

Cambridge Books Online

<http://ebooks.cambridge.org/>



Advances in Economics and Econometrics

Tenth World Congress

Edited by Daron Acemoglu, Manuel Arellano, Eddie Dekel

Book DOI: <http://dx.doi.org/10.1017/CBO9781139060035>

Online ISBN: 9781139060035

Hardback ISBN: 9781107016064

Paperback ISBN: 9781107627314

Chapter

11 - ExtrapoLATE-ing: External Validity and Overidentification in the
LATE Framework pp. 401-434

Chapter DOI: <http://dx.doi.org/10.1017/CBO9781139060035.012>

Cambridge University Press

ELEVEN

Extrapolating: External Validity and Overidentification in the LATE Framework

Joshua D. Angrist and Iván Fernández-Val

1.0 Introduction

Local Average Treatment Effects (LATE) capture the causal effect of an instrument-induced shift in treatment. This effect is necessarily tied to the instrument that generates the shift. The interpretation of instrumental variable (IV) estimates as instrument-specific should not be surprising or troubling – when this point is cast in terms of specific examples, we wonder how it could be otherwise. Quarter-of-birth instruments for a wage equation reveal the payoff to schooling induced by compulsory-attendance laws and not the value of a bought-and-paid-for MBA. Still, a clear statement of the nature of causal effects revealed by any instrument raises questions about the external validity of this estimate. Can we use a given IV estimate to identify the effects induced by another source of variation? What about an unconditional average effect? Can we go from average effects on compliers to average effects on the entire treated population?

The usual answer to these questions is “no.” Except in special cases, we cannot go farther, at least not without additional assumptions. As described by Angrist, Imbens, and Rubin (1996), the treated population includes two groups: (1) compliers whose behavior is affected by the instrument at hand, and (2) always-takers who are treated irrespective of whether a Bernoulli instrument is switched off or on. The nontreated are likewise composed of compliers and never-takers; the latter group avoids treatment

Our thanks to Alberto Abadie, Manuel Arellano, Gary Chamberlain, Victor Chernozhukov, Guido Imbens, Frank Vella, participants in the 2010 Econometric Society World Congress, and seminar participants at Boston College, Boston University, Carnegie Mellon, Georgetown, Michigan, and the Harvard–MIT Econometrics Workshops for helpful discussions and comments.

no matter what. In the absence of strong homogeneity or distributional assumptions, the data are uninformative for always-takers and never-takers. Moreover, each instrument typically generates its own compliant sub-population. Effects for one group of compliers need not generalize to another.

One route to external validity is structural. A latent-index choice framework sometimes allows us to fill gaps in the data. For example, Heckman, Tobias, and Vytlacil (2001, 2003) and Angrist (2004) used parametric latent-index models to identify and compare alternative causal effects. Chamberlain (2010) developed a Bayesian semiparametric procedure for extrapolation that relies on models for variation in outcome distributions as a function of the first stage. In a recent paper, Heckman (2010) summarized a literature on IV that establishes theoretical links between parameters such as LATE and effects on the treated. There is no free lunch, however; these links provide *nonparametric* identification of effects other than LATE only with instruments that drive the probability of treatment over a wide range (in fact, from zero to one if we hope to recover the population average causal effect). Such “super-instruments” are rare if not unknown in applied work. In practice, most instruments that identify causal effects are discrete with finite support, and many are Bernoulli.

These theoretical challenges notwithstanding, the predictive value of a particular set of IV estimates may be revealed empirically when a researcher succeeds in isolating multiple instruments for the same underlying causal relation. A pioneering effort in this direction is the Oreopoulos (2006) study of the economic returns to schooling. Oreopoulos compared IV estimates of the returns to schooling across instruments of different strengths. Some of Oreopoulos’s instruments are derived from compulsory-attendance policies that had modest effects on schooling. However, two of his policy experiments generate instruments with large first-stage effects, close to a half-year increase in schooling. Moreover, in these examples, there are few never-takers, so LATE is the average effect of treatment on the nontreated. As it turns out, Oreopoulos’s IV estimates of the returns to schooling using marginal and full-bore instruments are similar, suggesting a robust causal effect that is likely to have considerable predictive value. Angrist, Lavy, and Schlosser (2010) made a similar homogeneity argument for IV estimates of the causal effects of family size on human capital (a relationship known as the “quantity–quality trade-off”). Sex-composition instruments, which have a modest first stage, generate causal effects similar to those found

using twins instruments, for which the first stage is larger by an order of magnitude, with no never-takers.¹

These examples are encouraging because they suggest that in a number of important applications, IV estimates are reasonably stable across instruments. In many applications, however, heterogeneous effects are likely to be important. For example, Ichino and Winter-Ebmer (1999) documented substantial heterogeneity in alternative IV estimates of the returns to schooling in Germany. The range of variation in this case far exceeds that which can be attributed to sampling variance.

In this chapter, we ask whether instrumental variable estimates using different instruments, possibly with very different compliant subpopulations, can be reconciled solely by differences in the observed characteristics of compliers. An important consequence of the Abadie (2003) weighting theorem is that the distribution of complier characteristics is identified and easy to describe empirically. A natural first step when comparing alternative IV estimates is to compare and contrast the observed characteristics of compliers. Assuming that treatment-effect heterogeneity is limited in a way that we make precise below, we can use the distribution of complier characteristics to construct an estimator that converts covariate-specific LATEs into broader parameters such as the effect of treatment on the treated and LATE using alternative instruments. Our first contribution is to explain and illustrate this approach to external validity.²

The second item on our agenda is the use of overidentification tests in pursuit of external validity. In the classical simultaneous equation framework, statistically significant differences between alternative IV estimates signal a failure of internal validity, perhaps due to violations of the exclusion restriction. In the LATE framework, by contrast, different (internally valid) instruments capture different causal effects. At the same time, covariate-specific overidentification tests and summary conditional tests weighted across covariate cells indicate whether differences in the observed characteristics of compliant subpopulations are enough to explain differences in

¹ See also Ebenstein (2009), who compared LATEs generated by first stages of varying strength for the effect of fertility on labor supply in the United States and Taiwan; and Cruces (2005) and Cruces and Galiani (2007), both of which explored the external validity of the fertility estimates generated by the Angrist and Evans (1998) sex composition instruments using data from Latin America.

² Other efforts in this direction include Angrist (2004) and Aronow and Sovey (2010), both of which focused on the possibility of using LATE to learn about unconditional average treatment effects.

unconditional effects. If so, it seems fair to say that the IV estimates at hand meet an empirically useful standard of external validity.

In practice, the question of whether covariates explain the difference between two sets of IV estimates need not have a simple answer. For some covariate values, or perhaps over a certain range, there may be a good match. In other cases, the match will be poor and the underlying estimates essentially unreconciled. We therefore use overidentification test statistics to design a hybrid testing-and-weighting scheme that isolates covariate-defined subsamples for which alternative IV estimates can be reconciled. We think of these samples as coming from a subpopulation for which heterogeneity in treatment effects is a function solely of observed characteristics. For this subpopulation, the predictive value of IV estimates is likely to be especially high.

The ideas in this chapter are illustrated through a comparison of alternative IV estimates of the labor supply consequences of childbearing. As in the study by Angrist and Evans (1998), the instruments are constructed from twin births and sibling sex composition. These instruments have very different first stages and produce significantly different estimates of the causal effect of a third birth. We show here that differences in the characteristics of instrument-specific complier subpopulations can account for most of the difference between the two sets of IV estimates.

2.0 Framework

We imagine that each individual is associated with two potential outcomes, y_0 and y_1 , which describe outcomes realized under alternative assignments of a Bernoulli treatment, D . The observed outcome, y , is linked to potential outcomes as follows:

$$y = y_0 + (y_1 - y_0)D. \quad (1)$$

A random-coefficients notation for this is:

$$y = \alpha + rD + \eta$$

where $\alpha = \mathbb{E}[y_0]$, $\eta = y_0 - \alpha$, and $r = y_1 - y_0$ is an individual-level causal effect.

We also define potential treatment status indexed against a Bernoulli instrument, z . Potential treatment status is D_1^z when the instrument is switched on and D_0^z when the instrument is switched off. The variables D_1^z and D_0^z are superscripted to signal the fact that they are tied to z . Observed

treatment status is:

$$D = D_0^z + (D_1^z - D_0^z)z$$

or, in random-coefficients notation,

$$D = \gamma + pz + v$$

where $\gamma = \mathbb{E}[D_0^z]$, $v = D_0^z - \gamma$, and $p = D_1^z - D_0^z$.

IV using z as an instrument for the effect of D on Y with no covariates is the Wald (1940) estimator. The Wald estimand can be interpreted as the effect of D on those whose treatment status can be changed by the instrument. Assuming, as we do here, that the instrument can make treatment move only in one direction (e.g., make treatment more likely), those whose treatment status is changed by z have $D_1^z = 1$ and $D_0^z = 0$. The causal effect on this group is LATE (Imbens and Angrist 1994). Formally, we have:

Assumption 1 (LATE):

- (a) *Independence and Exclusion:* $(Y_1, Y_0, D_1^z, D_0^z) \perp\!\!\!\perp z$.
- (b) *First-stage:* $\mathbb{E}[D_1^z - D_0^z] \neq 0$ and $0 < \mathbb{P}[z = 1] < 1$.
- (d) *Monotonicity:* $D_1^z \geq D_0^z$ a.s., or vice versa.

Theorem 1 (LATE): Under Assumption 1:

$$\frac{\mathbb{E}[Y | z = 1] - \mathbb{E}[Y | z = 0]}{\mathbb{E}[D | z = 1] - \mathbb{E}[D | z = 0]} = \mathbb{E}[Y_1 - Y_0 | D_1^z > D_0^z] = \mathbb{E}[r | p > 0] =: \Delta^z.$$

Proof: See Imbens and Angrist (1994). ■

As noted by Angrist, Imbens, and Rubin (1996), the LATE framework partitions the population exposed to an instrument into three instrument-dependent groups. These groups are defined by the way people react to the instrument.

Definition 1 (Groups Defined by Instrument z):

- (a) *z-Compliers:* The subpopulation with $D_1^z = 1$ and $D_0^z = 0$.
- (b) *z-Always-takers:* The subpopulation with $D_1^z = D_0^z = 1$.
- (c) *z-Never-takers:* The subpopulation with $D_1^z = D_0^z = 0$.

LATE using z as an instrument is the effect of treatment on the group of z -compliers.

Table 1. *Wald Estimates of the Effects of Family Size on Labor Supply*

Dependent Variable	Mean	OLS (1)	Twins Instrument		Same-Sex Instrument		Both 2SLS Estimates (6)
			First Stage (2)	Wald Estimates (3)	First Stage (4)	Wald Estimates (5)	
Weeks worked	20.83	-8.98 (0.072)	0.603 (0.008)	-3.28 (0.634)	0.060 (0.002)	-6.36 (1.18)	-3.97 (0.558)
	Overid: $\chi^2(1)$ (p-value)	-	-	-	-	-	5.3 (0.02)
Employment	0.565	-0.176 (0.002)		-0.076 (0.014)		-0.132 (0.026)	-0.088 (0.012)
	Overid: $\chi^2(1)$ (p-value)	-	-	-	-	-	3.5 (0.06)

Note: The table reports OLS, Wald, and 2SLS estimates of the effects of a third birth on labor supply using twins and sex composition instruments. Data are from the Angrist and Evans (1998) extract from the 1980 U.S. census 5 percent sample, including women aged 21–35 with at least two children. OLS models include controls for mother's age, age at first birth, ages of the first two children, and dummies for race. The sample size is 394,840.

3.0 Covariate-Driven Heterogeneity

3.1 Two Instruments for One Effect

The case for omitted variables bias in the relationship between fertility and labor supply is clear: Mothers with weak labor force attachment or low earnings potential may be more likely to have children than mothers with strong labor force attachment or high earnings potential. This makes the observed association between family size and employment difficult to interpret because mothers who have big families work less anyway. Angrist and Evans (1998) solved this omitted-variables problem using two instruments, both of which lend themselves to Wald-type estimation strategies.

The first Wald estimator uses twin births, an instrument for the effects of family size introduced by Rosenzweig and Wolpin (1980). The twins instrument in Angrist and Evans (1998) is a dummy indicating multiple second births in a sample of mothers with at least two children. The twins first stage is about 0.6, an estimate reported in Column 2 of Table 1. This means that 40 percent of mothers with two or more children would have had

a third birth anyway; a multiple third birth increases this proportion to 1. The twins instrument rests on the idea that the occurrence of a multiple birth is essentially random – unrelated to potential outcomes or demographic characteristics – and that a multiple birth affects labor supply solely by increasing fertility.

The second Wald estimator in Table 1 uses a dummy for same-sex sibling pairs as an instrument. This is motivated by the fact that American parents with two children are more likely to have a third child if the first two are same-sex than if sex composition is mixed. This is illustrated in Column 4 of Table 1, which shows that parents of same-sex siblings are about 6 percentage points more likely to have a third birth than those with a mixed-sex sibship (the probability of a third birth among parents with a mixed-sex sibship is 0.38). Internal validity of the same-sex instrument rests on the claim that sibling sex composition is essentially random and affects mothers' labor supply solely by increasing fertility.

Twins and sex-composition IV estimates both suggest that the birth of a third child substantially reduces weeks worked and employment. Wald estimates using twins instruments show a precisely estimated reduction in weeks worked of a little more than three weeks, with an employment reduction of about 0.8. These results, which appear in Column 3 of Table 1, are smaller in absolute value than the corresponding ordinary least squares (OLS) estimates reported in Column 1 (the latter include a set of controls listed in the table). This suggests that the OLS estimates are exaggerated by selection bias. It is interesting and perhaps surprising that Wald estimates constructed using a same-sex dummy, reported in Column 5, are larger in magnitude than the twins estimates. The juxtaposition of twins and sex-composition instruments suggests that different instruments need not generate similar estimates of causal effects even if both instruments are valid.

The last column of Table 1 reports 2SLS estimates of childbearing using both twins and same-sex instruments, along with the associated overidentification test statistic. The overidentification test statistic generates p-values of 0.02 and 0.06, implying that the twins and sex-composition IV estimates are at least marginally significantly different from one another.

Twins and same-sex IV estimates reflect behavior in different compliant subpopulations. To show this, we let x be a Bernoulli-distributed characteristic – for example, a dummy indicating college graduates. Are sex-composition compliers more or less likely to be college graduates than other women with two children? This question is answered by the following

complier characteristics ratio:

$$\begin{aligned} \frac{\mathbb{P}[x = 1 \mid D_1^z > D_0^z]}{\mathbb{P}[x = 1]} &= \frac{\mathbb{P}[D_1^z > D_0^z \mid x = 1]}{\mathbb{P}[D_1^z > D_0^z]} \\ &= \frac{\mathbb{E}[D \mid z = 1, x = 1] - \mathbb{E}[D \mid z = 0, x = 1]}{\mathbb{E}[D \mid z = 1] - \mathbb{E}[D \mid z = 0]} \end{aligned}$$

where the first equality follows by Bayes rule. This second equality shows that the relative likelihood that a z -complier is a college graduate is given by the ratio of the first stage for college graduates to the overall first stage.

This calculation is illustrated in Table 2, which reports compliers' characteristics ratios for the age of the second-born and mother's schooling as described by dummies for high school graduates, some college, and college graduates. Twins compliers have younger second-born children, reflecting the fact that few women who recently had their second child will have had time to have a third child. The birth of a third child in this group is therefore especially likely to have been caused by a multiple pregnancy. (Among second-borns who are less than about a year old, the only explanation for a third birth in the family is a multiple birth.) This is important because the birth of a third child may matter less if the second child is young, helping explain the finding that twins-IV estimates are smaller than same-sex estimates (Gelbach 2002 showed that the presence of a child younger than age 5 in the household is a key labor-supply mediator).

Twins compliers are more likely to be college graduates than the average mother, whereas sex-composition compliers are less educated. This fact also helps explain the smaller Wald estimates generated by twins instruments because Angrist and Evans (1998) showed that the labor-supply consequences of childbearing decline with mother's schooling.

A general method for constructing the mean or other features of the distribution of covariates for compliers uses Abadie's (2003) kappa-weighting scheme. A consequence of Theorem 3.1 in Abadie (2003) is that:

$$\mathbb{E}[x \mid D_1 > D_0] = \frac{\mathbb{E}[\kappa^z(x) x]}{\mathbb{E}[\kappa^z(x)]} \quad (2)$$

where:

$$\kappa^z(x) = 1 - \frac{D(1-z)}{1 - \mathbb{P}[z = 1 \mid x = x]} - \frac{(1-D)z}{\mathbb{P}[z = 1 \mid x = x]}.$$

Table 2. *Complier Characteristics for Twins and Sex Composition Instruments*

Variable	Population Mean	Mean for Twins Compliers		Mean for Same-Sex Compliers	
	$E[x_{1i}]$ (1)	$E[x_{1i} D_{1i} > D_{0i}]$ (2)	$E[x_{1i} D_{1i} > D_{0i}]/E[x_{1i}]$ (3)	$E[x_{1i} D_{1i} > D_{0i}]$ (4)	$E[x_{1i} D_{1i} > D_{0i}]/E[x_{1i}]$ (5)
A. Dummy characteristics					
Age of second child is less than or equal to 4 years	0.343	0.449	1.31	0.194	0.565
High school graduate	0.488	0.498	1.02	0.515	1.06
Some college	0.202	0.212	1.05	0.212	1.05
College graduate	0.132	0.151	1.14	0.092	0.702
B. Discrete, ordered characteristics					
Age of second child	6.59	5.51	0.835	7.14	1.08
Mother's schooling	12.13	12.43	1.03	12.09	1.00

Note: The table reports an analysis of complier characteristics for twins and sex composition instruments. The ratios in columns 3 and 5 in Panel A give the relative likelihood that compliers have the characteristic indicated at left. The values in columns 2 and 4 in Panel B represent Abadie's (2003) kappa-weighted means. Data are from the 1980 census 5 percent sample including mothers aged 21–35 with at least two children, as in Angrist and Evans (1998). The sample size is 394,840.

Intuitively, this works because, as Abadie showed, the weighting function, $\kappa^z(x)$, “finds compliers,” even though it is not a simple indicator for compliers. Estimates of $\mathbb{E}[x \mid D_1 > D_0]$ for age of second child and mother’s education are reported in the last two rows of Table 2. These estimates show a marked difference in the average second-child age and a smaller difference in schooling. The main difference between the schooling of twins and same-sex compliers is in the proportion of college graduates.

3.2 Covariates and Extrapolation

Covariates play two roles in our analysis. First, they may be necessary for identification. For example, we might want to control for race and maternal age when using twins instruments because the probability of multiple births varies by race and increases with maternal age. Second, we use covariates for extrapolation. Specifically, we argue that in some cases, including the twins and same-sex comparisons, variation in causal effects across covariate cells is sufficient to explain differences between IV estimates.

The foundation of our analysis with covariates is a *conditional* independence assumption. This assumption expresses the idea that we think of the instruments as being “as good as randomly assigned,” conditional on covariates, x . Generalizing Assumption 1, we have:

Assumption 2 (Conditional LATE):

- (a) *Independence and Exclusion:* $(Y_1, Y_0, D_1^z, D_0^z) \perp\!\!\!\perp Z \mid x$ a.s.:
- (b) *First-stage:* $\mathbb{E}[D_1^z - D_0^z \mid x] \neq 0$ and $0 < \mathbb{P}[Z = 1 \mid x] < 1$ a.s.
- (c) *Monotonicity:* $\mathbb{P}[D_1^z \geq D_0^z \mid x] = 1$ a.s., or $\mathbb{P}[D_1^z \leq D_0^z \mid x] = 1$ a.s.

For each value of x , we define covariate-specific LATE using z as an instrument:

$$\Delta^z(x) := \mathbb{E}[Y_1 - Y_0 \mid D_1^z > D_0^z, x = x]. \tag{3}$$

As noted by Frolich (2007), when conditioning is required for identification, unconditional LATE can be constructed by averaging $\Delta^z(x)$:

$$\begin{aligned} \Delta^z &= \mathbb{E}[\Delta^z(x) \mid D_1^z > D_0^z] = \int \Delta^z(x) dF_x(x \mid D_1^z > D_0^z) \\ &= \int \Delta^z(x) \frac{\mathbb{E}[D \mid z = 1, x = x] - \mathbb{E}[D \mid z = 0, x = x]}{\mathbb{E}[D \mid z = 1] - \mathbb{E}[D \mid z = 0]} dF_x(x) \end{aligned} \tag{4}$$

where $F_x(\cdot | D_1^z > D_0^z)$ is the distribution of x for z -compliers and $F_x(\cdot)$ is the distribution of x in the population.

We first show how to construct average causal effects such as the effect on the treated, $\mathbb{E}[Y_1 - Y_0 | D = 1]$, from $\Delta^z(x)$. This is possible because we assume that heterogeneity in causal effects across instruments is entirely due to changes in the observable characteristics of compliers. Specifically, we start with:

Assumption 3 (CEI): *Conditional Effect Ignorability for an instrument z :*

$$\mathbb{E}[Y_1 - Y_0 | D_1^z, D_0^z, x] = \mathbb{E}[Y_1 - Y_0 | x] \quad a.s.$$

A sufficient condition for CEI is:

$$Y_1 = Y_0 + g(x) + \nu$$

where $g(x)$ is any function and ν is mean-independent of (D_1^z, D_0^z) conditional on x . In other words, heterogeneity in average causal effects is solely due to observed covariates.³

To see what this assumption means in a latent-index specification, suppose that:

$$D = 1[h(x, z) > \eta]$$

where η is a random factor involving unobserved costs and benefits of D assumed to be independent of z conditional on x . This latent-index model characterizes potential treatment assignments as:

$$D_0^z = 1[h(x, 0) > \eta] \text{ and } D_1^z = 1[h(x, 1) > \eta].$$

The associated model for potential outcomes is:

$$Y_0 = g_0(x) + \epsilon_0$$

$$Y_1 = g_1(x) + \epsilon_1$$

where the errors here have mean zero conditional on the covariates and instrument. Assuming $h(x, 1) \geq h(x, 0)$ a.s., conditional LATE can be written:

$$\begin{aligned} \Delta^z(x) &= \mathbb{E}[Y_1 - Y_0 | x, h(x, 1) > \eta > h(x, 0)] \\ &= g_1(x) - g_0(x) + \mathbb{E}[\epsilon_1 - \epsilon_0 | x, D_1^z > D_0^z] \end{aligned} \quad (5)$$

³ This is a conditional-on-covariates version of Restriction 2 in Angrist (2004) and is similar to the Frangakis and Rubin (2002) notion of principal stratification, which isolates covariate-defined subpopulations in which selection bias is likely to be minimal.

The CEI assumption implies that $\epsilon_1 - \epsilon_0$ is mean independent of (D_1^z, D_0^z) conditional on \mathbf{x} , so that Equation (5) simplifies to:

$$\Delta^z(\mathbf{x}) = \mathbb{E}[y_1 - y_0 \mid \mathbf{x}] = g_1(\mathbf{x}) - g_0(\mathbf{x}).$$

In the language of Rubin (1977), CEI is a type of ignorability assumption for treatment effects. Given this ignorability, we might wonder whether we need to be concerned about selection bias in the first place. Under CEI, selection bias arises due to correlation between η and ϵ_0 . For example, a latent-index specification compatible with the CEI sets $\epsilon_j = \theta + \xi_j$, $j = 0, 1$, where θ is correlated with η but ξ_j is not.⁴ CEI rules out Roy (1951)-type selection into treatment. In other words, η is assumed to be independent of unobserved gains, denoted by $\epsilon_1 - \epsilon_0$ in the latent-index specification. CEI does not rule out selection bias, but it eliminates an important source of heterogeneity in average causal effects.⁵ Although the empirical importance of Roy selection has yet to be established, the Roy model is an important econometric benchmark. Here, however, we focus on an effort to manage the treatment-effect heterogeneity driven by observable characteristics.

The latent-index specification can be used to formulate a structural justification for the CEI assumption in our empirical application. We start by combining elements of selection models in Olsen (1980) and Vytlacil (2002):

$$\mathbb{E}[\epsilon_1 - \epsilon_0 \mid \eta, \mathbf{x}] = \rho(\mathbf{x})(\eta - 1/2) \text{ and } \eta \mid \mathbf{x}, z \sim U(0, 1).$$

Conditional LATE then becomes:

$$\begin{aligned} \Delta^z(\mathbf{x}) &= g_1(\mathbf{x}) - g_0(\mathbf{x}) + \rho(\mathbf{x})\mathbb{E}[\eta - 1/2 \mid h(\mathbf{x}, 1) > \eta > h(\mathbf{x}, 0), \mathbf{x}] \\ &= g_1(\mathbf{x}) - g_0(\mathbf{x}) + \rho(\mathbf{x})[h(\mathbf{x}, 1) + h(\mathbf{x}, 0) - 1]/2. \end{aligned}$$

For each $\mathbf{x} = \mathbf{x}$, CEI turns on whether $\rho(\mathbf{x}) = 0$.⁶ Following Imbens and Newey (2009), we assume that treatment (i.e., fertility) decisions are based on a comparison of predicted benefits and costs of childbearing. Specifically, women choose to have a third child if:

$$1[h(\mathbf{x}, z) > \eta] = 1\{\lambda(\mathbf{x})\mathbb{E}[y_1 - y_0 \mid \mathbf{x}, \eta] > c(\mathbf{x}, z)\}$$

⁴ See Appendix A for an illustrative limited dependent-variable model with selection bias that satisfies CEI.

⁵ This has been noted by others working with models of this type; see, e.g., Vella and Verbeek (1999).

⁶ Interestingly, this example has the property that LATE is the unconditional average treatment effect when the first stage is symmetric (i.e., $h(\mathbf{x}, 1) = 1 - h(\mathbf{x}, 0)$). See Angrist (2004) for other selection models with this property.

where $\lambda(x)$ is the weight given to outcome gaps; $c(x, z)$ is the expected cost of having a third child; the instrument, z , is a cost-shifter independent of potential outcomes; and η is private information about $y_1 - y_0$, orthogonal to x . Then, $\rho(x)$ is close to zero when either $\lambda(x)$ is small (i.e., labor-supply consequences are of little import), or η matters little given x (e.g., for women with a young second-born who are already at home or for relatively educated women who more easily can afford to pay for child care). Variation in $\rho(x)$ also offers a possible explanation for why CEI might be satisfied for some women but not for others. We consider econometric models in which CEI is partially satisfied in Section 3.4.

3.3 Reweighting LATE

Our covariate-based strategy reweights conditional treatment effects across covariate cells. This is similar to matching estimators designed to control for selection bias; such estimators reweight the conditional-mean function for outcome variables, capturing causal effects when identification is based on a selection-on-observables story (see, e.g., Hahn 1998). In this case, however, we rely on instrumental variables to control for selection bias while using covariates to manage treatment-effect heterogeneity.

Theorem 2 (LATE-Reweight): *Let z be an instrument that satisfies Assumption 3 and let $s^z = s(D_0^z, D_1^z, z)$ be an indicator for any group defined by z . For example, for z -compliers, we have $s^z = D_1^z - D_0^z$; for the treated, $s^z = (1 - z)D_0^z + zD_1^z = D$; for the nontreated, $s^z = (1 - z)(1 - D_0^z) + z(1 - D_1^z) = 1 - D$; and for the entire population, $s^z = 1$. Under Assumption 2 and $\mathbb{E}[|Y|] < \infty$:*

$$\begin{aligned} \mathbb{E}[Y_1 - Y_0 \mid s^z = 1] &= \mathbb{E}[\Delta^z(x) \mid s^z = 1] = \int \Delta^z(x) dF_x(x \mid s^z = 1) \\ &= \int \Delta^z(x) \omega_s^z(x) dF_x(x) \end{aligned}$$

where $\omega_s^z(x) = \mathbb{P}[s^z = 1 \mid x = x] / \mathbb{P}[s^z = 1]$ and $\int \omega_s^z(x) dF_x(x) = 1$.

Proof: By the law of iterated expectations:

$$\begin{aligned} \mathbb{E}[Y_1 - Y_0 \mid s^z = 1] &= \mathbb{E}[\mathbb{E}[Y_1 - Y_0 \mid s^z = 1, x] \mid s^z = 1] \\ &= \mathbb{E}[\mathbb{E}[Y_1 - Y_0 \mid x] \mid s^z = 1] = \mathbb{E}[\Delta^z(x) \mid s^z = 1] \end{aligned}$$

where the second and third equalities follow from CEI. By Bayes rule

$$\int \Delta^z(x) dF_x(x | s^z = 1) = \int \Delta^z(x) \omega_s^z(x) dF_x(x)$$

where $\omega_s^z(x) = \mathbb{P}[s^z = 1 | x = x] / \mathbb{P}[s^z = 1]$ and

$$\int \omega_s^z(x) dF_x(x) = \mathbb{E}\{\mathbb{P}[s^z = 1 | x = x]\} / \mathbb{P}[s^z = 1] = 1$$

by the law of iterated expectations. ■

The LATE-Reweight theorem allows us to go from LATE to the population average treatment effect (ATE), the effect of treatment on the treated (TOT), and the effect of treatment on the nontreated (TNT). The relevant weighting functions, $\omega_\Delta^z(x)$, can be written in terms of observed variables as:

$$\omega_\Delta^z(x) = \frac{\mathbb{E}[D | z = 1, x = x] - \mathbb{E}[D | z = 0, x = x]}{\mathbb{E}[D | z = 1] - \mathbb{E}[D | z = 0]} \tag{6}$$

for effects on z-compliers;

$$\omega_{TOT}^z(x) = \mathbb{E}[D | x = x] / \mathbb{E}[D] \tag{7}$$

for effects on the treated;

$$\omega_{TNT}^z(x) = \mathbb{E}[1 - D | x = x] / \mathbb{E}[1 - D] \tag{8}$$

for effects on the nontreated; and

$$\omega_{ATE}^z(x) = 1 \tag{9}$$

for the population.

3.4 Overidentification

Differences in the observable characteristics of twins and same-sex compliers may explain the difference between the Wald estimates constructed using these two instruments. If so, we can reweight covariate-specific LATEs to go from one to the other. We show this for two instruments, z and w . The difference between the LATE generated by each can be decomposed as:

$$\begin{aligned} \Delta^z - \Delta^w &= \int [\Delta^z(x) - \Delta^w(x)] \omega_\Delta^z(x) dF_x(x) \\ &\quad + \int \Delta^w(x) [\omega_\Delta^z(x) - \omega_\Delta^w(x)] dF_x(x). \end{aligned} \tag{10}$$

The first term reflects differences in conditional LATEs between z -compliers and w -compliers, while the second term captures differences in complier characteristics. If z and w satisfy CEI, the following *compatibility condition* holds:

$$\Delta^w(x) = \Delta^z(x) \text{ a.s.} \tag{11}$$

and the first term of the decomposition (10) is zero.⁷ In other words, CEI for both instruments implies that the instruments w and z satisfy overidentifying restrictions conditional on x . The following theorem shows that in such cases, we can use the distribution of compliers for a hypothetical instrument to construct the treatment effects that might be generated by instruments other than those we have.

Theorem 3 (LATE-Overid): *Let z and w be two instruments that satisfy Assumptions 2 and 3. Let $\Delta^w(x)$ be defined as in Equation (3) using the instrument w . Let $s^w = s(D_0^w, D_1^w, w)$ be an indicator for any group defined by the instrument w with corresponding potential treatment assignments (D_0^w, D_1^w) . If $\mathbb{E}[|Y|] < \infty$:*

$$\mathbb{E}[Y_1 - Y_0 \mid s^w = 1] = \int \Delta^w(x)\omega_s^w(x)dF_x(x) = \int \Delta^z(x)\omega_s^w(x)dF_x(x). \tag{12}$$

Proof: The first equality in Equation (12) follows from Theorem 2 applied to w . To establish the second equality, note that the law of iterated expectations and CEI for z give

$$\mathbb{E}[Y_1 - Y_0 \mid s^w = 1] = \mathbb{E}[\mathbb{E}[Y_1 - Y_0 \mid x] \mid s^w = 1] = \mathbb{E}[\Delta^z(x) \mid s^w = 1].$$

The result then follows by Bayes rule. ■

The LATE-Overid theorem allows us to determine whether differences in the distribution of complier covariates are enough to explain differences in IV estimates across instruments. If so, it seems fair to say that the underlying covariate-specific results have predictive value for subpopulations defined by these covariate values and therefore some claim to external validity.

⁷ A weaker alternative to CEI that implies the compatibility condition is $\mathbb{E}[Y_1 - Y_0 \mid D_1^z > D_0^z, x] = \mathbb{E}[Y_1 - Y_0 \mid D_1^w > D_0^w, x]$ a.s.; in other words, equality of conditional LATEs only for compliers. We find the stronger CEI assumption more appealing because it seems hard to imagine a mechanism or model that generates this more limited sort of independence without generating full CEI.

CEI also implies that the conditional LATEs, $\Delta^z(x)$, are overidentified if we observe both z and w . We can therefore construct more precise estimators for conditional LATEs using cell-by-cell GMM procedures that use both z and w to form moment conditions. Moreover, GMM overidentification tests can be used to statistically assess the LATE compatibility condition, (11). Comparisons of estimates of the two expressions on the right-hand side of Equation (12) also serve as a test of compatibility.

In practice, of course, covariates need not account fully for the difference between two LATEs. For some covariate values, there may be a good match, while for others, CEI fails. A rationale for partial fulfillment of CEI emerges in our latent-index example. In this context, partial CEI is like saying that some types of women select on gains while others do not. The values for which CEI is satisfied define a subpopulation for which heterogeneous treatment effects can be understood to be solely a function of observable characteristics. For this subpopulation, we can define an average causal effect for which the predictive value of IV estimates is likely to be especially high.

Definition 2 (CATE): The Compatible Average Treatment effect is:

$$\Delta^{z,w} := \int \Delta^z(x) dF_x(x \mid \Delta^z(x) = \Delta^w(x)).$$

If the compatibility condition, $\Delta^z(x) = \Delta^w(x)$, holds for all values of x , CATE is ATE. However, if compatibility holds for only some values of x , CATE is ATE for the subpopulation defined by these values.

The compatible subpopulation may be of interest for a number of reasons. First, we might be looking for the largest and most representative subpopulation for which a given set of IV estimates has predictive value. This might be, for example, a compatible subset of the treated. Second, we may be interested in constructing a precise estimate of covariate-specific treatment effects. CATE defines a subpopulation where it is possible to use two instruments to construct more precise estimates.

4.0 Estimation and Inference

We assume that the effects of interest are to be estimated in a random sample of size n .

Assumption 4 (Sampling): $\{R_i = (Y_i, D_i, X_i, Z_i, W_i), i = 1, \dots, n\}$ are *i.i.d.* observations of the random vector $R = (Y, D, X, Z, W)$.

We also assume that the covariates, x , take on a finite and fixed number of values. The education and age covariates in the empirical example satisfy this condition. Generalization to continuous covariates seems straightforward but requires additional technical machinery. For example, it seems likely that with continuous covariates, we should allow for a gradual failure of CEI as opposed to discrete cutoffs. We therefore leave this extension for future work.

Assumption 5 (Discrete Covariates): For a finite set, $\mathcal{X} = \{x_1, \dots, x_K\}$, $\mathbb{P}[x \in \mathcal{X}] = 1$.

The effects in Theorems 2 and 3 can be written as follows:

$$\Delta_{s^u}^L = \mathbb{E}[\Delta^L(x)\omega_s^U(x)], \quad \omega_s^U(x) = \mathbb{P}[s^U = 1 \mid x = x] / \mathbb{P}[s^U = 1]$$

where $L = U = z$ for Theorem 2 and $L = z$ and $U = w$ or vice versa for Theorem 3. More generally, superscript L indexes the group where conditional LATEs are obtained and s^U is an indicator for the group with the covariate distribution of interest, defined using instrument U .

Estimation is straightforward in our finite-dimensional setting. We replace expectations \mathbb{E} and probabilities \mathbb{P} by empirical analogs \mathbb{E}_n and \mathbb{P}_n , where $\mathbb{E}_n[g(\mathbf{R})] = n^{-1} \sum_{i=1}^n g(\mathbf{R}_i)$ and $\mathbb{P}_n[g(\mathbf{R}) \in C] = n^{-1} \sum_{i=1}^n 1[g(\mathbf{R}_i) \in C]$ for any function g and set C . For conditional expectations and probabilities, we let $\mathbb{E}_n[\cdot \mid x = x, U = u]$ and $\mathbb{P}_n[\cdot \mid x = x, U = u]$ denote empirical analogs in the covariate cell, where $x = x$ and $U = u$ for $U \in \{z, w\}$ and $u \in \{z, w\}$. This gives:

$$\hat{\Delta}_{s^u}^L = \mathbb{E}_n[\hat{\Delta}^L(x)\hat{\omega}_s^U(x)], \quad \hat{\omega}_s^U(x) = \mathbb{P}_n[s^U = 1 \mid x = x] / \mathbb{P}_n[s^U = 1] \quad (13)$$

where $\hat{\Delta}^L(x)$ is any consistent estimator of $\Delta^L(x)$. For example, $\hat{\Delta}^L(x)$ can be the Wald estimator with instrument $L \in \{z, w\}$ in cell $x = x$, or 2SLS using both z and w as instruments in cell $x = x$. For treated, nontreated, and the entire population, the indicator s^U is observed, so construction of the empirical $\hat{\omega}_s^U(x)$ is straightforward. For compliers, we can estimate $\omega_s^U(x)$ using the sample analog of Equation (6):

$$\hat{\omega}_{\Delta}^U(x) = \frac{\mathbb{E}_n[D \mid U = 1, x = x] - \mathbb{E}_n[D \mid U = 0, x = x]}{\mathbb{E}_n[D \mid U = 1] - \mathbb{E}_n[D \mid U = 0]}, \quad U \in \{z, w\}. \quad (14)$$

Consistency of $\hat{\Delta}_{s^u}^L$ follows from the law of large numbers and the Slutsky theorem.

Theorem 4 (Consistency): Let z and w be two instruments that satisfy Assumptions 2 and 3. Under Assumptions 4 and 5, and $\mathbb{E}[|Y|] < \infty$:

$$\hat{\Delta}_{s^u}^L = \mathbb{E}_n[\hat{\Delta}^L(x)\hat{\omega}_s^u(x)] \rightarrow_p \Delta_{s^u}^L = \mathbb{E}[Y_1 - Y_0 \mid s^w = 1]; L, U \in \{z, w\}$$

where $\hat{\omega}_s^u(x)$ and $\hat{\Delta}^L(x)$ are any consistent estimators of $\omega_s^u(x)$ and $\Delta^L(x)$, for all $x \in \mathcal{X}$.

The estimators developed from Theorems 2 and 3 are smooth functions of GMM-type estimators and are therefore asymptotically normal under general conditions. The following result uses the delta method to characterize the relevant limiting distributions. In particular, we show that $(\hat{\Delta}_{s^u}^z, \hat{\Delta}_{s^u}^w)$ are asymptotically jointly normal, allowing us to draw inferences about the effects of interest and to test some of the implications of CEI. Let $p_k = \mathbb{P}(x = x_k)$, $\vartheta_s^u(x) = \mathbb{P}[s^u = 1 \mid x = x]$ (i.e., the numerator of $\omega_s^u(x)$), and let $\hat{\vartheta}_s^u(x)$ be an estimator of $\vartheta_s^u(x)$.

Theorem 5 (Asymptotic Distribution): Let z and w be two instruments that satisfy Assumptions 2 and 3. Assume that $p_k > \varepsilon > 0$ for all $k \in \{1, \dots, K\}$, and for $k \in \{1, \dots, K\}$ and $u \in \{z, w\}$:

$$\begin{aligned} & \sqrt{n} \begin{pmatrix} \hat{\vartheta}_s^u(x_k) - \vartheta_s^u(x_k) \\ \hat{\Delta}^z(x_k) - \Delta^z(x_k) \\ \hat{\Delta}^w(x_k) - \Delta^w(x_k) \end{pmatrix} \\ & \rightarrow_d Z_k^u \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} V_s^u(x_k)/p_k & C_{s\Delta}^{uz}(x_k)/p_k & C_{s\Delta}^{uw}(x_k)/p_k \\ C_{s\Delta}^{uz}(x_k)/p_k & V_{\Delta}^z(x_k)/p_k & C_{\Delta\Delta}^{zw}(x_k)/p_k \\ C_{s\Delta}^{uw}(x_k)/p_k & C_{\Delta\Delta}^{zw}(x_k)/p_k & V_{\Delta}^w(x_k)/p_k \end{bmatrix} \right) \end{aligned}$$

where (Z_1^u, \dots, Z_K^u) are independent. Under Assumptions 9 and 10, for $u \in \{z, w\}$:

$$\sqrt{n} \begin{pmatrix} \hat{\Delta}_{s^u}^z - \Delta_{s^u}^z \\ \hat{\Delta}_{s^u}^w - \Delta_{s^u}^w \end{pmatrix} \rightarrow_d \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} V_{s^u}^z & C_{s^u}^{zw} \\ C_{s^u}^{zw} & V_{s^u}^w \end{bmatrix} \right)$$

where:

$$\begin{aligned} V_{s^u}^L &= \sum_{k=1}^K p_k [\omega_s^u(x_k)^2 + \tilde{V}_s^u(x_k) + 2\omega_s^u(x_k)\tilde{C}_{s\Delta}^{uL}(x_k)][\Delta^L(x_k) - \Delta_{s^u}^L]^2 \\ &+ \sum_{k=1}^K p_k \omega_s^u(x_k)^2 V_{\Delta}^L(x_k), \end{aligned}$$

for $L \in \{z, w\}$, $\tilde{V}_s^U(x_k) = V_s^U(x_k) / [\sum_{k=1}^K p_k \vartheta_s^U(x_k)]^2$, $\tilde{C}_{s\Delta}^{UL}(x_k) = C_{s\Delta}^{UL}(x_k) / \sum_{k=1}^K p_k \vartheta_s^U(x_k)$, and

$$C_{s^v}^{z^w} = \sum_{k=1}^K p_k [\omega_s^U(x_k)^2 + \tilde{V}_s^U(x_k) + \omega_s^U(x_k) \{ \tilde{C}_{s\Delta}^{UZ}(x_k) + \tilde{C}_{s\Delta}^{UW}(x_k) \}] [\Delta^Z(x_k) - \Delta_{s^v}^Z] [\Delta^W(x_k) - \Delta_{s^v}^W] + \sum_{k=1}^K p_k \omega_s^U(x_k)^2 C_{\Delta\Delta}^{z^w}(x_k).$$

Proof: Let $\hat{p}_k = \mathbb{P}_n(x = x_k)$. By a standard Central Limit Theorem for multinomial sequences, $\sqrt{n}(\hat{p}_1 - p_1, \dots, \hat{p}_K - p_K)$ converges in distribution to a multivariate normal with zero mean, variances $p_k(1 - p_k)$ and covariances $-p_k p_j$, for $k, j = 1, \dots, K$, $k \neq j$. Write, for $L \in \{z, w\}$:

$$\hat{\Delta}_{s^v}^L = \frac{\sum_{k=1}^K \hat{p}_k \hat{\vartheta}_s^U(x_k) \hat{\Delta}^L(x_k)}{\sum_{k=1}^K \hat{p}_k \hat{\vartheta}_s^U(x_k)} \text{ and } \Delta_{s^v}^L = \frac{\sum_{k=1}^K p_k \vartheta_s^U(x_k) \Delta^L(x_k)}{\sum_{k=1}^K p_k \vartheta_s^U(x_k)}.$$

Let $\vec{\pi} = (\pi_1, \dots, \pi_K)$, $\vec{v} = (v_1, \dots, v_K)$, and $\vec{\delta} = (\delta_1, \dots, \delta_K)$. If $\sum_k \pi_k v_k \neq 0$, the function $f(\vec{\pi}, \vec{v}, \vec{\delta}) = \sum_k \pi_k v_k \delta_k / \sum_k \pi_k v_k$ is continuously differentiable in $(\vec{\pi}, \vec{v}, \vec{\delta})$ with partial derivatives:

$$\frac{\partial f(\vec{\pi}, \vec{v}, \vec{\delta})}{\partial \pi_k} = \tilde{v}_k \tilde{\delta}_k, \quad \frac{\partial f(\vec{\pi}, \vec{v}, \vec{\delta})}{\partial v_k} = \frac{\pi_k \tilde{\delta}_k}{\sum_k \pi_k v_k}, \quad \frac{\partial f(\vec{\pi}, \vec{v}, \vec{\delta})}{\partial \delta_k} = \pi_k \tilde{v}_k$$

for $\tilde{v}_k = v_k / \sum_k \pi_k v_k$ and $\tilde{\delta}_k = \delta_k - \sum_k \pi_k \tilde{v}_k \delta_k$.

Set $\pi_k = p_k$, $v_k = \vartheta_s^U(x_k)$, and $\delta_k = \Delta^L(x \vartheta_k)$. The result then follows using the delta method. ■

The joint normality assumption for the components of our reweighting estimators holds under standard regularity conditions. In particular, this follows for the estimators of the weighting function by the Central Limit Theorem for binary sequences. For IV, GMM, and other moment-based estimators of the conditional LATEs, such as generalized empirical likelihood, existence of second moments (i.e., $\mathbb{E}[Y^2] < \infty$) is sufficient. The first term in the expressions for $V_{s^v}^z$, $V_{s^v}^w$, and $C_{s^v}^{z^w}$ reflects sampling variation due to the estimation of the covariate-cell probabilities, p_k , and the weighting functions, whereas the second term arises from the estimation of conditional LATEs. The first term is zero if, for example, the conditional LATEs are constant (i.e., $\Delta^L(x) = \Delta_{s^v}^L$ for all $x \in \mathcal{X}$).

In practice, there are two routes to asymptotic inference. We can estimate asymptotic distributions analytically using sample analogs or approximate them numerically by resampling or simulation. We use bootstrap methods

in the empirical application. Resampling methods are convenient and save us from having to estimate complicated analytical formulas for asymptotic variances and covariances. Consistency of a bootstrap approximation to the distributions of our reweighting estimators follows from Hall and Horowitz (1996), Hahn (1996), and Brown and Newey (2002) theorems for GMM, and application of the delta method for the bootstrap (see, e.g., Theorem 23.5 in van der Vaart 1998).

There are many ways to bootstrap. We use the empirical likelihood (EL) bootstrap proposed by Brown and Newey (2002) for GMM estimators. This method resamples from the empirical likelihood distribution (ELD) that imposes the moment conditions in the sample, instead of from the empirical distribution (ED). Let $\mathbb{E}[g(\mathbf{r}, \theta)] = 0$ be the moment conditions that define the conditional LATEs, weighting functions, and effects of interest, where θ includes all unknown parameters. The ELD $(\hat{\pi}_1, \dots, \hat{\pi}_n)$ is the solution to:

$$\max_{\pi_1, \dots, \pi_n} \sum_{i=1}^n \ln(\pi_i), \text{ s.t. } \sum_{i=1}^n \pi_i g(\mathbf{r}_i, \hat{\theta}) = 0, \sum_{i=1}^n \pi_i = 1, \pi_i \geq 0$$

where $\hat{\theta}$ is the EL or another consistent estimator of θ .⁸ ELD therefore is the closest to ED in terms of Kullback–Leibler distance. ELD and ED are equal in exactly identified models, but they generally differ under overidentification. In practice, the difference between the EL bootstrap and the standard nonparametric bootstrap is that the former resamples from the data with probabilities $\hat{\pi}_i$ instead of $1/n$.

Consistent estimation of CATE requires that we condition on the unobservable events $\{\Delta^z(x) = \Delta^w(x)\}$. In finite samples, we never have $\hat{\Delta}^z(x) = \hat{\Delta}^w(x)$. We therefore use cell-by-cell overidentification tests to find compatible values of x . Instead of discarding cells that fail the identification test for some small significance level, we reweight estimates of conditional LATE by a decreasing function of the overidentification test statistic. Letting $J(x)$ denote the overidentification test statistic for the instruments z and w in the cell $x = x$, the resulting estimator of CATE is:

$$\hat{\Delta}^{z,w} = \mathbb{E}_n[\hat{\Delta}^{z,w}(x)\hat{\omega}_{CATE}(x)],$$

$$\hat{\omega}_{CATE}(x) = \exp\{-J(x)/a_n(x)\}/\mathbb{E}_n[\exp\{-J(x)/a_n(x)\}] \quad (15)$$

where $\hat{\Delta}^{z,w}(x)$ is the GMM estimate in cell $x = x$ that uses z and w as instruments or any other moment estimator (e.g., 2SLS), and $a_n(x)$ is a sequence such that $a_n(x) \rightarrow \infty$ and $a_n(x) = o(n)$, for $x \in \mathcal{X}$.

⁸ In the empirical application, we use the EL estimator of θ to obtain the ELD.

The sequences $a_n(x)$ guarantee the consistency of the reweighting estimator (15) for CATE. These sequences play a role similar to the penalty terms used in Andrews (1999) to obtain consistent model-selection procedures for GMM estimators. To formally establish consistency, it is convenient to introduce additional notation. Let $\mathcal{X}_0 = \{x \in \mathcal{X} : \Delta^z(x) = \Delta^w(x)\}$ denote the set of covariate values that satisfy the compatibility condition, with complement $\bar{\mathcal{X}}_0 = \{x \in \mathcal{X} : \Delta^z(x) \neq \Delta^w(x)\}$.

Theorem 6 (CATE Consistency): *Let z and w be two instruments that satisfy Assumption 2. Let $a_n(x)$ be sequences such that $a_n(x) \rightarrow \infty$ and $a_n(x) = o(n)$ for all $x \in \mathcal{X}$. Assume that $\hat{\Delta}^{z,w}(x) \rightarrow_p \Delta^z(x)$ and $J(x) = O_p(1)$ for all $x \in \mathcal{X}_0$, and $\hat{\Delta}^{z,w}(x) = O_p(1)$ and $J(x) = O_p(n)$ for all $x \in \bar{\mathcal{X}}_0$. Under Assumptions 4 and 5, $\mathbb{P}\{x \in \mathcal{X}_0\} > 0$, and $\mathbb{E}[|Y|] < \infty$:*

$$\hat{\Delta}^{z,w} = \mathbb{E}_n[\hat{\Delta}^{z,w}(x)\hat{\omega}_{CATE}(x)] \rightarrow_p \Delta^{z,w} = \mathbb{E}[\Delta^z(x) \mid x \in \mathcal{X}_0]$$

where $\hat{\omega}_{CATE}(x) = \exp\{-J(x)/a_n(x)\}/\mathbb{E}_n[\exp\{-J(x)/a_n(x)\}]$.

Proof: We write:

$$\begin{aligned} \hat{\Delta}^{z,w} &= \sum_{k=1}^K \hat{p}_k \hat{\Delta}^{z,w}(x_k) \hat{\omega}_{CATE}(x_k) \text{ and} \\ \Delta^{z,w} &= \sum_{k=1}^K p_k \Delta^z(x_k) 1\{x_k \in \mathcal{X}_0\} / \sum_{k=1}^K p_k 1\{x_k \in \mathcal{X}_0\}. \end{aligned}$$

By the Law of Large Numbers, $\hat{p}_k \rightarrow_p p_k$. For $x_k \in \mathcal{X}_0$, $\hat{\Delta}^{z,w}(x_k) \rightarrow_p \Delta^z(x_k)$, $J(x_k) = O_p(1)$ and $\exp\{-J(x_k)/a_n(x_k)\} \rightarrow_p 1$. For $x_k \in \bar{\mathcal{X}}_0$, $\hat{\Delta}^{z,w}(x_k) = O_p(1)$, $J(x_k) = O_p(n)$ and $\exp\{-J(x_k)/a_n(x_k)\} \rightarrow_p 0$. Hence, $\hat{\Delta}^{z,w}(x_k) \exp\{-J(x_k)/a_n(x_k)\} \rightarrow_p \Delta^z(x) 1\{x_k \in \mathcal{X}_0\}$.

The result follows by the Slutsky Theorem, noting that:

$$\begin{aligned} \mathbb{E}_n[\exp\{-J(x)/a_n(x)\}] &= \sum_{k=1}^K \hat{p}_k \exp\{-J(x_k)/a_n(x_k)\} \\ &\rightarrow_p \sum_{k=1}^K p_k 1\{x_k \in \mathcal{X}_0\} = \mathbb{P}\{x \in \mathcal{X}_0\} > 0. \end{aligned}$$

For both $x \in \mathcal{X}_0$ and $x \in \bar{\mathcal{X}}_0$, the convergence-rate assumptions for $\hat{\Delta}^{z,w}(x)$ and $J(x)$ in the statement of the theorem are satisfied by GMM-type

estimators under standard regularity conditions for asymptotic normality (see, e.g., Assumption 1 in Andrews 1999).

Although consistency of our CATE estimator is relatively easy to show, inference is challenging. As with the reweighting estimators in Theorem 5, the limiting distribution of $\hat{\Delta}^{z,w}$ depends on the limiting distributions of the weighting functions. Here, however, the limiting distribution of $\hat{\omega}_{CATE}(x)$ converges at nonstandard rates, complicating the analysis. Slow convergence in this case is a by-product of the need for a term like $a_n(x)$ in the weighting function to ensure consistency and the fact that, in practice, we choose this to grow more slowly than \sqrt{n} .⁹ Moreover, convergence to the limiting distribution cannot be uniform in the data-generating process because the CATE estimator implicitly conditions on a pretest (see, e.g., Leeb and Pötscher 2008). In a related setting, Andrews and Guggenberger (2009) addressed a pretest problem using subsampling. Subsampling is imprecise in our application due to small cell sizes. We have no easy solution in this case other than to caution that convergence to the relevant limiting distribution may be slow and to conjecture the pointwise validity of bootstrap methods.

A second and less serious inference complication arises from the fact that the EL bootstrap imposes CEI at all covariate values, while the purpose of CATE is to allow deviations from CEI. In the empirical application, therefore, we supplement the standard errors obtained from the EL bootstrap with standard errors obtained by a nonparametric bootstrap that does not impose the compatibility condition on the bootstrap data generating process.

5.0 Results

Our empirical exploration of covariate-reweighting focuses on a categorical representation of second-child age and mother's schooling. As shown in Table 2, these covariates are strongly related to compliance probabilities. A younger second child in the household reduces the likelihood of a third birth, if only because less time has passed since the birth of the second. Education matters because college-educated women are less likely to choose to have a third child than less-educated women. Twins compliers therefore are relatively more likely to have a young second-born and to be highly educated. Sex-composition compliers, by contrast, are relatively unlikely

⁹ As in Crump, Hotz, Imbens, and Mitnik (2009), we could simplify here by doing inference conditional on the sample. We do not take this route because we are interested in predictive population inference as opposed to sample-specific causal inference.

to be college graduates or to have a third child soon after the birth of the second. Labor supply effects also are likely to vary with second-child age and mother's schooling. The birth of a third child has little effect on the work behavior of a woman with a young second-born who is at home anyway. Likewise, a relatively educated woman should be affected less by the birth of a child than other women because for women earning higher wages, it makes sense to pay for child care in the market. Differences in twins and same-sex complier subpopulations therefore might account for the fact that twins instruments generate a smaller effect of childbearing than sex-composition instruments.

In an effort to see whether this conjecture is substantiated empirically, Table 3 reports IV estimates in each of 12 cells defined by second-child age and mother's schooling. The age categories are less than or equal to 4 years, greater than 4 and less than or equal to 8, and greater than 8. The schooling categories are high school dropout, high school graduate, some college, and college graduate. The first column of the table reports the probability mass function in the contingency table generated by these categories. The next two columns describe the distribution of covariates for the treated and untreated, relative to the entire population. Women who do and do not have a third child are clearly very different.

Cell-specific IV estimates are fairly imprecise, as shown in Columns 4 and 6 of Table 3. From these noisy cell-by-cell estimates alone, it is difficult to see how the causal effect of a third birth varies with individual characteristics. A clearer pattern emerges, however, once the cells are "weighted-up," a point we return to in Table 4.

Table 3 also reports estimates of:

$$\omega_s^z(x) = \frac{\mathbb{P}[s^z = 1 \mid x = x]}{\mathbb{P}[s^z = 1]}$$

the weighting function for twins compliers (Column 5) and same-sex compliers (Column 7). This is the ratio of the relevant first stage in the cell to the overall first stage of Equation (6). The distribution of twins and same-sex compliers over cells clearly is different from the cell distribution in the random sample. As suggested by the summary comparisons in Table 2, the complier distributions for the two instruments also are very different from one another. Specifically, twins compliers are much more likely to have a young second-born child, while few same-sex compliers are in this group. Twins compliers also are relatively educated, while the schooling gradient in compliance probabilities for same-sex is less pronounced.

Table 3. *LATE Decompositions*

Covariate		Covariate pmf			IV Estimates and Weighting Functions							
					Twins Instrument		Same-Sex Instrument		Both Instruments			
Age	Education	$P(X)$ (1)	$P(X D = 1)/P(X)$ (2)	$P(X D = 0)/P(X)$ (3)	$\Delta^z(X)$ (4)	$\omega_{\Delta}^z(X)$ (5)	$\Delta^w(X)$ (6)	$\omega_{\Delta}^w(X)$ (7)	$\Delta^{z,w}(X)$ (8)	J-p value (9)	$\omega_{CATE}(X)$ (10)	
A. Weeks worked												
[0,4]	HS drop	0.06	0.65	1.24	-4.63 (1.84)	1.23	-1.26 (5.49)	0.77	-4.35 (1.60)	0.55	1.45	
	HS grad	0.15	0.46	1.36	-4.24 (1.15)	1.36	-2.66 (5.00)	0.59	-4.17 (1.04)	0.76	1.99	
	Some col	0.06	0.46	1.36	-3.79 (1.69)	1.36	-4.93 (6.86)	0.66	-3.86 (1.59)	0.87	2.14	
	Col grad	0.05	0.41	1.40	-5.36 (1.99)	1.40	-2.21 (9.77)	0.54	-5.24 (1.84)	0.75	1.95	
(4,8]	HS drop	0.07	1.29	0.81	-4.43 (2.65)	0.81	-9.35 (3.36)	1.20	-6.29 (2.06)	0.25	0.47	
	HS grad	0.17	0.93	1.05	-3.07 (1.52)	1.05	-5.59 (2.09)	1.38	-3.94 (1.23)	0.33	0.79	
	Some col	0.07	0.91	1.06	-1.02 (2.22)	1.06	-7.39 (3.91)	1.13	-2.59 (1.93)	0.16	0.21	
	Col grad	0.04	0.89	1.08	-1.52 (2.63)	1.08	-2.88 (5.97)	0.94	-1.72 (2.32)	0.83	2.09	
(8+]	HS drop	0.10	1.79	0.47	0.29 (4.47)	0.47	-12.04 (4.74)	0.80	-5.53 (3.24)	0.06	0.04	
	HS grad	0.17	1.40	0.73	-2.41 (2.22)	0.73	-9.76 (2.25)	1.30	-6.15 (1.60)	0.02	0.01	
	Some col	0.06	1.33	0.78	-4.40 (3.55)	0.78	-4.72 (4.54)	1.10	-4.52 (2.82)	0.96	2.20	
	Col grad	0.02	1.15	0.90	6.78 (4.99)	0.90	17.90 (9.01)	1.12	9.48 (4.17)	0.28	0.44	

		B. Employment									
[0,4]	HS drop	0.06	0.65	1.24	-0.154 (0.048)	1.23	-0.035 (0.143)	0.77	-0.143 (0.043)	0.43	1.15
	HS grad	0.15	0.46	1.36	-0.081 (0.027)	1.36	-0.123 (0.117)	0.59	-0.083 (0.026)	0.73	2.14
	Some col	0.06	0.46	1.36	-0.089 (0.038)	1.36	-0.035 (0.157)	0.66	-0.086 (0.037)	0.74	2.15
	Col grad	0.05	0.41	1.40	-0.130 (0.047)	1.40	-0.023 (0.231)	0.54	-0.125 (0.046)	0.65	1.90
(4.8]	HS drop	0.07	1.29	0.81	-0.140 (0.064)	0.81	-0.150 (0.082)	1.20	-0.144 (0.051)	0.92	2.42
	HS grad	0.17	0.93	1.05	-0.081 (0.034)	1.05	-0.156 (0.046)	1.38	-0.108 (0.028)	0.19	0.39
	Some col	0.07	0.91	1.06	-0.034 (0.048)	1.06	-0.161 (0.084)	1.13	-0.066 (0.042)	0.19	0.34
	Col grad	0.04	0.89	1.08	0.075 (0.061)	1.08	-0.202 (0.134)	0.94	0.028 (0.054)	0.06	0.03
(8+]	HS drop	0.10	1.79	0.47	0.047 (0.101)	0.47	-0.188 (0.106)	0.80	-0.064 (0.073)	0.11	0.13
	HS grad	0.17	1.40	0.73	-0.066 (0.046)	0.73	-0.167 (0.047)	1.30	-0.117 (0.033)	0.13	0.20
	Some col	0.06	1.33	0.78	-0.110 (0.071)	0.78	-0.023 (0.092)	1.10	-0.075 (0.058)	0.47	1.32
	Col grad	0.02	1.15	0.90	0.034 (0.098)	0.90	0.224 (0.171)	1.12	0.082 (0.081)	0.33	0.68

Note: Standard errors for estimates in parentheses. The p -value for the joint J-statistic for all covariate values is 0.25 for weeks and 0.29 for LFP. The sample size is 394,840.

Table 4. *Reweighting LATE*

Population Represented (effect)	Instrument Used for $\Delta(X)$	Weighting Function $\omega(X)$	Weeks Worked		Employment	
			Estimate (1)	t for Diff (2)	Estimate (3)	t for Diff (4)
Twins compliers (LATE)	twins	twins	-3.15 (0.62)	1.32	-0.075 (0.014)	0.93
	same-sex		-2.71 (0.81)		-0.068 (0.018)	
Same-sex compliers (LATE)	same-sex	same-sex	-6.30 (1.15)	1.44	-0.131 (0.026)	0.77
	twins		-5.08 (1.58)		-0.115 (0.037)	
Everyone (ATE)	twins	1	-2.84 (0.76)	1.99	-0.067 (0.017)	1.58
	same-sex		-5.88 (1.35)		-0.123 (0.031)	
Treated (TOT)	twins	$P(X D = 1)/P(X)$	-2.38 (1.07)	2.85	-0.056 (0.024)	2.14
	same-sex		-7.08 (1.28)		-0.136 (0.029)	
Nontreated (TNT)	twins	$P(X D = 0)/P(X)$	-3.15 (0.62)	1.15	-0.075 (0.014)	1.02
	same-sex		-5.08 (1.58)		-0.115 (0.037)	
Compatible (CATE)	twins	$\exp[-12^* J(X)/n(X)]$	-3.80 (0.80)	0.18 [0.09]	-0.099 (0.018)	0.17 [0.08]
	same-sex		-3.66 (1.01)		-0.095 (0.023)	
			[0.96]	[0.021]		
	twins, same-sex		-4.00 (0.77)	-0.101 (0.017)		
		[0.62]	[0.014]			

Note: Standard errors for estimates in parentheses. *T*-statistics are for the difference between same-sex and twins estimates. Standard errors and *t*-statistics obtained by Brown and Newey (2002) GMM bootstrap with 1,000 repetitions. Standard errors and *t*-statistics reported in brackets were obtained by nonparametric bootstrap with 1,000 repetitions. The sample size is 394,840.

Table 4 reports LATE-reweighted estimates combining conditional LATEs and weighting functions for twins compliers and same-sex compliers. The first and fifth row weight conditional LATE estimates to produce an overall unconditional effect using the sample analog of Equation (4). For example, in the first row, cell-by-cell twins estimates are weighted across cells using the probability of twins compliance in each cell. This is close to the Wald estimate using twins instruments in Table 1; compare -3.15 to -3.28 and -0.075 to -0.076 . Likewise, the marginalized same-sex estimates using the probability of same-sex compliance are -6.30 for weeks worked and -0.132 for employment, close to the corresponding Wald estimates. The fact that weighted conditional LATE comes out close to Wald estimates that ignore covariates suggests that covariates are unnecessary for identification of LATE, although they are helpful for extrapolation.

Reweighting conditional same-sex estimates using twins weights brings these estimates remarkably close to the marginalized twins estimates. Compare, for example, the estimate of the effect on weeks worked of -3.15 using twins in each cell to the estimate of -2.71 using same-sex instruments in each cell. Reweighting conditional twins and same-sex estimates using same-sex weights also produces a good match. In this case, the estimate of -6.30 using same-sex instruments in each cell can be compared to an estimate of -5.08 using twins instruments in each cell in the fifth and seventh rows of Table 4. In other words, differing LATE estimates generated using twins and same-sex instruments can be reconciled by reweighting covariate-specific estimates using a set of common compliance weights.

This reconciliation is an encouraging finding which suggests that external validity is an attainable goal in this context. On the other hand, Table 4 also shows that ATE and TOT using twins and same-sex instruments are not well matched. This is disappointing because, by the same argument that reconciles twins and same-sex estimates of LATE, we should be able to generate similar estimates of ATE and TOT using either instrument to construct cell-specific IV estimates. At the same time, the match for TNT is not as bad as that for ATE and TOT. The fact that some parameters can be matched more easily across instruments than others suggests that a few poorly matched cells are what drives the cross-instrument imbalance in estimates of ATE and TOT – that is, a handful of cells generate estimates that depend on the instrument.

CATE, shown at the bottom of Table 4, solves this problem, and produces a good match for the average treatment effect in compatible cells by down-weighting cells where CEI is most at odds with the data. CATE estimates

are -3.80 to -3.66 for effects on weeks worked and -0.099 to -0.095 for effects on employment. The estimates of CATE reported in Table 4 set $a_n(x) = \log n(x)/K$, where $n(x)$ is the number of individuals with $x = x$ and K is the number of covariate cells. This relatively slow normalization works well in our application because the cell-level overidentification test is never very large. Still, by applying the most weight to cells in which CEI appears to be satisfied, CATE generates remarkably similar estimates of the population ATE, whether the underlying cell-level IV estimates use twins or same-sex instruments.

The match generated by CATE is for a subpopulation that need not be representative of the entire population of interest. Figure 1 describes the compatible subpopulation by plotting the weighting functions used by different estimators. CATE essentially discards the two low-education cells for women with an older second-born, with weights something like the histogram for twins compliers. Because twins instruments induce one-sided noncompliance, the group of twins compliers is the same as the nontreated population (Angrist, Lavy, and Schlosser 2010). Thus, the compatible subpopulation is similar to the population of nontreated women, although not exactly the same (Table 3 shows that reweighting to estimate effects on TNT is not very successful). Still, like the nontreated population, the subpopulation for which IV estimates of the effects of childbearing appear to have strong predictive value consists mainly of mothers who have a young second-born child and are somewhat more likely to have gone to college.

6.0 Summary and Directions for Further Work

In the LATE framework, differences in IV estimates need not signal a failure of the exclusion restriction. Rather, these differences may be attributable to differences in the type of people who are affected by the underlying experiments implicit in any IV identification strategy. At the same time, we often hope to use one set of IV estimates to predict causal effects in settings other than the one generating the estimates; the question of external validity turns on our ability to do this reliably.

Here, we begin with the idea that differences in IV estimates of LATE for the same causal relation might be driven by a combination of treatment-effect heterogeneity across covariate cells and differences in covariate distributions for instrument-specific-compliant subpopulations. Limiting heterogeneity across cells and instruments to be a function solely of observed characteristics, we can reweight one set of IV estimates to generate

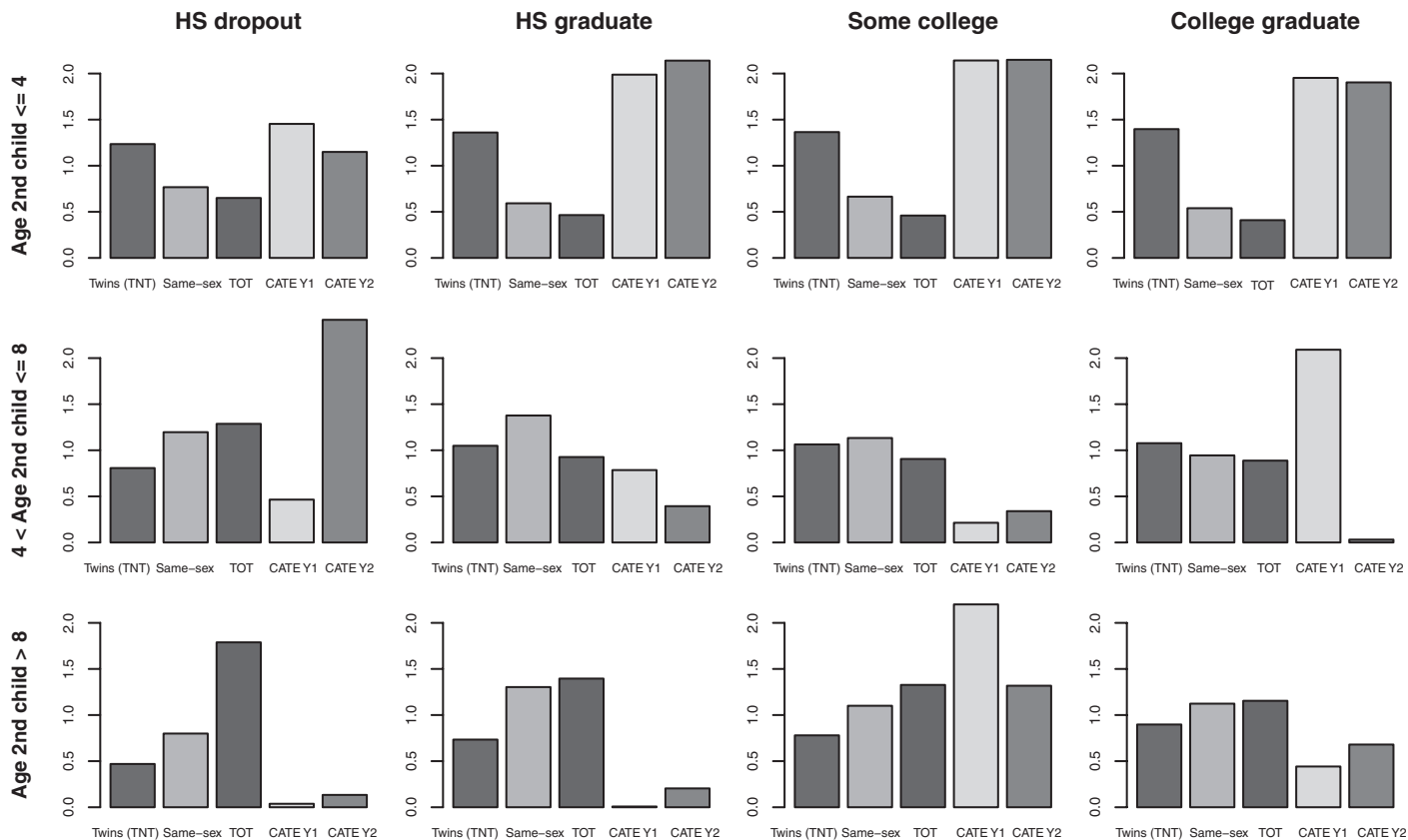


Figure 1. The distribution of compliers and related subpopulations across covariate cells: Y1 = weeks worked, Y2 = employment.

effects for compliant subpopulations other than the one defined by the instrument at hand. This approach turns out to do a good job of explaining why twins instruments produce smaller estimates of the labor supply consequences of childbearing than sex-composition instruments in the Angrist and Evans (1998) dataset.

The CEI assumption that lies at the heart of our approach rules out Roy (1951)–type selection into treatment on the basis of outcome gains. “No Roy selection” is unlikely to be compelling in many settings; gain-driven selection motivates a wide range of theoretical discussions of causal effects in labor economics and other applied fields (see, e.g., Rosen and Willis 1979 for a Roy model of schooling). At the same time, it seems hard to argue with the idea that any analysis of treatment-effect heterogeneity should *begin* with effect variation that is associated with the characteristics we observe.

In an effort to bridge the gap between heterogeneity associated with observed characteristics and latent gains, we also have explored an approach that allows some covariate values to satisfy our CEI assumption, while others – perhaps only a few – do not. This idea seems to work well in our application, generating, for example, remarkably similar estimates of population ATE using twins and same-sex instruments when the sample is reweighted toward cells that appear to satisfy CEI. At the same time, we acknowledge that brazen pretesting induces a complicated limiting distribution that we have not yet succeeded in characterizing and that may not always be useful for applied work. The development of robust and convenient inference procedures for CATE-type estimators seems a natural direction for further work on the external validity of IV estimates.

APPENDIX A: CEI IN AN ILLUSTRATIVE LDV MODEL

Because the text illustration of CEI uses a linear additive structure for potential outcomes, this example shows how we might have CEI in a nonadditive model with nonrandom selection into treatment and a Bernoulli outcome. To simplify, we drop covariates.

Suppose that treatment is determined by:

$$D = 1[h_0(1 - z) + h_1z > \eta], \quad \eta \mid z \sim U(0, 1/2), \quad 0 \leq h_0 < h_1 \leq 1/2$$

so that potential assignments are:

$$D_0^z = 1[h_0 > \eta] \text{ and } D_1^z = 1[h_1 > \eta].$$

The associated model for potential outcomes also has the latent-index representation:

$$y_j = 1[\eta + g_j > u_j], \quad u_j \mid \eta, z \sim U(0, 1), \quad 0 \leq g_j \leq 1/2, \quad j \in \{0, 1\}.$$

Iterating expectations, we can show that $\mathbb{E}[y_0] = g_0 + 1/4$ and

$$\mathbb{E}[y_0 \mid D = 1] = g_0 + \frac{1}{2} \frac{h_0^2 + (h_1^2 - h_0^2)\mathbb{P}[z = 1]}{h_0 + (h_1 - h_0)\mathbb{P}[z = 1]}.$$

Selection bias arises because the conditional mean of y_0 is a function of D . However, CEI holds because:

$$\Delta^z = \mathbb{E}[y_1 - y_0 \mid D_1^z, D_0^z] = \mathbb{E}[\mathbb{E}[y_1 - y_0 \mid D_1^z, D_0^z, \eta, z] \mid D_1^z, D_0^z] = g_1 - g_0.$$

References

- Abadie, A. (2003), “Semiparametric Instrumental Variable Estimation of Treatment Response Models.” *Journal of Econometrics*, 113(2), 231–63.
- Andrews, D. W. K. (1999), “Estimation When a Parameter Is on a Boundary.” *Econometrica*, 67(6), 1341–83.
- Andrews, D. W. K., and P. Guggenberger (2009), “Hybrid and Size-Corrected Subsampling Methods,” *Econometrica*, 77(3), 721–62.
- Angrist, J. D. (2004), “Treatment Effect Heterogeneity in Theory and Practice,” *Economic Journal*, 114(494), C52–C83.
- Angrist, J. D., and W. N. Evans (1998), “Children and Their Parents’ Labor Supply: Evidence from Exogenous Variation in Family Size,” *American Economic Review*, 88(3), 450–77.
- Angrist, J. D., G. W. Imbens, and D. B. Rubin (1996), “Identification of Causal Effects Using Instrumental Variables,” *Journal of the American Statistical Association*, 91, 444–72.
- Angrist, J. D., V. Lavy, and A. Schlosser (2010), “Multiple Experiments for the Causal Link between the Quantity and Quality of Children,” *Journal of Labor Economics*, 28, 773–824.
- Aronow, P., and A. J. Sovey (2010), “Beyond LATE: A Simple Method for Recovering Sample Average Treatment Effects,” Yale University, discussion paper.
- Brown, B. W., and W. K. Newey (2002), “Generalized Method of Moments, Efficient Bootstrapping, and Improved Inference,” *Journal of Business & Economic Statistics*, 20(4), 507–17.
- Chamberlain, G. (2010), “Bayesian Aspects of Treatment Choice,” Harvard University, discussion paper.
- Cruces, G. (2005), “Poverty, Income Fluctuations and Work: Argentina 1991–2002,” London School of Economics and Political Science, Ph.D. dissertation.
- Cruces, G., and S. Galiani (2007), “Fertility and Female Labor Supply in Latin America: New Causal Evidence,” *Labour Economics*, 14, 565–73.

- Crump, R. K., V. J. Hotz, G. W. Imbens, and O. A. Mitnik (2009), "Dealing with Limited Overlap in Estimation of Average Treatment Effects," *Biometrika*, 96(1), 187–99.
- Ebenstein, A. (2009), "When Is the Local Average Treatment Close to the Average? Evidence from Fertility and Labor Supply," *Journal of Human Resources*, 44(4), 955–75.
- Frangakis, C. E., and D. B. Rubin (2002), "Principal Stratification in Causal Inference," *Biometrics*, 58(1), 21–9.
- Frölich, M. (2007), "Nonparametric IV Estimation of Local Average Treatment Effects with Covariates," *Journal of Econometrics*, 139(1), 35–75.
- Gelbach, J. B. (2002), "Public Schooling for Young Children and Maternal Labor Supply," *American Economic Review*, 92(1), 307–22.
- Hahn, J. (1996), "A Note on Bootstrapping Generalized Method of Moments Estimators," *Econometric Theory*, 12, 187–97.
- Hahn, J. (1998), "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects," *Econometrica*, 66(2), 315–32.
- Hall, P., and J. L. Horowitz (1996), "Bootstrap Critical Values for Tests Based on Generalized-Method-of-Moments Estimators," *Econometrica*, 64(4), 891–916.
- Heckman, J. J. (2010), "Building Bridges between Structural and Program Evaluation Approaches to Evaluating Policy," *Journal of Economic Literature*, 48(2), 356–98.
- Heckman, J. J., J. L. Tobias, and E. Vytlacil (2001), "Four Parameters of Interest in the Evaluation of Social Programs," *Southern Economic Journal*, 68(2), 210–23.
- Heckman, J. J., J. L. Tobias, and E. Vytlacil (2003), "Simple Estimators for Treatment Parameters in a Latent-Variable Framework," *Review of Economics and Statistics*, 85(3), 748–55.
- Ichino, A., and R. Winter-Ebmer (1999), "Lower and Upper Bounds of Returns to Schooling: An Exercise in IV Estimation with Different Instruments," *European Economic Review*, 43, 889–901.
- Imbens, G. W., and J. D. Angrist (1994), "Identification and Estimation of Local Average Treatment Effects," *Econometrica*, 62(2), 467–75.
- Imbens, G. W., and W. K. Newey (2009), "Identification and Estimation of Triangular Simultaneous Equations Models without Additivity," *Econometrica*, 77(5), 1481–512.
- Leeb, H., and M. M. Pötscher (2008), "Model Selection," in T. G. Andersen, R. A. Davis, J.-P. Kreiss, and T. Mikosch (eds.), *Handbook of Financial Time Series*, Berlin Heidelberg, Springer-Verlag.
- Olsen, R. J. (1980), "A Least Squares Correction for Selectivity Bias," *Econometrica*, 48(7), 1815–20.
- Oreopoulos, P. (2006), "Estimating Average and Local Average Treatment Effects of Education When Compulsory Schooling Laws Really Matter," *American Economic Review*, 96(1), 152–75.
- Rosenzweig, M. R., and K. I. Wolpin (1980), "Testing the Quantity–Quality Fertility Model: The Use of Twins as a Natural Experiment," *Econometrica*, 48(1), 227–40.
- Roy, A. D. (1951), "Some Thoughts on the Distribution of Earnings," *Oxford Economic Papers*, 3(2), 135–46.
- Rubin, D. B. (1977), "Assignment to a Treatment Group on the Basis of a Covariate," *Journal of Educational Statistics*, 2(1), 1–26.
- Van der Vaart, A. W. (1998), *Asymptotic Statistics*, New York: Cambridge University Press.

- Vella, F., and M. Verbeek (1999), "Estimating and Interpreting Models with Endogenous Treatment Effects," *Journal of Business & Economic Statistics*, 17(4), 473–78.
- Vytlacil, E. (2002), "Independence, Monotonicity, and Latent Index Models: An Equivalence Result," *Econometrica*, 70(1), 331–41.
- Wald, A. (1940), "The Fitting of Straight Lines if Both Variables Are Subject to Error," *Annals of Mathematical Statistics*, 11(3), 284–300.
- Willis, R. J., and S. Rosen (1979), "Education and Self-Selection," *The Journal of Political Economy*, 87(5), S7–S36.

