



HCEO WORKING PAPER SERIES

Working Paper



HUMAN CAPITAL AND
ECONOMIC OPPORTUNITY
GLOBAL WORKING GROUP

The University of Chicago
1126 E. 59th Street Box 107
Chicago IL 60637

www.hceconomics.org

Simple Tests for Selection: Learning More from Instrumental Variables*

Dan A. Black

University of Chicago

Joonhwi Joo

University of Texas at Dallas

Robert LaLonde

University of Chicago

Jeffrey A. Smith

University of Wisconsin at Madison

Evan J. Taylor

University of Arizona

June 13, 2020

Abstract

We provide simple tests for selection on unobserved variables in the Vytlacil-Imbens-Angrist framework for Local Average Treatment Effects (LATEs). Our setup allows researchers not only to test for selection on either or both of the treated and untreated outcomes, but also to assess the magnitude of the selection effect. We show that it applies to the standard binary instrument case, as well as to experiments with imperfect compliance and fuzzy regression discontinuity designs, and we link it to broader discussions regarding instrumental variables. We illustrate the substantive value-added by our framework with three empirical applications drawn from the literature.

Keywords: instrumental variable, local average treatment effect, selection, test

*First draft: December 2013. We thank seminar participants at the CESifo Education Group Meetings (especially Georg Graetz and Edwin Leuven), Chicago, UIC, Indiana, IZA, Kentucky, UNC, UC-Riverside, and the Hitotsubashi Institute for Advanced Study, along with Josh Angrist, Chris Berry, Ben Feigenberg, Anthony Fisher, Anthony Fowler, Martin Huber, David Jaeger, Darren Lubotsky, Lois Miller, Nikolas Mittag, Douglas Staiger, and Gerard van den Berg for helpful comments. We thank Bill Evans for providing us with the data from Angrist and Evans (1998). Sadly, Bob LaLonde passed away while we were working on a revision of this paper. He is greatly missed.

1 Introduction

In the years since the publication of Imbens and Angrist (1994), applied researchers have embraced the interpretation of Instrumental Variables (IV) estimators as measuring the impact of treatment on the subset of respondents who comply with the instrument. Imbens and Angrist call this the Local Average Treatment Effect, or LATE. The LATE framework allows researchers to consistently estimate interpretable causal parameters in the presence of treatment effect heterogeneity.

The LATE framework, however, comes with some costs. First, the LATE approach requires the assumption that instruments have a monotonic impact on behavior. Put differently, the instruments must induce all agents to behave in a weakly uniform manner when subjected to a change in the value of the instrument. Informally, if the instrument induces some agents to enter the treatment, then the instrument must not induce any agent to leave the treatment. Second, relative to approaches based on conditional independence, LATE identifies a treatment effect only for a very specific sub-population. Third, treatment effect heterogeneity complicates the interpretation of the traditional Durbin-Wu-Hausman test for the equivalence of the IV and Ordinary Least Squares (OLS) estimands.¹ More broadly, the relationship between the OLS and IV estimates becomes uninformative about the existence of selection within the LATE framework. Thus, researchers face the paradox of using IV to address potential selection on unobserved variables, but with no clear evidence that such selection exists.

Consider the binary treatment framework of Angrist et al. (1996), with a binary instrument $Z_i \in \{0, 1\}$. Without loss of generality, let $Z_i = 1$ increase the likelihood of treatment. They show that we may divide agents into three mutually exclusive sets: the always-takers (A), defined as the sub-population with $D_i(1) = D_i(0) = 1$, the never-takers (N), defined as the sub-population with $D_i(1) = D_i(0) = 0$, and the compliers (C), defined as the sub-

¹The “OLS estimand” equals the coefficient on D in a parametric linear model with Y as the dependent variable and X and D as independent variables. The “IV estimand” equals the coefficient on D in the same linear model estimated by two-stage least squares with Z as the instrument.

population with $D_i(1) = 1$ and $D_i(0) = 0$, where $D_i = D_i(Z_i)$ denotes the treatment choice of agent i as a function of the binary instrument Z_i , with $D_i = 1$ for treatment and $D_i = 0$ for no treatment. In this framework, the Wald estimand corresponds to a LATE or

$$\Delta^w = \frac{E(Y_i|Z_i = 1) - E(Y_i|Z_i = 0)}{E(D_i|Z_i = 1) - E(D_i|Z_i = 0)} = E(Y_{1i} - Y_{0i}|C) \quad (1)$$

where Y_{1i} denotes the treated potential outcome of agent i , Y_{0i} denotes the untreated potential outcome of agent i , and $Y_i = D_i Y_{1i} - (1 - D_i) Y_{0i}$ denotes the observed outcome.

Selection on unobserved variables means that one or more of the following four conditions fails:

$$E(Y_{0i}|N, X_i) = E(Y_{0i}|C, X_i) \text{ and } E(Y_{0i}|C, X_i) = E(Y_{0i}|A, X_i) \quad (2)$$

$$E(Y_{1i}|N, X_i) = E(Y_{1i}|C, X_i) \text{ and } E(Y_{1i}|C, X_i) = E(Y_{1i}|A, X_i) \quad (3)$$

where X_i denotes a vector of conditioning variables not affected by the treatment. These conditions do not imply the equivalence of the OLS and IV estimands: By construction, the set of compliers differs from the set of people receiving treatment. Nor does equivalence of the estimands imply the conditions.² How then do we test for such selection?

In this paper, we provide a set of simple tests for the presence of selection on the treated and untreated outcomes and derive measures of the magnitude of the selection in each case. We care about the incidence and magnitude of selection because (i) they inform our understanding of the underlying behavior, (ii) they affect the reasonableness of generalizing impact estimates obtained via instrumental variables beyond the compliers, and (iii) they affect the relative attractiveness of estimates based on IV and on conditional independence in particular substantive contexts. We illustrate (i) in our empirical applications and expand

²To see this, consider the following example: Suppose that $Pr(C) = Pr(A) = Pr(N) = 1/3$ and that $Pr(Z_i = 1) = 1/2$. Further, let $E(Y_{1i}|A) = 1$, $E(Y_{1i}|C) = 0$, $E(Y_{0i}|C) = 0$, and $E(Y_{0i}|N) = 1$. In this case, the OLS estimand equals zero. The IV estimand also equals zero but $E(Y_{1i}|A) > E(Y_{1i}|C)$ and $E(Y_{0i}|N) > E(Y_{0i}|C)$ so we clearly have selection on Y_{1i} and Y_{0i} .

on (ii) and (iii) in the penultimate section of the paper.

Our tests consider two of the four conditions in equations (2) and (3) (or their distributional analogues) as the data do not provide information about the other two. As such, we test necessary (but not sufficient) conditions for the absence of selection. We develop and apply tests of both conditional mean independence and conditional independence of distributions. Relative to the traditional Durbin-Wu-Hausman test that compares the IV and OLS estimates, our tests reveal substantively relevant information regarding whether selection occurs on the treated outcome, the untreated outcome, or both.

Drawing on the work of Black, Sanders, Taylor, and Taylor (2015), our tests of conditional mean independence come in two forms. First, we compare the conditional mean outcomes of agents who comply with the instrument when not treated to those of agents who never take treatment. Second, we compare the conditional mean outcomes of agents who comply with the instrument when treated to those of agents who always take treatment. Mechanically, we implement these tests by estimating outcome equations for those who are untreated, or treated, as a function of the covariates and the instruments (or the probability of selection). With a simple adjustment, our tests allow researchers to assess the economic magnitude of any selection effects as well. Our distributional tests proceed along similar lines, but replace the dependent variable with other features of the outcome distribution.

Our tests resemble those in Heckman's (1979) seminal paper on the bivariate normal selection model. In the two-step estimator for the normal selection model with a common treatment effect, the inverse Mills ratio represents the control function, and the coefficient on the inverse Mills ratio identifies the correlation between the errors of the outcome equation and the selection equation. Under the null hypothesis of no selection on unobserved variables, a simple test for selection asks if the coefficient on the inverse Mills ratio equals zero. Allowing selection on both the treated and untreated outcomes in a trivariate normal framework, as in Björklund and Moffitt (1987), makes the parallel to our tests an exact one; see Blundell, Dearden, and Sianesi, (2005) for implementation. In more general selection models that

retain the common effect assumption but dispense with normality, the exact form of the control function is unknown, and the control function is estimated semiparametrically as in the estimators examined in Newey, Powell, and Walker (1990). Yet the nature of the test remains the same.

Our paper is closely related to Heckman, Schmieder, and Urzua (2010), hereinafter HSU, who derive both parametric and non-parametric tests for the correlated random coefficient model. Formally, HSU develop a test for the independence of treatment status and the idiosyncratic effect of treatment conditional on covariates. Note that selection on the idiosyncratic component of the treatment effect implies selection on one or both of Y_{1i} and Y_{0i} , while selection on outcomes need not imply selection on the idiosyncratic component of the treatment effect.

Building on the work of Heckman and Vytlačil (2005, 2007a,b), who show that conditional independence of Y_{0i} and Y_{1i} implies constant marginal treatment effects, HSU (2010) propose parametric and nonparametric tests that regress the realizations of the dependent variable against the estimated propensity score (which includes the instruments) to see if the realizations of the outcome variables are linear functions of the propensity score. But as HSU note, their non-parametric tests suffer from low power in sample sizes common in empirical studies. In addition, our tests are considerably easier to implement than their non-parametric tests, which generally require the use of the bootstrap procedures of Romano and Wolf (2005) as well as the step-down method of multiple hypothesis testing in Romano and Shaikh (2006). Our tests also provide more insight into the precise nature of the selection problem because we allow for selection on one or both of Y_{1i} and Y_{0i} .

Similarly, in the context of a Marginal Treatment Effects (MTE) model, Brinch, Mogstad, and Wiswall (2017) propose testing for a constant MTE by regressing $Y_i = Y_{0i} - D_i(Y_{1i} - Y_{0i})$ against D_i , Z_i and their interactions. As they note, the extension of their test to a model with covariates is straightforward, but for a linear parametric model the test could involve the estimation of a large number of interaction terms. Kowalski (2018a) describes alternative

tests that relax the linearity assumptions in Brinch et al. (2017) and provides links to the literature in health economics that seeks to empirically distinguish moral hazard and adverse selection; Kowalski (2018b) applies her tests in the context of the Oregon Health Insurance experiment.

Bertanha and Imbens (2019) consider closely related tests in the context of fuzzy regression discontinuity designs and also provide an extensive discussion of the Durbin-Wu-Hausman test in the context of treatment effect heterogeneity. Huber (2013) provides a Wald test for the exogeneity of non-compliance in experiments closely related to ours, but does not extend his analysis to other IV settings. Angrist (2004) proposes a test that compares the estimated treatment effect for compliers to an estimate obtained using the always-takers and the never-takers. His test does not distinguish among selection on one or both of Y_{1i} and Y_{0i} and assumes the magnitude of the treatment effect does not vary with covariates. Guo, Cheng, Lorch, and Smallet (2014) provide a substantially more complex test than ours in the IV context. Battistin and Rettore (2008) and Costa Dias, Ichimura, and van den Berg (2013) consider related tests that exploit particular empirical contexts. Donald, Hsu, and Lieli (2014) derive a test that uses differences between estimates of the Average Treatment Effect on the Treated (ATET) and LATE parameters in the context of one-sided non-compliance.³

We complement the existing literature by generalizing the simple tests in Huber (2013) and Bertanha and Imbens (2019) within a common framework that encompasses standard binary instruments, experiments with imperfect compliance, and fuzzy regression discontinuity designs and by separately considering selection on Y_{1i} and Y_{0i} , which illuminates the underlying economics. We extend the literature by (a) considering tests of conditional independence of distributions as well as conditional mean independence; (b) providing measures of the amount of selection that complement our tests; (c) linking our tests to the classical literature on selection models; (d) relating our tests to ongoing discussions about instrumental

³More distant relations include Kitagawa (2015) and Huber and Mellace (2015), both of which propose tests of instrument validity in the LATE context.

variables methods and the trade-offs involved between approaches that build on instruments and those that build on conditional independence; and (e) providing three empirical applications, one each for standard instrumental variables, fuzzy regression discontinuity, and experiments with imperfect compliance.

We find it peculiar that while the LATE revolution has led to a more sophisticated interpretation of IV estimates, researchers rarely make an empirical case via testing for the use of instrumental variables methods. Heckman, Ichimura, Smith, and Todd (1998) find that most of the difference between non-experimental and experimental estimates of the treatment effect in the Job Training Partnership Act (JTPA) data results from lack of common support and from differences in the distributions of covariates, leaving selection on unobserved variables to account for *only about seven percent* of the difference. Blundell et al. (2005) find little evidence that their matching estimates suffer from selection bias when estimating their single treatment model using the very rich National Child Development Survey data. Similarly, when using comparable outcome measures for the treated and untreated units, Diaz and Handa (2006)'s propensity score matching estimates are quite similar to the experimental evidence from the famous PROGRESA experiment with their set of conditioning variables.

While by no means conclusive, matching on a rich set of covariates motivated by theory and the institutional context and limiting the analysis to the region of common support between treated and untreated units substantially reduces bias in many empirical contexts. Put differently, sometimes unconditional differences in mean outcomes between treated and untreated units arise mainly from selection on observed variables and lack of common support rather than from selection on unobserved variables. Indeed, given the necessary, and often empirically quite large, increase in the variance of estimates when using instrumental variables methods, a researcher might well prefer a precise but modestly biased OLS estimate to a consistent but imprecise IV estimate, much as in non-parametric estimation, a researcher trades off bias and variance via the choice of a bandwidth or other tuning parameter. We

develop a simple estimator of the magnitude of the bias implicit in the OLS estimate which allows a quantitative comparison of the costs and benefits of IV relative to OLS.

In the next section, we outline the restrictions necessary for matching and OLS. In section 3, we outline the necessary assumptions for Imbens and Angrist’s (1994) IV estimation and the latent index approach of Vytlačil (2002). In section 4, we outline our simple tests of conditional independence and provide formulas for the magnitude of the selection. Section 5 reports on our empirical applications. In section 6 we relate our tests to the broader literatures on instrumental variables and to the choice between IV and OLS estimation of the linear model. The final section concludes.

2 Matching, OLS, and Selection on Observed Variables

Using the notation introduced above, we define the causal impact of the treatment on agent i as $\delta_i = Y_{1i} - Y_{0i}$. The fundamental problem of evaluation is that we observe only one of the two potential outcomes; we must estimate the other, which the literature calls the missing counterfactual. Matching estimators represent one intuitive class of estimators for generating the missing counterfactuals. These estimators rely on the assumption that researchers have sufficiently rich covariates that any differences in the treatment decisions of the agents are independent of the agents’ potential outcomes conditional on the covariates. Let X_i denote those covariates. Formally, matching estimators rely on two assumptions:

First, and most vexing, matching estimators require a Conditional Independence Assumption (CIA). Various “flavors” of CIA correspond to different parameters of interest. The strongest flavor demands that:

$$(Y_{0i}, Y_{1i}) \perp\!\!\!\perp D_i | X_i \tag{CIA}$$

where $\perp\!\!\!\perp$ denotes statistical independence. This version of the CIA applies to the Average

Treatment Effect (ATE), or $\Delta^{ATE} = E(Y_{1i} - Y_{0i})$. The CIA for Y_{0i} assumes

$$Y_{0i} \perp\!\!\!\perp D_i | X_i \tag{CIA^0}$$

This version of the CIA allows the estimation of $\Delta^{ATE} = E(Y_{1i} - Y_{0i} | D_i = 1)$. Because researchers observe Y_{1i} for those who are treated, estimation of the Δ^{ATE} only requires the weaker (CIA⁰) rather than the (CIA). Similarly, when estimating the Average Treatment Effect on the Non-treated (ATEN), researchers need only assume

$$Y_{1i} \perp\!\!\!\perp D_i | X_i \tag{CIA^1}$$

which allows the estimation of $\Delta^{ATEN} = E(Y_{1i} - Y_{0i} | D_i = 0)$. Of course, the (CIA) implies that both (CIA¹) and (CIA⁰) hold.

Technically, the ATE, ATET or ATEN do not require full conditional independence for identification. Rather, they “only” require conditional mean independence, which comes in three parallel flavors:

$$E(Y_{ji} | X_i, D_i = 1) = E(Y_{ji} | X_i, D_i = 0) \quad j \in \{0, 1\} \tag{CMIA}$$

$$E(Y_{0i} | X_i, D_i = 1) = E(Y_{0i} | X_i, D_i = 0) \tag{CMIA^0}$$

$$E(Y_{1i} | X_i, D_i = 1) = E(Y_{1i} | X_i, D_i = 0) \tag{CMIA^1}$$

We mainly consider tests of (CMIA⁰) and (CMIA¹), but also develop and implement tests of (CIA⁰) and (CIA¹).

Second, matching estimators also require the Common Support Assumption (CSA) or

$$0 < Pr(D_i = 1 | X_i) < 1. \tag{CSA}.$$

The (CSA) requires an untreated comparison unit with approximately the same realization

of the covariates as each treated unit. Unlike the (CIA), researchers may test the (CSA) in observational data. When the (CSA) fails in practice, as it sometimes does, researchers generally change the definition of the relevant population to that over which the (CSA) holds, reflecting the limited variation that the data provide.⁴

When applying semi-parametric or non-parametric matching methods, researchers commonly specify the potential outcome functions as

$$Y_{1i} = g_1(X_i) + \epsilon_{1i} \tag{4}$$

$$Y_{0i} = g_0(X_i) + \epsilon_{0i} \tag{5}$$

where $(g_0(\cdot), g_1(\cdot))$ denote the unknown conditional mean functions and $(\epsilon_{0i}, \epsilon_{1i})$ summarize the residual uncertainty associated with unobserved variables. With the (CSA) and the appropriate version of the (CIA) or (CMIA), researchers may use a variety of methods to estimate the unknown conditional mean functions. A common alternative to matching uses OLS to estimate parametric linear models. In this approach, researchers specify the functional form of the conditional mean function by replacing $g_1(X_i)$ in (4) with $X_i'\beta_1$ and replacing $g_0(X_i)$ in (5) with $X_i'\beta_0$, or by pooling the data and estimating a common β . With these substitutions, the researcher avoids invoking the (CSA), but not the (CMIA).⁵

Estimates identified via conditional independence often get criticized because such assumptions appear implausible in many substantive contexts given the available conditioning variables. To avoid a CIA, applied researchers often turn to IV estimation. While traditional IV assumptions presume a common treatment effect, Imbens and Angrist (1994) demonstrate that under different assumptions IV estimation allows for heterogeneous treatment effects. Researchers now routinely invoke their LATE framework when applying IV methods with

⁴See the discussions in Black and Smith (2004) and Crump, Hotz, Imbens, and Mitnik (2009).

⁵Other estimators that build on conditional independence include Inverse Propensity Weighting (IPW), so-called “double robust” estimators that combine IPW with a linear model, and various machine learning estimators. See e.g. Heckman et al. (1999), Imbens (2004), Smith and Todd (2005), Huber et al. (2013), and Busso et al. (2014) for variations on matching as well as introductions to these other estimators.

binary instruments. It is difficult to overemphasize the importance of this advance. Models that omit selection into treatment based upon (possibly very partial) knowledge of heterogeneous treatment effects seem incapable of capturing the complexity of human behavior. Incorporating such treatment effect heterogeneity allows researchers to consider and estimate far more plausible and interesting models, including the justifiably famous Roy (1951) model. Indeed, Heckman Urzua, and Vytlačil (2006) term such heterogeneous impacts “essential heterogeneity.”

3 The IV and Control Function Approaches

Consider a binary instrument $Z_i \in \{0, 1\}$. We may state the assumptions of the LATE estimator as the Existence of Instruments (EI) and Monotonicity (M). Formally,

$$(i) (Y_{0i}, Y_{1i}, D_i(0), D_i(1)) \perp\!\!\!\perp Z_i | X_i \tag{EI}$$

$$(ii) Pr(D_i = 1 | X_i, Z_i) \text{ is a non-trivial function of } Z_i$$

$$\forall z^0, z^1, \text{ either } D_i(z^0) \geq D_i(z^1) \forall i \text{ or } D_i(z^0) \leq D_i(z^1) \forall i \tag{M}$$

The (EI) assumption requires that, conditional on X_i , the instrument only affects the outcome through its effect on treatment $D_i(Z_i)$. The (M) assumption requires that all agents respond to the instrument in the same direction, not that the function $Pr(D_i = 1 | X_i, Z_i)$ be monotone in Z_i ; this led Heckman et al. (2006) to rename the condition uniformity, although the somewhat confusing term monotonicity was too well-established to be displaced. The (M) assumption is restrictive; it need not hold in all substantive contexts. Should the (M) assumption fail while the (EI) assumption holds, IV estimation provides a mixture of treatment effects associated with agents who both enter and leave the treatment as the instrument varies. To keep the notation simple, we continue to assume $Z_i \in \{0, 1\}$; our arguments, however, generalize to continuous instruments.

Imbens and Angrist note that the latent index models pioneered by Heckman and various co-authors imply the (EI) and (M) conditions. In an important paper, Vytlacil (2002) shows the equivalence of the two approaches. In our notation, one may define the expectations of the errors in equations (4) and (5) as zero, or $E(\epsilon_{1i}|X_i) = E(\epsilon_{0i}|X_i) = 0$. This is, of course, a convenient normalization with any nonzero mean absorbed into the conditional mean functions. When we observe only one of the two potential outcomes, we do not know that the conditional expectations $E(\epsilon_{1i}|D_i = 1)$ and $E(\epsilon_{0i}|D_i = 0)$ equal zero because of the missing data problem. To formalize this argument a bit, we follow Vytlacil (2002) and let

$$D_i = 1(h(Z_i, X_i) - U_i \geq 0) \text{ and } h(Z_i, X_i) \text{ be a nontrivial function of } Z_i \quad (\text{V1})$$

$$Z_i \perp\!\!\!\perp (Y_{1i}, Y_{0i}, U_i) | X_i \quad (\text{V2})$$

where $1(\cdot)$ is the logical indicator function, U_i is a unobserved random variable that affects treatment choice, and $h(Z_i, X_i)$ is the index function. Selection on unobservables means that (Y_{1i}, Y_{0i}) are not independent of U_i conditional on X_i .

With assumptions (EI) and (M) (or the equivalent assumptions (V1) and (V2) for latent index models), we may write

$$E(\epsilon_{1i}|X_i, Z_i, D_i = 1) = c_1(X_i, Pr(X_i, Z_i)) \quad (6)$$

$$E(\epsilon_{0i}|X_i, Z_i, D_i = 0) = c_0(X_i, Pr(X_i, Z_i)) \quad (7)$$

where $Pr(Z_i, X_i) = Pr(D_i = 1|X_i, Z_i)$ is the conditional probability of treatment or propensity score⁶ and where we denote the control functions that embody the conditional means of ϵ_{1i} and ϵ_{0i} by $c_1(\cdot)$ and $c_0(\cdot)$. We can obtain $E(\epsilon_{1i}|D_i = 1)$ and $E(\epsilon_{0i}|D_i = 0)$ by integrating out over X_i and Z_i in (6) and (7). The key to (6) and (7) is that Z_i enters only through the probability of treatment. In the absence of selection the control function equals zero

⁶Unlike the propensity score used in matching estimators under the (CIA), this propensity score also includes at least one instrument; see Heckman and Navarro-Lazano (2004) for further discussion.

everywhere; in the presence of selection, including the control function solves the selection problem.

The control function approach allows an easier interpretation of the independence assumption embedded in (EI); it simply requires that Z_i be independent of (U_i, Y_{1i}, Y_{0i}) conditional on X_i . Given the equivalence of the LATE and control function assumptions, we refer to the (EI) and (M) assumptions, or (V1) and (V2), as the Vytlacil-Imbens-Angrist (VIA) assumptions. We turn now to our tests, which replace the control function with a simple function of the instrument.

4 Testing Conditional Independence

4.1 Instrumental Variables

As noted above, the various versions of the (CIA) and (CMIA) allow researchers to ignore the possibility of selection on unobserved variables, although they typically invoke them without looking for evidence of such selection. In contrast, the VIA assumptions allow researchers to consistently estimate LATEs for compliers. In the case of a single instrument we augment the linear parametric versions of equations (4) and (5) to obtain

$$E(Y_{1i}|X_i, Z_i) = X_i'\beta_1 + \alpha_1 Z_i \tag{8}$$

$$E(Y_{0i}|X_i, Z_i) = X_i'\beta_0 + \alpha_0 Z_i \tag{9}$$

With non-binary instruments researchers may wish to add higher order terms (i.e., replace Z_i with $f(Z_i)$) though this raises subtle but important issues of model selection that lie outside the scope of this paper. With multiple instruments, researchers could replace Z_i with the estimated propensity score, $\hat{Pr}(Z_i, X_i)$ and adjust the standard errors for generated regressors as in Murphy and Topel (2002).⁷

⁷Joo and LaLonde (2014) present a control function version of our test along these lines.

The model behind equations (8) and (9) represents an important departure from the canonical model used in IV applications, given by

$$Y_i = X_i'\beta + \phi D_i + \epsilon_i \quad (10)$$

In contrast, the model underlying equations (8) and (9) consists of

$$Y_{1i} = X_i'\beta_1 + \epsilon_{1i} \quad (11)$$

$$Y_{0i} = X_i'\beta_0 + \epsilon_{0i}. \quad (12)$$

Equations (11) and (12) allow estimation of heterogeneous treatment effects, $\Delta(X_i) = (\beta_1 - \beta_0)X_i + (\epsilon_{1i} - \epsilon_{0i})$, that differ with the realization of the covariates while still maintaining the (CMIA). There is generally no theoretical reason to prefer equation (10) to equations (11) and (12), but the demands on instrument strength usually dissuade researchers from using the model described by equations (11) and (12) when they resort to IV estimation. Avoiding any interactions involving the treatment indicator yields a partially linear model and represents a strong, substantive restriction on the model (and on our test).

In the case of matching estimators, we augment equations (4) and (5) and specify the conditional mean functions as

$$E(Y_{1i}|X_i, Z_i, D_i = 1) = g_1(X_i) + \alpha_1 Z_i \quad (13)$$

$$E(Y_{0i}|X_i, Z_i, D_i = 0) = g_0(X_i) + \alpha_0 Z_i \quad (14)$$

To clarify the relationship among the various forms of the (CMIA) and our test, it is useful to outline the samples used and hypotheses involved when estimating these auxiliary regressions. Formally, we estimate (9) or (14) using the sample of untreated observations to test the null that the (CMIA⁰) holds, corresponding to $\alpha_0 = 0$, against the alternative that it does not hold, or $\alpha_0 \neq 0$. Similarly, we estimate equation (8) or (13) using the sample of

treated observations to test the null that (CMIA¹) holds, corresponding to $\alpha_1 = 0$, against the alternative that it does not hold, or $\alpha_1 \neq 0$.

To develop some intuition for the tests, assume that $D_i(1) \geq D_i(0)$ and divide agents under VIA into the three types defined in the introduction: always-takers, never-takers, and compliers. The test given in either equation (8) or equation (13) simply compares $E(Y_{1i}|X, A)$ to $E(Y_{1i}|X, C)$. As Black et al. (2015) note, this is easily done because $E(Y_{1i}|X_i, Z_i = 0) = E(Y_{1i}|X_i, A)$ and $E(Y_{1i}|X_i, Z_i = 1) = E(Y_{1i}|X_i, A \cup C)$. Thus, at $X_i = x$ we have that

$$\alpha_1(x) \equiv \frac{Pr(C|x)}{Pr(C|x) + Pr(A|x)} [E(Y_{1i}|x, C) - E(Y_{1i}|x, A)] \quad (15)$$

The regression coefficient in either equation (8) or equation (13) then simply integrates over the realizations of X_i , or $\alpha_1 = \int \alpha_1(x) w_1(x) dF(x)$, where $F(x)$ is the distribution of the covariates X_i and $w_1(x)$ is a weighting function that depends on the joint distribution of X and Z . Put differently, the tests look for evidence of a non-constant control function in equation (6), which constitutes evidence that unobserved variables affect the outcomes. A parallel argument applies to α_0 . In both cases, if the nature of the selection varies with X_i , trouble may ensue; see Słoczyński (2020) for more on these weighting functions.

The finding that either $\alpha_0 \neq 0$ or $\alpha_1 \neq 0$ constitutes evidence of either selection or violation of the exclusion restrictions (i.e., the failure of EI) or both. Assuming the validity of the exclusion restriction, rejection of one or both of the null hypotheses provides simple and compelling evidence for violation of the (CMIA), and thus of the (CIA) as well.

The tests also allow researchers to assess whether any selection arises on Y_{0i} , which represents a violation of (CMIA⁰), or on Y_{1i} , which represents a violation of (CMIA¹), or both. In addition, as with the tests for selection in the bivariate and trivariate normal models, our tests allow researchers to determine the signs of the relevant selection effects and their magnitudes, which we discuss in section 4.4. This allows researchers to provide a much more nuanced discussion of the nature of agents' underlying choice behavior. Given the equivalence that Vytlacil (2000, 2002) demonstrates, it is perhaps not surprising that we

may learn more about the selection problem using IV methods than we learn from current practices.

4.2 Fuzzy regression discontinuity

In regression discontinuity designs, treatment depends on a running variable S_i and has the feature that the probability of treatment jumps (i.e., has a discontinuity) at some particular value of S_i . We assume that the jump occurs at $S_i = 0$. To use both of our tests we require fuzziness on both sides of the discontinuity; formally, we require

$$1 > \lim_{S_i \downarrow 0} Pr(D_i = 1 | X_i, S_i) > \lim_{S_i \uparrow 0} Pr(D_i = 1 | X_i, S_i) > 0$$

or the same condition but with the two limits reversed. With both treated and untreated units on only one side of the discontinuity, a researcher can apply our test for one of (Y_{1i}, Y_{0i}) but not both. As emphasized by Imbens and Lemieux (2008) and Lee and Lemieux (2010), when faced with a fuzzy RD, researchers who use the discontinuity at $S_i = 0$ as an instrument for treatment estimate a LATE at $S_i = 0$.

Because of the discrete change in the treatment probability at $S_i = 0$, under selection on unobserved variables we would expect a jump in the control function at the same point. More formally, selection on unobserved variables implies a jump in the value of $E(Y_{0i} | X_i, S_i, D_i = 0)$ as S_i crosses zero, while the (CMIA⁰) assumption implies a smooth $E(Y_{0i} | X_i, S_i, D_i = 0)$ function around zero. This suggests a simple test based on a model of the form

$$E(Y_{0i} | X_i, S_i, D_i = 0) = g_0(X_i, S_i) + \alpha_0 1(S_i \geq 0) \tag{16}$$

with the null hypothesis being that $\alpha_0 = 0$, or the corresponding version for testing (CMIA¹)

$$E(Y_{1i} | X_i, S_i, D_i = 1) = g_1(X_i, S_i) + \alpha_1 1(S_i \geq 0) \tag{17}$$

with the null hypothesis being that $\alpha_1 = 0$. Estimation of the sample analogue of (16) makes use only of untreated observations, which constitute a mixture of compliers and never-takers. Similarly, estimation of the sample analogue of (17) uses only treated observations, which constitute a mixture of compliers and always-takers.

As noted in the introduction, Bertanha and Imbens (2019) consider closely related tests for fuzzy regression discontinuity designs. Indeed, they state, “As a matter of routine, we recommend that researchers present graphs with estimates of these two conditional expectations in addition to graphs with estimates of the expected outcome conditional on the forcing variable alone.” We concur.

4.3 Experiments with Imperfect Compliance

As Heckman (1996) emphasizes, random assignment creates an instrument for treatment. As documented in Heckman, Hohmann, Smith, and Khoo (2000), many social experiments have imperfect compliance, with both treatment group dropout and control group substitution into similar treatments provided elsewhere. For instance, their Table II reports that, among those recommended for classroom training prior to random assignment, somewhere between 49 and 59 percent of the treatment group in the Job Training Partnership Act (JTPA) experiment received services, depending on the demographic group, while between 27 and 40 percent of the control group received services.

In the presence of dropout and substitution, applied researchers will often rely on the Bloom (1984) estimator. To use the Bloom estimator, the researcher need only use random assignment to the treatment group as an instrument for the receipt of treatment. As random assignment provides a binary instrument, the Wald estimator recovers the LATE for those who comply with the experimental protocol. One could easily implement our tests to check for selection on Y_{1i} or Y_{0i} in experiments such as these. As noted above, Huber (2013) provides a Wald test of the exogeneity of non-compliance in experiments closely related to ours.

4.4 Recovering Estimates of the Magnitude of the Selection Effect

To recover estimates of the magnitude of the selection effect, continue to assume that $Z_i = 1$ encourages treatment, and ignore covariates for notational simplicity. We have

$$E(Y_{0i}|Z_i = 0, D_i = 0) = \frac{Pr(C)}{Pr(C) + Pr(N)}E(Y_{0i}|C) + \frac{Pr(N)}{Pr(C) + Pr(N)}E(Y_{0i}|N) \quad (18)$$

while

$$E(Y_{0i}|Z_i = 1, D_i = 0) = E(Y_{0i}|N) \quad (19)$$

so that

$$\alpha_0 = E(Y_{0i}|Z_i = 1, D_i = 0) - E(Y_{0i}|Z_i = 0, D_i = 0) = \frac{Pr(C)}{Pr(C) + Pr(N)}(E(Y_{0i}|N) - E(Y_{0i}|C)). \quad (20)$$

Thus, a measure of the selection effect for Y_{0i} , which we denote B_{0i} , is simply

$$B_{0i} = E(Y_{0i}|N) - E(Y_{0i}|C) = \frac{Pr(C) + Pr(N)}{Pr(C)}\alpha_0. \quad (21)$$

An eerily similar derivation yields

$$B_{1i} = E(Y_{1i}|C) - E(Y_{1i}|A) = \frac{Pr(C) + Pr(A)}{Pr(C)}\alpha_1. \quad (22)$$

To implement these measures empirically, we may use the OLS estimates of (α_0, α_1) . We know that $Pr(A) = Pr(D_i = 1|Z_i = 0)$, $Pr(N) = Pr(D_i = 0|Z_i = 1)$, and $Pr(C) = 1 - Pr(N) - Pr(A)$ under assumption (M), so we have sample analogues of all the terms on the right-hand sides of equations (21) and (22).

4.5 Full Independence

Up to this point we have concentrated on testing (CMIA⁰) and (CMIA¹). As noted above, the most common causal estimands require only one or both of these conditional mean independence assumptions. In this subsection we develop additional tests of other implications of full conditional independence, i.e., tests of other implications of (CIA⁰) and (CIA¹).⁸ We have two motivations: first, some researchers lust after causal estimands that require full conditional independence; second, even researchers interested in the ATET, ATEN or ATE may value the additional information provided by these tests, especially in cases where our tests of conditional mean independence only marginally fail to reject the null. We continue to emphasize ease of implementation in the tests we propose.

One category of tests compares higher moments rather than means. For example, in the IV case with a binary instrument considered in Section 3.1, we can test the nulls

$$F(Y_{0i}|N, X_i) = F(Y_{0i}|C, X_i) \quad (23)$$

$$F(Y_{1i}|C, X_i) = F(Y_{1i}|A, X_i) \quad (24)$$

where $F(\cdot)$ denotes a cumulative distribution function, by estimating versions of (8) and (9) or of (13) and (14) with higher moments of Y_{1i} and Y_{0i} as the dependent variables.

A second category of test builds on non-linear transformations of continuous outcomes. For example, a very simple version of this strategy replaces the levels of Y_{0i} and Y_{1i} with their natural logarithms in the tests already described. The variant we pursue builds on distribution regressions of the form

$$1(Y_{1i} \leq q|X_i, Z_i, D_i = 1) = g_1^q(X_i) + \alpha_1^q Z_i \quad (25)$$

$$1(Y_{0i} \leq q|X_i, Z_i, D_i = 0) = g_0^q(X_i) + \alpha_0^q Z_i \quad (26)$$

⁸Of course, for binary outcomes, conditional mean independence implies full independence.

where q denotes some quantile of the unconditional distribution of Y_i and $g_j^q(X_i)$ $j \in \{0, 1\}$ are index functions. Under (CIA¹) the null implies $\alpha_1^q = 0$ while under (CIA⁰) it implies $\alpha_0^q = 0$, in both cases for all $q \in (0, 1)$.

We implement both the higher moment tests and the distribution regression tests in our applications in Sections 5.1 and 5.2 where we have continuous outcomes. For the higher moment tests, we consider second, third and fourth moments about the mean.⁹ For the distributional regressions, we look at either deciles or ventiles, taking care with the mass point at zero for the earnings outcome in Section 5.1. Using deciles or ventiles rather than, say, percentiles, lessens concerns about non-independence of adjacent tests at some cost in statistical power. Both approaches to testing full independence raise multiple testing issues (as does our consideration of multiple outcomes in our tests of conditional mean independence). Given the illustrative nature of our empirical applications, we view these issues as beyond our scope; Schochet (2008) and List, Shaikh, and Xu (2019) provide valuable overviews of the multiple testing literature aimed at applied readers.

5 Empirical Applications

5.1 Angrist and Evans (1998) data

Our first application draws on Angrist and Evans (1998). This paper uses data from the 1980 and 1990 US Censuses to measure the causal impact of children on maternal labor supply. Because fertility is likely endogenous with respect to women’s labor supply decisions, Angrist and Evans devise an ingenious instrumental variables strategy. Limiting their sample to women who have at least two children, Angrist and Evans noticed that women whose first two children are the same sex are more likely to have additional children than women whose first two children are of opposite sexes. For instance, in the 1980 Census, married women

⁹For numerical reasons, we divide the outcome by its standard deviation when calculating the third and fourth moments.

whose first two children are of the same sex are about six percentage points more likely to have additional children than women whose first two children are of opposite sexes. For our analysis, we focus on the labor supply decisions of women in the 1980 Census.

In many ways, this design is ideal. Because of the random nature of child sex determination, the sample is split approximately equally between families whose first two children are of the same sex and those whose children are of opposite sexes. In these data, 51.1% of the children born are male, and in 50.6% of families the first two children are of the same sex. Formally, the system that Angrist and Evans estimate is:

$$Y_i = X_i'\beta + \phi M_i + \epsilon_i \quad (27)$$

$$M_i = X_i'b + \gamma Z_i + \eta_i \quad (28)$$

where M_i is an indicator for having more than two children. The covariates include the age of the mother, the age of the mother at first birth, indicators for whether the mother is black or whether the mother is non-black and non-white (white is the omitted category), an indicator for whether the mother is Hispanic, an indicator for whether the first child was a boy, and an indicator for whether the second child was a boy. The instrument, Z_i , is an indicator for whether the first two children were either two boys or two girls. For dependent variables, we use a subset of those explored by Angrist and Evans: whether the mother worked in the previous year, the number of weeks worked in that year, typical hours worked in that year, and her income from working. We set all of these variables to zero for women who did not work in the previous year. The sample is limited to women 21 to 35 years of age; see Angrist and Evans (1998) for more details.

In Table 1, we replicate the Angrist and Evans results in the 1980 Census; see their Table

7, columns (1) and (2). We also use a semiparametric approach and estimate

$$Y_i = g(X_i) + \phi M_i + v_i \tag{29}$$

$$M_i = h(X_i) + \gamma Z_i + u_i \tag{30}$$

where $g(\cdot)$ and $h(\cdot)$ are unknown functions. As we have only discrete conditioning variables, we estimate $g(X_i)$ and $h(X_i)$ using fully saturated regressions. Our parametric results, both the OLS and two-stage least squares estimates, exactly match Angrist and Evans' findings. Moreover, the semi-parametric estimates turn out virtually identical to the parametric estimates of Angrist and Evans, which is not too surprising given that nature makes the sex of women's offspring independent of all of our observed characteristics.

Of course, to interpret the IV estimand as a LATE we need to assume the VIA conditions. Angrist and Evans document that the instrument does indeed raise fertility. In addition, we need to assume that the instrument provides an exclusion restriction in the sense that having the first two children of the same sex does not directly affect women's labor supply decisions, and we need to assume the monotonicity (or uniformity) condition so that having two children of the same sex reduces no one's fertility. With these (strong) assumptions, we may now implement our parametric tests of the CIAs using

$$Y_{0i} = X_i' \beta_0 + \alpha_0 Z_i + \epsilon_{0i} \tag{31}$$

$$Y_{1i} = X_i' \beta_1 + \alpha_1 Z_i + \epsilon_{1i} \tag{32}$$

and our semi-parametric tests by repeating the estimation replacing $X_i' \beta_0$ with $g_0(X_i)$ in (31) and $X_i' \beta_1$ with $g_1(X_i)$ in (32). For our semiparametric analysis we drop the indicator for having a boy as the second child in order to avoid making the Z_i variable perfectly collinear. We estimate equation (31) on the sample of 236,092 women who have exactly two children, and equation (32) using the sample of 158,743 who have three or more children.

Table 2 presents the results of our tests of conditional mean independence. The data

strongly reject the null that $\alpha_0 = 0$, thereby providing strong evidence against the (CMIA⁰). For each of the four outcomes, we reject the null hypothesis at a five-percent confidence level. In each case, we estimate a positive coefficient α_0 on Z_i , where $Z_i = 1$ among the non-treated corresponds to the never-takers. Thus, we find that the never-takers have higher earnings, hours worked, and weeks worked, and are more likely to work at all conditional on our covariates relative to the compliers who do not have a third child.

In contrast, we find little evidence against the (CMIA¹). Unlike the estimates of α_0 , our estimates of α_1 are statistically insignificant and economically very small. Thus, we find no evidence of selection when estimating the missing counterfactual Y_{1i} . Frankly, we find this result stunning. There appear to be no meaningful differences in the labor supply behavior of the always-takers and the compliers when we condition on the limited observed characteristics available in the US Census data. Before undertaking this analysis, we fully expected to find a two-sided selection problem. The data disagreed.

Table 3 adds the results from our tests of independence of distributions for the three continuous outcomes. The upper panel shows the results from estimating equations (25) and (26) for the mass point at zero and for each decile above the mass point. Both sets of distributional tests strongly corroborate the lessons from our tests of conditional mean independence: substantively large and statistically significant rejections of the null for Y_{0i} and substantively small and rather imprecise failures to reject the null for Y_{1i} . The lower panel shows the results from tests based on (23) and (24) with the 2nd, 3rd, and 4th (central) moments replacing the mean. For completeness, we repeat the results for the 1st moment (i.e., the mean) from Table 2. The tests based on higher moments shed much less light, as only two of 18 tests reject their null, one for Y_{0i} with the hours worked outcome and one for Y_{1i} with the income outcome.

To describe the magnitude of the selection effects we use the semi-parametric estimates from Table 2. Compared to the compliers, we find that the never-takers are five percentage points more likely to have worked last year, worked about three weeks more, worked about

two hours more per week, and earned \$1,965 more per year. Comparing the compliers to the always-takers, we find that compliers were one percentage point more likely to work last year, they worked about 0.4 extra weeks per year, they worked a tenth of an hour more per week, and earned \$38 dollars less per year than the always-takers. Obviously, the compliers represent a poor comparison group for the never-takers while differing only very little from the always-takers.

5.2 Angrist and Lavy (1999) data

Angrist and Lavy (1999) address the perennial economics of education topic of the causal effect of class size on student outcomes using data from Israel. They exploit the use of Maimonides' rule that limits class sizes to 40 in Israeli public schools to construct regression discontinuity estimates of class size effects. Under Maimonides' rule enrollment in a particular combination of school, grade and academic year constitutes the running variable. When the running variable crosses from 40 to 41, average class size drops from 40 to 20.5. Similar things happen when the running variable moves from 80 to 81 or from 120 to 121. In practice, for a variety of institutional reasons (such as extra funds provided to schools serving relatively disadvantaged students and the fact that enrollment is measured in the fall of the academic year while class size is measured in the spring) the data present a fuzzy regression discontinuity rather than the strict one that would arise if Maimonide's rule alone determined class size.

For their primary estimates Angrist and Lavy (1999) operationalize class size as a continuous treatment, and do two-stage least squares with the first stage instrument consisting of the predicted enrollment based solely on Maimonide's rule. In keeping with our focus on the discontinuity framing of their analysis, we replicate the estimates in columns (5) and (11) of their Table V. These estimates refer to fourth graders in 1991 and use only observations on school-grade combinations with enrollments within five students (above or below) of a multiple of 40.¹⁰ The outcomes consist of reading and math achievement tests scores in

¹⁰Their paper provides a startling reminder that the now-ubiquitous applied econometric toolkit for re-

“natural” units rather than standard deviation units. Our replicate estimates appear in the top row of Table 4. Except for the third digit in one of the standard errors, we match their numbers exactly. In a common effect world, we can interpret the -0.098 estimate for reading comprehension as indicating that a one student increase in class size decreases the number of correct answers by about 0.1 in a context where, following their Table I, the mean math score equals 28.1 with a standard deviation of 6.5.¹¹ In a heterogeneous treatment effects world, the two-stage least squares estimator that Angrist and Lavy (1999) employ implicitly yields a weighted average of the LATEs at the multiple discontinuities.

To make their setup comparable to our other examples, we discretize both their continuous treatment and their continuous instrument. In particular, our binary treatment consists of having a class size of less than 32, while our binary instrument consists of having a predicted class size of less than 32. The bottom row of Table 4 shows the estimates we obtain using the same sample as in the top row, but with the binary treatment and binary instrument.¹² To interpret these estimates, first note that the mean change in class size for the compliers equals 3.18 students. Thus, under a LATE interpretation, a change in class size of that magnitude decreases expected math test scores by -1.33 and increases expected reading comprehension test scores by 1.37 within that group, although neither of these estimates approach statistical significance.

Table 5 presents the results from our tests of (CMIA⁰) and (CMIA¹) using the Angrist and Lavy (1999) data. We find little evidence of selection in either case but note that, consistent with the standard errors on our estimates in Table 4, we do not have nearly as much power as we would like. The tests of (CIA⁰) and (CIA¹) in Table 6 tell a modestly different story. In particular, our distributional regressions provide meaningful but not overwhelming evidence

gression discontinuity designs, including sophisticated bandwidth selection tools, did not yet exist in 1999.

¹¹The value of 6.5 refers to the standard deviation of class-level mean math scores. The individual-level variation in math scores is surely much larger. Thus, the estimated causal effect, expressed in terms of the common metric of standard deviations of individual scores, is rather small. See their Section V.B for more about magnitudes.

¹²Angrist and Lavy (1999) undertakes a similar exercise, see their Table VI. Ours differs in that we do not obtain separate estimates for each of the three discontinuities, preferring instead a more precisely estimated weighted average, at some cost in simplicity of interpretation.

against (CIA⁰) for the reading comprehension test and against (CIA¹) for the math test in the form of multiple (two out of four ventiles) rejections of the null in each case. As we would expect a large common component underlying test performance in the two subjects, these findings make us hesitant to extrapolate from the LATE interpretation to an ATET or ATEN interpretation in either subject; put differently, background knowledge leads us to think that the rejections should “spill over” across subjects.

5.3 Eberwein, Ham, and LaLonde (1997) data

When facing control group substitution and treatment group dropout (i.e. two-sided non-compliance) in an experiment, researchers will often estimate two treatment parameters: the intent-to-treat parameter, estimated as the mean difference in a dependent variable between the treatment and control groups, and the impact of treatment for those who comply with the treatment protocol, estimated using Bloom’s (1984) estimator.

We examine the impact of training for a sample of adult women who took part in the Job Training Partnership Act (JTPA) experiment; see Bloom et al. (1997) for a discussion of the experiment and analysis of the results. Our sample, the same one used by Eberwein, Ham, and LaLonde (1997), consists of women recommended for classroom training (the “CT-OS treatment stream” in the jargon of the experiment) prior to random assignment.¹³ We measure training as the onset of self-reported classroom training within nine months of randomization. We focus on classroom training and ignore other (usually much less intensive) services, such as job search assistance, received by some members of both the treatment and control groups for simplicity. We rely on the self-reported training data for both groups, rather than the self-reports for the controls and the JTPA administrative data for the treatment group, for comparability.¹⁴ The administrative data from the experiment provide our treatment indicator. For our outcome variable, we use an indicator for self-

¹³The reader should not compare our results directly to those in Donald et al. (2014) as their data on women also includes those recommended for services other than classroom training.

¹⁴Smith and Whalley (2020) offer a depressing exploration of the concordance of administrative and self-reported measures of service receipt in the JTPA study data.

reported employment in the eighteenth month after random assignment.¹⁵

There was much non-compliance in this sample. Only about 65 percent of the treatment group reports receiving classroom training in the first nine months after random assignment. There was much control group substitution as well: about 34 percent of the control group reports receiving classroom training in the first nine months after random assignment.

In Table 7, we provide two sets of estimates of the intent-to-treat parameter. In column (1) we provide the simple mean difference estimates given by

$$Y_i = \beta + \phi R_i + \epsilon_i, \tag{33}$$

where Y_i is the outcome variable, R_i is an indicator for whether the participant was assigned to the treatment group during random assignment, ϵ_i is the error term, and (β, ϕ) are parameters to be estimated. The estimated intent-to-treat parameter, $\hat{\phi}$, equals 0.041 and statistically differs from zero at the five-percent level. This relatively modest effect, however, hides a larger impact of treatment for people who complied with the treatment protocol, which equals 0.136 and again achieves statistical significance at the five-percent level. The differential arises, of course, because random assignment only increases the rate of treatment by about 0.305, the coefficient on the indicator for random assignment to the treatment group from the first stage.

Nothing in this analysis, however, informs the researcher regarding the presence or absence of selection into treatment based on unobserved variables. Toward that end, we next estimate the following equations

$$Y_{0i} = \beta_0 + \alpha_0 R_i + e_{0i} \tag{34}$$

$$Y_{1i} = \beta_1 + \alpha_1 R_i + e_{1i} \tag{35}$$

¹⁵As noted in Section 4.5, with a binary outcome the conditional mean fully characterizes the distribution, leaving no scope for additional tests of conditional independence of distributions.

where (Y_{0i}, Y_{1i}) denote the outcomes of those not receiving training and receiving training. We estimate Equation (34) using the 1,233 adult women who do not receive training and estimate equation (35) using the 1,501 who do receive training. We find little evidence that the compliers have different Y_{0i} than those who never take training. The coefficient on the indicator for assignment to the treatment group in equation (34) is small, 0.006, and statistically insignificant at the five-percent level. In contrast, the coefficient on the indicator for assignment to the treatment group in equation (35) is large, 0.068, and statistically significant. These estimates imply that while the always-takers have a mean employment rate of 0.50, the compliers when treated have a mean employment rate of 0.65. Thus, the always-takers are adversely selected with respect to the likelihood of employment in the treated state.

A finding of substantively large selection on unobserved variables in the absence of covariates will hardly surprise most readers. Thus, we augment our equations with a vector of variables measuring the educational and demographic characteristics of those randomly assigned, as well as their pre-random assignment labor market activity and transfer payment receipt; see the notes to Table 7 for a complete list. Their inclusion (as expected) has only modest effects on the intent-to-treat and LATE estimates, although some may be dismayed that the estimates no longer clear the five-percent hurdle. More surprisingly, the results of our tests for selection on unobserved variables also change very little when we add covariates. In the Y_{1i} regression, the coefficient on the assignment indicator for those receiving treatment falls from 0.068 to 0.062 and remains significant at the five percent level. Despite the inclusion of detailed information on labor supply in the 12 months prior to random assignment and other controls, the coefficient on the treatment group indicator falls by only about nine percent. The observed variables examined here account for little of the selection. In contrast, our test consistently fails to detect unobserved differences between the compliers and the never-takers.

6 Discussion

6.1 Standard IV issues and our test

First, weak instruments imply a relatively small number of compliers which in turn implies relatively low power for our test. Put differently, with a weak instrument, comparing the conditional means of, say, always-takers and compliers will provide only noisy evidence regarding the null of no selection in the absence of a large selection effect, a large sample, or both.

Second, the case for the monotonicity (M) assumption typically rests on some combination of institutional knowledge and economic theory specific to a particular empirical context. Our test provides no help in detecting failures of the (M) assumption. When it does fail, the untreated units include what Angrist et al. (1996) call defiers, agents who change treatment status when the value of the instrument changes but in an unexpected way, in addition to compliers and never-takers. Similarly, the treated units now comprise always-takers, compliers and defiers. The presence of defiers undoes the LATE interpretation of the IV estimand. Klein (2010) provides a gateway into the literature that relaxes the monotonicity assumption.

Finally, because our test implicitly represents a joint test of the (EI) assumption and the null of no selection bias, failure of the (EI) assumption can lead to incorrect inferences regarding the presence or absence of selection. When failure of the (EI) assumption leads the test to reject the joint null, researchers may proceed to place heavy weight on the IV estimates, when in fact they converge to an unknown mixture of the population LATE and the bias associated with the invalid instrument.

6.2 Bias, variance, and mean-squared error

Statistics reminds us of the trade-off between bias and variance. Consider the canonical model in (10) and imagine for the moment either a common effect world or a world where the (CMIA) holds. In that world, as Cameron and Trivedi (2005, p. 107) note, with one

instrument we may compare the variance of the parameter of interest in (10) estimated using instrumental variables to the variance of the same parameter estimated using OLS using the equation:

$$SE(\hat{\delta}^{IV}) = \frac{SE(\hat{\delta}^{OLS})}{\hat{\rho}(\tilde{D}, \tilde{Z})} \quad (36)$$

where $\hat{\rho}(\tilde{D}, \tilde{Z})$ denotes the partial correlation coefficient after removing the variation associated with the other covariates X , what Black and Smith (2006) term the Yulized residuals in honor of Yule’s (1907) brilliant paper.

Many empirical contexts yield quite modest partial correlations, 0.10 or even 0.05 depending on the strength of the instrument, implying dramatic decreases in the precision of the estimates when using IV methods. For example, Black et al. (2015) estimate $\hat{\rho}(\tilde{D}, \tilde{Z}) = 0.059$ and Angrist and Evans (1998) estimate a value of 0.064. Others contexts yield larger values: for instance, the ratio of the standard errors for the OLS and IV estimates in Table 2 of Bertanha and Imbens (2019) implies a large partial correlation. In many situations, serious researchers will trade off increases in bias for reductions in variance. Our method provides researchers with a means of assessing this bias and so allows them to make a quantitatively informed decision regarding whether or not the IV cure for the OLS bias disease is worse than the disease itself.

To see how nefarious the IV cure may be, recall the canonical common effects model given in (10). Let $\hat{\delta}^{OLS} = \delta + e^{OLS}$ and $\hat{\delta}^{IV} = \delta + e^{IV}$, where e^{OLS} and e^{IV} denote the errors in the two estimates. Researchers concerned only with the distance of their estimate from the true parameter value (i.e., with the absolute values of e^{OLS} and e^{IV}) should prefer the IV estimator only when

$$(|\hat{\rho}(\tilde{D}, \tilde{\epsilon})| |\hat{\rho}(\tilde{D}, \tilde{Z})| - |\hat{\rho}(\tilde{Z}, \tilde{\epsilon})|) > 0, \quad (37)$$

where $\hat{\rho}(\cdot, \cdot)$ again denotes the sample correlation between its arguments.¹⁶ Any correlation between the treatment indicator, D , and the error term ϵ due to omitted variables and

¹⁶This is just the finite sample analogue of equation (7) in Bound, Baker and Jaeger (1995).

sampling variation gets down-weighted by the relative strength of the instrument. As the estimated correlation $\hat{\rho}(\tilde{D}, \tilde{Z})$ equals 0.064 in our Angrist and Evans’ application, the correlation $\hat{\rho}(\tilde{D}, \tilde{\epsilon})$ must be an order of magnitude larger than $\hat{\rho}(\tilde{Z}, \tilde{\epsilon})$ before we prefer the IV estimates.

We might turn to formal decision theory to determine which estimator to use. An oft-used metric in decision theory is mean squared error. The expression for

$$\hat{\delta}_{IV} = \frac{\text{cov}(\tilde{Y}, \tilde{Z})}{\text{cov}(\tilde{D}, \tilde{Z})} \tag{38}$$

illustrates a peculiar feature about IV estimation: Because of the $\text{cov}(\tilde{D}, \tilde{Z})$ term in the denominator of the estimator, the expectation of $\hat{\delta}_{IV}$ does not exist in the case with only one instrument.¹⁷ In just-identified models, both the bias and the variance of the IV estimator are not defined.¹⁸ In contrast, while the OLS estimator is biased and inconsistent, the bias and variance of the estimator are finite. In the terminology of decision theory, the IV estimator is inadmissible: we should always pick OLS estimation. In our view, this conclusion may say more about the use of mean squared error as a criterion than the choice of estimators.

6.3 Different estimators, different estimands

In a heterogeneous treatment effects (i.e., realistic) context, the IV and OLS estimands correspond to different causal parameters. Among others, Deaton (2010), Heckman and Urzua (2010), and Imbens (2010) debate the value of these estimands at a relatively high level of generality. In our view, the relevance of the LATE parameter depends on the substantive context and on the composition of the complier population. Our tests and the related measures of bias developed above provide a framework within which to think about

¹⁷This result dates back to the 1960s; see Nelson and Startz (1990) for an excellent discussion and some additional results.

¹⁸As noted by Nelson and Startz (1990), the failure of the existence of the first moment $\hat{\delta}_{IV}$ is the result of the estimator performing poorly when $\widehat{\text{cov}}(\tilde{D}, \tilde{Z}) \approx 0$, as the estimates fluctuate widely depending on the magnitude and sign of $\widehat{\text{cov}}(\tilde{Z}, \tilde{\epsilon})$.

generalizing from the LATE to the ATET without additional assumptions.¹⁹

Imbens and Rubin (1997) and Abadie (2003) demonstrate that one can recover estimates of the marginal distributions of Y_{1i} and Y_{0i} for the compliers and thus $E(Y_{1i}|C)$ and $E(Y_{0i}|C)$. In Table 8, we produce estimates of $E(Y_{0i}|N)$, $E(Y_{0i}|C)$, $E(Y_{0i}|N \cup C)$, $E(Y_{1i}|A)$, $E(Y_{1i}|C)$, and $E(Y_{1i}|A \cup C)$ for each of the four outcomes in the Angrist and Evans data. To facilitate comparisons, we use IPW to reweight the data so that each set (i.e., A , N , or C) has the same distribution of the covariates as the aggregated data.²⁰ Two features stand out: First, we find little substantive difference between the always-takers and compliers in the treated state, but, second, we find that the compliers and the never-takers have substantially larger differences in the untreated state, consistent with Table 2.

Recall that we have no evidence for selection on unobserved variables for ϵ_{1i} , and, hence OLS estimation of the regression

$$Y_{1i} = X_i' \beta_1 + \epsilon_{1i} \quad (39)$$

using the treated units produces unbiased and consistent estimates of β_1 . We can then use the estimated $\hat{\beta}_1$ to predict a treated outcome for each untreated unit and then estimate $\hat{\Delta}^{ATEN}$ using the observed untreated outcomes and the predictions. But given the estimates in Table 8, this appears unwise. While we have no direct evidence against the hypothesis that $E(Y_{1i}|N, X_i) = E(Y_{1i}|A \cup C, X_i)$, we do have evidence that the never-takers and compliers differ in unobserved ways in the untreated state.

Given that we have evidence of the similarity of the always-takers and the compliers in the treated state, we may wish to use $E(Y_{0i}|X_i, C)$ to form the missing counterfactual for $E(Y_{0i}|X_i, A)$ and estimate Δ^{ATE} for the always-takers and compliers. This requires us to assume that $E(Y_{0i}|X_i, C) = E(Y_{0i}|X_i, A)$ as well as maintaining the VIA assumptions (so we may use the test for selection bias). In contrast, the IV estimator simply maintains the VIA

¹⁹Brinch et al. (2017) and Kowalski (2018a,b) consider the value of additional assumptions in addressing this same issue.

²⁰Doing so requires us to drop 3,104 observations (less than one percent of the data) to impose the common support condition.

assumptions. But the estimand of the IV estimator in this case lacks intrinsic interest: It is the impact of having a third child for the set of women who would have stopped with two children had they had a boy and a girl for their first two children. In contrast, the Δ^{ATET} , even with the potential bias, retains its substantive interest.

This argument applies specifically to the Angrist-Evans application. In the JTPA experiment, the compliers behaved in accordance with the treatment protocol in the experiment. We care about their LATE because it tells us about the effect of classroom training on those induced by the program to sign up for it. This constitutes an economically interesting parameter, albeit one estimated with large standard errors in the available data. Of course, interest in the LATE in a particular context does not imply lack of interest in Δ^{ATET} or Δ^{ATEN} .

6.4 Pre-testing and an alternative

A researcher who follows the methodology outlined in the paper for conditional mean independence will end up with an IV estimate, an OLS estimate, two test statistics (with p-values), and two bias estimates, one each for selection on Y_{0i} and Y_{1i} . What should the researcher do with these statistics? One strategy, inspired by the way researchers used the Durbin-Wu-Hausman (DWH) test back in the 1970s and 1980s, would view each test as a pre-test and report only the estimates indicated by the test, interpreted in the appropriate way. Thus, the researcher would report the IV estimate, interpreted as a LATE, given evidence of selection on Y_{0i} or Y_{1i} , and would report an OLS (or other CIA-based) estimate, interpreted as an ATET (or even an ATE), in the absence of such evidence. This strategy has the virtue of simplicity, and performs better than one might expect in the Monte Carlo analyses in Donald et al. (2014).

We see two main reasons not to adopt this strategy. First, the estimated variances do not incorporate the additional variance component induced by the pre-test. Guggenberger (2010) shows that incorporating this component in the traditional common effect DWH world really

matters. Now, of course, one might argue that in producing the typical empirical economics paper researchers examine many estimates that they do not report, yet which influence their choice of which ones to report, which implies that the typical empirical economics paper is already up to its eyeballs in pre-test issues. Yet no one complains, so why complain here? More importantly in our view, the traditional approach lacks nuance and hides valuable information from the reader. In many contexts, the test will reject or fail to reject only marginally; in many contexts, the data will imply only a modest bias, whether statistically significant or not. In such contexts, we think most readers will want to assign non-zero posterior weights to both estimates.

In addition, to aid in this process, the researcher might even formally combine the IV and OLS estimates in some reasonable way, along the lines suggested in Donald et al. (2014), who combine to minimize variance, or Hansen (2017), which combines based on the strength of the DWH test result. We think such combinations represent a promising avenue for future research.

7 Conclusion

In this paper, we derived simple tests for selection on unobserved variables when using instrumental variables. Our tests of conditional mean independence generalize various existing tests in the literature while our tests of independence of distributions go beyond what the existing literature provides. Using a Wald-like estimator, one can use the estimates generated by our test to assess the magnitude of the selection effect as well and thereby gain a much better understanding of the precise nature of any selection on unobserved variables.

Since the publication of Vytlacil (2002), we have understood the equivalence between the assumptions necessary for the LATE interpretation of IV estimates and models of selection into treatment based on latent indices. But IV estimation has always seemed to provide less information about the nature of the selection effect than control function estimation in the

context of the bivariate or trivariate normal model. In this paper, we showed that simple auxiliary regressions will produce rich insights into the nature and magnitude of the selection effect when using IV estimation. Our three empirical applications nicely demonstrate the knowledge our framework produces, as does its application in Azzam, Bates, and Fairris (2018).

8 Bibliography

Abadie, Alberto, 2003. “Semiparametric Instrumental Variable Estimation of Treatment Response Models” *Journal of Econometrics* 113 231-263.

Angrist, Joshua, 2004. “Treatment Effects Heterogeneity in Theory and Practice” *Economic Journal* 114(494) C52-C83.

Angrist, Joshua and William Evans, 1998. “Children and Their Parents Labor Supply: Evidence from Exogenous Variation in Family Size” *American Economic Review* 88(3) 450-77.

Angrist, Joshua, Guido Imbens, and Donald Rubin, 1996. “Identification of Causal Effects using Instrumental Variables” (with discussion) *Journal of the American Statistical Association* 91 444-72.

Angrist, Joshua, and Victor Lavy, 1999. “Using Maimonides’ Rule to Estimate the Effect of Class Size on Scholastic Achievement” *Quarterly Journal of Economics* 114(3) 533-575.

Azzam, Tarek, Michael Bates, and David Fairris, 2018. “Do Learning Communities Increase First Year College Retention? Testing the External Validity of Randomized Control Trials” Unpublished manuscript, University of California at Riverside.

Battistin, Eric and Enrico Rettore, 2008. “Ineligible and Eligible Non-Participants as a Double Comparison Group in Regression Discontinuity Designs” *Journal of Econometrics* 142 715-730.

Bertanha, Marinho and Guido Imbens, 2019. “External Validity in Fuzzy Regression Discontinuity Designs” *Journal of Business and Economic Statistics* forthcoming.

Björklund, Anders and Robert Moffitt, 1987. “The Estimation of Wage Gains and Welfare Gains in Self-Selection Models” *Review of Economics and Statistics* 69(1) 42-49.

Black, Dan, Seth Sanders, Evan Taylor, and Lowell Taylor, 2015. “The Impact of the Great

Migration on the Mortality of African-Americans: Evidence from Deep South” *American Economic Review* 105(2) 477-503.

Black, Dan and Jeffrey Smith, 2004. “How Robust is the Evidence on the Effects of College Quality? Evidence from Matching” *Journal of Econometrics* 121(1-2) 99-121.

Black, Dan and Jeffrey Smith. 2006. “Estimating the Returns to College Quality with Multiple Proxies for Quality” *Journal of Labor Economics* 24(3) 701-28.

Bloom, Howard, 1984. “Accounting for No-Shows in Experimental Evaluation Designs” *Evaluation Review* 8(2) 225-46.

Bloom, Howard, Larry Orr, Stephen Bell, George Cave, Fred Doolittle, Winston Lin, Johannes Bos, 1997. “The Benefits and Costs of JTPA Title II-A Programs: Key Findings from the National Job Training Partnership Study” *Journal of Human Resources* 32(3) 549-76.

Blundell, Richard, Lorrain Dearden, and Barbara Sianesi, 2005. “Evaluating the Effect of Education on Earnings: Models, Methods and Results from the National Child Development Survey” *Journal of the Royal Statistical Society, Series A* 167(3) 473-512.

Bound, John, David Jaeger, and Regina Baker, 1995. “Problems with Instrumental Variables Estimation when the Correlation Between the Instruments and the Endogenous Explanatory Variable is Weak” *Journal of the American Statistical Association* 90(430) 443-450.

Brinch, Christian, Magne Mogstad, and Matthew Wiswall, 2017. “Beyond LATE with a Discrete Instrument.” *Journal of Political Economy* 125(4) 985-1039.

Busso, Matias, John DiNardo, Justin McCrary, 2014. “New Evidence on the Finite Sample Properties of Propensity Score Reweighting and Matching Estimators” *Review of Economics and Statistics* 96(5) 885-897.

Cameron, Colin and Pravin Trivedi, 2005. *Microeconometrics: Methods and Applications* Cambridge, UK: Cambridge University Press.

Costa Dias, Monica, Hidehiko Ichimura, and Gerard van den Berg. 2013. “Treatment Evaluation with Selective Participation and Ineligibles.” *Journal of the American Statistical Association* 142 715-730.

Crump, Richard, V. Joseph Hotz, Guido Imbens, and Oscar Mitnik, 2009. “Dealing with Limited Overlap in Estimation of Average Treatment Effects” *Biometrika* 96(1) 187-99.

Deaton, Angus, 2010. “Instruments, Randomization, and Learning about Development” *Journal of Economic Literature* 48(2) 424-455.

- Diaz, Juan Jose and Sudhanshu Handa, 2006. "An Assessment of Propensity Score Matching as a Nonexperimental Estimator" *Journal of Human Resources* 41(2) 319-45.
- Donald, Stephen, Yu-Chin Hsu, and Robert Lieli, 2014. "Testing the Unconfoundedness Assumption via Inverse Probability Weighted Estimators of (L)ATT" *Journal of Business and Economic Statistics* 23(3) 395-415.
- Eberwein, Curtis, John Ham, and Robert LaLonde, 1997. "The Impact of Being Offered and Receiving Classroom Training on the Employment Histories of Disadvantaged Women: Evidence from Experimental Data" *Review of Economic Studies* 64(4) 655-82.
- Guggenberger, Patrick, 2010. "The Impact of a Hausman Pretest on the Asymptotic Size of a Hypothesis Test" *Econometric Theory* 26(2) 369-382.
- Guo, Zijian, Jing Cheng, Scott Lorch, and Dylan Small, 2014. "Using an Instrumental Variable to Test for Unmeasured Confounding" *Statistics in Medicine* 33(20) 3528-3546.
- Hansen, Bruce, 2017. "Stein-like 2SLS Estimator" *Econometric Reviews* 36(6-9) 840-852.
- Heckman, James, 1979. "Sample Selection Bias as a Specification Error" *Econometrica* 47(1) 206-248.
- Heckman, James, 1996. "Randomization as an Instrumental Variable" *Review of Economics and Statistics* 78(2) 336-340.
- Heckman, James, Neil Hohmann, Jeffrey Smith, and Michael Khoo, 2000. "Substitution and Drop Out Bias in Social Experiments: A Study of an Influential Social Experiment" *Quarterly Journal of Economics* 115(2) 651-94.
- Heckman, James, Hidehiko Ichimura, Jeffrey Smith, and Petra Todd, 1998. "Characterizing Selection Bias Using Experimental Data" *Econometrica* 66(5) 1017-98.
- Heckman, James, Robert Lalonde and Jeffrey Smith, 1999. "The Economics and Econometrics of Active Labor Market Programs," in *Handbook of Labor Economics, Volume 3*, eds. Orley Ashenfelter and David Card. Amsterdam: North-Holland, 1865-2097.
- Heckman, James and Salvador Navarro-Lazano, 2004. "Using Matching, Instrumental Variables, and Control Functions to Estimate Economic Choice Models" *Review of Economics and Statistics* 86(1) 30-57.
- Heckman, James, Daniel Schmieder, and Sergio Urzua, 2010. "Testing the Correlated Random Coefficients Model" *Journal of Econometrics* 158(2) 177-203.
- Heckman, James and Sergio Urzua, 2010. "Comparing IV with Structural Models: What Simple IV Can and Cannot Identify" *Journal of Econometrics* 156(1) 27-37.

Heckman, James, Sergio Urzua, and Edward Vytlacil, 2006. “Understanding Instrumental Variables in Models with Essential Heterogeneity” *Review of Economics and Statistics* 88(3) 389-432.

Heckman, James and Edward Vytlacil, 2005. “Structural Equations, Treatment Effects, and Econometric Policy Evaluation” *Econometrica* 73(3) 669-738.

Heckman, James and Edward Vytlacil, 2007a. “Econometric Evaluation of Social Programs, Part1: Causal Models, Structural Models and Econometric Policy Evaluation” in *Handbook of Econometrics, Volume 6B*, eds. James Heckman and Edward Leamer, 4780-4873.

Heckman, James and Edward Vytlacil, 2007b. “Econometric Evaluation of Social Programs, Part 2: Using Marginal Treatment Effect to Organize Alternative Estimators to Evaluate Social Programs and to Forecast their Effects in New Environments” in *Handbook of Econometrics, Volume 6B*, eds. James Heckman and Edward Leamer, 4875-5143.

Huber, Martin, 2013. “A Simple Test for the Ignorability of Non-compliance in Experiments” *Economic Letters* 120(3) 389-91.

Huber, Martin, Michael Lechner, and Connie Wunsch, 2013. “The Performance of Estimators Based on the Propensity Score” *Journal of Econometrics* 175(1) 1-21.

Huber, Martin and Giovanni Mellace, 2015. “Testing Instrument Validity for LATE Identification Based on Inequality Moment Constraints” *Review of Economics and Statistics* 97(2) 398-411.

Imbens, Guido, 2004. “Nonparametric Estimation of Average Treatment Effects Under Exogeneity: A Review” *Review of Economics and Statistics* 86(1) 4-29.

Imbens, Guido, 2010. “Better LATE Than Nothing: Some Comments on Deaton (2009) and Heckman and Urzua (2009)” *Journal of Economic Literature* 48(2) 399-423.

Imbens, Guido and Joshua Angrist, 1994. “Identification and Estimation of Local Average Treatment Effects” *Econometrica* 62(2) 467-475.

Imbens, Guido, and Thomas Lemieux, 2008. “Regression Discontinuity Designs: A Guide to Practice” *Journal of Econometrics* 142(2) 615-635.

Imbens, Guido, and Donald Rubin, 1997. “Estimating Outcome Distributions for Compliers in Instrumental Variables Models” *Review of Economic Studies* 64 555-574.

Joo, Joonhwi, and Robert LaLonde, 2014. “Testing for Selection Bias” IZA Discussion Paper No. 8455.

- Kitagawa, Toru, 2015. “A Test for Instrument Validity” *Econometrica* 83(5) 2043-2063.
- Klein, Tobias, 2010. “Heterogeneous treatment Effects: Instrumental Variables without Monotonicity?” *Journal of Econometrics* 155 99-116.
- Kowalski, Amanda, 2018a. “How to Examine External Validity Within an Experiment.” NBER Working Paper No. 24834.
- Kowalski, Amanda, 2018b. “Reconciling Seemingly Contradictory Results from the Oregon Health Insurance Experiment and the Massachusetts Health Reform.” NBER Working Paper No. 24647.
- Lee, David and Thomas Lemieux, 2010. “Regression Discontinuity Designs in Economics” *Journal of Economic Literature* 48(2) 281-355.
- List, John, Azeem Shaikh, and Yang Xu, 2019. “Multiple Testing in Experimental Economics” *Journal of Experimental Economics* 22(4) 773-793.
- Murphy, Kevin and Robert Topel, 2002. “Estimation and Inference in Two-step Models” *Journal of Business and Economic Statistics* 20(1) 88-97.
- Nelson, Charles and Richard Startz, 1990. “Some Further Results on the Exact Small Sample Properties of the Instrumental Variables Estimator” *Econometrica* 58(4) 967-976.
- Newey, Whitney, James Powell, and James Walker. 1990. “Semiparametric Estimation of Selection Models: Some Empirical Results” *American Economic Review* 80(2): 324-328.
- Romano, Joseph, and Azeem Shaikh, 2006. “Stepup Procedures for Control of Generalizations of the Familywise Error Rate” *Annals of Statistics* 34(4) 1850-73.
- Romano, Joseph, and Michael Wolf, 2005. “Exact and Approximate Stepdown Methods for Multiple Hypothesis Testing” *Journal of the American Statistical Association* 100(469) 94-108.
- Roy, Andrew, 1951. “Some Thoughts on the Distribution of Earnings” *Oxford Economics Papers* 3(2) 135-146.
- Schochet, Peter, 2008. *Technical Methods Report: Guidelines for Multiple Testing in Impact Evaluations* NCEE 2008-4018. Washington, DC: National Center for Education Evaluation and Regional Assistance, Institute of Education Sciences, U.S. Department of Education.
- Sloczyński, Tymon, 2020. “Interpreting OLS Estimands When Treatment Effects Are Heterogeneous: Smaller Groups Get Larger Weights” IZA Discussion Paper No. 13283.

Smith, Jeffrey and Petra Todd, 2005. "Does Matching Overcome LaLondes Critique of Non-experimental Estimators?" *Journal of Econometrics* 125(1) 305-53.

Smith, Jeffrey and Alexander Whalley. 2020. "How Well Do We Measure Public Job Training?" Unpublished manuscript, University of Wisconsin.

Vytlacil, Edward, 2000. *Three Essays on the Nonparametric Evaluation of Treatment Effects* Dissertation, University of Chicago.

Vytlacil, Edward, 2002. "Independence, Monotonicity, and Latent Index Models: An Equivalence Result" *Econometrica* 70(1) 331-41.

Yule, Udry, 1907. "On the Theory of Correlation for Any Number of Variables, Treated by a New System of Notation" *Proceedings of the Royal Society* 79(529) 181-93.

Table 1: Causal Impact of Having More than Three Children on Mother's Labor Supply, Angrist and Evans (1998)

	Parametric OLS model	Parametric IV model	Semiparametric model	Semiparametric IV model
Worked last year	-0.176*** (0.002)	-0.120*** (0.025)	-0.174*** (0.0016)	-0.117*** (0.025)
Weeks worked	-8.97*** (0.071)	-5.66*** (1.108)	-8.90*** (0.073)	-5.53*** (1.109)
Hours worked	-6.66*** (0.061)	-4.59*** (0.945)	-6.59*** (0.062)	-4.45*** (0.946)
Income	-3,768*** (33.45)	-1,960*** (541.5)	-3,739 (35.47)	-1,915*** (542.0)
First Stage: Same sex coefficient	---	0.062*** (0.0015)	----	0.062*** (0.0015)
N	394,835	394,835	394,835	394,835

*5 percent significance level, ** 1 percent significant level, *** 0.1 percent significance level

Notes: Covariates in the parametric model include the age of the mother, the age of the mother at first birth, indicators for whether the mother is black or non-black and non-white, an indicator for whether the mother is Hispanic, an indicator for whether the first child was a boy, and an indicator for whether the second child was a boy. For the semiparametric model, we drop the indicator for the second child being a boy to avoid perfect colinearity with the instrument, an indicator that both of the first two children are the same sex. The semiparametric IV regression model uses a fully saturated model in the covariates and an additively separable term for having more children. The F-statistic on the instrument for the parametric model equals 1,711. For the semiparametric model it equals 1,702. For the semiparametric model, 72 cases have predicted values of one for the probability of having more children and 168 have predicted probabilities of zero. Our parametric estimates exactly match those of Angrist and Evans, Table 7, columns (1) and (2).

Table 2: Test of CIA using Means, Angrist and Evans (1998) Data

	OLS		Semiparametric	
	(CMIA ⁰) test	(CMIA ¹) test	(CMIA ⁰) test	(CMIA ¹) test
Worked last year				
Dependent variable	Y ₀	Y ₁	Y ₀	Y ₁
Coefficient on instrument (standard error) [selection effect]	0.0046* (0.002) [0.047]	0.0015 (0.003) [0.011]	0.0051** (0.002) [0.052]	0.0017 (0.003) [0.012]
N	236,092	158,743	236,092	158,743
Weeks worked				
Dependent variable	Y ₀	Y ₁	Y ₀	Y ₁
Coefficient on instrument (standard error) [selection effect]	0.297*** (0.090) [3.01]	0.053 (0.104) [0.37]	0.315*** (0.090) [3.19]	0.063 (0.105) [0.44]
N	236,092	158,743	236,092	158,743
Hours worked				
Dependent variable	Y ₀	Y ₁	Y ₀	Y ₁
Coefficient on instrument (standard error) [selection effect]	0.205** (0.075) [2.08]	-0.0004 (0.093) [0.00]	0.221** (0.075) [2.24]	0.016 (0.093) [0.11]
N	236,092	158,743	236,092	158,743
Income				
Dependent variable	Y ₀	Y ₁	Y ₀	Y ₁
Coefficient on instrument (standard error) [selection effect]	188*** (45.43) [1,904]	-7.01 (48.28) [-49]	194*** (45.40) [1,965]	-5.49 (48.41) [-38]
N	236,092	158,743	236,092	158,743

*5 percent significance level, ** 1 percent significant level, *** 0.1 percent significance level

Notes: Covariates in the parametric model include the age of the mother, the age of the mother at first birth, indicators for whether the mother is black or non-black and non-white, an indicator for whether the mother is Hispanic, an indicator for whether the first child was a boy, and an indicator for whether the second child was a boy. For the semiparametric model, we drop the indicator for the second child being a boy to avoid perfect colinearity with the instrument, an indicator that both of the first two children are the same sex. The semiparametric model uses a fully saturated model in the covariates.

Table 3: Test of CIA using Percentiles and Moments, Angrist and Evans (1998) Data

Dependent Variable	Weeks worked		Hours worked		Income	
	(CIA ⁰) test	(CIA ¹) test	(CIA ⁰) test	(CIA ¹) test	(CIA ¹) test	(CIA ⁰) test
	Y ₀	Y ₁	Y ₀	Y ₁	Y ₀	Y ₁
<i>Percentiles</i>						
Zero	-0.005** (0.002)	-0.001 (0.003)	-0.005** (0.002)	-0.001 (0.003)	-0.006** (0.002)	-0.001 (0.003)
50 th	-0.006*** (0.002)	-0.001 (0.003)	-0.007*** (0.002)	-0.002 (0.002)	-0.007** (0.002)	-0.002 (0.002)
60 th	-0.007*** (0.002)	-0.001 (0.002)	-0.007** (0.002)	-0.001 (0.002)	-0.005*** (0.002)	0.001 (0.002)
70 th	-0.006*** (0.002)	-0.002 (0.002)	-0.005** (0.002)	-0.001 (0.002)	-0.006** (0.002)	0.001 (0.002)
80 th	-0.004* (0.002)	-0.001 (0.002)	-0.002** (0.001)	0.001 (0.001)	-0.007*** (0.002)	-0.001 (0.002)
90 th	---	---	---	---	-0.002 (0.001)	-0.00003 (0.001)
<i>Moments</i>						
1 st moment	0.297*** (0.090)	0.053 (0.104)	0.205** (0.075)	-0.0004 (0.093)	188*** (45.43)	-7.01 (48.28)
2 nd moment	0.788 (1.067)	0.860 (1.932)	-0.498 (1.304)	-3.54 (1.933)	0.572* (0.236)	-0.304 (2.946)
3 rd moment	0.017 (0.520)	0.005 (0.009)	0.004 (0.013)	-0.039 (0.020)	0.381 (0.228)	0.075 (0.433)
4 th moment	0.005 (0.004)	0.007 (0.013)	-0.048 (0.049)	-0.163* (0.078)	4.27 (3.529)	2.20 (6.923)
N	236,092	158,743	236,092	158,743	236,092	158,743

* 5 percent significance level **1 percent significance level *** 0.1 percent significance level

Notes: Percentile regressions use a parametric model where the dependent variable is an indicator for being equal to or below given percentile. All higher moments are centered about the mean. The third and fourth moments are normalized by the standard deviation. The second moment of income we divide coefficient and standard error by a million. Covariates include the age of the mother, the age of the mother at first birth, indicators for whether the mother is black or non-black and non-white, an indicator for whether the mother is Hispanic, and an indicator for whether the first child was a boy, and an indicator for whether the second child was a boy.

Table 4: Causal Impact of Class Size on Test Scores, Angrist and Lavy (1999)

	Reading comprehension	Math
Class size (continuous measure)	-0.098 (0.090)	0.095 (0.113)
N	415	415

	Reading comprehension	Math
Class size less than 35	1.365 (1.888)	-1.335 (2.382)
n	415	415

*5 percent significance level, ** 1 percent significant level, *** 0.1 percent significance level

Notes: Results from 2SLS regressions. Dependent variable is the average score of a class. All regressions control for percentage disadvantaged. Our replication of Angrist-Lavy uses their Discontinuity Sample and match results from their Table 5, columns (5) and (11) precisely, and using their continuous instrument. Our treatment indicator is defined as having a class size of less than 35. Our instrument is defined as having the Angrist-Lavy instrument taking on a value of less than 35. The F-statistic on the binary instrument equals 49.4. If we regress our instrument against class size, controlling percentage disadvantaged, the point estimate indicates our instrument results in a reduction of class size of 7.58

Table 5: Test of CMIA, Angrist and Lavy (1999) Data

	(CMIA ⁰) test	(CMIA ¹) test
Reading comprehension	Y_0	Y_1
Coefficient on instrument (standard error)	-0.673 (1.236)	1.959 (1.348)
Math		
Coefficient on instrument (standard error)	-1.934 (1.762)	0.786 (1.567)
N	177	238

*5 percent significance level, ** 1 percent significant level, *** 0.1 percent significance level

Notes: All regressions control for percentage disadvantaged. Our treatment indicator is defined as having a class size of less than 35. Our instrument is defined as having the Angrist-Lavy instrument taking on a value of less than 35.

Table 6: Test of CIA using Percentiles and Moments, Angrist and Lavy (1999) Data

Dependent Variable	Reading Comprehension		Math	
	(CIA ⁰) test	(CIA ¹) test	(CIA ⁰) test	(CIA ¹) test
	Y ₀	Y ₁	Y ₀	Y ₁
<i>Percentiles</i>				
20 th	-0.042 (0.056)	-0.087 (0.053)	-0.003 (0.059)	0.013 (0.056)
40 th	-0.047 (0.073)	-0.127* (0.062)	0.192* (0.084)	-0.038 (0.064)
60 th	0.131 (0.084)	-0.134* (0.061)	0.178* (0.084)	-0.093 (0.064)
80 th	0.026 (0.071)	-0.067 (0.054)	0.102 (0.075)	-0.075 (0.052)
<i>Moments</i>				
1 st moment	-0.673 (1.236)	1.959 (1.348)	-1.934 (1.762)	0.786 (1.567)
2 nd moment	1.78 (14.04)	11.11 (13.12)	-8.77 (17.17)	29.81 (18.36)
3 rd moment	0.027 (0.540)	0.097 (0.489)	0.347 (0.521)	-0.276 (0.524)
4 th moment	0.815 (1.219)	1.806 (1.111)	-0.672 (1.022)	2.237 (1.239)
N	177	238	177	238

* 5 percent significance level **1 percent significance level *** 0.1 percent significance level

Notes: Percentile regressions use a parametric model where the dependent variable is an indicator for being equal to or below given percentile. All higher moments are centered about the mean. The third and fourth moments are normalized by the standard deviation. All regressions control for percentage disadvantaged

Table 7: Impact of Training for Adult Women Recommended for Classroom Training in the National JTPA Study

Covariates	No	Yes
Mean of employment for control group	0.505	0.505
Mean training for control group	0.344	0.344
Intent to treat (standard error) (n=2,374)	0.041** (0.020)	0.037* (0.020)
Bloom estimator		
First-stage treatment indicator (standard error) [F-statistic on instrument] (n=2,374)	0.305*** (0.019) [246]	0.305*** (0.019) [246]
Impact of classroom training on compliers (standard error) [n=2,374]	0.136** (0.067)	0.122* (0.065)
Treatment group indicator for Y_0 regression (standard error) [selection effect] (n=1,233)	0.006 (0.029) [0.013]	0.005 (0.027) [0.011]
Treatment group indicator for Y_1 regression (standard error) [selection effect] (n=1,501)	0.068** (0.032) [0.145]	0.062** (0.031) [0.132]

* 5 percent significance level **1 percent significance level *** 0.1 percent significance level

Notes: The dependent variable is an indicator variable for whether the participant is employed in the 18th month after random assignment. The treatment indicator equals one when the participant is assigned to the treatment group. The classroom training variable is an indicator for whether the participant received classroom training in the first nine months after random assignment. For the specification with covariates, the set of covariates include age and the square of age and a vector of indicator variables. The indicator variables indicate whether the participant has never been married, whether the participant is currently married, whether the participant is a non-Hispanic black, whether the participant is Hispanic, whether the participant is another race/ethnicity (white, non-Hispanic is the excluded category), whether the participant has less than a high school degree, whether the participant has a General Education Development degree, whether the participant has more than a high school degree (high school degree is the excluded category), whether the participant was on Aid to Families with Dependent Children (AFDC) at the time of random assignment, whether the participant was on AFDC for two years or more, whether the participant was on food stamps at the time of random assignment, whether the participant had children under five years of age in the household, whether the participant had children under 18 in the household, whether the participant reported problems with her English skills, whether the participant reported never working for pay, whether the participant reported never working full time, whether the participant worked in the 12 months prior to random assignment, a cubic in the fraction of the year that the participant worked prior to random assignment, and 15 indicators for the experimental sites. To avoid dropping observations, if a variable was missing, we set its value to zero and added an indicator variable equal to one when the variable was missing.

Table 8: Moments from the Angrist and Evans (1998) Data

	Worked last year	Weeks worked	Hours worked	Income
$\bar{Y}_{0,N}$	0.638	24.62	21.57	\$8,819
$\bar{Y}_{0,NUC}$	0.633	24.32	21.34	\$8,616
$\bar{Y}_{1,A}$	0.462	15.67	14.96	\$5,051
$\bar{Y}_{1,AUC}$	0.463	15.74	14.95	\$5,049
$\bar{Y}_{0,C}$	0.599	22.29	19.79	\$7,258
$\bar{Y}_{1,C}$	0.472	16.23	14.90	\$5,037
N	391,731	391,731	391,731	391,731

Notes: We reweight the covariates – the age of the mother (individual year indicators), the age of the mother at first birth (individual year indicators), indicators for whether the mother is black or non-black and non-white, an indicator for whether the mother is Hispanic, an indicator for whether the first child was a boy, and an indicator for whether the second child was a boy – to have the same CDF for distribution of the aggregate data using inverse probability weighting. We drop 3,104 observations from Table 1 and 2 that violate the common support assumption. There are 1,418 cells.