathcal{B} is compact, and by arguments as in equations (8.10) and (8.11), we can show that $\omega_n = o_p(A^1)$, so that we focus on the leading term A^1 . After change of variable, this is

$$\begin{aligned} &\sqrt{n} \int [D_z m_{X|Z}(\zeta)']^- \{ D_z m_{X|Z}(\zeta)' - \\ &\quad \times \int \frac{h^{-1} \nabla_\psi \mathcal{K}(\psi)}{f_Z(\zeta)} m_{X|Z}(\psi h + \zeta)' f_Z(\psi h + \zeta) d\psi - \int \frac{\nabla_z f_Z(\zeta) \mathcal{K}(\psi)}{(f_Z(\zeta))^2} m_{X|Z}(\psi h + \zeta)' f_Z(\psi h + \zeta) d\psi \} \\ &\quad \times [D_z m_{X|Z}(\zeta)']^- \nabla_z m_{Y|Z}(\zeta) f_Z(\zeta) d\zeta. \end{aligned}$$

Then make use of partial integration and apply a standard Taylor expansion, to obtain that

$$A^1 = \sqrt{n} \int [D_z m_{X|Z}(\zeta)']^- BX(\zeta) [D_z m_{X|Z}(\zeta)']^- \nabla_z m_{Y|Z}(\zeta) f_Z(\zeta) d\zeta + O(\sqrt{nh}^r). \quad (8.13)$$

where denotes higher order bias terms, i.e. $BX(\zeta) = \sum_{l=1\dots r} \mu_k h^l BX_l(\zeta)$, and $BX_l(\zeta)$ contains sums of products of all higher order derivatives of $m_{X|Z}$ and f_Z , where the order of the product of derivatives combined is at most of order $l + 1$. The expectations of these terms exist due to assumption 9, and provided that $r = 2L$ in connection with assumption 11. Consequently, $\sqrt{n}S_{2n}^2 = o_p(1)$. Under similar conditions on $BY(\zeta)$ (cf. assumption 9), and by similar arguments $\sqrt{n}S_{3n}^2 = o_p(1)$ and hence the bias expression proves asymptotically negligible under our assumptions.

Step 1d: Finally, the first terms provide us with the variance. Since $\mathbb{E}[(V'_i, Q'_i)' | Z_i] = 0$, σ_{1i} is uncorrelated with σ_{2i} and σ_{3i} . The result follows by application of a standard central limit theorem. *Q.E.D.*

Step 2: The Behavior of Higher Order Terms

The characteristic feature of all terms in the expansion is that they involve higher powers in $G_i - \hat{G}_i$ or $\hat{B}_i - B_i$. Intuitively, what happens is that these terms will add a factor that tends to zero faster as the variance terms cancel, and the term is of the order of the squared bias terms. To fix ideas, recall that $B_i = \nabla_z m_{Y|Z}(Z_i)$ and consider

$$\begin{aligned} & n^{-1/2} \sum_i G_i^- (G_i - \hat{G}_i) (\hat{B}_i - B_i) \\ &= n^{1/2} \int G_i^- (m_{X|Z}(z)' - D_z \hat{m}_{X|Z}(z)') (\nabla_z \hat{m}_{Y|Z}(z) - \nabla_z m_{Y|Z}(z)) \hat{F}_Z(dz). \end{aligned}$$

The expression on the right hand side is in absolute value bounded by

$$n^{1/2} \sup_{z \in \mathcal{B}} |D_z m_{X|Z}(z)' - D_z \hat{m}_{X|Z}(z)'| \sup_{z \in \mathcal{B}} |\nabla_z \hat{m}_{Y|Z}(z) - \nabla_z m_{Y|Z}(z)| \underbrace{n^{1/2} \int |G_i^-| \hat{F}_Z(dz)}_{C_n}$$

Since $\sup_{z \in \mathcal{B}} |D_z m_{X|Z}(z)' - D_z \hat{m}_{X|Z}(z)'| \sup_{z \in \mathcal{B}} |\nabla_z \hat{m}_{Y|Z}(z) - \nabla_z m_{Y|Z}(z)| = O_p(h^{2r} + (nh^{L+2})^{-1} \ln(n))$ by an extensions to a theorem of Masry (1994), and C_n converges to a nondegenerate random variable, provided that the second moment of G_i^- is finite (which is implied by assumption 3), this term is $o_p(1)$ under general conditions. Materially similar, yet more involved arguments can be used to establish the assertion for the other higher order terms, using assumptions 3 and 4. *Q.E.D.*

Step 3: Modifications with Generated Dependent Variables - Theorem 3

To have an idea why T_{2n} and T_{3n} vanish, consider first T_{2n} in Step 3a, and then T_{3n} in Step 3b.

Step 3a: Recall that

$$T_{2n} = n^{-1} \sum_i [D_z \widehat{m}_{X|Z}(Z_i)']^{-} \sum_{j \neq i} \nabla_z W_j(Z_i) [\mathbb{K}\{(P_j - 0.5)/h\} - \mathbb{I}\{P_j < 0.5\}] B(Z_j).$$

As before, we analyze this expression in several steps. We start by considering

$$T_{2n}^* = n^{-1} \sum_i [D_z m_{X|Z}(Z_i)']^{-} \sum_{j \neq i} \nabla_z W_j(Z_i) [\mathbb{K}\{(P_j - 0.5)/h\} - \mathbb{I}\{P_j < 0.5\}] B(Z_j),$$

and note that $T_{2n} = T_{2n}^* + R_n$, where R_n contains the difference $[D_z \widehat{m}_{X|Z}(Z_i)']^{-} - [D_z m_{X|Z}(Z_i)']^{-}$ instead of $[D_z m_{X|Z}(Z_i)']^{-}$. As is easy to see (given the discussion above), R_n produces a faster vanishing higher order bias term. Quite obviously, this expression has a similar structure as the one analyzed in Step 1b above, save for the fact that Y_j is replaced by $\mathbb{K}\{(P_j - 0.5)/h\} - \mathbb{I}\{P_j < 0.5\}$. Following the same argumentation as the one in Step 1c, we arrive at the crucial decomposition $T_{2n}^* = T_{2n}^{**} + \varrho_n$, where ϱ_n are terms that converge faster by Glivenko-Cantelli and compact support \mathcal{B} , and T_{2n}^{**} is defined as follows:

$$\begin{aligned} & T_{2n}^{**} \\ &= \int \int \int h^{-(L+1)} (D_z m_{X|Z}(\zeta)^- f_Z(\zeta)^{-1} (\nabla_\zeta \mathcal{K}((z - \zeta)/h) - \nabla_z f_Z(\zeta) f_Z(\zeta)^{-1} h \mathcal{K}((\zeta - z)/h)) \\ & \quad \times [\mathbb{K}\{(p - 0.5)/h\} - \mathbb{I}\{p < 0.5\}] D_z m_{X|Z}(\zeta)^- \nabla_z m_{Y|Z}(\zeta) F_{PZ}(dp, dz) F_Z(d\zeta) \\ &= \int g(\zeta) \int \int h^{-1} \nabla_\psi \mathcal{K}(\psi) [\mathbb{K}\{(p - 0.5)/h\} - \mathbb{I}\{p < 0.5\}] f_Z(\zeta + \psi h) F_{P|Z}(dp; \zeta + \psi h) d\psi \times \\ & \quad D_z m_{X|Z}(\zeta)^- \nabla_z m_{Y|Z}(\zeta) - g(\zeta) \nabla_z f_Z(\zeta) \int \int \mathcal{K}(\psi) [\mathbb{K}\{(p - 0.5)/h\} - \mathbb{I}\{p < 0.5\}] \times \\ & \quad f_Z(\zeta + \psi h) F_{P|Z}(dp; \zeta + \psi h) d\psi g(\zeta) \nabla_z m_{Y|Z}(\zeta) F_Z(d\zeta) \\ &= Q_{1n} - Q_{2n}, \end{aligned}$$

where $g(\zeta) = D_z m_{X|Z}(\zeta)^- f_Z(\zeta)^{-1}$. Next, consider the inner integral in Q_{1n} :

$$\begin{aligned} & h^{-1} \int \int \nabla_\psi \mathcal{K}(\psi) \mathbb{K}\{(p - 0.5)/h\} f_Z(\zeta + \psi h) dF_{P|Z}(dp; \zeta + \psi h) d\psi \\ & - h^{-1} \int \int \nabla_\psi \mathcal{K}(\psi) \mathbb{I}\{p < 0.5\} f_Z(\zeta + \psi h) dF_{P|Z}(dp; \zeta + \psi h) d\psi \\ &= h^{-1} \int \int \nabla_\psi \mathcal{K}(\psi) K(\tau) F_{P|Z}(0.5 + h\tau; \zeta + \psi h) f_Z(\zeta + \psi h) d\tau d\psi \quad (8.14) \\ & - h^{-1} \int \nabla_\psi \mathcal{K}(\psi) F_{P|Z}(0.5; \zeta + \psi h) f_Z(\zeta + \psi h) d\psi, \end{aligned}$$

where we made use of Fubini's theorem in connection with standard arguments for integrals of kernels. Next, use integration by parts to obtain that the rhs of (8.14) equals

$$\begin{aligned} & \int \mathcal{K}(\psi) \nabla_{\psi} [F_{P|Z}(0.5; \zeta + \psi h) f_Z(\zeta + \psi h)] d\psi \\ & - \int \mathcal{K}(\psi) \int K(\tau) \nabla_{\psi} [F_{P|Z}(0.5 + h\tau; \zeta + \psi h) f_Z(\zeta + \psi h)] d\tau d\psi \end{aligned} \quad (8.15)$$

Inserting $F_{P|Z}(0.5 + h\tau; \zeta + \psi h) = F_{P|Z}(0.5; \zeta + \psi h) + h\tau f_{P|Z}(0.5; \zeta + \psi h) + \dots + (r_1!)^{-1} h^{r_1} \tau^{r_1} \partial_p^{r_1-1} f_{P|Z}(0.5 + \lambda h\tau; \zeta + \psi h)$, where $\lambda \in (0, 1)$, we obtain that (8.15) reduces, under the familiar assumption on all moments of the kernel up to order r_1 to be zero to

$$- \int \mathcal{K}(\psi) (r_1!)^{-1} h^{r_1} \mu_{r_1} \nabla_{\psi} \partial_p^{r_1-1} f_{PZ}(0.5; \zeta + \psi h) d\psi,$$

plus a term of smaller order. Applying standard arguments, in particular expand $\partial_p^{r_1-1} f_{P|Z}(0.5; \zeta + \psi h)$ in ψ , we obtain that $\sqrt{n} T_{2n}^{**} = o_p(1)$, provided that $\sqrt{nh^{r_1}} h^r = o(1)$. The same argumentation holds for Q_{2n} .

Step 3b: Next, consider

$$T_{3n} = n^{-1} \sum_i [D_z \widehat{m}_{X|Z}(Z_i)']^{-} \sum_{j \neq i} \nabla_z W_j(Z_i) \left[\mathbb{K} \left\{ (\widehat{P}_j - 0.5)/h \right\} - \mathbb{K} \left\{ (P_j - 0.5)/h \right\} \right] B(Z_j),$$

we can rewrite the last term on the right hand side as:

$$\begin{aligned} T_{3n} &= n^{-1} \sum_i [D_z \widehat{m}_{X|Z}(Z_i)']^{-} \sum_{j \neq i} \nabla_z W_j(Z_i) h^{-1} K \left\{ (P_j - 0.5)/h \right\} \left(\widehat{P}_j - P_j \right) B(Z_j) + R_{1n} \\ &= T_{4n} + R_{1n}, \end{aligned}$$

where R_{1n} denotes higher order terms in a mean value expansion, and $R_n = o_p(T_{4n})$. Using again $[D_z \widehat{m}_{X|Z}(Z_i)']^{-} = [D_z m_{X|Z}(Z_i)']^{-} + \left[[D_z \widehat{m}_{X|Z}(Z_i)']^{-} - [D_z m_{X|Z}(Z_i)']^{-} \right]$, which produces a leading term T_{5n} and again a faster converging remainder, we find that

$$\begin{aligned} & \sqrt{n} T_{5n} \\ &= n^{-1/2} \sum_i [D_z m_{X|Z}(Z_i)']^{-} n^{-1} \sum_{j \neq i} \nabla_z \frac{h^{-L-2} \mathcal{K}((Z_j - Z_i)/h)}{\widehat{f}_Z(Z_i)} K((P_j - 0.5)/h) \left(\widehat{P}_j - P_j \right) \\ &= \sqrt{n} \int \int [D_z m_{X|Z}(z)']^{-} h^{-2} \int \nabla_{\psi} \frac{\mathcal{K}(\psi)}{\widehat{f}_Z(z)} K_1((p(z + \psi h_1, \varpi) - 0.5)/h) \\ & \quad \times (\widehat{p}(z + \psi h_1, \varpi) - p(z + \psi h_1, \varpi)) f_{ZV}(z + \psi h, \varpi) d\psi d\varpi F_Z(dz) \\ &= \sqrt{n} \int \int [D_z m_{X|Z}(z)']^{-} h^{-2} \widehat{f}_Z(z)^{-1} \times \\ & \quad \nabla_z [K_1((p(z, \varpi) - 0.5)/h) (\widehat{p}(z, \varpi) - p(z, \varpi)) f_{ZV}(z, \varpi)] d\varpi F_Z(dz) + \rho_n \\ &= T_{6n} + \rho_n, \end{aligned}$$

where $\rho_n = o_p(T_{6n})$. Hence, $\sqrt{n}T_{5n}$ is bounded in absolute value by

$$c_1 \sup_{z,v \in \mathcal{B} \times \mathcal{V}} |\nabla_z \widehat{p}(z, v) - \nabla_z p(z, v)| b_{1n} + c_2 \sup_{z,v \in \mathcal{B} \times \mathcal{V}} |\widehat{p}(z, v) - p(z, v)| b_{2n}, \quad (8.16)$$

where $b_{1n} = n^{-1/2} \sum_i \left| [D_z m_{X|Z}(Z_i)']^- K_1((p(Z_i, V_i) - 0.5) / h_1) f_{ZV}(Z_i, V_i) \right|$ and $b_{2n} = n^{-1/2} \sum_i \left| [D_z m_{X|Z}(Z_i)']^- \nabla_z [K_1((p(Z_i, V_i) - 0.5) / h_1) f_{ZV}(Z_i, V_i)] \right|$ converge to nondegenerate distributions. To see this, pick $b_{1n} = h_1^{1/2} (nh_1)^{-1/2} \sum_i \left| [D_z m_{X|Z}(Z_i)']^- \right| f_{ZV}(Z_i, V_i) \times K_1((P_i - 0.5) / h_1) = h_1^{1/2} b_{3n}$, where b_{3n} is a nonparametric estimator of

$$\mathbb{E} \left[\left| [D_z m_{X|Z}(Z_i)']^- \right| f_{ZV}(Z_i, V_i) | P_i = 0.5 \right] f_P(0.5).$$

Observe that b_{3n} converges to a nondegenerate limiting distribution provided that the second moment of $\left| [D_z m_{X|Z}(Z_i)']^- \right| f_{ZV}(Z_i, V_i)$ exist. But this follows by elementwise square integrability in assumption 3, together with the boundedness assumption 14. Hence,

$$c_1 \sup_{z,v \in \mathcal{B} \times \mathcal{V}} |\nabla_z \widehat{p}(z, v) - \nabla_z p(z, v)| b_{1n} = o_p(1).$$

Similar arguments can be made for the second summand in (8.16), using the boundedness of the derivatives in assumption 14. Consequently, $\sqrt{n}T_{5n} = o_p(1)$, implying that $\sqrt{n}T_{3n} = o_p(1)$. *Q.E.D.*

Proof of the Identification Theorems

Proof of Theorem 4

To start out with, rewrite

$$Y = \mathbb{I}\{X'\beta(A_1) + A_2'\gamma(A_1) > 0\} = \mathbb{I}\{X'\beta + \underbrace{X'(\beta(A_1) - \beta) + A_2'\gamma(A_1)}_U > 0\}. \quad (8.17)$$

and observe first that $k_{U|ZV}^{0.5}(Z, V) = g(V)$ implies that $k_{U|XV}^{0.5}(X, V) = g(V)$. To see this, note that by the definition of the α -quantile to obtain

$$\mathbb{P}(U \leq k_{U|Z,V}^\alpha(Z, V) | Z, V) = \alpha = \mathbb{P}(U \leq k_{U|X,V}^\alpha(X, V) | X, V),$$

for any $\alpha \in (0, 1)$. Taking conditional expectations with respect to (X, V) on both sides produces

$$\mathbb{E} [\mathbb{E} \{ \mathbb{I}(U \leq k_{U|Z,V}^\alpha(Z, V)) | Z, V \} | X, V] = \mathbb{P}(U \leq k_{U|X,V}^\alpha(X, V) | X, V).$$

But due to $k_{U|Z,V}^\alpha(Z, V) = g(V)$, and the law of iterated expectations, we have that

$$\mathbb{E} [\mathbb{I}(U \leq g(V)) | X, V] = \mathbb{P}(U \leq k_{U|X,V}^\alpha(X, V) | X, V),$$

implying that $k_{U|X,V}^\alpha(X, V) = g(V)$, provided U is continuously distributed.

Hence, if we assume the median exclusion restriction $k_{U|ZV}^{0.5}(Z, V) = g(V)$, we obtain that $\nabla_x k_{Y^*|XV}^{0.5}(X, V) = \beta$. Since $Y^* = X'\beta(A_1) + A_2'\gamma(A_1) = \phi(X, A)$, and $X \perp A|V$, we can apply Hoderlein and Mammen's (2007) theorem to obtain that the (constant) derivative has the following interpretation:

$$\beta = \mathbb{E} [\beta(A_1)|X = x, V = v, Y^* = k_{Y^*|XV}^{0.5}(x, v)], \quad (8.18)$$

for all $(x, v) \in \text{supp}(X) \times \text{supp}(V)$. *Q.E.D.*

8.0.1 Proof of Theorem 5

Next, consider the case defined by assumptions 16: W.l.o.g, we consider two subsets of the support of V , denoted S_1 and S_2 . Then, let $l \nearrow$ on $S_1 = (-\infty, a)$, \searrow on $S_2 = (a, \infty)$ with inverses l_1, l_2 . Let $\vartheta(z)'\beta \leq \max_{v \in S_1} l(v) = l(a)$. Then,

$$\begin{aligned} m_{\bar{Y}|Z}(z) &= \mathbb{P} [m_{X|Z}(Z)'\beta > l(V)|Z = z] \\ &= \mathbb{P} [m_{X|Z}(Z)'\beta > l(V), V \in S_1|Z = z] + \mathbb{P} [m_{X|Z}(Z)'\beta > l(V), V \in S_2|Z = z] \\ &= \mathbb{P} [l_1(m_{X|Z}(Z)'\beta) > V \wedge a|Z = z] + \mathbb{P} [l_2(m_{X|Z}(Z)'\beta) > V \vee a|Z = z] \\ &= \int_{-\infty}^{l_1(m_{X|Z}(z)'\beta)} f_{V|Z}(v|z)dv + \int_{l_2(m_{X|Z}(z)'\beta)}^{\infty} f_{V|Z}(v|z)dv \end{aligned} \quad (8.19)$$

Taking derivatives by applying Leibnitz' rule produces

$$\begin{aligned} \nabla_z m_{\bar{Y}|Z}(z) &= D_z m_{X|Z}(z)'\beta \left[\frac{\partial l_1}{\partial s} (m_{X|Z}(z)'\beta) f_{V|Z}(l_1(m_{X|Z}(z)'\beta)|z) - \frac{\partial l_2}{\partial s} (m_{X|Z}(z)'\beta) f_{V|Z}(l_2(m_{X|Z}(z)'\beta)|z) \right] \\ &\quad + \int_{-\infty}^{l_1(m_{X|Z}(z)'\beta)} \nabla_z [\log f_{V|Z}(v|z)] f_{V|Z}(v|z)dv \\ &\quad + \int_{l_2(m_{X|Z}(z)'\beta)}^{\infty} \nabla_z [\log f_{V|Z}(v|z)] f_{V|Z}(v|z)dv \\ &= D_z m_{X|Z}(z)'\beta [\dots] + \mathbb{E} \left\{ \tilde{Y} Q_z(V; Z) | Z = z \right\}. \end{aligned} \quad (8.)$$

where $Q_z(V; Z) = \nabla_z [\log f_{V|Z}(V; Z)]$, and all the integrals on the right hand side of the first and second equality exist by assumption 16

Appendix 2: Graphs and Tables

Table A.1: Variables in Data Set

1	dma	dma code, code for television market
2	income	household income in \$
3	owncable	does household have cable TV
4	ownsat	does household have satellite TV
5	cableco	cable company
6	age	what range best describes your age
7	hhsiz	household size
8	hhcomp	household composition
9	educ	education
10	hisp	hispanic or not
11	single	single or couple
12	state	
13	rent	renter status (do they rent or own the house)
14	typeres	type of residence (house, apartment, condominium)
15	angle	dish angle
16	avgpbi	instrument, average price of basic cable across other cable franchises
17	avgppi	same for premium
18	tvself	tv choice (1: basic cable, 2: premium cable, 0: nohighTV, 3 or 4: satellite)
19	yearst	year established (satellite dish)
20	chancap	channel capacity
21	airchan	number of over the air-channels available
22	paychan	number of pay channels available
23	othchan	other channels
24	ppv	pay per view available
25	cityff	city fixed fee (tax)
26	pricebe	price of basic cable
27	gender	gender
28	varelev	variance of the local terrain and the average elevation

Table A.1: Variables in Data Set(cont.)

29	mild	local weather index
30	bright	local weather index
31	stable	local weather index
32	climate	local weather index
33	twoway	cable franchise char - probably whether signals can be sent both ways
34	hboprice	HBO price
35	density	population density in an area (city density)
36	cnts	number of sampled households in that cable franchise market
37	poprank	city code (market area: necessary to merge with damachers, cable98)

Table A.2 Summary Statistics for Forrester Data

	Mean	Std. Dev.	25%	50%	75%
Satellite	0.10	0.30	0.00	0.00	0.00
Cable	0.72	0.45	0.00	1.00	1.00
Household income in \$	57,366	28,642	32,500	55,000	87,500
Rent	0.22	0.42	0.00	0.00	0.00
Single unit dwelling	0.78	0.41	1.00	1.00	1.00
Household size	2.16	1.88	1.00	1.00	3.00
Single	0.18	0.38	0.00	0.00	0.00
Age of HH	50.59	15.42	39.00	49.00	61.00
Education in years	14.06	2.69	12.00	13.00	16.00

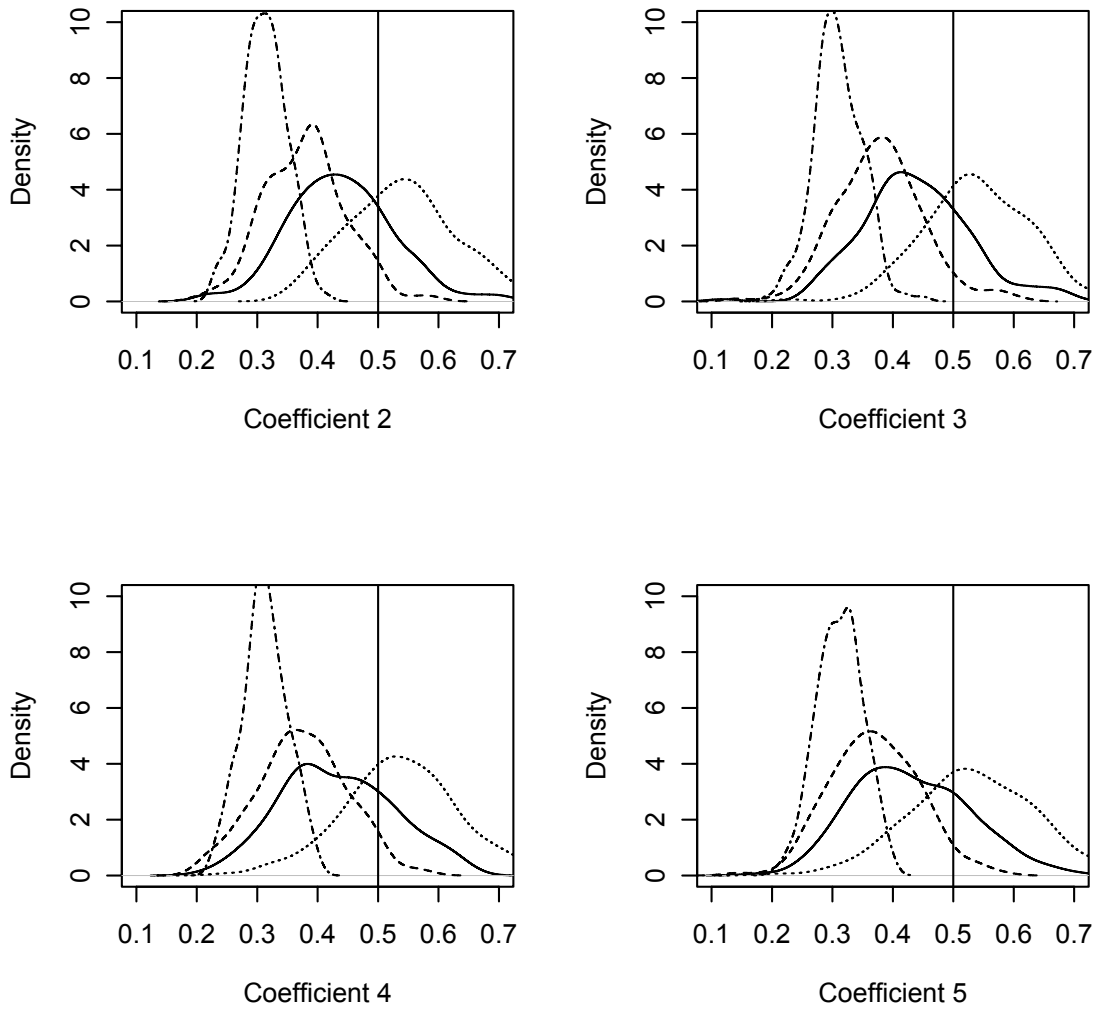
The education level corresponds to the mean education in a non-single household.

Table A.3 Summary Statistics for Warren's Factbook Data

	Mean	Std. Dev	25%	50%	75%
Monthly cable price in \$	25.45	8.39	20.88	24.43	29.95
HBO price in \$	11.13	1.51	9.95	10.95	12.45
Channel capacity	65.36	17.44	54.00	62.00	78.00
Pay-per-view available	0.92	0.26	1.00	1.00	1.00
Year franchise began	1974.94	9.82	1971	1976	1982
City franchise fee	4.06	1.55	3.00	5.00	5.00
Number of over-the-air channels	11.46	3.38	8.00	12.00	14.00
Observations	132				

Fig.1: Comparison of Distribution of Estimator for Centrality Parameter

Heterogeneity Robust (Solid Line)
Conditional Independence (Broken)
Rivers Young (Broken and Dotted)
Oracle (Dotted)

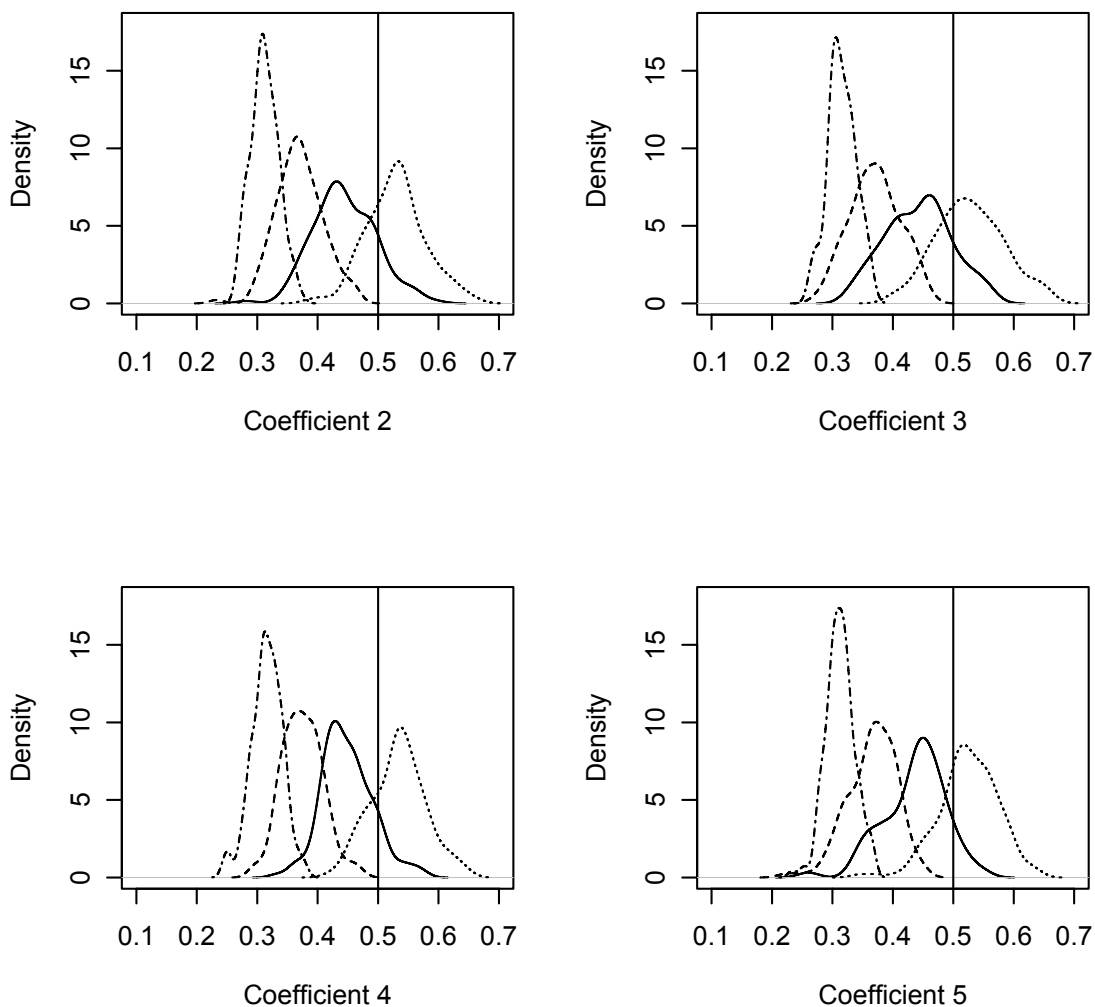


n = 2500

Kernel Density of Estimators in 400 Realizations of the DGP

Fig.2: Comparison of Distribution of Estimator for Centrality Parameter

Heterogeneity Robust (Solid Line)
Conditional Independence (Broken)
Rivers Young (Broken and Dotted)
Oracle (Dotted)

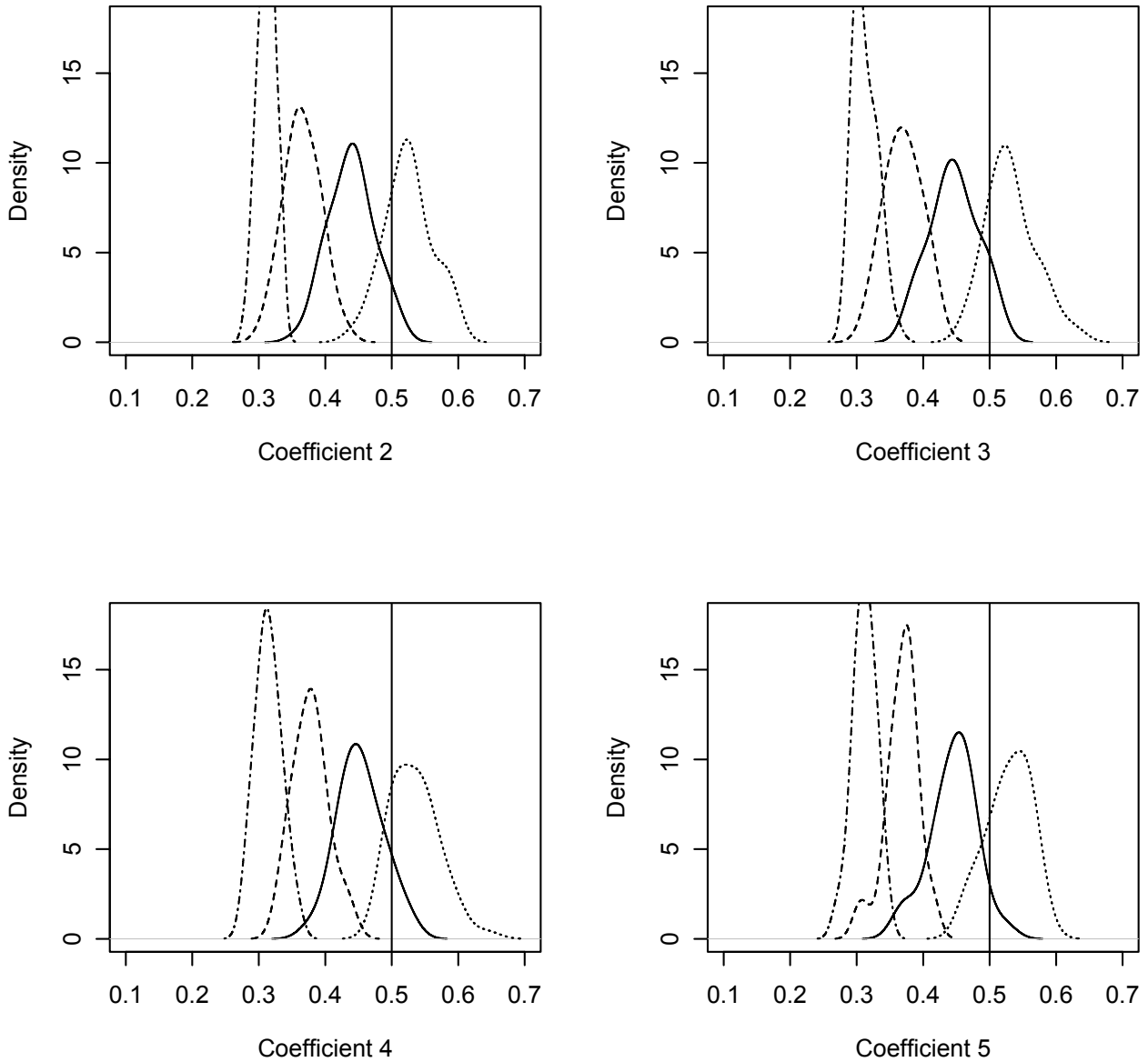


n = 7500

Kernel Density of Estimator in 200 Realizations of the DGP

Fig. 3: Comparison of Distribution of Estimators for Centrality Parameter

Heterogeneity Robust (Solid Line)
Conditional Independence (Broken)
Rivers Young (Broken and Dotted)
Oracle (Dotted)



n = 15000

Kernel Density of Estimators in 100 Realizations of DGP

Fig. 4: Finite Sample Distribution of the Bootstrap Standard Errors

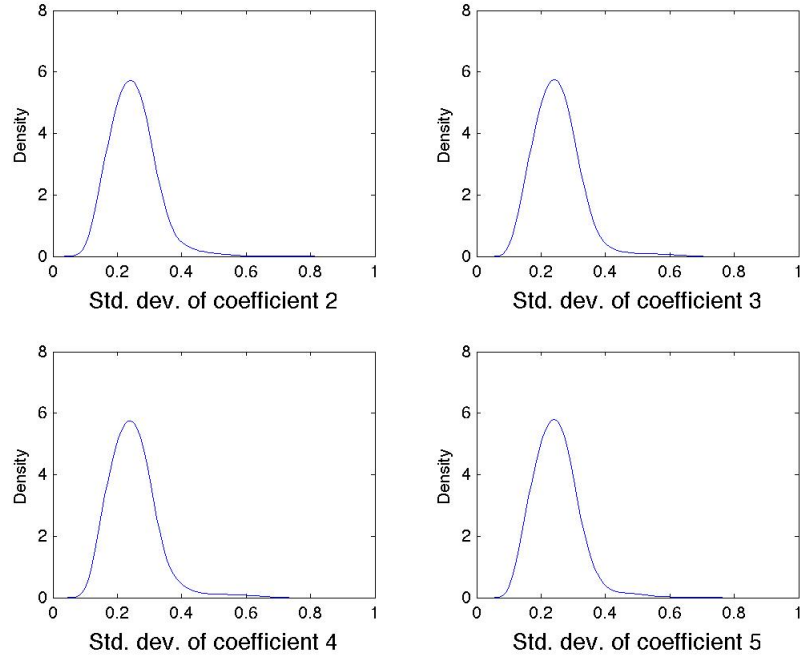


Fig. 5: Finite Sample Distribution of the Asymptotic Standard Errors

