

When is TSLS *Actually* LATE?*

Christine Blandhol[†]

John Bonney[‡]

Magne Mogstad[§]

Alexander Torgovitsky[¶]

August 25, 2025

Abstract

Linear instrumental variable estimators, such as two-stage least squares (TSLS), are commonly interpreted as estimating non-negatively weighted averages of causal effects, referred to as local average treatment effects (LATEs). We examine whether the LATE interpretation actually applies to the types of TSLS specifications that are used in practice. We show that if the specification includes covariates—which most empirical work does—then the LATE interpretation does not apply in general. Instead, the TSLS estimator will, in general, reflect treatment effects for both compliers *and* always/never-takers, and some treatment effects for the always/never-takers will *necessarily* be negatively weighted. We show that the only specifications that have a LATE interpretation are “saturated” specifications that control for covariates nonparametrically, implying that such specifications are both sufficient and *necessary* for TSLS to have a LATE interpretation, at least without additional parametric assumptions. This result is concerning because, as we document, empirical researchers almost never control for covariates nonparametrically, and rarely discuss or justify parametric specifications of covariates. We apply our results to thirteen empirical studies and find strong evidence that the LATE interpretation of TSLS is far from accurate for the types of specifications actually used in practice. We offer concrete recommendations for practice motivated by our theoretical and empirical results.

*We thank Richard Boylan, Deniz Dutz, Sergio Firpo, Mikkel Gandil, Tymon Słoczyński, Chris Walters, and Thomas Wiemann for helpful comments.

[†]Department of Economics, Princeton University.

[‡]Department of Economics, Stanford University. Research supported by National Science Foundation Graduate Research Fellowship under grant DGE-1656518.

[§]Kenneth C. Griffin Department of Economics, University of Chicago; Statistics Norway; NBER.

[¶]Kenneth C. Griffin Department of Economics, University of Chicago. Research supported by National Science Foundation grant SES-1846832.

1 Introduction

Instrumental variable (IV) strategies are widely used for causal inference in economics, political science, sociology, epidemiology, and other fields. Since the work of [Imbens and Angrist \(1994\)](#), it has been increasingly common to interpret linear IV estimators as estimating a local average treatment effect (LATE), or at least a non-negatively weighted average of LATEs.

The LATE interpretation is most commonly derived for simplified IV specifications that do not include covariates. We examine whether the LATE interpretation extends to the types of linear IV specifications that are used in practice. We show that if the IV specification includes covariates—which most empirical work does—then the LATE interpretation does not apply in general. Instead, the linear IV estimand with covariates is generally composed of treatment effects for both compliers and always-takers, and some always-taker treatment effects are necessarily negatively weighted.

Our finding challenges the claim by [Angrist and Pischke \(2009, pg. 173\)](#) that

2SLS with covariates produces an average of covariate-specific LATEs. . . These results provide a simple casual [typo in original] interpretation for 2SLS in most empirically relevant settings.

The formal justification that [Angrist and Pischke \(2009\)](#) provide for this assertion is based on a saturated two stage least squares (TSLS) specification that controls for covariates nonparametrically, described by the authors as the “saturate and weight approach” (Theorem 4.5.1; originally Theorem 3 in [Angrist and Imbens, 1995](#)). Drawing on this justification, they continue on pg. 178 by suggesting

It seems reasonable to imagine that models with fewer parameters . . . nevertheless approximate some kind of covariate-averaged LATE.

Our results show that this imagining is incorrect: saturated specifications are *necessary* for TSLS with covariates to be interpretable as an average of covariate-specific LATEs, at least without additional parametric assumptions.

Are saturated specifications “empirically relevant?” In Section 2, we report the results of a survey on the specification of linear IV estimators in published empirical papers in economics. Of the 99 papers in our survey that use a linear IV estimator with covariates, we found only five papers that used a saturated specification at least once and only a single paper that exclusively used saturated specifications. The implication of our results for the 98 other papers is that they may not be estimating an average of covariate-specific LATEs. In fact, they may be estimating a quantity that doesn’t even satisfy the minimal requirement of being a non-negatively weighted average of subgroup-specific treatment effects, a property we describe as weakly causal.

Section 2 also contains an exposition of our main findings in the special case of a binary treatment and binary instrument. This case exposes the central intuition: if the covariates are not specified flexibly, then the TSLS estimand depends not only on treatment *effects*, but also on potential outcome *levels*. We call this phenomenon level dependence. Because the TSLS estimand is generally level dependent, it does not necessarily have a unique decomposition into a weighted average of subgroup treatment effects. Consequently, analyzing whether the TSLS estimand is weakly causal is more complicated than simply checking for non-negative weights.¹

In Section 3, we tackle this challenge by providing a conceptual definition of a weakly causal estimand that is separated from the form that the estimand takes. We then provide sufficient and necessary conditions for an estimand to be weakly causal. The characterization has two components. First, a weakly causal estimand cannot be level dependent. Second, a weakly causal estimand should not apply negative weight to the treatment effects for any subgroup.

In Section 4, we specialize this definition to TSLS estimands. We show that a *necessary* condition for the TSLS estimand to be weakly causal is that the TSLS specification has rich covariates, meaning that it exactly reproduces the conditional mean of the instrument. Specifications that are saturated in covariates, such as the Angrist and Pischke (2009) “saturate and weight” specification, will always have rich covariates. But a non-saturated TSLS specification only has rich covariates if an implicit parametric functional form assumption happens to be correct. Saturated specifications can be extremely data hungry, which may explain why they were so seldom used in our survey of empirical papers.

Kolesár (2013) provided the most general sufficient conditions for the TSLS estimand to be equal to a non-negatively weighted average of LATEs. His conditions maintain rich covariates as an *assumption*. Our results show that rich covariates is also *necessary* for the TSLS estimand to have even a weakly causal interpretation, let alone an interpretation as a non-negatively weighted average of LATEs.

The implication of our results is that the Angrist and Pischke (2009) interpretation of TSLS as a non-negatively weighted average of LATEs is fragile. In particular, it depends on rich covariates, which is an implicit parametric functional form assumption that appears to always be left unstated in empirical work. Although our survey turned up only a single paper that used a TSLS specification guaranteed to satisfy rich

¹In this sense our results and analysis are quite different from the recent literature on two-way fixed effects models (e.g. Goodman-Bacon, 2021; Sun and Abraham, 2021), which point out interpretation problems that arise in event studies if there are heterogeneous treatment effects due to observables (in particular, cohorts). When analyzed without covariates, these estimands are not level dependent, but may have negative weights. A consequence is that the problems we point to remain even with constant treatment effects (Section 4.5), unlike in the two-way fixed effects literature.

covariates, we found numerous papers that nevertheless invoked the widespread LATE interpretation. Our results draw this interpretation into question.

In Section 5, we consider alternatives to TSLS. One alternative is to change the TSLS specification to be saturated. However, as our empirical survey suggested, and as our simulations confirm, this is often impractical due to the large number of regressors produced by saturating. An alternative is to use double/debiased machine learning (Chernozhukov et al., 2018, “DDML”) to estimate a partially linear IV (PLIV) modification of TSLS that controls for covariates in an additive but nonparametric way. This frees the researcher of the need to choose a parameterization of the covariates, but comes at a computational cost. It also estimates a quantity that, while weakly causal, might still be hard to interpret.

When the instrument is binary, a related and potentially more attractive alternative is to estimate an unconditional average causal response (ACR), which reduces to an unconditional LATE when the treatment is also binary. This can be done either non-parametrically with DDML or semi-parametrically using instrument propensity score weighting (e.g. Tan, 2006; Uysal, 2011; Słoczyński et al., 2024). A third potential alternative for the binary instrument, binary treatment case is Abadie’s (2003) κ -weighting approach. The implementation of κ -weighting requires explicitly parameterizing the conditional mean of the instrument given the covariates, the same object that we show needs to be implicitly assumed to be correctly specified for TSLS to be weakly causal. However, we show that κ -weighting too will only be weakly causal if rich covariates is satisfied, the same necessary condition as for TSLS.

In Section 6, we compare the TSLS estimator with these alternatives in thirteen empirical papers. We find strong evidence that rich covariates is often not satisfied in practice. DDML PLIV estimates can be nearly as different from TSLS as TSLS is different from its comparable OLS estimate. The Ramsey (1969) RESET test tends to do a good job detecting when failures of rich covariates lead to sizable differences between TSLS and a DDML PLIV estimate. DDML PLIV estimates can still be dramatically different from DDML or semi-parametric estimates of the unconditional ACR/LATE.

In Section 7, we provide some concluding remarks and recommendations for practice, all of which can be implemented in Stata or R with mature software packages. These recommendations show that it is still possible to estimate an unconditional ACR/LATE or a statistically-weighted average of conditional LATEs in the presence of covariates. But not with the types of TSLS specifications that are currently being used in practice.

Słoczyński (2020, 2024) has recently made a different critique of the interpretation of TSLS estimators. He maintains rich covariates as an assumption and shows that the TSLS estimand can still fail to be weakly causal if the direction of monotonicity varies with covariates but the TSLS specification does not include instrument-covariate

interactions in the first stage. In contrast, our analysis focuses on the necessity of the rich covariates condition under a stronger, unconditional monotonicity condition. [Słoczyński \(2024\)](#) also makes the important theoretical point that even when rich covariates is satisfied, the resulting linear IV estimand may be quite different from the type of unconditional LATE that practitioners might expect. We do not discuss any theory about this point, although we do illustrate it our empirical applications.

Rich covariates remains necessary under the weaker monotonicity condition considered by [Słoczyński \(2020, 2024\)](#). Taken together, our paper and [Słoczyński \(2024\)](#) show that two conditions are necessary for the TSLS estimand to be interpretable as a non-negatively weighted average of LATEs: (i) rich covariates, and (ii) a first stage equation flexible enough to capture changes in the direction of monotonicity across covariate values. The necessity of these conditions provides a definitive answer to the question: “When is TSLS *actually* LATE?” That answer: probably not often.

2 Overview

In this section, we demonstrate our main results in the special case of a binary treatment and a binary instrument.

2.1 Linear IV with covariates is not LATE

Let $T \in \{0, 1\}$ be a binary treatment and $Z \in \{0, 1\}$ be a binary instrument. The outcome is Y with potential outcomes $Y(0)$ and $Y(1)$ related via $Y = (1 - T)Y(0) + TY(1)$. Potential treatment states are $T(0)$ and $T(1)$ with $T = (1 - Z)T(0) + ZT(1)$. The vector of covariates is X .

Assume that Z is conditionally exogenous in the sense of being independent of $(Y(0), Y(1), T(0), T(1))$ conditional on X . Suppose that the [Imbens and Angrist \(1994\)](#) monotonicity condition holds so that $\mathbb{P}[T(1) \geq T(0)] = 1$. The monotonicity condition implies that the group variable $G \equiv (T(0), T(1))$ can take three values with non-zero probability: $G = (0, 0) \equiv \text{NT}$ are the never-takers, $G = (0, 1) \equiv \text{CP}$ are the compliers, and $G = (1, 1) \equiv \text{AT}$ are the always-takers.

Consider a linear IV regression with outcome variable Y , endogenous variable T , excluded instrument Z , and a vector of control variables X that includes a constant. The Frisch-Waugh-Lovell Theorem can be used to show that the IV estimand (the population coefficient on T) is given by

$$\beta_{\text{iv}} = \frac{\mathbb{E}[Y\tilde{Z}]}{\mathbb{E}[T\tilde{Z}]}, \quad \text{where} \quad \tilde{Z} \equiv Z - \mathbb{E}[Z|X] \quad (1)$$

are the residuals from a regression of Z on X , and

$$\mathbb{L}[Z|X] \equiv X' \mathbb{E}[XX']^{-1} \mathbb{E}[XZ]$$

are the population fitted values from regressing (linearly projecting) Z onto X . The IV estimand, β_{iv} , is often interpreted as reflecting a non-negatively weighted average of treatment effects for only the compliers. The following proposition shows that this is not true in general.

Proposition 1. Suppose that $\mathbb{E}[Y(t)|X] = \eta'_t X$ for some unknown parameters η_t , $t = 0, 1$.² Let $\Delta(\text{CP}, x) \equiv \mathbb{E}[Y(1) - Y(0)|G = \text{CP}, X = x]$ and $\Delta(\text{AT}, x) \equiv \mathbb{E}[Y(1) - Y(0)|G = \text{AT}, X = x]$ denote the conditional average treatment effects for the compliers and always-takers, respectively. Then

$$\beta_{iv} = \mathbb{E}[\omega(\text{CP}, X)\Delta(\text{CP}, X)] + \mathbb{E}[\omega(\text{AT}, X)\Delta(\text{AT}, X)], \quad (2)$$

$$\begin{aligned} \text{where } \omega(\text{CP}, X) &\equiv \mathbb{E}[Z|X] (1 - \mathbb{L}[Z|X]) \mathbb{P}[G = \text{CP}|X] \mathbb{E}[\tilde{Z}T]^{-1} \\ \text{and } \omega(\text{AT}, X) &\equiv \mathbb{E}[\tilde{Z}|X] \mathbb{P}[G = \text{AT}|X] \mathbb{E}[\tilde{Z}T]^{-1}. \end{aligned}$$

If $\mathbb{E}[\tilde{Z}T] > 0$, then the complier weights $\omega(\text{CP}, X)$ are negative if and only if $\mathbb{L}[Z|X] > 1$. The always-taker weights $\omega(\text{AT}, X)$ are strictly negative with positive probability unless $\mathbb{E}[\tilde{Z}|X] = 0$ deterministically.

Proposition 1 shows that, in general, β_{iv} reflects not only the compliers, but also the always-takers. The monotonicity condition implies that the first stage coefficient is positive, so $\mathbb{E}[\tilde{Z}T] > 0$. The weights on the always-takers therefore have the same sign as the random variable $\mathbb{E}[\tilde{Z}|X] = \mathbb{E}[Z|X] - \mathbb{L}[Z|X]$. Because X contains a constant, $\mathbb{E}[\tilde{Z}] = \mathbb{E}[\mathbb{E}[\tilde{Z}|X]] = 0$, implying that $\mathbb{E}[\tilde{Z}|X]$ is either always equal to zero, or else it has positive probability of taking both positive and negative values. As a consequence, whenever $\mathbb{L}[Z|X] \neq \mathbb{E}[Z|X]$, the IV estimand incorporates negatively weighted treatment effects for some groups, which means that it fails to satisfy even a minimal condition for “being causal.”

This reasoning shows that in order for the LATE interpretation to hold, it is necessary that $\mathbb{L}[Z|X] = \mathbb{E}[Z|X]$, a condition we call rich covariates. Specifications that are saturated in covariates, such as “saturate and weight” (Angrist and Pischke, 2009), have rich covariates. If Z and X are independent, as can be the case in some controlled and natural experiments, then any specification with a constant will have rich covari-

²This additional linearity assumption is made in order to simplify the weights. Removing the assumption only *amplifies* the negative interpretation issues exposed by Proposition 1. Our general results in Section 4 do not maintain this assumption.

ates.³ Outside these two cases, having rich covariates is a parametric assumption. If it fails, then the IV estimand β_{iv} reflects not just compliers, but also negatively-weighted always-takers.

There is no reason to expect, a priori, that the weights on the always-taker treatment effects in (2) will be small in magnitude. In many applications, the proportion of always-takers, $\mathbb{P}[G = AT|X]$, is considerably larger than the proportion of compliers, $\mathbb{P}[G = CP|X]$. As a consequence, even negative values of $\mathbb{E}[\tilde{Z}|X]$ that are small in magnitude can produce large negative weights on the always-taker treatment effects.

Decomposition (2) is not the only one possible. Instead of interpreting β_{iv} as a weighted average of compliers and always-takers, one can interpret it as a weighted average of compliers and never-takers, or of all three groups, as shown in the next proposition.

Proposition 2. Suppose that $\mathbb{E}[Y(t)|X] = \eta'_t X$ for some (unknown) parameters η_t , and both $t = 0, 1$. Let $\Delta(NT, x) \equiv \mathbb{E}[Y(1) - Y(0)|G = NT, X = x]$ denote the conditional average treatment effect for the never-takers. Then for any real number ϵ ,

$$\begin{aligned}\beta_{iv} &= \mathbb{E}[\omega_\epsilon(CP, X)\Delta(CP, X)] + \mathbb{E}[\omega_\epsilon(AT, X)\Delta(AT, X)] + \mathbb{E}[\omega_\epsilon(NT, X)\Delta(NT, X)], \\ \text{where } \omega_\epsilon(CP, X) &\equiv \left(\epsilon \mathbb{E}[\tilde{Z}|X] + \mathbb{L}[Z|X](1 - \mathbb{E}[Z|X])\right) \mathbb{P}[G = CP|X] \mathbb{E}[\tilde{Z}T]^{-1}, \\ \omega_\epsilon(AT, X) &\equiv \epsilon \mathbb{E}[\tilde{Z}|X] \mathbb{P}[G = AT|X] \mathbb{E}[\tilde{Z}T]^{-1}, \\ \text{and } \omega_\epsilon(NT, X) &\equiv (\epsilon - 1) \mathbb{E}[\tilde{Z}|X] \mathbb{P}[G = NT|X] \mathbb{E}[\tilde{Z}T]^{-1}.\end{aligned}$$

Each choice of ϵ in Proposition 2 provides a different interpretation of β_{iv} , with Proposition 1 corresponding to $\epsilon = 1$. However, unless rich covariates holds, so that $\mathbb{E}[\tilde{Z}|X] = \mathbb{E}[Z|X] - \mathbb{L}[Z|X] = 0$, any choice of ϵ still leads to an interpretation that involves either the always-takers or the never-takers, or both, with negative weights for some values of X , as well as potentially negative weights for the compliers. Only in specifications with rich covariates is β_{iv} a non-negatively weighted average among compliers alone.

Proposition 2 shows that a causal interpretation can be partially salvaged if there is one-sided non-compliance. For example, if there are no always-takers, so that $\mathbb{P}[G = AT|X] = 0$, then one can take $\epsilon = 1$, so that β_{iv} is a weighted average among compliers alone. The same is true if there are no never-takers by taking $\epsilon = 0$. The complier weights can still be negative in these cases if $\mathbb{L}[Z|X]$ does not lie in $[0, 1]$, but rich covariates is stronger than necessary to rule this out. However, these conclusions about one-sided non-compliance depend on the simplifying linearity assumption that $\mathbb{E}[Y(t)|X] = \eta'_t X$, which we do not maintain in our general results in Section 4.

³In Appendix SA.1, we discuss the case in which Z is randomly assigned conditional on a subset of X , as might occur in a stratified experiment.

2.2 Intuition

The intuition behind Propositions 1 and 2 can be seen by writing the numerator of β_{iv} as

$$\mathbb{E}[Y\tilde{Z}] = \mathbb{E}\left[\mathbb{E}\left[Y\tilde{Z}|X\right]\right] = \mathbb{E}\left[\overbrace{\mathbb{C}[Y, \tilde{Z}|X]}^{\text{only contains complier treatment effects}}\right] + \underbrace{\mathbb{E}\left[\mathbb{E}[Y|X]\mathbb{E}[\tilde{Z}|X]\right]}_{\text{contains all three groups}}, \quad (3)$$

where \mathbb{C} denotes covariance. The first term in (3) is the average of the numerator of a nonparametric IV specification that *conditions* on X . The argument in [Imbens and Angrist \(1994\)](#) shows that this term is equal to an average of scaled LATEs, which only reflects treatment effects for the compliers. It is the second term of (3) that causes problems. This term reflects the difference between nonparametric conditioning and linear projection.⁴

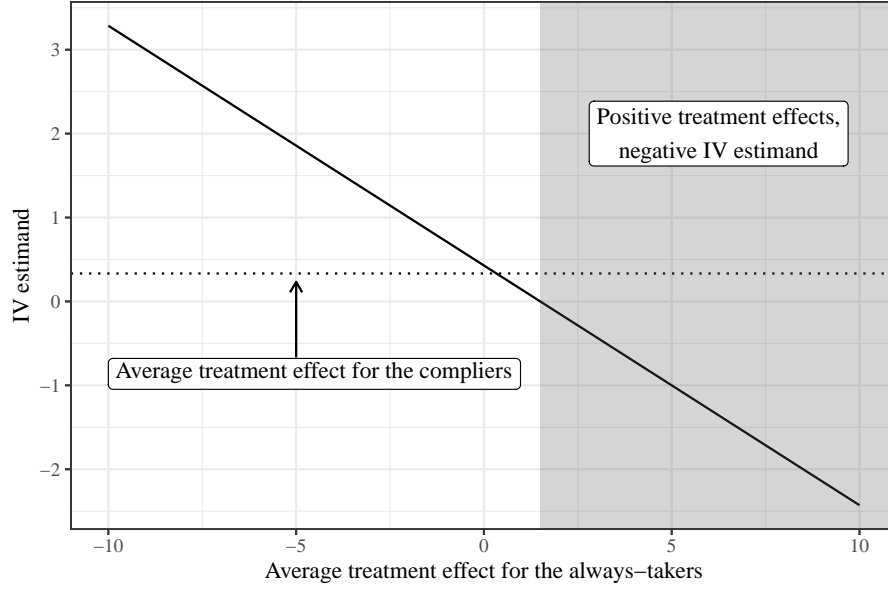
When covariates are not rich, so that $\mathbb{E}[\tilde{Z}|X] \neq 0$, the second term in (3) generally depends on $\mathbb{E}[Y|X]$, a quantity which is determined not only by compliers, but also by always-takers and never-takers. This creates level dependence in β_{iv} because the always-takers always have $Y = Y(1)$ and the never-takers always have $Y = Y(0)$: β_{iv} depends on the levels of the always-taker and never-taker potential outcomes. As we show in Section 3, level dependent estimands do not have a causal interpretation because the levels can lead β_{iv} to have the “wrong sign.”

The expression in Proposition 1 arises from centering the term $\mathbb{E}[Y|X]$ in (3) around $\mathbb{E}[Y(0)|X]$. The simplifying linearity assumption implies that $\mathbb{E}[Y(0)|X] = \eta'_0 X$ is uncorrelated with $\mathbb{E}[\tilde{Z}|X]$. Since never-takers always have $Y = Y(0)$, the centering removes the average untreated outcome for the never-takers, leaving only a weighted average of the complier and always-taker treatment effects. Alternatively, we can center around $\mathbb{E}[Y(1)|X] = \eta'_1 X$, which leaves a weighted average of the complier and never-taker treatment effects. Both decompositions are equally valid ways to rewrite a single number, β_{iv} , as a weighted average of $\Delta(\text{CP}, X)$, $\Delta(\text{AT}, X)$, and $\Delta(\text{NT}, X)$. Taking an ϵ -weighted average of these two decompositions yields the expression in Proposition 2, which creates a family of equally-valid decompositions.

The theory we develop in Section 3 is designed to handle this type of non-uniqueness in decomposition and determine, in a general setting, necessary conditions for the existence of *some* “good” decomposition. For the simplified case considered here, with a binary treatment, a binary instrument, and the linearity assumption $\mathbb{E}[Y(t)|X] = \eta'_t X$, this type of analysis can be done directly, as in Proposition 2. Our analysis of more general TSLS specifications in Section 4 shows that the necessity of rich covariates for a causal interpretation is a conclusion that applies more broadly.

⁴[Firpo et al. \(2020\)](#) make a similar point in the context of balance tests for stratified experiments.

Figure 1: IV with covariates is not LATE



2.3 Numerical illustration

As a simple illustration of these results, suppose that $X \in \{(1, -1), (1, 0), (1, 1)\}$ with equal probability, where the first component corresponds to a constant. Then suppose that

$$\mathbb{E}[Z|X = x] = \mathbb{P}[Z = 1|X = (1, x)] = \begin{cases} 4/5 & \text{if } x \in \{-1, 1\} \\ 2/5 & \text{if } x = 0 \end{cases}.$$

Regressing Z onto X yields the constant regression line:

$$\mathbb{L}[Z|X] = X' \mathbb{E}[XX']^{-1} \mathbb{E}[XZ] = 2/3,$$

so that $\mathbb{E}[\tilde{Z}|X] = \mathbb{E}[Z|X] - \mathbb{L}[Z|X] \neq 0$ and is both positive and negative with non-zero probability.

Suppose that the conditional group share probabilities are given by:

$$\begin{aligned} (\text{never-takers}) \quad & \mathbb{P}[G = \text{NT}|X = (1, x)] = 1/3 \\ (\text{compliers}) \quad & \mathbb{P}[G = \text{CP}|X = (1, x)] = 1/6 + |x|/6 \\ (\text{always-takers}) \quad & \mathbb{P}[G = \text{AT}|X = (1, x)] = 1/2 - |x|/6. \end{aligned}$$

Simplifying the algebra in Proposition 1 yields

$$\omega(\text{CP}, (1, x)) = \begin{cases} 12/7, & \text{if } |x| = 1 \\ 3/7, & \text{if } x = 0 \end{cases} \quad \text{and} \quad \omega(\text{AT}, (1, x)) = \begin{cases} 6/7, & \text{if } |x| = 1 \\ -18/7, & \text{if } x = 0 \end{cases}.$$

For simplicity, assume that $Y(0) = 0$, so that treatment effects are determined solely by $Y(1)$, and that $\mathbb{E}[Y(1)|G = \text{CP}, X = x] \equiv \mu(\text{CP})$ and $\mathbb{E}[Y(1)|G = \text{AT}, X = x] = \mu(\text{AT})$ do not depend on x . Then Proposition 1 shows that

$$\beta_{\text{iv}} = \frac{9}{7}\mu(\text{CP}) - \frac{2}{7}\mu(\text{AT}).$$

Figure 1 shows the value of β_{iv} as a function of $\mu(\text{AT})$, keeping $\mu(\text{CP}) = 1/3$. If it were true that LATE only reflects the compliers, then we would expect to see a flat line, so that the IV estimand doesn't depend on the treatment effect for the always-takers. Not only is the line not flat, it slopes down. This means that the IV estimand can be negative even when both the compliers and the always-takers have positive treatment effects.

2.4 Survey on IV specifications used in empirical work

Propositions 1 and 2 show that using an IV specification that is saturated in covariates is needed for the LATE interpretation asserted by Angrist and Pischke (2009). To get a sense of how common it is to saturate in covariates, we surveyed the specifications used in the empirical economics literature.

Our sample was constructed by searching the Web of Science Database for articles published between January 2000 and October 2018 containing the words “instrument” or “instrumental variable” in the abstract, title, or topic words. We restricted the search to the following five journals: *Journal of Political Economy*, *American Economic Review*, *Quarterly Journal of Economics*, *Review of Economic Studies*, and *Econometrica*. In total, 266 articles matched our search criteria.

We restricted our attention to papers that use at least one IV specification in an empirical application. This produced 122 papers; the other 144 papers not included were either methodological papers without an empirical application, or were papers that used the word “instrument” in a different context, such as to describe a policy or financial instrument. Column (1) of Table 1 tabulates the papers used in our survey by the journal in which they were published.

Column (2) shows that over 92% of the papers in our survey use TSLS (including exactly identified linear IV) for at least some of their results. Column (3) counts the subset of the papers in column (2) for which *all* TSLS specifications in the main body of the paper include at least one covariate, or the authors explicitly state the exogeneity

Table 1: IV papers by journal and type

	(1)	(2)	(3)	(4)
	All papers	Papers using TSLS	Papers using TSLS with covariates	Papers using TSLS with covariates, referring to LATE
American Economic Review	100% 44	95% 42	82% 36	27% 12
Quarterly Journal of Economics	100% 28	93% 26	86% 24	14% 4
Journal of Political Economy	100% 23	91% 21	83% 19	30% 7
Econometrica	100% 15	73% 11	73% 11	27% 4
Review of Economic Studies	100% 12	100% 12	75% 9	25% 3
All	100 % 122	92% 112	81% 99	25% 30

assumption for the instrument as conditional on covariates.⁵ Comparing columns (2) and (3) shows that using covariates in TSLS is extremely common practice; only 13 out of the 112 papers that use TSLS include any specifications without covariates. Column (4) shows that almost a third of the papers that use TSLS with covariates also explicitly use the phrases “compliers,” “local average treatment effect,” or “LATE” to describe their results.

In Table 2, we categorize the papers in column (3) of Table 1 by the TSLS specifications they use. Column (2) shows that only 5% of these papers use any specification that is saturated in covariates. These are typically preliminary specifications with only a set of fixed effects. Column (3) shows that all of these papers use at least one specification that is *not* saturated in covariates, with only one exception. The one exception is [Chamberlain and Imbens \(2004\)](#). Column (4) shows that those authors also saturate the first stage in both the covariates and the instruments, as prescribed by [Angrist and Pischke’s \(2009\)](#) “saturate and weight approach.”

2.5 Implications for empirical practice

Avoiding the conclusion of Propositions 1 and 2 requires choosing a specification with rich covariates, that is, one that ensures $\mathbb{L}[Z|X] = \mathbb{E}[Z|X]$.

⁵Another possible justification for including covariates is to improve statistical precision. This motivation was rarely stated explicitly in the papers in our survey. While it is difficult to infer researchers’ unstated reasons for choosing particular specifications, it seems unlikely that they would *only* use specifications with covariates if covariates were only being used to improve precision.

Table 2: TSLS papers with covariates by journal and empirical specification

	(1)	(2)	(3)	(4)
			At least one specification	
	Papers using TSLS with covariates	Saturated in covariates	Not saturated in covariates	Saturated in instruments and covariates
American Economic Review	100% 36	0% 0	100% 36	0% 0
Quarterly Journal of Economics	100% 24	4% 1	100% 24	0% 0
Journal of Political Economy	100% 19	16% 3	100% 19	0% 0
Econometrica	100% 11	9% 1	91% 10	9% 1
Review of Economic Studies	100% 9	0% 0	100% 9	0% 0
All	100 % 99	5% 5	99% 98	1% 1

Notes: This table classifies the papers from column (3) of Table 1 by TSLS specification.

The saturate and weight (SW) specification ([Angrist and Pischke, 2009](#)) is saturated in covariates, so has rich covariates. However, it also uses a first stage that is fully saturated in both the covariates *and* the instruments, meaning that the regressors are indicators for all possible instrument-covariate combinations. This results in many excluded variables and potential many instruments bias, which may explain why the SW specification was used by only a single paper in the survey. In fact, that one paper ([Chamberlain and Imbens, 2004](#)) is a methodological consideration of many instruments bias.

However, the interactions between covariates and instruments used in the SW specification may not be necessary for the LATE interpretation. Excluded interactions were not used in (1) and yet Propositions 1 and 2 show that if covariates are rich, then β_{iv} will be composed of only non-negatively weighted complier effects. The reason is that we assumed that the [Angrist and Imbens \(1995\)](#) monotonicity condition operates in the same direction for every covariate group. In contrast, the SW specification is premised on a weaker version of the monotonicity assumption that allows the direction of monotonicity to vary with covariates. [Słoczyński \(2020, 2024\)](#) shows that including the instrument-covariate interaction terms used in SW is necessary when considering this weaker monotonicity condition.

Our results show that flexibly controlling for covariates is important for ensuring that TSLS has a causal interpretation. If a flexible covariate specification cannot be used, then another response is to test the null hypothesis that $\mathbb{L}[Z|X] = \mathbb{E}[Z|X]$. The

most well-known test is Ramsey’s (1969) RESET test (e.g. Wooldridge, 2010, pp. 137–138), which is straightforward to implement in either Stata or R. No papers in our survey reported such a test. If Z is binary, then it is also sensible to check that the fitted values $\mathbb{L}[Z|X]$ lie between 0 and 1, which is necessary for $\mathbb{L}[Z|X] = \mathbb{E}[Z|X]$. Alternatively, researchers can consider using an estimator other than TSLS. We discuss alternative estimators in Section 5 and apply them in Section 6.

3 Definition and characterization of weakly causal estimands

In this section we define a weak property that an estimand should satisfy in order to “be causal.” We do this because, as Proposition 2 showed, if rich covariates fails, then the TSLS estimand might have multiple equally valid decompositions. Alternatively, if the simplifying linearity assumption maintained in Proposition 2 is dropped, the TSLS estimand might not have any decomposition in terms of only treatment effects. These complications motivate a more abstract definition of a weakly causal estimand that is separated from the functional form that the estimand takes. We develop the weakly causal property in the context of a nonparametric IV model using potential outcomes notation (e.g. Angrist et al., 1996) with an ordered treatment and a multivalued instrument. The results generalize the special case of a binary treatment and binary instrument discussed in Section 2.

3.1 The nonparametric instrumental variables model

A discrete, ordered treatment variable T takes values in $\mathcal{T} \equiv \{t_0, t_1, \dots, t_J\}$, listed in increasing order. We are interested in the causal effects that T has on an outcome variable, Y . We observe a scalar- or vector-valued instrumental variable (IV) Z that takes values in a set $\mathcal{Z} \equiv \{z_0, z_1, \dots, z_K\}$. The case in Section 2 corresponds to $\mathcal{T} = \{0, 1\}$ and $\mathcal{Z} = \{0, 1\}$. There is a vector of covariates X with support \mathcal{X} .

Associated with each level of the IV is a potential treatment choice, $T(z)$. Associated with each level of the treatment is a potential outcome, $Y(t)$, which does not directly depend on the instrument due to the usual exclusion restriction. The potential and actual treatments and outcomes are related through

$$T = \sum_{z \in \mathcal{Z}} \mathbb{1}[Z = z]T(z) \quad \text{and} \quad Y = \sum_{t \in \mathcal{T}} \mathbb{1}[T = t]Y(t).$$

We maintain the following standard nonparametric exogeneity condition throughout our analysis.

Assumption EXO. (Exogeneity) $(\{T(z)\}_{z \in \mathcal{Z}}, \{Y(t)\}_{t \in \mathcal{T}}) \perp\!\!\!\perp Z|X$.

We assume that each of T , Z , and X are discretely distributed with finite support. This is just for mathematical simplicity. Our theoretical results can be extended to allow for T to be a continuous scalar, and both X and Z to be vectors with continuous components. The changes required essentially involve replacing sums with integrals and finite indices with function arguments. We also assume throughout that the expectation of Y exists.

Our analysis uses a partition of individuals into mutually exclusive and exhaustive groups based on their potential treatment choices. Let $G \equiv (T(z_0), T(z_1), \dots, T(z_K))$ denote an individual's choice group, that is, their configuration of potential treatment choices under each of the instrument values. Let \mathcal{G} denote the values that G can take. In the binary treatment, binary instrument case, G takes values in $\mathcal{G} = \{(0, 0), (1, 1), (0, 1), (1, 0)\}$, corresponding to the groups Angrist et al. (1996, Table 1) called the never-takers, always-takers, compliers, and defiers, respectively. Using the group notation, Assumption EXO can be equivalently written as follows.

Assumption EXO. (Exogeneity, group form) $(G, \{Y(t)\}_{t \in \mathcal{T}}) \perp\!\!\!\perp Z | X$.

3.2 Definition of a weakly causal estimand

Consider the *group treatment responses* (GTRs)

$$\mu_j(g, x) \equiv \mathbb{E}[Y(t_j) | G = g, X = x],$$

which are the expected potential outcomes across choice and covariate groups.⁶ We collect the GTRs as $\mu \equiv \{\mu_j(g, x) : j = 0, 1, \dots, J, g \in \mathcal{G}, x \in \mathcal{X}\}$, which takes values in \mathbb{R}^{d_μ} . Let β be a quantity whose value depends on μ . We use the following definition as a minimal requirement for β to be interpreted as “causal.”

Definition WC. β is **weakly causal** if both of the following statements are true for any μ :

If $\mu_j(g, x) - \mu_{j-1}(g, x) \geq 0$ for all $j \geq 1$, all $g \in \mathcal{G}$, and every $x \in \mathcal{X}$, then $\beta \geq 0$.

If $\mu_j(g, x) - \mu_{j-1}(g, x) \leq 0$ for all $j \geq 1$, all $g \in \mathcal{G}$, and every $x \in \mathcal{X}$, then $\beta \leq 0$. (4)

Definition WC is a natural requirement to place on an estimand. The requirement is merely that *if* the causal effect of the treatment has the same sign for every treatment contrast, and every choice and covariate subgroup, then the summary estimand β also has that sign. That is, β is weakly causal if it is not a systematically misleading measure of the sign of the underlying group- and covariate-specific treatment effects.

⁶As a minor abuse of notation, we assume that $\mu_j(g, x)$ is well-defined for all (g, x) , even if g is not in the support of G given X , so that $\mathbb{P}[G = g, X = x] = 0$. This convention has no impact on our results.

Definition WC is intended to be an extremely weak criterion. An estimand can be weakly causal and still be completely uninteresting. For example, the trivial estimand $\beta = 0$ is weakly causal. However, it seems unlikely that an estimand that fails to be weakly causal could still reasonably be described as reflecting the causal effect of T on Y , since it may not even have the right sign. As minimal as Definition WC is, we have already seen in Figure 1 that a linear IV estimand can fail to satisfy it, even if the instrument satisfies exclusion and exogeneity (Assumption EXO).

3.3 Weak causality and non-negatively weighted averages

We consider estimands that can be written as

$$\beta = \mathbb{E}[b(T, X, Z)Y] \quad (5)$$

for some function b . For example, β_{iv} in Section 2 satisfies (5) with $b(T, X, Z) = \tilde{Z} / \mathbb{E}[T\tilde{Z}] = (Z - \mathbb{L}[Z|X]) / \mathbb{E}[T(Z - \mathbb{L}[Z|X])]$. The following proposition decomposes these estimands into GTRs.

Proposition 3. Suppose that β has form (5), and that Assumption EXO holds. Then

$$\beta = \sum_{g,x} \omega_0(g, x) \mu_0(g, x) + \sum_{g,x} \sum_{j=1}^J \omega_j(g, x) (\mu_j(g, x) - \mu_{j-1}(g, x)), \quad (6)$$

where $\omega_j(g, x) \equiv \mathbb{E}[\mathbb{1}[T \geq t_j] b(t_j, x, Z) | G = g, X = x] \mathbb{P}[G = g, X = x]$ for all $j \geq 0$.

The next proposition shows that an estimand that is weakly causal can be written as a non-negatively weighted average of subgroup-specific treatment effects. This criterion is widely-used (e.g. Angrist, 1998; Lee, 2008; Angrist and Pischke, 2009; Card et al., 2015; Goodman-Bacon, 2021; Sun and Abraham, 2021; Goldsmith-Pinkham et al., 2024). The proposition also shows that there are two reasons that an estimand can fail to be weakly causal: either it places negative weights on treatment effects or it depends on the levels of the potential outcomes (or both).

Proposition 4. Suppose that β has the form (5) and that Assumption EXO holds. Then β is weakly causal if and only if:

- **(Non-negative weights)** $\omega_j(g, x) \geq 0$ for all $j \geq 1$, and all g and x .
- **(Level independence)** $\omega_0(g, x) = 0$ for all g and x .

If these conditions are satisfied, then

$$\beta = \sum_{g,x} \sum_{j=1}^J \omega_j(g, x) (\mu_j(g, x) - \mu_{j-1}(g, x)) \quad (7)$$

for non-negative weights $\omega_j(g, x) \geq 0$.

Proposition 3 shows that if β has form (5), then β can always be written as (6). Proposition 4 uses that representation to show that if β cannot also be written like (7) with weights that are non-negative, then one of two things must be true: either β only reflects treatment effects, but some of these effects are negatively weighted, or else β reflects not just treatment effects but also the levels of potential outcomes. The first situation violates the non-negative weights requirement, which is naturally necessary for β to be weakly causal (recall Figure 1). The second situation violates the level independence requirement. Level independence is necessary for β to be weakly causal because it prevents the possibility that all treatment effects are positive, even while the levels of the GTRs are such that $\beta < 0$.

4 When is TSLS weakly causal?

In this section we specialize the general results of the previous section to a class of TSLS estimands.

4.1 TSLS specifications and estimands

A TSLS specification is characterized by four components: (i) the outcome variable; (ii) the endogenous variables that are included in the second stage and are regressands in the first stage; (iii) the excluded variables that are excluded from the second stage but are regressors in the first stage; and (iv) the included variables that are regressors in both stages. The nonparametric IV model specifies the outcome variable, Y , but not which combinations of T , Z , and X go in the first and second stages. A TSLS specification produces a TSLS estimator, the probability limit of which is called the TSLS estimand.

We consider TSLS specifications where there is a single endogenous variable, T , a single scalar excluded variable, Z , and a vector of included variables, X . For this case, the TSLS estimand is the same as the linear IV estimand because Z and T have the same dimension. We consider more general TSLS specifications in Appendix SA.2. In what follows, we reserve the phrase TSLS for specifications with strictly more excluded variables than endogenous variables.

The coefficient on T for the linear IV (née TSLS) estimand with a single endogenous variable and a single excluded variable is given by

$$\beta_{\text{iv}} = \frac{\mathbb{E}[\tilde{Z}Y]}{\mathbb{E}[\tilde{Z}T]} = \mathbb{E} \left[\left(\frac{\tilde{Z}}{\mathbb{E}[\tilde{Z}T]} \right) Y \right]. \quad (8)$$

Proposition 3 shows that β_{iv} can be written as (6) with

$$\omega_j(g, x) = \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E} \left[\mathbf{1}[T \geq t_j] \tilde{Z} | G = g, X = x \right] \mathbb{P}[G = g, X = x]. \quad (9)$$

Proposition 4 shows that whether β_{iv} is weakly causal is determined by $\omega_j(g, x)$.

4.2 Main result

Monotonicity conditions are essential for TSLS estimands to have weakly causal interpretations. We maintain the following monotonicity condition in the main text.

Assumption MON. (Monotonicity) Label the values of Z in increasing order as $z_0 \leq z_1 \leq \dots \leq z_K$. Then

$$\mathbb{P}[T(z_0) \leq T(z_1) \leq \dots \leq T(z_K) | X = x] = 1 \quad \text{for all } x.$$

Assumption MON means that increasing the instrument weakly increases treatment for all individuals. This is a strong form of monotonicity because it operates in the same direction conditional on $X = x$ for all values of x . Results under weaker forms of monotonicity can be found in Appendix SA.2.

Our main result is Theorem 1, which uses the following definition.

Definition RC. Let $\mathbb{L}[Z|X = x] \equiv \mathbb{E}[ZX'] \mathbb{E}[XX']^{-1}x$ be the population fitted value at $X = x$ from regressing Z onto X . An IV specification has **rich covariates** if $\mathbb{E}[Z|X = x] = \mathbb{L}[Z|X = x]$ for every $x \in \mathcal{X}$.

Theorem 1. Suppose that Assumptions EXO and MON are satisfied. Then β_{iv} is weakly causal if *and only if* the IV specification has rich covariates.

Theorem 1 shows that rich covariates is both sufficient and necessary for the linear IV estimand to have a weakly causal interpretation. The sufficient direction is a special case of Kolesár (2013, Theorem 1), who explicitly maintained rich covariates as an *assumption* (Kolesár, 2013, Assumption L). The necessary direction shown here is novel. It shows that rich covariates is an essential assumption.

As Kolesár (2013, pp. 10–11) notes, there are two important special cases in which an IV specification will have rich covariates. One is when X represents a saturated specification consisting of a vector of indicators for a set of exclusive and mutually exhaustive events. The other is when Z is mean independent of X so that $\mathbb{E}[Z|X = x] = \mathbb{E}[Z]$ is constant in x , which implies that

$$\mathbb{L}[Z|X = x] \equiv \mathbb{E}[ZX'] \mathbb{E}[XX']^{-1}x = \mathbb{E}[Z] \mathbb{E}[1X'] \mathbb{E}[XX']^{-1}x = \mathbb{E}[Z],$$

because X contains a constant. Outside these two special cases, the claim that an IV specification has rich covariates is an implicit parametric assumption. Theorem 1 shows that this parametric assumption must be defended in order to argue that β_{iv} has a causal interpretation.

4.3 Implications for OLS under selection on observables

Theorem 1 also applies to selection on observables by taking $Z = T$, under which Assumption MON is immediately satisfied.

Corollary 1. Suppose that $Z = T$ and that Assumption EXO is satisfied. Let β_{ols} denote the coefficient on T for the OLS estimand generated by the ordinary least squares regression of Y on T and X . Then β_{ols} is weakly causal if and only if $\mathbb{L}[T|X] = \mathbb{E}[T|X]$.

[Angrist \(1998\)](#) proposed implementing a selection on observables strategy using the OLS estimator described in Corollary 1 with a saturated specification of covariates. He described the difference between this regression coefficient and nonparametric matching as “partly cosmetic” ([Angrist, 1998](#), pg. 255). Based on these results, [Angrist and Pischke \(2009, Section 3.3.1\)](#) argue that “the differences between regression and matching are unlikely to be of major empirical importance.”

However, Corollary 1 shows that [Angrist’s \(1998\)](#) argument cannot be extrapolated beyond the saturated case that he considered. The result implies that any deviation from full saturation will mean that the OLS estimand fails to be weakly causal unless one assumes that the propensity score $\mathbb{P}[T = 1|X = x] = \mathbb{E}[T|X] = \mathbb{L}[T|X]$ is actually linear in X . Moreover, whenever [Angrist’s \(1998\)](#) saturated specification can actually be implemented, the overlap condition $\mathbb{P}[T = 1|X = x] \in (0, 1)$ must hold for every x , or else there would be perfect collinearity.

The implication of Corollary 1 then is that there are only two situations in which [Angrist’s \(1998\)](#) linear regression implementation of selection on observables will be weakly causal. First, when the propensity score is implicitly assumed to be linear. Second, when it is also possible to nonparametrically estimate conditional average treatment effects x -by- x . The first case involves a parametric assumption, while in the second case one could just as well weight the x -by- x treatment effects into a parameter such as the average treatment effect that is not only weakly causal but also has a clear counterfactual interpretation. Outside these two cases, β_{ols} is not weakly causal.

4.4 Specifications with more general excluded variables

[Śłoczyński \(2020, 2024\)](#) considers the interpretation of TSLS estimands with a binary treatment and a binary instrument under the assumption that the specification has rich covariates. He considers both Assumption MON, which he calls “strong” monotonicity,

and a “weak monotonicity” counterpart in which the direction of monotonicity can vary with x . [Słoczyński \(2020, 2024\)](#) shows that if Assumption MON fails, but weak monotonicity is satisfied, then β_{iv} will not be weakly causal even if the specification has rich covariates. The problem is that the IV specification includes only a single excluded variable, Z , so is not flexible enough to pick up changes in monotonicity in the first stage. [Słoczyński \(2020, 2024\)](#) shows that this problem can be resolved by using the “saturate and weight” TSLS specification in [Angrist and Pischke \(2009\)](#), which includes interactions between X and Z as excluded variables.

[Kolesár \(2013\)](#) provided general sufficient conditions for the TSLS estimand to be interpreted as a non-negatively weighted average of treatment effects under weak monotonicity. As noted above, [Kolesár \(2013\)](#) maintained rich covariates as an assumption, whereas we show that rich covariates is a necessary condition. [Kolesár \(2013\)](#) showed that given rich covariates, the TSLS estimand can be written as a weighted average of treatment effects. Whether the weights are non-negative depends on whether the first stage equation is able to sufficiently well approximate the nonparametric propensity score. In Appendix SA.2, we provide a lower-level sufficient-and-necessary characterization of when the weights are non-negative in terms of the first stage specification being “monotonicity-correct.” The takeaway from that characterization reinforces [Słoczyński](#)’s findings that even when rich covariates is satisfied, an additional necessary condition for TSLS to be weakly causal is that the first stage is sufficiently flexible to reproduce the direction of monotonicity across covariate groups.

The rich covariates condition extends readily to more general types of excluded variables. Suppose that the excluded variables are a vector $i(Z, X)$ with population first stage coefficient vector γ . Let $\dot{Z} \equiv \gamma' i(Z, X)$. In Appendix SA.2 we show that a necessary condition for the resulting TSLS estimand to be weakly causal is that $\mathbb{E}[\dot{Z}|X = x] = \mathbb{L}[\dot{Z}|X = x]$ for all x , a condition that naturally generalizes the case considered here with $i(Z, X) = Z$ scalar. This condition has basically the same content as Definition RC, but involves the aggregate \dot{Z} instead of just Z itself.

4.5 Constant, linear treatment effects

Suppose we assume that treatment effects are constant and linear.

Assumption CLE. (Constant, linear effects) There exists a constant Δ such that $\mu_j(g, x) - \mu_{j-1}(g, x) = \Delta(t_j - t_{j-1})$ for every $j \geq 1$, $g \in \mathcal{G}$ and $x \in \mathcal{X}$.

Theorem 1 continues to hold under Assumption CLE, except Assumption MON no longer needs to be maintained.

Proposition 5. Suppose that Assumptions EXO and CLE are satisfied. Then β_{iv} is weakly causal if *and only if* the IV specification has rich covariates.

The assumptions of Proposition 5 allow for a simple illustration of the level dependence phenomenon. Assumption CLE implies that $Y = Y(t_0) + \Delta T$, so

$$\beta_{\text{iv}} = \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E}[\tilde{Z}(Y(t_0) + \Delta T)] = \Delta + \mathbb{E}[\tilde{Z}T]^{-1} \overbrace{\mathbb{E}[\tilde{Z}Y(t_0)]}^{\text{depends on } Y(t_0)}. \quad (10)$$

Using Assumption EXO, the potentially level dependent term can be written as

$$\mathbb{E}[\tilde{Z}Y(t_0)] = \mathbb{E} \left[\mathbb{E}[\tilde{Z}|X] \mathbb{E}[Y(t_0)|X] \right]. \quad (11)$$

The nonparametric IV model does not restrict $\mathbb{E}[Y(t_0)|X]$ at all. Level dependence will therefore happen whenever $\mathbb{E}[\tilde{Z}|X] \neq 0$ with positive probability, which in turn happens whenever the IV specification does not have rich covariates.

Proposition 5 shows that the necessity of rich covariates does not have to do with heterogeneous or nonlinear treatment effects per se. Rather, it is a fundamental consequence of the exercise started by [Imbens and Angrist \(1994\)](#) of interpreting a linear IV estimand through a nonparametric IV model. The linear IV estimator was designed for the linear IV model; giving it a causal interpretation within a nonparametric IV model requires additional parametric assumptions.

Instead of that additional parametric assumption being rich covariates, one can maintain a parametric assumption on a conditional mean of the potential outcomes.

Assumption LIN. (Linear potential outcome mean) $\mathbb{E}[Y(t_j)|X = x] = \eta'x$ for some η and some j .

Proposition 6. Suppose that Assumptions EXO, CLE, and LIN are satisfied. Then $\beta_{\text{iv}} = \Delta$, so β_{iv} is weakly causal.

Assumption LIN—or something similar—is explicitly stated in classical and textbook treatments of IV models, e.g. [Heckman and Robb \(1985, pp. 184–186\)](#) or [Wooldridge \(2010, pg. 939\)](#). But it is not part of the nonparametric IV model that is used to justify the widely-invoked “LATE interpretation” of the linear IV estimator ([Angrist and Imbens, 1995](#)). As [Abadie \(2003, pg. 247\)](#) points out, an undesirable implication of Assumption LIN is that one can have $\beta_{\text{iv}} = \Delta$ even if the excluded “instrument” Z is actually some nonlinear function of X alone, an example of what [Angrist and Pischke \(2009, pg. 191\)](#) describe as “back-door identification.”

A higher-level alternative to having rich covariates or imposing Assumption LIN is to directly assume that the left-hand side of (11) is zero. This assumption appears in [Wooldridge’s \(2010, pg. 937\)](#) discussion of the binary treatment case as the assumption that $\mathbb{L}[Y(t_0)|X, Z]$ does not depend on Z . If we put aside knife-edge balancing cases, (11) shows that this assumption either requires rich covariates or Assumption

LIN. However, considering the high-level assumption usefully exposes the fundamental problem with using the nonparametric IV model to justify linear IV: Assumption EXO by itself *does not* imply that a hypothetical linear regression of $Y(t_0)$ onto X and Z would yield a zero coefficient on Z , even though this condition is essential for giving the linear IV estimand a causal interpretation.

Assumption LIN was also maintained in Propositions 1 and 2, which showed that β_{iv} is not weakly causal without rich covariates. This does not contradict Proposition 6 because of the addition of constant, linear treatment effects (Assumption CLE). When Assumption CLE is removed to allow for heterogeneous treatment effects, Assumption LIN no longer suffices as a substitute for rich covariates.

5 Alternatives to linear IV

5.1 Partially linear IV

Theorem 1 shows that β_{iv} is weakly causal if and only if the IV specification has rich covariates. If rich covariates is satisfied, it follows from (8) that $\beta_{iv} = \beta_{rich}$, where

$$\beta_{rich} \equiv \frac{\mathbb{E}[Y(Z - \mathbb{E}[Z|X])]}{\mathbb{E}[T(Z - \mathbb{E}[Z|X])]} = \frac{\mathbb{E}[\mathbf{C}[Y, Z|X]]}{\mathbb{E}[\mathbf{C}[T, Z|X]]}. \quad (12)$$

If rich covariates is not satisfied, then it may be that $\beta_{iv} \neq \beta_{rich}$, however we can still consider β_{rich} as what the IV estimand would have been had rich covariates actually been satisfied. Given Assumptions EXO and MON, β_{rich} always satisfies the minimal requirement of being weakly causal.

One way to estimate β_{rich} is to use a richer linear IV specification that controls for covariates so flexibly that rich covariates must be satisfied. If X is discrete, then this is the same as using a saturated specification with one dummy variable for each discrete value of X . These types of specifications are discussed in Angrist (1998) and Angrist and Pischke (2009), but were rarely used in the IV papers in our survey (Table 2). They take an extreme position on the bias-variance trade-off that is difficult to defend for settings in which X takes many values.

Chernozhukov et al. (2018) show how machine learning (ML) methods can be used to estimate a modification of the classical linear IV model where the linear function of covariates has been replaced by an unknown function. They describe the model as partially linear IV (PLIV). It is straightforward to show that the coefficient on treatment in their model is equal to β_{rich} . Chernozhukov et al. (2018) show how to construct the Neyman orthogonal score for the PLIV model, which depends on the

treatment coefficient as well as the functions

$$\nu(x) \equiv (\mathbb{E}[Y|X = x], \mathbb{E}[T|X = x], \mathbb{E}[Z|X = x]). \quad (13)$$

They then show how to use the orthogonality of the score in conjunction with cross-fitting to construct consistent and asymptotically normal estimators of the treatment coefficient under nonparametric assumptions about the unknown functions that comprise ν . They suggest estimating ν using supervised ML algorithms such as random forests, gradient boosted trees, and neural networks.

5.2 Unconditional average causal response

Proposition 4 showed that any weakly causal estimand, such as β_{rich} , can be written as a non-negatively weighted average of subgroup treatment effects. Given rich covariates, the general form of the weights in (37) becomes

$$\omega_j(g, x) = \mathbb{E}[\tilde{Z}T]^{-1} \mathbf{C} [\mathbb{1}[T \geq t_j], Z|G = g, X = x] \mathbb{P}[G = g, X = x]. \quad (14)$$

While (14) has a reasonable statistical interpretation—larger groups and contrasts that covary more with the instrument get more weight—it does not appear to have a more concrete counterfactual interpretation. One obstacle is that the instrument can be multivalued, which even without covariates turns the linear IV estimand into a statistically-weighted average of treatment effects across different complier groups (Imbens and Angrist, 1994). If the instrument $Z \in \{0, 1\}$ is binary, then a parameter that does have a concrete counterfactual interpretation is the unconditional average causal response (ACR) (Angrist and Imbens, 1995):

$$\beta_{\text{acr}} \equiv \mathbb{E}[Y(T(1)) - Y(T(0))|T(1) > T(0)]. \quad (15)$$

Note that β_{acr} is the LATE when $T \in \{0, 1\}$ is binary, so that $T(1) = 1$ and $T(0) = 0$.

As Słoczyński (2020, 2024) points out, the difference between β_{rich} and β_{acr} can be large. To see why, let $\beta_{\text{acr}}(x) \equiv \mathbb{E}[Y(T(1)) - Y(T(0))|T(1) > T(0), X = x]$ be the ACR conditional on $X = x$. Then iterating expectations shows that

$$\beta_{\text{acr}} = \mathbb{E} \left[\beta_{\text{acr}}(X) \frac{\mathbb{P}[T(1) > T(0)|X]}{\mathbb{P}[T(1) > T(0)]} \right], \quad (16)$$

whereas, with a bit of algebra, it can also be shown that

$$\beta_{\text{rich}} = \mathbb{E} \left[\beta_{\text{acr}}(X) \frac{\mathbb{P}[T(1) > T(0)|X] \mathbb{V}[Z|X]}{\mathbb{E}[\mathbb{P}[T(1) > T(0)|X] \mathbb{V}[Z|X]]} \right]. \quad (17)$$

The difference between β_{acr} and β_{rich} arises because the latter puts extra weight on values of X with more variation in Z . [Słoczyński \(2020, 2024\)](#) argues that β_{acr} is likely what empirical researchers have in mind, and he shows that the difference in weights can make β_{rich} misleading. So, even if rich covariates holds, β_{tsls} may not be what an empirical researcher expects. In Section 6, we find empirical evidence that β_{rich} and β_{acr} can be quite different.

We produce this evidence by directly estimating β_{acr} using two different approaches. Both approaches are based on a finding due to [Tan \(2006\)](#) and [Frölich \(2007\)](#) that

$$\beta_{\text{acr}} = \frac{\mathbb{E}[\mathbb{E}[Y|Z = 1, X] - \mathbb{E}[Y|Z = 0, X]]}{\mathbb{E}[\mathbb{E}[T|Z = 1, X] - \mathbb{E}[T|Z = 0, X]]}. \quad (18)$$

The first approach comes from [Chernozhukov et al. \(2018\)](#), who show how to estimate β_{acr} using DDML. The orthogonal score that they derive involves the five functions

$$\nu(x) \equiv (\mathbb{E}[Y|Z = 0, X = x], \mathbb{E}[Y|Z = 1, X = x], \mathbb{E}[T|Z = 0, X = x], \\ \mathbb{E}[T|Z = 1, X = x], \mathbb{E}[Z|X = x]),$$

all of which can be estimated nonparametrically using ML algorithms. The second approach comes from [Uysal \(2011\)](#), [Heiler \(2022\)](#), and [Słoczyński et al. \(2024\)](#), who exploit the connection that (18) has with propensity score weighting: the numerator looks like the average treatment effect of Z on Y and the denominator like the average treatment effect of Z on T . They propose a weight-normalized inverse propensity score estimator and derive its asymptotic properties. Implementing the estimator requires parameterizing the instrument propensity score, $\mathbb{E}[Z|X]$.⁷

5.3 Abadie’s (2003) κ

[Abadie \(2003, Section 4.2.1\)](#) and [Angrist and Pischke \(2009, pp. 179–180\)](#) suggest using a weighted regression to control for covariates when both $T \in \{0, 1\}$ and $Z \in \{0, 1\}$ are binary. As [Abadie \(2003\)](#) showed, a weighted linear regression of Y on T and X with weights given by

$$\kappa \equiv 1 - \frac{T(1 - Z)}{1 - \mathbb{E}[Z|X]} - \frac{(1 - T)Z}{\mathbb{E}[Z|X]} \quad (19)$$

is, in the population, the same as an unweighted linear regression of Y on T and X among the subpopulation of compliers. [Abadie \(2003\)](#) showed that if rich covariates hold, then the κ -weighted estimate of the coefficient on T is numerically equivalent to

⁷See also [MaCurdy et al. \(2011\)](#), [Donald et al. \(2014\)](#), [Ogburn et al. \(2015\)](#), [Sun and Tan \(2022\)](#), and [Singh and Sun \(2024\)](#) for related estimators.

the linear IV estimate, so estimates a weakly causal estimand.

The next proposition shows that rich covariates turns out to also be *necessary* for the κ -weighted estimand to be weakly causal.

Proposition 7. Suppose that Assumptions EXO and MON are satisfied and that both T and Z are binary. Let β_{abadie} be the population estimand for a weighted linear regression of Y on T and X with weights given by κ . Then β_{abadie} is weakly causal if and only if the linear IV specification has rich covariates, so that $\mathbb{E}[Z|X] = \mathbb{L}[Z|X]$. When this is true, $\beta_{\text{abadie}} = \beta_{\text{iv}} = \beta_{\text{rich}}$.

Proposition 7 shows that when rich covariates holds, the κ -weighting estimand is equal to the linear IV estimand, so there is no reason to prefer it, as the linear IV estimand is simpler.⁸ When rich covariates does not hold, both the linear IV and κ -weighting estimands are not weakly causal.

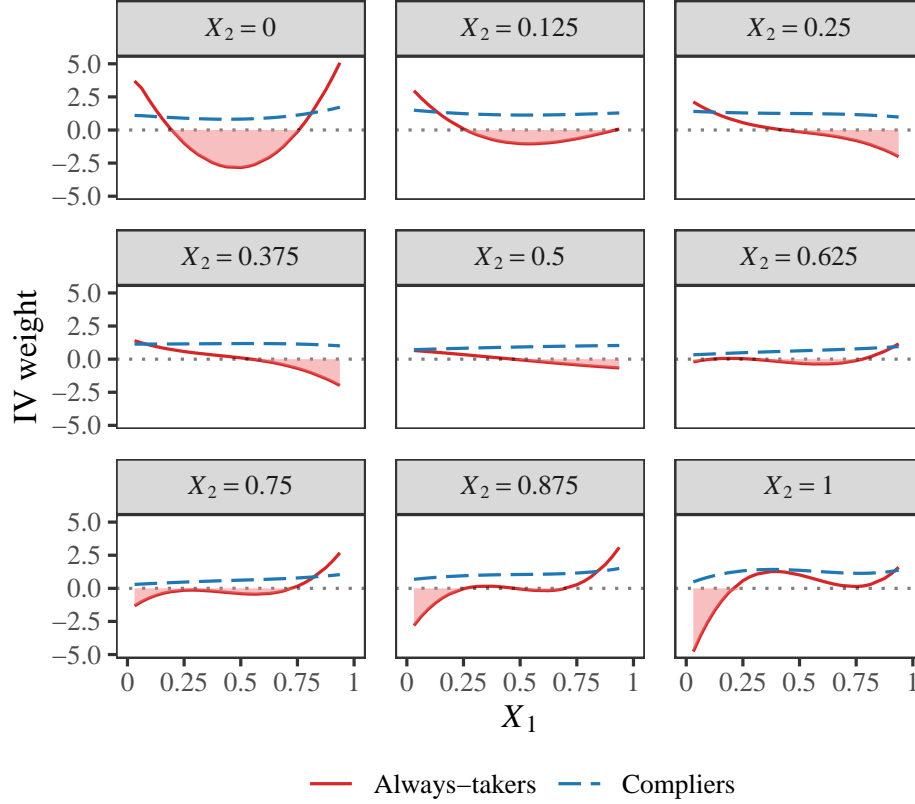
Angrist and Pischke (2009, pp. 180–181) use Angrist’s (2001) reanalysis of Angrist and Evans (1998) as an example to dismiss the relevance of Abadie’s (2003) approach. Yet Angrist (2001, pg. 12) also reports that “the covariates are not highly correlated with the twins instruments...” Our findings show why it is misleading to extrapolate the Angrist and Pischke (2009) argument to other empirical settings: the case when Z is mean independent of X is one where *any* covariate specification is rich. If Z and X are dependent—as is often the case when covariates are used in an IV analysis—then the linear IV estimand will not have a complier interpretation unless $\mathbb{E}[Z|X = x]$ is modeled correctly. At the same time, Proposition 7 also implies that the implementation of Abadie’s κ proposed by Angrist and Pischke (2009, pp. 180–181) *only* has a causal interpretation when the IV specification has rich covariates.

5.4 Monte Carlo simulation

In this section, we report the results of a Monte Carlo simulation based on a data generating process (DGP) calibrated to Card’s (1995) data on the returns to schooling, which we reanalyze in the next section. We use covariates $X \equiv (X_1, X_2)$, where X_1 takes a number of values that we vary across simulations, while X_2 always takes nine values. We generate the binary instrument Z —presence of a nearby college in Card’s application—by specifying $\mathbb{E}[Z|X = x]$ to be an interacted cubic polynomial fit to the Card data with X_1 as experience and X_2 as region indicators (Figure SA.1). We generate the binary treatment T (college attendance) so that $\mathbb{P}[T = 1|Z = z]$ matches its estimated counterpart in Card’s data. Then, we generate the outcome Y (log wages) using an optimization procedure that matches several estimates in Card’s data while

⁸These considerations are about the estimand; they do not take into account differences in the statistical properties of the linear IV and κ -weighting estimators.

Figure 2: Weights for β_{iv} in the simulation DGP



Notes: The figure shows the weights in Proposition 1 for the linear IV estimand when X_1 takes 24 values. The weights vary by both choice group and $X = (X_1, X_2)$. The weights for the compliers are always non-negative, but the weights for the always-takers are often negative, as shown in shaded red. The decomposition underlying the figure is not unique (Proposition 2). Figure SA.2 shows the analogous figure for a decomposition involving only compliers and never-takers.

also ensuring that Assumption LIN is satisfied, as in the simplified case discussed in Section 2. See Appendix SA.3 for more details.

We use this DGP to compare the performance of five estimators.

The first is a linear IV estimator that controls for X_1 linearly while including a full set of indicators for X_2 , but omits any nonlinear or interaction terms. This specification does not satisfy rich covariates, so is not weakly causal. Figure 2 illustrates the weights for the estimand β_{iv} of this estimator using the Proposition 1 decomposition into compliers and always-takers, for the case when X_1 takes 24 values. All complier groups are positively-weighted. However, always-takers receive considerable weight, both positive and negative. The overall value of β_{iv} is .660, which reflects the sum of .391 from positively-weighted compliers, .614 from positively-weighted always-takers, and $-.345$

from negatively-weighted always-takers. Figure SA.2 shows that writing β_{iv} in terms of complier and never-takers instead of always-takers also leads to negative weights, as shown in Proposition 2. The second estimator is Abadie’s κ -weighting estimator using the same covariate specification, which Proposition 7 showed is also not weakly causal. The estimand for the κ -weighting estimator, β_{abadie} , is very similar to that for β_{iv} regardless of how many values X_1 takes (Figure SA.3).

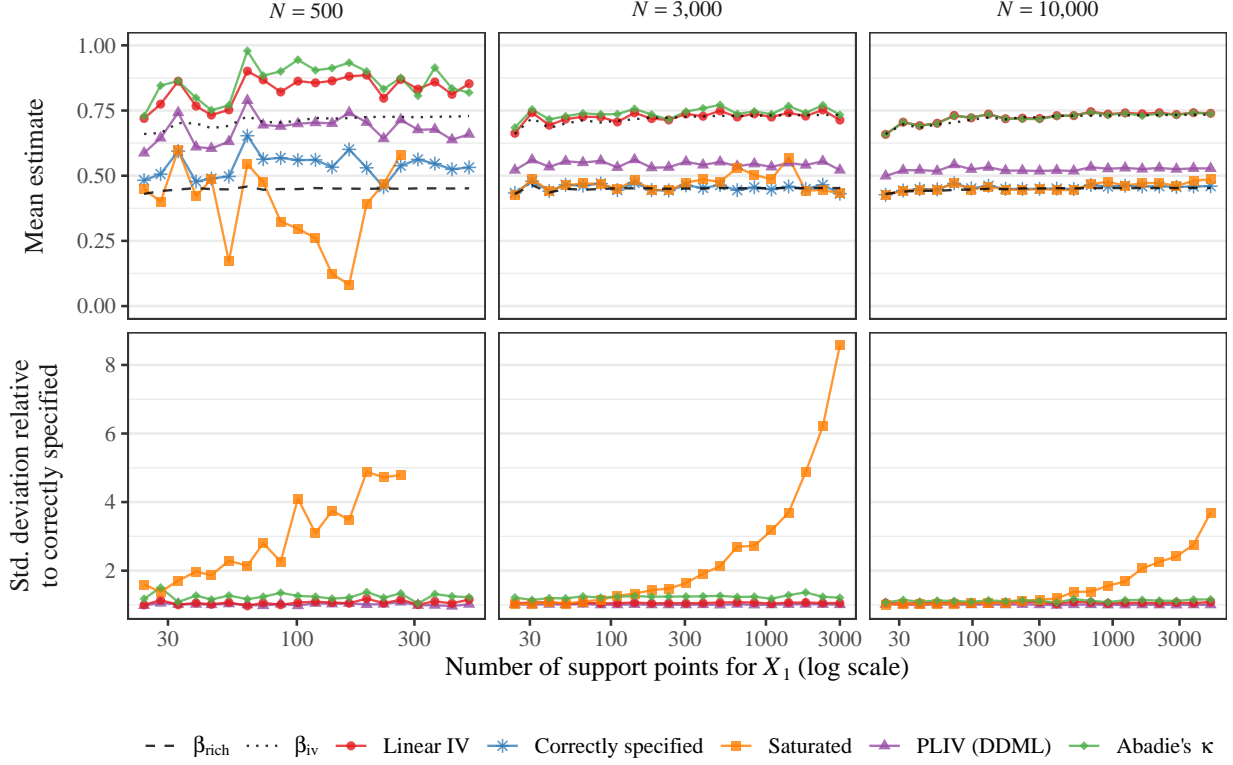
The third estimator is a linear IV estimator that includes nonlinear and interaction terms, so that rich covariates is satisfied. We call this estimator “correctly specified,” as it assumes that X has been chosen so that $\mathbb{L}[Z|X] = \mathbb{E}[Z|X]$. The fourth estimator is a linear IV estimator that is saturated in X . This estimator also satisfies rich covariates, but the number of regressors it uses increases with the support of X_1 , which we will vary in the simulation. The fifth estimator is a DDML estimator for the PLIV model using an ensemble of a random forest with 1000 trees and three variables at each split, gradient boosted trees with 1000 stages, and a neural network with two neurons.⁹ Each of these three estimators can be viewed as estimating $\beta_{rich} = .430$, which is a weakly causal estimand comprised of only non-negatively weighted complier effects.

The top row of Figure 3 compares the means of the five estimators, with some more detailed results reported in Table SA.1. The facet columns of Figure 3 vary the sample size while the x-axis varies the size of the support of X_1 . The linear IV estimator converges to the negatively-weighted estimand β_{iv} , so it exhibits a bias for β_{rich} that does not decrease with the sample size. The correctly specified and saturated estimators both converge to β_{rich} , however when the sample size is small relative to the number of covariate values, the saturated estimator exhibits substantial bias. The bias of the DDML-PLIV estimator for β_{rich} is larger than the correctly specified estimator, but decreases as the sample size increases. Table SA.2 shows that using more expressive algorithms by themselves (without an ensemble) can eliminate the bias, but this comes with the risk of using a very poor-performing algorithm, especially with smaller sample sizes.

The bottom row of Figure 3 compares the standard deviations of the five estimators. The comparison is taken relative to the correctly-specified estimator to keep the magnitude comparable across sample sizes. The linear, κ -weighting, correctly specified, and DDML-PLIV estimators all exhibit broadly similar standard deviations across sample

⁹ The ensemble is formed by short-stacking with convex weights chosen through non-negative least squares, as advocated by Ahrens et al. (2023, 2024b). The DDML estimates are random due to the sample splits used in cross-fitting. We repeat each estimate across five sample splits and report the resulting median, as recommended by Chernozhukov et al. (2018). We implemented our simulations in R using the `ddml` package (Ahrens et al., 2024a). We used the following packages on the back-end of `ddml` to implement machine learning algorithms: `gbm` for gradient-boosting (Ridgeway, 2007), `nnet` for neural networks (Ripley and Venables, 2016), and `ranger` for random forests (Wright and Ziegler, 2017). In Table SA.2, we also report some results using lasso as implemented by the `glmnet` package (Friedman et al., 2008).

Figure 3: Simulation results: bias and standard deviation

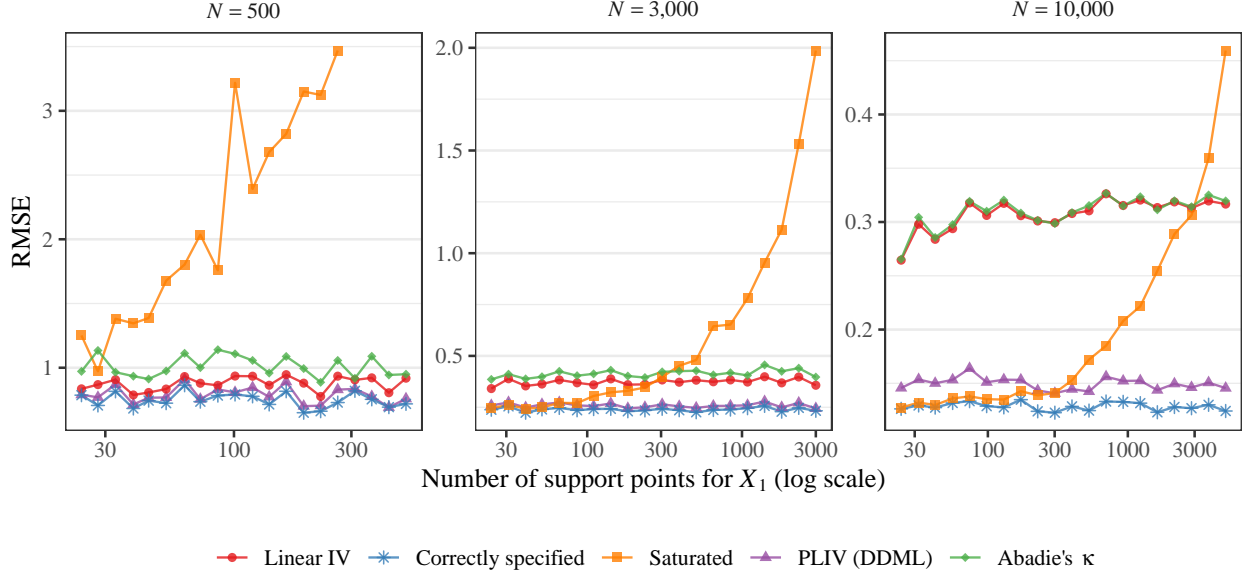


Notes: Each point is constructed from 500 draws from the DGP discussed in Appendix SA.3. Means and standard deviations are computed after trimming the top and bottom .025 quantiles of the distribution. Results for the saturated estimator for $N = 500$ and more than 300 support points have many undefined draws so are excluded.

sizes and number of covariate values. The flexibility of these estimators does not depend on the number of values that the covariates take, so increasing the support of X_1 does not have a large impact on their standard deviations. In contrast, the standard deviation of the saturated linear IV estimator explodes as the number of covariates increases.

Figure 4 summarizes these findings by reporting the root mean-squared error (RMSE) of the five estimators. The saturated specification performs well when the number of covariate values is small relative to the sample size, but poorly when the number of covariate values is moderate or large. This is likely the reason that saturated specifications are used so infrequently in the literature (Section 2.4), as even a small number of distinct discrete covariates leads to a saturated specification with a large number of covariate values due to the curse of dimensionality. The ideal, but infeasible, solution would be to use a linear IV estimator in which $\mathbb{E}[Z|X]$ is known to be correctly specified, so that rich covariates is satisfied. Without such knowledge, this solution en-

Figure 4: Simulation results: root mean-squared error



Notes: See notes for Figure 3. The bias component of the root mean-squared error (RMSE) is calculated relative to β_{rich} . Note that the scale of the y-axis changes across the panels.

tails an assumption that the specification is in fact correct. If assuming rich covariates is unattractive, then the DDML-PLIV estimator provides a feasible alternative that can be viewed as nonparametrically estimating the weakly causal quantity β_{rich} . Figure 4 shows that in our simulation the DDML-PLIV estimator is comparable in terms of RMSE to the (infeasible) correctly specified linear IV estimator.

6 Applications

In this section, we use our findings to reanalyze several empirical studies. We begin with [Card's \(1995\)](#) estimates of the returns to education as a classic and familiar example. Next, we turn to the papers by [Nunn and Wantchekon \(2011\)](#) and [Dube and Harish \(2020\)](#) as more modern examples of how linear IV is applied and interpreted in practice. Finally, we reanalyze the main estimates for ten studies from our survey in Section 2.4.

For all studies, we reproduce the original linear IV estimates alongside their comparable OLS estimates. We conduct a [Ramsey \(1969\)](#) RESET test of the null hypothesis that $\mathbb{E}[Z|X]$ is linear in X , that is, of rich covariates.¹⁰ We then estimate β_{rich} with DDML using the same ensemble as in Section 5.4. Standard errors for all estimators are heteroskedasticity and/or cluster-robust depending on the original application.¹¹

¹⁰We implement the RESET test using the second and third orders of the fitted values.

¹¹For the DDML estimates, we still report the median estimate, but we use 100 random sample splits instead of five as in the simulations. The DDML standard errors include an adjustment for uncertainty over

6.1 Card (1995)

Card (1995) used a sample of 24–34-year-old men from the 1976 interview of the NLSY to estimate the returns to education. The outcome Y is log hourly wage. The treatment T is years of education. The instrument Z is a binary indicator for the presence of an accredited four-year college in the local labor market when the respondent was 14 years old. In his main results, Card (1995, Table 3A, column (5)) includes the following covariates as X : a quadratic in years of potential experience, a race indicator for Black, geography indicators for living in the South and in an urban area, a set of indicators for region of residence in 1966, and an indicator for residence in an urban area in 1966. All of these terms enter additively, so the covariate specification is not saturated and might not satisfy rich covariates.

Column (1) of Table 3 reproduces Card’s IV estimate of the returns to education and OLS estimates of the comparable OLS estimand, β_{ols} , that instruments for T with itself. The original linear IV estimate of .132 uses covariates and increases by about 30% if the covariates are omitted. The RESET test overwhelmingly rejects the null hypothesis that the specification has rich covariates. By Theorem 1, this is the same as rejecting the null hypothesis that β_{iv} has a weakly causal interpretation. It doesn’t necessarily imply that β_{iv} is not equal to β_{rich} , a quantity which does have a weakly causal interpretation. However, it is important to keep in mind the simple point that two estimands can be equal even if one has a causal interpretation and the other does not; the estimates themselves say nothing without an underlying theory to justify their interpretation.

The DDML estimate of β_{rich} reported in the fourth row is modestly smaller than the IV estimate of β_{iv} , with a similar standard error. Some perspective on the magnitude of this difference is given in the row titled relative specification bias, where we report an estimate of $|\beta_{\text{iv}} - \beta_{\text{rich}}|/|\beta_{\text{iv}}|$ at about .076, or roughly an 8% difference. The subsequent row reports an estimate of $|\beta_{\text{iv}} - \beta_{\text{rich}}|/|\beta_{\text{ols}} - \beta_{\text{rich}}|$, which at roughly 21% shows that the difference between β_{iv} and β_{rich} represents a sizable fraction of the “selection bias” between OLS and the DDML estimate.

The sixth row reports DDML estimates of β_{acr} . While both β_{rich} and β_{acr} are weakly causal, the DDML estimate of β_{acr} is roughly half the size of the DDML estimate of β_{rich} , with a comparable standard error. This difference likely reflects the difference in weights discussed in Section 5.2, providing an empirical illustration of a critique made by Słoczyński (2024). In the sixth row, we report an alternative estimate of β_{acr} that uses the normalized instrument-propensity score weighting (IPSW) estimator proposed

splits (Chernozhukov et al., 2018, pg. C30). In the applications, we used the Stata `ddml` package (Ahrens et al., 2023) together with `pystacked` (Ahrens et al., 2022) to implement machine learning algorithms from `scikit-learn` (Pedregosa et al., 2011).

Table 3: Comparison of IV estimates for three applications

	(1)	(2)	(3)
	Card (1995)	Nunn & Wantchekon (2011)	Dube & Harish (2020)
OLS	0.075 (0.004)	−0.203 (0.033)	0.115 (0.035)
IV, no covariates	0.188 (0.026)	−0.190 (0.111)	1.011 (0.522)
IV, with covariates	0.132 (0.054)	−0.271 (0.088)	0.400 (0.211)
PLIV (DDML)	0.122 (0.053)	−0.071 (0.091)	0.318 (0.240)
Abadie’s κ	—	—	−0.404 (4.711)
LATE/ACR (DDML)	0.067 (0.046)	—	0.203 (0.141)
LATE/ACR (IPSW)	0.085 (0.052)	—	0.573 (2.453)
Ramsey RESET test p -val. (H_0 : rich covariates)	0.000	0.000	0.000
Relative specification bias	0.076	0.738	0.204
Specification vs. selection bias	0.213	1.515	0.400
Outcome variable	log(hourly wage)	Trust in neighbors	At war
Outcome variable, mean	6.262	1.732	0.296
Treatment variable	Years of education	log(1 + slave exports)	Queen ruling
Treatment variable, mean	13.263	0.621	0.160
Included variables	14	99	66
Sample size	3,010	16,679	3,586

Notes: Heteroskedasticity- or cluster-robust standard errors are reported in parentheses. Standard errors for Abadie’s κ are block bootstrapped with the top and bottom .5% of the bootstrap distributed trimmed. Standard errors for IPSW are also computed through block bootstrap. LATE estimates are not reported in column (2) because the instrument is not binary. Relative specification bias is an estimate of $|\beta_{iv} - \beta_{rich}|/|\beta_{iv}|$ and specification vs. selection bias is an estimate of $|\beta_{iv} - \beta_{rich}|/|\beta_{ols} - \beta_{rich}|$. Estimates of $\mathbb{E}[Z|X]$ are trimmed to [.01, .99] for the [Dube and Harish \(2020\)](#) application when using Abadie’s κ and IPSW.

by [Uysal \(2011\)](#), [Heiler \(2022\)](#), and [Śloczyński et al. \(2024\)](#). The estimate and standard error are quite similar to those for DDML.

6.2 Two modern applications

We now turn to two more recent examples. These examples are explicit in their use of an extensive set of covariates to justify the exogeneity of the instrument.

We first consider an influential paper by [Nunn and Wantchekon \(2011\)](#), who estimate the effect of the slave trade on modern day measures of trust in Africa using data from the 2005 Afrobarometer survey. The outcome Y is the respondent’s reported level of

trust in their neighbors. The treatment T is the natural log of (one plus) total historical slave exports for the respondent’s ethnic group. The instrument Z is the historical distance of the respondent’s ethnic group from the nearest coast. In their main results, [Nunn and Wantchekon \(2011\)](#) include as covariates X a set of country fixed effects, a set of demographic controls for the respondent, measures of ethnic homogeneity for the respondent’s district, a set of variables intended to proxy for the amount of European influence, distance of the ethnic group’s historical homeland to the Saharan slave trade, and a historical measure of the ethnic group’s reliance on fishing. In total, there are 93 covariates. The authors are explicit that their motivation for incorporating these covariates is to help ensure the exogeneity of their instrument ([Nunn and Wantchekon, 2011](#), pg. 3239).

Column (2) of Table 3 reproduces the IV estimates from column (2) of Table 6 in [Nunn and Wantchekon \(2011\)](#) alongside the comparable OLS estimate. The RESET test again overwhelmingly rejects the null of rich covariates. In this case, the IV estimate of β_{iv} is almost four times as large as the DDML estimate of β_{rich} , representing one and a half times the difference in magnitude between the IV and OLS estimates. The DDML estimate has a similar standard error to the IV estimate. Based on the DDML estimate, the null hypothesis that the slave trade had no impact on levels of trust would not be rejected at conventional significance levels, contrary to the central finding of [Nunn and Wantchekon \(2011\)](#). Note that unconditional ACR estimates are not defined for this application because the instrument is not binary.

Next, we consider a paper by [Dube and Harish \(2020\)](#), who estimate the effect of queen rule on war using panel data on the polities of Europe covering the years 1480 to 1913. The outcome Y is a binary indicator for whether a polity-year observation was at war. The treatment T is a binary indicator for whether a queen ruled in that polity-year. The instrument Z is an indicator for whether the previous monarch had a legitimate firstborn male child. The covariates X in their main results ([Dube and Harish, 2020](#), Table 3, column (3)) are polity and decade identifiers, whether the previous monarchs were corulers unrelated to one another, whether they had any legitimate children (with and without missing birth years), and whether the gender of the previous firstborn child is missing.

[Dube and Harish \(2020\)](#) justify most of their covariates with concerns about exogeneity of the instrument. For example, they argue that controlling for whether the previous monarch had any legitimate children is necessary because the firstborn son instrument is mechanically zero whenever the previous monarch had no children ([Dube and Harish, 2020](#), pp. 2601–2602). In Table SA.3 we show that without polity fixed effects their IV estimates are implausibly large, sometimes exceeding the logical value of 1, albeit with large standard errors. With both polity and decade fixed effects, but

without the previous monarch controls, their estimates are close to half as large in magnitude. Covariates apparently matter substantially for their conclusions.

Dube and Harish (2020, pg. 2605) explicitly invoke a LATE interpretation for their estimates:

If there are heterogeneous treatment effects, the IV estimate will be the LATE (Imbens and Angrist 1994). It will tell us the effect for the specific group of women who were eligible to rule and induced into ruling because of the presence of a firstborn female or sister among previous monarchs (i.e., the set of women who were compliers).

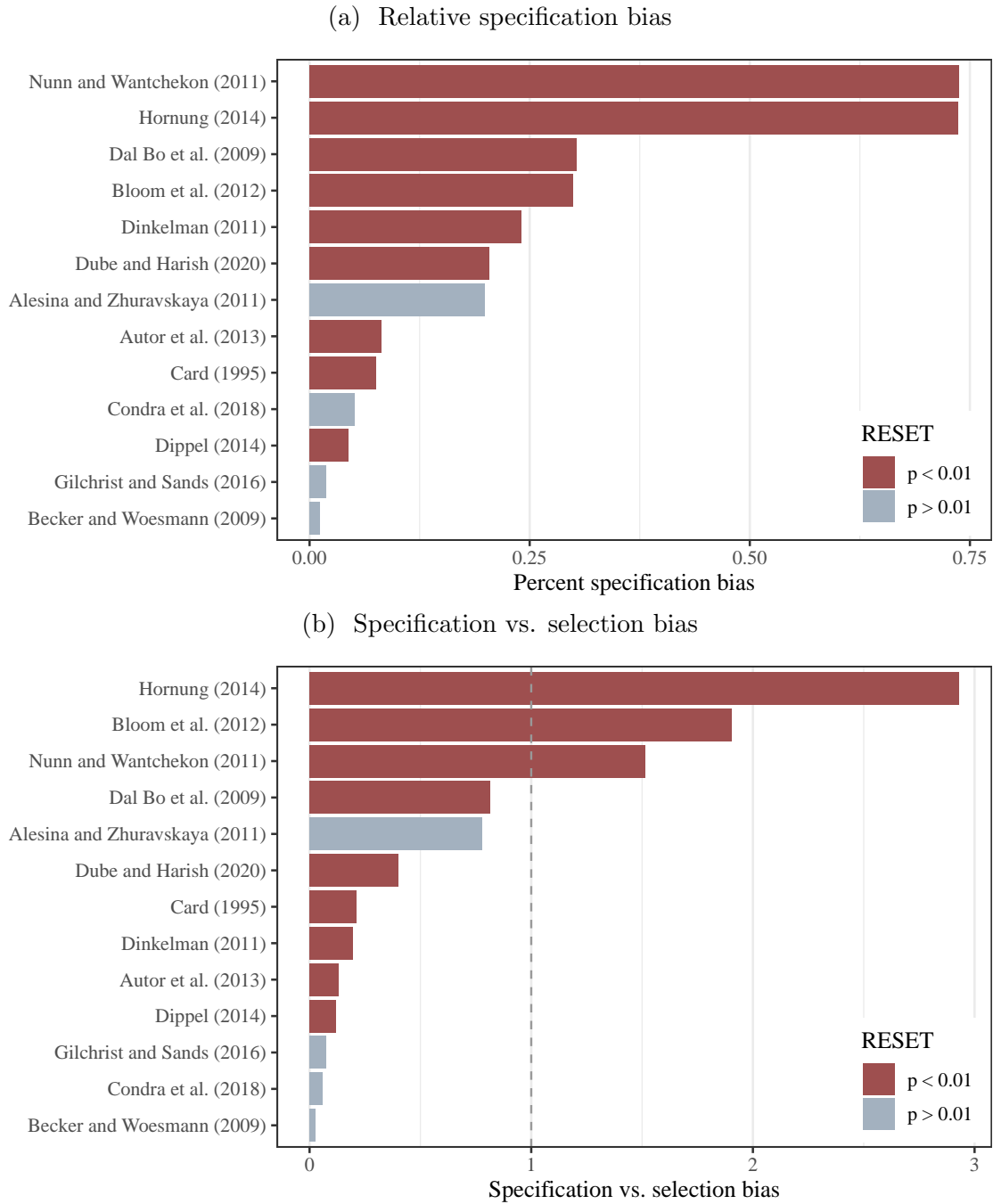
Both the treatment and instrument are binary, so the idea of a single LATE as introduced in Imbens and Angrist (1994) is well-defined. However, even if covariates are rich, linear IV estimates not an unconditional LATE, but β_{rich} , which is a statistically-weighted average of different covariate-specific LATEs. If covariates are not rich, then β_{iv} is not even weakly causal, let alone an estimate of the LATE.

Column (3) of Table 3 replicates Table 3, column (3) of Dube and Harish (2020) along with the comparable OLS estimates. The RESET test once again overwhelmingly rejects the null hypothesis that β_{iv} is weakly causal. The DDML estimate of β_{rich} is about 20% smaller than the original estimate, or 40% of the difference between the original IV and OLS estimates. While estimated with similar precision as linear IV, it is no longer significantly different from zero at conventional levels.

In the fifth row of column (3), we report the κ -weighting estimator discussed in Section 5.3, which we can apply here because both the instrument and treatment are binary. The estimates of κ use a logit to estimate $\mathbb{E}[Z|X]$. Proposition 7 showed that this estimator will converge to $\beta_{\text{iv}} = \beta_{\text{rich}}$ if rich covariates hold. In this application, we find that the resulting estimate is has the opposite sign and is extremely noisy. This may be because of the large number of fixed effects, which makes it difficult to estimate a logit.

Because the instrument is binary, we can also estimate β_{acr} using either DDML or IPSW. The treatment is also binary, so $\beta_{\text{acr}} = \beta_{\text{late}} \equiv \mathbb{E}[Y(1) - Y(0)|D(1) > D(0)]$. The DDML estimate of β_{late} is about half the size of the original IV estimate and about two thirds the size of the DDML estimate of β_{rich} . Although it is estimated more precisely, it is not statistically different from zero at conventional levels. The IPSW estimate of β_{late} is larger than the original IV estimate, but extremely noisy. As with the κ -weighting estimator, this may be because of the large number of fixed effects.

Figure 5: Relative magnitude across several applications



Notes: These figures present the absolute difference between linear IV and DDML estimates relative to the linear IV estimate (panel A) and relative to the difference between the OLS and DDML estimates (panel B). The definitions of the two relative biases are as in Table 3. Details on the sample size and number of included variables for each specification are provided in Table SA.5.

6.3 Patterns from multiple studies

To conduct a more systematic evaluation, we return to our survey of IV papers from Section 2.4. We consider all papers from the survey for which data was available and for which we were able to replicate the main IV estimates. We limit our attention to papers that fit the framework of Section 4 with a single endogenous variable and a single instrument used as the sole excluded variable. We omit papers that use panel data with two-way fixed effects, as their implementation with DDML requires more complex methods (e.g. [Semenova et al., 2023](#)). Imposing these restrictions leaves us with ten studies. For comparison, we also include [Card \(1995\)](#), [Nunn and Wantchekon \(2011\)](#), and [Dube and Harish \(2020\)](#), bringing the total to thirteen.

Figure 5 summarizes the differences between the IV estimate of β_{iv} and the DDML estimate of β_{rich} for the main specification in each of these thirteen studies. Panel (a) measures the differences relative to the original IV estimates, while panel (b) measures them relative to the difference between the original IV and comparable OLS estimates. The bars are shaded according to whether the RESET test rejects the null that β_{iv} is weakly causal at the 1% level. Table SA.4 provides tabular results for each study and shows that standard errors for the IV and DDML estimates are generally similar.

The RESET test rejects in nine out of the thirteen studies, implying that for most of these studies β_{iv} is not weakly causal. The magnitude of the difference between the IV estimate of β_{iv} and the DDML estimate of β_{rich} varies across studies, but is often large measured either relative to the original estimates or relative to the difference between the OLS and IV estimates. Cases when this difference is large are reliably detected by the RESET test. The one exception, [Alesina and Zhuravskaya \(2011\)](#), has only 97 observations, so is likely underpowered. Conversely, the studies for which the RESET test does not reject also tend to exhibit small differences between the IV and DDML estimates, suggesting that $\beta_{iv} = \beta_{rich}$, and that the original linear IV estimate is indeed weakly causal. For some studies, such as [Dippel \(2014\)](#), the RESET test rejects, but the difference between the IV and DDML estimates is fairly small. This is not a contradiction: two numbers can be equal even if one is a non-negative weighted average and the other is not. Someone who finds this possibility troubling is expressing dissatisfaction with only imposing the extremely weak requirement of weak causality.

7 Conclusion and recommendations for practice

In discussing the LATE interpretation of linear IV estimates, [Angrist and Krueger \(1999, pg. 1326\)](#) conjectured:

That is, IV estimates in models with covariates can be thought of as producing a weighted average of covariate-specific Wald estimates as long as the model

for covariates is saturated In other cases it seems reasonable to assume that some sort of approximate weighted average is being generated, but we are unaware of a precise causal interpretation that fits all cases.

We have shown that this seemingly-reasonable assumption is false. Unless rich covariates is satisfied, the linear IV estimand cannot be interpreted as “weakly causal,” and so cannot be interpreted as a non-negatively weighted average of LATEs. We tested the null hypothesis of rich covariates in several empirical studies and found that it was commonly rejected.

Based on our theoretical results and empirical applications, we recommend that researchers using linear IV estimators take the following steps.

1. Consider the role of covariates in the IV analysis. If covariates are not essential for justifying instrument exogeneity, then report estimates without covariates. Estimates with covariates can still be reported if the covariates are helpful for precision. If covariates play an important role in justifying exogeneity, then think carefully about which covariates ought to be included and why. Using a “kitchen sink” approach to controlling for covariates makes it less likely that rich covariates is satisfied and so more likely that the resulting IV estimate is not weakly causal.
2. Report the [Ramsey \(1969\)](#) RESET test for a regression of the instrument on the covariates. The null hypothesis of this test is equivalent to the null hypothesis that rich covariates holds, which our results show is also equivalent to the null hypothesis the IV estimand is weakly causal for the types of IV specifications considered in the main text. The RESET test can be implemented in Stata with the command `estat ovtest` and in R through the `resettest` function in the `lmtest` package ([Zeileis and Hothorn, 2002](#)).¹² If the RESET test rejects, then proceed to the next step. Otherwise, proceed to step four.
3. Estimate β_{rich} and report the result alongside the linear IV estimate of β_{iv} . The DDML estimator of β_{rich} can be viewed as a nonparametric estimator and seems to perform well in our simulations. A Stata implementation of DDML has been developed by [Ahrens et al. \(2023\)](#). There are at least two R packages ([Ahrens et al., 2024a](#); [Bach et al., 2024](#)).
4. If the instrument is binary, then estimate the unconditional ACR, β_{acr} , which is equal to the unconditional LATE, β_{lato} , if the treatment is also binary. This can be implemented with DDML in either Stata or R. It can also be implemented in

¹²Testing whether a TSLS estimand with multiple excluded variables is weakly causal is less straightforward because the rich covariates condition now concerns the conditional mean of the aggregated excluded variables (see Section 4.4). A bootstrapped version of the RESET test may be an adequate solution in this case. See our working paper ([Blandhol et al., 2022](#)) for more details and an example.

Stata with IPSW (Słoczyński et al., 2024). We are not aware of an IPSW package for R, although it is straightforward to construct point estimates by fitting a binary response model and then constructing four weighted means.

It is important to emphasize that the criterion of “weakly causal” used throughout the analysis is an extremely weak one. Being weakly causal may be necessary for a quantity to represent an interesting causal effect, but it is not sufficient. Even if rich covariates is satisfied, β_{rich} may be hard to interpret. As we showed empirically, it can also be quite different from a more interpretable quantity, such as the unconditional ACR, β_{acr} , or unconditional LATE, β_{late} .

These interpretation difficulties were already reason to explore alternative IV methods designed to estimate quantities, such as an unconditional LATE or the average treatment on the treated, that are not only weakly causal but also have clear counterfactual interpretations. Such methods rely on explicitly stated parametric assumptions (e.g. Heckman, 1976; Imbens and Rubin, 1997; Heckman et al., 2003) or are semiparametric (e.g. Carneiro et al., 2011; Brinch et al., 2017; Mogstad et al., 2018; Słoczyński et al., 2024) or nonparametric (e.g. Heckman and Vytlacil, 1999; Manski and Pepper, 2000; Chernozhukov et al., 2018). By showing that common interpretations of linear IV estimands also rely on either parametric assumptions or nonparametric implementations, our findings provide another reason to pursue such approaches.

A Proofs

Proof of Proposition 1. The expression for β_{iv} is a special case of Proposition 2 with $\epsilon = 1$.

If $\mathbb{E}[T\tilde{Z}] > 0$, then because $\mathbb{E}[Z|X] \in [0, 1]$ for binary Z , the sign of $\omega(\text{CP}, X)$ depends on the sign of $1 - \mathbb{L}[Z|X]$, which is negative if and only if $\mathbb{L}[Z|X] > 1$. The sign of $\omega(\text{AT}, X)$ varies with X according to the sign of $\mathbb{E}[\tilde{Z}|X]$. Because X contains a constant, $\mathbb{E}[\mathbb{E}[\tilde{Z}|X]] = \mathbb{E}[\tilde{Z}] = 0$, so $\mathbb{E}[\tilde{Z}|X]$ is either zero with probability 1, or else it has positive probability of taking both positive and negative values. In the latter case, the sign of $\omega(\text{AT}, X)$ is negative for some values of X regardless of whether $\mathbb{E}[T\tilde{Z}]$ is positive or negative. *Q.E.D.*

Proof of Proposition 2. The numerator of β_{iv} can be written as

$$\mathbb{E}[Y\tilde{Z}] = \mathbb{E}\left[\mathbb{E}\left[Y\tilde{Z}|X\right]\right] = \mathbb{E}\left[\mathbb{C}[Y, \tilde{Z}|X]\right] + \mathbb{E}\left[\mathbb{E}[Y|X]\mathbb{E}[\tilde{Z}|X]\right], \quad (20)$$

where \mathbb{C} denotes covariance. The same argument as in Imbens and Angrist (1994)

applied conditional-on-covariates yields

$$\mathbf{C}[Y, Z|X] = \Delta(\text{CP}, X) \mathbf{C}[T, Z|X] = \Delta(\text{CP}, X) \mathbb{P}[G = \text{CP}|X] \mathbb{E}[Z|X](1 - \mathbb{E}[Z|X]). \quad (21)$$

As for the second term of (20),

$$\begin{aligned} \mathbb{E}[Y|X] &= \mathbb{E}[Y|G = \text{AT}, X] \mathbb{P}[G = \text{AT}|X] + \mathbb{E}[Y|G = \text{NT}, X] \mathbb{P}[G = \text{NT}|X] \\ &\quad + \mathbb{E}[Y|G = \text{CP}, X] \mathbb{P}[G = \text{CP}|X] \\ &= \mathbb{E}[Y(1)|G = \text{AT}, X] \mathbb{P}[G = \text{AT}|X] + \mathbb{E}[Y(0)|G = \text{NT}, X] \mathbb{P}[G = \text{NT}|X] \\ &\quad + \mathbb{E}[(1 - Z)Y(0) + ZY(1)|G = \text{CP}, X] \mathbb{P}[G = \text{CP}|X]. \end{aligned} \quad (22)$$

Adding and subtracting $\mathbb{E}[Y(0)|G = \text{AT}, X] \mathbb{P}[G = \text{AT}|X]$ gives

$$\mathbb{E}[Y|X] = \Delta(\text{CP}, X) \mathbb{P}[G = \text{CP}|X] \mathbb{E}[Z|X] + \Delta(\text{AT}, X) \mathbb{P}[G = \text{AT}|X] + \eta'_0 X \quad (23)$$

due to both the exogeneity of Z and the linearity assumption on $\mathbb{E}[Y(0)|X = x]$. Alternatively, adding and subtracting $\mathbb{E}[Y(1)|G = \text{NT}, X] \mathbb{P}[G = \text{NT}|X]$ to (22) gives

$$\mathbb{E}[Y|X] = \Delta(\text{CP}, X) \mathbb{P}[G = \text{CP}|X] (\mathbb{E}[Z|X] - 1) - \Delta(\text{NT}, X) \mathbb{P}[G = \text{NT}|X] + \eta'_1 X. \quad (24)$$

So multiplying (23) by ϵ and summing it with (24) multiplied by $1 - \epsilon$ gives

$$\begin{aligned} \mathbb{E}[Y|X] &= \Delta(\text{CP}, X) \mathbb{P}[G = \text{CP}|X] (\mathbb{E}[Z|X] + \epsilon - 1) + \Delta(\text{AT}, X) \epsilon \mathbb{P}[G = \text{AT}|X] \\ &\quad + \Delta(\text{NT}, X) (\epsilon - 1) \mathbb{P}[G = \text{NT}|X] + \epsilon \eta'_0 X + (1 - \epsilon) \eta'_1 X. \end{aligned}$$

Because X and \tilde{Z} are orthogonal,

$$\begin{aligned} \mathbb{E}[\mathbb{E}[Y|X] \mathbb{E}[\tilde{Z}|X]] &= \mathbb{E} \left[\Delta(\text{CP}, X) \mathbb{P}[G = \text{CP}|X] (\mathbb{E}[Z|X] + \epsilon - 1) \mathbb{E}[\tilde{Z}|X] \right. \\ &\quad + \Delta(\text{AT}, X) \epsilon \mathbb{P}[G = \text{AT}|X] \mathbb{E}[\tilde{Z}|X] \\ &\quad \left. + \Delta(\text{NT}, X) (\epsilon - 1) \mathbb{P}[G = \text{NT}|X] \mathbb{E}[\tilde{Z}|X] \right]. \end{aligned} \quad (25)$$

Summing (21) and (25), and noting that

$$\begin{aligned} &\mathbb{E}[Z|X](1 - \mathbb{E}[Z|X]) + (\mathbb{E}[Z|X] + \epsilon - 1) \mathbb{E}[\tilde{Z}|X] \\ &= (\mathbb{E}[Z|X] - \mathbb{E}[\tilde{Z}|X]) (1 - \mathbb{E}[Z|X]) + \epsilon \mathbb{E}[\tilde{Z}|X] \\ &= \mathbb{L}[Z|X](1 - \mathbb{E}[Z|X]) + \epsilon \mathbb{E}[\tilde{Z}|X] \end{aligned}$$

yields a weighting expression with weights proportional to the claimed expression but missing a common multiple of $\mathbb{E}[\tilde{Z}T]^{-1}$, which comes from the denominator of β_{iv} .
Q.E.D.

Proof of Proposition 3. Note that T is only stochastic due to Z after conditioning on X and G , as a direct consequence of the definition of G . Assumption EXO then implies that T and $Y(t)$ are independent conditional on X and G . We use this observation to write

$$\begin{aligned}\beta &= \sum_{g,x} \mathbb{E} [b(T, x, Z)Y | G = g, X = x] \mathbb{P}[G = g, X = x] \\ &= \sum_{g,x,j} \mathbb{E} [\mathbb{1}[T = t_j]b(t_j, x, Z)Y(t_j) | G = g, X = x] \mathbb{P}[G = g, X = x] \\ &= \sum_{g,x,j} \mu_j(g, x) \mathbb{E} [\mathbb{1}[T = t_j]b(t_j, x, Z) | G = g, X = x] \mathbb{P}[G = g, X = x] \\ &\equiv \sum_{g,x,j} \mu_j(g, x)\psi_j(g, x),\end{aligned}$$

where all summations are taken over $g \in \mathcal{G}, x \in \mathcal{X}, j \in \{0, 1, \dots, J\}$, and

$$\psi_j(g, x) \equiv \mathbb{E} [\mathbb{1}[T = t_j]b(t_j, x, Z) | G = g, X = x] \mathbb{P}[G = g, X = x].$$

Notice that $\omega_j(g, x) = \sum_{k=j}^J \psi_k(g, x)$, so that (6) follows from Lemma 1.

Q.E.D.

Lemma 1. For any constants $\{a_j, c_j\}_{j=0}^J$,

$$\sum_{j=0}^J a_j c_j = a_0 \tilde{c}_0 + \sum_{j=1}^J (a_j - a_{j-1}) \tilde{c}_j,$$

where $\tilde{c}_j \equiv \sum_{k=j}^J c_k$.

Proof of Lemma 1. Since $c_j = \tilde{c}_j - \tilde{c}_{j+1}$,

$$\begin{aligned}\sum_{j=0}^J a_j c_j &= \sum_{j=0}^J a_j (\tilde{c}_j - \tilde{c}_{j+1}) \\ &= a_0 \tilde{c}_0 + \sum_{j=1}^J a_j \tilde{c}_j + \sum_{j=0}^{J-1} a_j \tilde{c}_{j+1} = a_0 \tilde{c}_0 + \sum_{j=1}^J (a_j - a_{j-1}) \tilde{c}_j,\end{aligned}$$

where the final equality used a change of variables in the second summand from j to $j + 1$.
Q.E.D.

Proof of Proposition 4. If $\omega_j(g, x) \geq 0$ and $\omega_0(g, x) = 0$ for all g and x , then it follows immediately from (6) that β satisfies Definition WC.

We will prove the converse by contraposition. That is, we will show that if either the non-negative weights or level independence condition is not satisfied, then there exists a μ such that $\mu_j(g, x) - \mu_{j-1}(g, x)$ has the same sign for every $j \geq 1$, and all g and x , and that this common sign is different from the sign of β . This shows that if the weights do not satisfy both the non-negative and level independence conditions, then β is not weakly causal. Or, by contraposition, if β is weakly causal, then the weights satisfy both conditions.

First, suppose that the level independence condition does not hold, but that the non-negative weights condition does hold. Then there exists a (g^*, x^*) such that $\omega_0(g^*, x^*) \neq 0$, but $\omega_j(g, x) \geq 0$ for all $j \geq 1$, g , and x . Set

$$\mu_j(g, x) = \begin{cases} \bar{\mu}, & \text{if } (g, x) \neq (g^*, x^*) \\ \mu^*, & \text{if } (g, x) = (g^*, x^*) \text{ and } j < j^* , \\ \mu^* + \Delta^*, & \text{if } (g, x) = (g^*, x^*) \text{ and } j \geq j^* \end{cases} \quad (26)$$

where $\bar{\mu}$, μ^* , and Δ^* are numbers we will choose, and $j^* \geq 1$ can be chosen arbitrarily. Then $\mu_j(g, x) - \mu_{j-1}(g, x)$ is zero for all $(g, x) \neq (g^*, x^*)$, while for $(g, x) = (g^*, x^*)$ it is Δ^* when $j = j^*$ and zero otherwise. In particular, the sign of $\mu_j(g, x) - \mu_{j-1}(g, x)$ is the sign of Δ^* for all $j \geq 1$ and all (g, x) , regardless of the values of $\bar{\mu}$ and μ^* . If μ is specified as in (26) with $\bar{\mu} = 0$ for simplicity, then (6) becomes

$$\beta = \omega_0(g^*, x^*)\mu^* + \omega_{j^*}(g^*, x^*)\Delta^*. \quad (27)$$

Fix any $\mu^* \neq 0$, so that $\omega_0(g^*, x^*)\mu^* \neq 0$. If $\omega_0(g^*, x^*)\mu^* > 0$, then choose a $\Delta^* < 0$ that is sufficiently small in magnitude so that $\omega_{j^*}(g^*, x^*)\Delta^* > -\omega_0(g^*, x^*)\mu^*$. Then from (27) we have $\beta = \omega_0(g^*, x^*)\mu^* + \omega_{j^*}(g^*, x^*)\Delta^* > 0$, so that these choices of μ^* and Δ^* produce a μ that violates the second condition of Definition WC. Similarly, if $\omega_0(g^*, x^*)\mu^* < 0$, then choose $\Delta^* > 0$ to be sufficiently small to ensure that $\omega_{j^*}(g^*, x^*)\Delta^* < -\omega_0(g^*, x^*)\mu^*$, so that $\beta < 0$, contradicting the first condition of Definition WC.

On the other hand, suppose that the non-negative weights condition does not hold, so there exist a j^* , g^* , and x^* such that $\omega_{j^*}(g^*, x^*) < 0$. Use the same construction as in (26) with these new values of j^* , g^* , and x^* , where j^* is no longer arbitrary. Set $\mu^* = 0$. Then (27) reduces to

$$\beta = \omega_{j^*}(g^*, x^*)\Delta^*.$$

Selecting any $\Delta^* > 0$ produces $\beta < 0$, establishing the existence of a μ that violates the first condition of Definition WC. *Q.E.D.*

Proof of Theorem 1. We evaluate the sufficient and necessary conditions in Proposition 4 using the expressions for $\omega_j(g, x)$ given in (37).

First, consider the level independence condition, which given (37) can be written as

$$\omega_0(g, x) = \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E}[\tilde{Z}|G = g, X = x] \mathbb{P}[G = g, X = x] = 0 \quad (28)$$

for all g and x . Assumption EXO implies that $\tilde{Z} \equiv Z - \mathbb{L}[Z|X]$ is independent of G given X , so

$$\mathbb{E}[\tilde{Z}|G = g, X = x] = \mathbb{E}[\tilde{Z}|X = x] = \mathbb{E}[Z|X = x] - \mathbb{L}[Z|X = x].$$

For every x there exists a $g \in \mathcal{G}$ such that $\mathbb{P}[G = g, X = x] > 0$, because G exhaustively partitions possible choice types. So (28) can hold for every g and x if and only if

$$\mathbb{E}[Z|X = x] = \mathbb{L}[Z|X = x]$$

for every x , that is, if and only if the specification has rich covariates. In particular, if the specification does not have rich covariates, then (28) is non-zero for some g and x , and so by Proposition 4, β_{iv} is not weakly causal.

To establish the sufficient direction, suppose that the specification has rich covariates and consider the non-negative weights condition in Proposition 4. Let $\mathcal{Z}_j(g)$ denote the set of instrument values for which individuals in choice group g would choose a treatment value t_j or larger. Then

$$\begin{aligned} \omega_j(g, x) &= \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E} \left[\tilde{Z} \mathbb{1}[T \geq t_j] \middle| G = g, X = x \right] \mathbb{P}[G = g, X = x] \\ &= \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E} \left[\tilde{Z} \mathbb{1}[Z \in \mathcal{Z}_j(g)] \middle| G = g, X = x \right] \mathbb{P}[G = g, X = x] \\ &= \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{C} \left[Z, \mathbb{1}[Z \in \mathcal{Z}_j(g)] \middle| X = x \right] \mathbb{P}[G = g|X = x] \mathbb{P}[X = x], \end{aligned} \quad (29)$$

where the third equality follows from Assumption EXO and the hypothesis of rich covariates. Given rich covariates,

$$\mathbb{E}[\tilde{Z}T] = \mathbb{E}[(Z - \mathbb{E}[Z|X]) \mathbb{E}[T|X, Z]] = \mathbb{E}[\mathbb{C}[Z, \mathbb{E}[T|X, Z]|X]], \quad (30)$$

which is non-negative because Assumptions EXO and MON imply that $\mathbb{E}[T|X, Z]$ is a weakly increasing function of Z (Angrist and Imbens, 1995; Vytlacil, 2002), and the covariance between two weakly increasing functions is non-negative (e.g. Thorisson, 1995, Section 2).

It remains to determine the sign of the covariance term in (29). Suppose that $\mathbb{P}[G = g|X = x] > 0$. Then the function $z \mapsto \mathbb{1}[z \in \mathcal{Z}_j(g)]$ must be weakly increasing. For otherwise, there would exist z, z' with $z < z'$ and $z \in \mathcal{Z}_j(g)$ but $z' \notin \mathcal{Z}_j(g)$, meaning that for group g , instrument value z leads to $T(z) \geq t_j$, while instrument value z' leads to $T(z') < t_j$. Given that $\mathbb{P}[G = g|X = x] > 0$, this would imply that

$$\mathbb{P}[T(z) \geq t_j > T(z')|X = x] \geq \mathbb{P}[G = g|X = x] > 0,$$

in contradiction with Assumption MON. It follows that the covariance term in (29) is non-negative, again because the covariance of two increasing functions of Z is non-negative. We conclude that $\omega_j(g, x) \geq 0$ for all j, g , and x , which by Proposition 4 shows that β_{iv} is weakly causal. Q.E.D.

Proof of Proposition 5. Write (10) as

$$\begin{aligned} \beta_{iv} &= \Delta + \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E}[\tilde{Z} \mathbb{E}[Y(t_0)|Z, G, X]] \\ &= \Delta + \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E}[\tilde{Z} \mathbb{E}[Y(t_0)|G, X]] \\ &= \Delta + \mathbb{E}[\tilde{Z}T]^{-1} \mathbb{E}[\mathbb{E}[\tilde{Z}|X] \mathbb{E}[Y(t_0)|G, X]] \equiv \Delta + \sum_{g,x} \omega_0(g, x) \mu_0(g, x), \end{aligned}$$

where $\omega_0(g, x)$ is as defined as (37), noting that $\mathbb{E}[\tilde{Z}|G, X] = \mathbb{E}[\tilde{Z}|X]$ due to Assumption EXO. If rich covariates holds, then $\omega_0(g, x) = 0$ for all g and x , so that $\beta_{iv} = \Delta$ is weakly causal by Proposition 4. Conversely, if rich covariates does not hold, then, as shown in the proof of Theorem 1, there exists a (g, x) such that $\omega_0(g, x) \neq 0$, so β_{iv} is not weakly causal, again by Proposition 4. Q.E.D.

Proof of Proposition 6. Given Assumption CLE, Assumption LIN also implies that

$$\mathbb{E}[Y(t_0)|X = x] = \mathbb{E}[Y(t_j) - Y(t_0)|X = x] + \mathbb{E}[Y(t_j)|X = x] = \Delta(t_j - t_0) + \eta'x,$$

so that $\mathbb{E}[Y(t_0)|X = x] = \eta'_0 x$, where η_0 is the same as η but has $\Delta(t_j - t_0)$ added to the coefficient on the constant component of x . Because \tilde{Z} is orthogonal to X ,

$$\mathbb{E}[\tilde{Z}Y(t_0)] = \mathbb{E}[\mathbb{E}[\tilde{Z}|X] \mathbb{E}[Y(t_0)|X]] = \mathbb{E}[\mathbb{E}[\tilde{Z}|X]X']\eta_0 = \mathbb{E}[\tilde{Z}X']\eta_0 = 0.$$

From (10), this implies that $\beta_{iv} = \Delta$, as claimed. Q.E.D.

Proof of Proposition 7. Abadie (2003, Proposition 5.1) showed that if rich covariates is satisfied, then the linear IV estimate is numerically equal to the κ -weighting estimate, implying that $\beta_{abadie} = \beta_{iv}$, with $\beta_{iv} = \beta_{rich}$ by definition in that case.

For the converse, consider the κ -weighting linear regression, which Abadie showed is the same as an unweighted linear regression of Y on T and X among the subpopulation

$G = \text{CP}$ of compliers. Assumption MON implies that $\mathbb{P}[T = Z|G = \text{CP}] = 1$. Assumption EXO then implies that $(Y(0), Y(1))$ is independent of T conditional on $G = \text{CP}$. Corollary 1 therefore implies that the population coefficient on T in the κ -weighting regression is weakly causal if and only if $\mathbb{E}[T|X, G = \text{CP}] = \mathbb{L}[T|X, G = \text{CP}] \equiv \gamma'X$ for some γ . However, Assumption EXO implies that

$$\gamma'X \equiv \mathbb{L}[T|X, G = \text{CP}] = \mathbb{E}[T|X, G = \text{CP}] = \mathbb{E}[Z|X, G = \text{CP}] = \mathbb{E}[Z|X], \quad (31)$$

which implies that $\mathbb{E}[Z|X]$ is linear in X , and therefore that $\mathbb{E}[Z|X] = \mathbb{L}[Z|X]$, so that rich covariates is satisfied. *Q.E.D.*

References

- ABADIE, A. (2003): “Semiparametric Instrumental Variable Estimation of Treatment Response Models,” *Journal of Econometrics*, 113, 231–263. 4, 20, 23, 24, 41
- ACKERBERG, D. A. AND P. J. DEVEREUX (2009): “Improved JIVE Estimators for Overidentified Linear Models with and without Heteroskedasticity,” *The Review of Economics and Statistics*, 91, 351–362. 53
- AHRENS, A., C. B. HANSEN, AND M. E. SCHAFFER (2022): “Pystacked: Stacking Generalization and Machine Learning in Stata,” . 29
- AHRENS, A., C. B. HANSEN, M. E. SCHAFFER, AND T. WIEMANN (2023): “Ddml: Double/Debiased Machine Learning in Stata,” . 26, 29, 35
- (2024a): *Ddml: Double/Debiased Machine Learning*. 26, 35
- (2024b): “Model Averaging and Double Machine Learning,” . 26
- ALESINA, A. AND E. ZHURAVSKAYA (2011): “Segregation and the Quality of Government in a Cross Section of Countries,” *American Economic Review*, 101, 1872–1911. 34
- ANGRIST, J. D. (1998): “Estimating the Labor Market Impact of Voluntary Military Service Using Social Security Data on Military Applicants,” *Econometrica*, 66, 249–288. 15, 18, 21
- (2001): “Estimation of Limited Dependent Variable Models With Dummy Endogenous Regressors,” *Journal of Business & Economic Statistics*, 19, 2–28. 24
- ANGRIST, J. D. AND W. N. EVANS (1998): “Children and Their Parents’ Labor Supply: Evidence from Exogenous Variation in Family Size,” *The American Economic Review*, 88, 450–477. 24
- ANGRIST, J. D. AND G. W. IMBENS (1995): “Two-Stage Least Squares Estimation of Average Causal Effects in Models with Variable Treatment Intensity,” *Journal of the American Statistical Association*, 90, 431–442. 2, 12, 20, 22, 40, 52, 53
- ANGRIST, J. D., G. W. IMBENS, AND A. B. KRUEGER (1999): “Jackknife Instrumental Variables Estimation,” *Journal of Applied Econometrics*, 14, 57–67. 53
- 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, 444–455. 13, 14
- ANGRIST, J. D. AND A. B. KRUEGER (1999): “Chapter 23 Empirical Strategies in Labor Economics,” Elsevier, vol. Volume 3, Part A, 1277–1366. 34
- ANGRIST, J. D. AND J.-S. PISCHKE (2009): *Mostly Harmless Econometrics: An Empiricist’s Companion*, Princeton University Press. 2, 3, 6, 10, 11, 12, 15, 18, 19, 20, 21, 23, 24, 52, 53
- AUTOR, D. H., D. DORN, AND G. H. HANSON (2013): “The China Syndrome: Local Labor Market Effects of Import Competition in the United States,” *American Economic Review*, 103, 2121–2168.
- BACH, P., M. S. KURZ, V. CHERNOZHUKOV, M. SPINDLER, AND S. KLAASSEN (2024):

- “DoubleML: An Object-Oriented Implementation of Double Machine Learning in R,” *Journal of Statistical Software*, 108, 1–56, arXiv:[2103.09603](#) [stat.ML]. 35
- BECKER, S. O. AND L. WOESSMANN (2009): “Was Weber Wrong? A Human Capital Theory of Protestant Economic History,” *Quarterly Journal of Economics*, 124, 531–596.
- BLANDHOL, C., J. BONNEY, M. MOGSTAD, AND A. TORGOVITSKY (2022): “When Is TSLS Actually LATE?” Tech. Rep. w29709, National Bureau of Economic Research, Cambridge, MA. 35, 53
- BLOOM, N., R. SADUN, AND J. VAN REENEN (2012): “The Organization of Firms Across Countries*,” *The Quarterly Journal of Economics*, 127, 1663–1705.
- BRINCH, C. N., M. MOGSTAD, AND M. WISWALL (2017): “Beyond LATE with a Discrete Instrument,” *Journal of Political Economy*, 125, 985–1039. 36
- CARD, D. (1995): “Using Geographic Variation in College Proximity to Estimate the Return to Schooling,” in *Aspects of Labour Market Behaviour: Essays in Honour of John Vanderkamp*, ed. by L. N. Christofides, K. E. Grant, and R. Swidinsky, Toronto: University of Toronto Press, 201–222. 24, 28, 29, 34, 57
- CARD, D., D. S. LEE, Z. PEI, AND A. WEBER (2015): “Inference on Causal Effects in a Generalized Regression Kink Design,” *Econometrica*, 83, 2453–2483. 15
- CARNEIRO, P., J. J. HECKMAN, AND E. J. VYTLACIL (2011): “Estimating Marginal Returns to Education,” *American Economic Review*, 101, 2754–81. 36
- CHAMBERLAIN, G. AND G. IMBENS (2004): “Random Effects Estimators with many Instrumental Variables,” *Econometrica*, 72, 295–306. 11, 12
- CHERNOZHUKOV, V., D. CHETVERIKOV, M. DEMIRER, E. DUFLO, C. HANSEN, W. NEWEY, AND J. ROBINS (2018): “Double/Debiased Machine Learning for Treatment and Structural Parameters,” *The Econometrics Journal*, 21, C1–C68. 4, 21, 23, 26, 29, 36
- CONDRA, L. N., J. D. LONG, A. C. SHAVER, AND A. L. WRIGHT (2018): “The Logic of Insurgent Electoral Violence,” *American Economic Review*, 108, 3199–3231.
- DAL BÓ, E., P. DAL BÓ, AND J. SNYDER (2009): “Political Dynasties,” *Review of Economic Studies*, 76, 115–142.
- DINKELMAN, T. (2011): “The Effects of Rural Electrification on Employment: New Evidence from South Africa,” *American Economic Review*, 101, 3078–3108.
- DIPPEL, C. (2014): “Forced Coexistence and Economic Development: Evidence from Native American Reservations,” *Econometrica*, 82, 2131–2165. 34
- DONALD, S. G., Y.-C. HSU, AND R. P. LIELI (2014): “Testing the Unconfoundedness Assumption via Inverse Probability Weighted Estimators of (L)ATT,” *Journal of Business & Economic Statistics*, 32, 395–415. 23
- DUBE, O. AND S. P. HARISH (2020): “Queens,” *Journal of Political Economy*, 128, 2579–2652. 28, 30, 31, 32, 34, 62
- EVDOKIMOV, K. S. AND M. KOLESÁR (2019): “Inference in Instrumental Variables Analysis with Heterogeneous Treatment Effects,” *Working paper*. 47
- FIRPO, S., M. N. FOGUEL, AND H. JALES (2020): “Balancing Tests in Stratified Randomized Controlled Trials: A Cautionary Note,” *Economics Letters*, 186, 108771. 8
- FRIEDMAN, J., T. HASTIE, R. TIBSHIRANI, B. NARASIMHAN, K. TAY, N. SIMON, AND J. YANG (2008): “Glmnet: Lasso and Elastic-Net Regularized Generalized Linear Models,” . 26
- FRÖLICH, M. (2007): “Nonparametric IV Estimation of Local Average Treatment Effects with Covariates,” *Journal of Econometrics*, 139, 35–75. 23
- GILCHRIST, D. S. AND E. G. SANDS (2016): “Something to Talk About: Social Spillovers in Movie Consumption,” *Journal of Political Economy*, 124, 1339–1382.
- GOLDSMITH-PINKHAM, P., P. HULL, AND M. KOLESÁR (2024): “Contamination Bias in Linear Regressions,” *American Economic Review*, 114, 4015–4051. 15
- GOODMAN-BACON, A. (2021): “Difference-in-Differences with Variation in Treatment Timing,” *Journal of Econometrics*, 225, 254–277. 3, 15
- HECKMAN, J., J. L. TOBIAS, AND E. VYTLACIL (2003): “Simple Estimators for Treatment

- Parameters in a Latent-Variable Framework,” *Review of Economics and Statistics*, 85, 748–755. 36
- HECKMAN, J. J. (1976): “The Common Structure of Statistical Models of Truncation, Sample Selection and Limited Dependent Variables and a Simple Estimator for Such Models,” *Annals of Economic and Social Measurement*. 36
- (2010): “Building Bridges between Structural and Program Evaluation Approaches to Evaluating Policy,” *Journal of Economic Literature*, 48, 356–98. 52
- HECKMAN, J. J. AND R. ROBB (1985): “Alternative Methods for Evaluating the Impact of Interventions,” in *Longitudinal Analysis of Labor Market Data*, ed. by J. J. Heckman and B. Singer, Cambridge University Press. 20
- HECKMAN, J. J. AND E. VYTLACIL (2005): “Structural Equations, Treatment Effects, and Econometric Policy Evaluation,” *Econometrica*, 73, 669–738. 52
- HECKMAN, J. J. AND E. J. VYTLACIL (1999): “Local Instrumental Variables and Latent Variable Models for Identifying and Bounding Treatment Effects,” *Proceedings of the National Academy of Sciences of the United States of America*, 96, 4730–4734. 36
- HEILER, P. (2022): “Efficient Covariate Balancing for the Local Average Treatment Effect,” *Journal of Business & Economic Statistics*, 40, 1569–1582. 23, 30
- HORNUNG, E. (2014): “Immigration and the Diffusion of Technology: The Huguenot Diaspora in Prussia,” *American Economic Review*, 104, 84–122.
- IMBENS, G. W. AND J. D. ANGRIST (1994): “Identification and Estimation of Local Average Treatment Effects,” *Econometrica*, 62, 467–475. 2, 5, 8, 20, 22, 32, 36, 49, 53
- IMBENS, G. W. AND D. B. RUBIN (1997): “Bayesian Inference for Causal Effects in Randomized Experiments with Noncompliance,” *The Annals of Statistics*, 25, 305–327. 36
- KOLESÁR, M. (2013): “Estimation in an Instrumental Variables Model with Treatment Effect Heterogeneity,” *Working paper*. 3, 17, 19, 47, 52, 53
- LEE, D. S. (2008): “Randomized Experiments from Non-Random Selection in U.S. House Elections,” *Journal of Econometrics*, 142, 675–697. 15
- MACURDY, T., X. CHEN, AND H. HONG (2011): “Flexible Estimation of Treatment Effect Parameters,” *American Economic Review*, 101, 544–551. 23
- MANSKI, C. F. AND J. V. PEPPER (2000): “Monotone Instrumental Variables: With an Application to the Returns to Schooling,” *Econometrica*, 68, 997–1010. 36
- MIKUSHEVA, A. AND L. SUN (2022): “Inference with Many Weak Instruments,” *The Review of Economic Studies*, 89, 2663–2686. 53
- MOGSTAD, M., A. SANTOS, AND A. TORGOVITSKY (2018): “Using Instrumental Variables for Inference About Policy Relevant Treatment Parameters,” *Econometrica*, 86, 1589–1619. 36, 57
- NUNN, N. AND L. WANTCHEKON (2011): “The Slave Trade and the Origins of Mistrust in Africa,” *American Economic Review*, 101, 3221–3252. 28, 30, 31, 34
- OGBURN, E. L., A. ROTNITZKY, AND J. M. ROBINS (2015): “Doubly Robust Estimation of the Local Average Treatment Effect Curve,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77, 373–396. 23
- PEDREGOSA, F., G. VAROQUAUX, A. GRAMFORT, V. MICHEL, B. THIRION, O. GRISEL, M. BLONDEL, P. PRETTENHOFER, R. WEISS, V. DUBOURG, J. VANDERPLAS, A. PASSOS, D. COURNAPEAU, M. BRUCHER, M. PERROT, AND E. DUCHESNAY (2011): “Scikit-Learn: Machine Learning in Python,” *Journal of Machine Learning Research*, 12, 2825–2830. 29
- RAMSEY, J. B. (1969): “Tests for Specification Errors in Classical Linear Least-Squares Regression Analysis,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 31, 350–371. 4, 13, 28, 35
- RIDGEWAY, G. (2007): “Generalized Boosted Models: A Guide to the Gbm Package,” *Update (Loma Linda University. Ethics Center)*, 1, 2007. 26
- RIPLEY, B. AND W. VENABLES (2016): “Package ‘Nnet’,” *R package version*, 7, 700. 26
- SEMENOVA, V., M. GOLDMAN, V. CHERNOZHUKOV, AND M. TADDY (2023): “Inference on Heterogeneous Treatment Effects in High-dimensional Dynamic Panels under Weak Depen-

- dence,” *Quantitative Economics*, 14, 471–510. 34
- SINGH, R. AND L. SUN (2024): “Double Robustness for Complier Parameters and a Semi-Parametric Test for Complier Characteristics,” *The Econometrics Journal*, 27, 1–20. 23
- SŁOCZYŃSKI, T. (2020): “When Should We (Not) Interpret Linear IV Estimands as LATE ?” . 4, 5, 12, 18, 19, 22, 23, 46, 49, 52
- (2024): “When Should We (Not) Interpret Linear IV Estimands as LATE?” . 4, 5, 12, 18, 19, 22, 23, 29, 46, 49, 52, 53
- SŁOCZYŃSKI, T., S. D. UYSAL, AND J. M. WOOLDRIDGE (2024): “Abadie’s Kappa and Weighting Estimators of the Local Average Treatment Effect,” *Journal of Business & Economic Statistics*, 1–14. 4, 23, 30, 36
- SUN, B. AND Z. TAN (2022): “High-Dimensional Model-Assisted Inference for Local Average Treatment Effects With Instrumental Variables,” *Journal of Business & Economic Statistics*, 40, 1732–1744. 23
- SUN, L. AND S. ABRAHAM (2021): “Estimating Dynamic Treatment Effects in Event Studies with Heterogeneous Treatment Effects,” *Journal of Econometrics*, 225, 175–199. 3, 15
- TAN, Z. (2006): “Regression and Weighting Methods for Causal Inference Using Instrumental Variables,” *Journal of the American Statistical Association*, 101, 1607–1618. 4, 23
- THORISSON, H. (1995): “Coupling Methods in Probability Theory,” *Scandinavian Journal of Statistics*, 22, 159–182. 40, 51
- UYSAL, S. D. (2011): “Three Essays on Doubly Robust Estimation Methods,” . 4, 23, 30
- VYTLACIL, E. (2002): “Independence, Monotonicity, and Latent Index Models: An Equivalence Result,” *Econometrica*, 70, 331–341. 40, 49, 55
- WOOLDRIDGE, J. M. (2010): *Econometric Analysis of Cross Section and Panel Data*, MIT press. 13, 20
- WRIGHT, M. N. AND A. ZIEGLER (2017): “Ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R,” *Journal of Statistical Software*, 77, 1–17. 26
- ZEILEIS, A. AND T. HOTHORN (2002): “Diagnostic Checking in Regression Relationships,” *R News*, 2, 7–10. 35

Supplemental Appendix

SA.1 Rich covariates under conditional random assignment

Suppose that $X = (X_1, X_2)$ has two components and that Z is randomly assigned conditional on X_1 . This could happen in a stratified experiment, where X_1 describes the strata. It could also happen in other settings, for example if judges Z are thought to be randomly assigned, but only conditionally on the day of the week, X_1 . The rich covariates condition is still the same in this case: $\mathbb{E}[Z|X_1, X_2] = \mathbb{L}[Z|X_1, X_2]$. However, random assignment implies that rich covariates reduces to the requirement that

$$\mathbb{E}[Z|X_1] = \mathbb{L}[Z|X_1, X_2], \quad (32)$$

because Z is independent of X_2 , given X_1 . In some situations, it will be natural to control for X_1 nonparametrically, for example using strata or day-of-week indicators. If this is done, then (32) reduces further to the requirement that

$$\mathbb{L}[Z|X_1] = \mathbb{L}[Z|X_1, X_2]. \quad (33)$$

Condition (33) can be evaluated by regressing Z onto X_1 and X_2 , then testing the null hypothesis that the coefficients on X_2 are zero.

SA.2 Generalizations

In this appendix, we generalize the discussion in the main text by considering specifications with more general forms of excluded variables. In doing so, we extend the results in [Słoczyński \(2020, 2024\)](#) about negative weighting under weaker forms of the monotonicity condition to settings with non-binary treatments and/or non-binary instruments.

SA.2.1 TSLS specification and estimand

We now consider specifications with a vector of excluded variables (excluded instruments), $I \equiv i(Z, X)$, where i is a known, vector-valued function. As in the main text, we continue to assume that the specification has a single endogenous variable, T . With a vector of excluded variables there is now a question of how to weight them. We consider the widely-used TSLS weighting with first stage variables $F \equiv [I', X']'$ and second stage variables $S \equiv [T, X']'$.

One way to interpret the first stage of TSLS is as a procedure for reducing the first stage variables F down to the same dimension as S by transforming I into a scalar.

That is, the first stage of TSLS replaces the vector of instruments I by the scalar “effective instrument”

$$\dot{Z} \equiv \gamma' I,$$

where γ is the vector of population coefficients on I in the first stage regression of T on I and X . The vector estimand of second stage coefficients, which we denote by α_{tsls} , can then be written as the standard IV estimand that uses $\dot{F} \equiv [\dot{Z}, X']'$ as instruments for $S \equiv [T, X']'$, that is

$$\alpha_{\text{tsls}} = \mathbb{E}[\dot{F}\dot{F}']^{-1} \mathbb{E}[\dot{F}Y]. \quad (34)$$

Alternatively, the first stage of TSLS can be viewed as constructing fitted values for the treatment,

$$\dot{T} \equiv \dot{Z} + \lambda' X \equiv \gamma' I + \lambda' X, \quad (35)$$

where λ is the vector of population coefficients on X in the first stage regression. The TSLS estimand is then the OLS estimand from a regression of Y onto \dot{T} and X .¹³

We assume that the standard rank condition holds, so that α_{tsls} exists. Our interest is in the component of α_{tsls} that corresponds to the coefficient on T , which we call β_{tsls} . The following proposition generalizes expression (8) for β_{iv} to β_{tsls} .

Proposition SA.1. Let $\tilde{Z} \equiv \dot{Z} - \mathbb{L}[\dot{Z}|X]$ denote the population residuals, where $\mathbb{L}[\dot{Z}|X] \equiv \mathbb{E}[\dot{Z}X'] \mathbb{E}[XX']^{-1}X$ are the population fitted values from regressing \dot{Z} onto X . Then

$$\beta_{\text{tsls}} = \frac{\mathbb{E}[\tilde{Z}Y]}{\mathbb{E}[\tilde{Z}^2]} = \frac{\mathbb{E}[\tilde{Z}Y]}{\mathbb{E}[\tilde{Z}T]} = \mathbb{E} \left[\left(\frac{\tilde{Z}}{\mathbb{E}[\tilde{Z}T]} \right) Y \right].$$

Proof of Proposition SA.1. The well-known two stage interpretation of α_{tsls} is

$$\alpha_{\text{tsls}} = \mathbb{E}[\dot{S}\dot{S}']^{-1} \mathbb{E}[\dot{S}Y], \quad (36)$$

where $\dot{S} \equiv \mathbb{L}[S|F] \equiv \mathbb{E}[SF'] \mathbb{E}[FF']^{-1}F$ are the fitted values from the population first stage regression. Because X is a subvector of both S and F , $\dot{S} = [\dot{T}', X']'$, where \dot{T} is the population fitted value from the first stage regression of T on I and X . Applying the Frisch-Waugh-Lovell Theorem to the second step regression (with full vector of

¹³Our definition of the TSLS estimand presumes the standard asymptotic framework where the number of observations is growing and the dimensions of I and X are fixed. [Kolesár \(2013\)](#) and [Evdokimov and Kolesár \(2019\)](#) consider alternative frameworks that allow for the dimensions of either or both of these vectors to also be growing.

coefficients (36)), the component of α_{tsls} corresponding to the coefficient on \dot{T} can be written as

$$\beta_{\text{tsls}} = \mathbb{E}[RY] / \mathbb{E}[R^2],$$

where $R \equiv \dot{T} - \mathbb{L}[\dot{T}|X]$ are the residuals from projecting the population fitted treatment variable, \dot{T} , onto the covariates, X . Using (35), these residuals can be written as

$$R \equiv \dot{T} - \mathbb{L}[\dot{T}|X] = \left(\dot{Z} + \lambda' X \right) - \mathbb{L} \left[\dot{Z} + \lambda' X | X \right] = \dot{Z} - \mathbb{L}[\dot{Z}|X] \equiv \tilde{Z}.$$

This shows that $\beta_{\text{tsls}} = \mathbb{E}[\tilde{Z}Y] / \mathbb{E}[\tilde{Z}^2]$. Because \tilde{Z} is a residual from a projection onto X , we can also use (35) to write

$$\mathbb{E}[\tilde{Z}^2] = \mathbb{E}[\tilde{Z}\dot{Z}] = \mathbb{E}[\tilde{Z}(\dot{T} - \lambda' X)] = \mathbb{E}[\tilde{Z}T] - \mathbb{E}[\tilde{Z}(T - \dot{T})] = \mathbb{E}[\tilde{Z}T],$$

where the final equality follows because \tilde{Z} is a linear function of I and X , and therefore orthogonal to the first stage residuals, $T - \dot{T}$. *Q.E.D.*

Applying Proposition 3 to Proposition SA.1 shows that β_{tsls} can be written as

$$\beta = \sum_{g,x} \omega_0(g,x) \mu_0(g,x) + \sum_{g,x} \sum_{j=1}^J \omega_j(g,x) (\mu_j(g,x) - \mu_{j-1}(g,x)), \quad (6)$$

with weights given by

$$\omega_j(g,x) = \mathbb{E}[\tilde{Z}^2]^{-1} \mathbb{E} \left[\mathbb{1}[T \geq t_j] \tilde{Z} | G = g, X = x \right] \mathbb{P}[G = g, X = x]. \quad (37)$$

As shown in Proposition 4, whether β_{tsls} is weakly causal is determined by $\omega_j(g,x)$, which is determined by the TSLS specification through \tilde{Z} .

SA.2.2 Rich covariates in more general TSLS specifications

The necessary direction of Theorem 1 extends immediately, just with a more general definition of \tilde{Z} .

Corollary SA.1. Suppose that Assumptions EXO is satisfied. If β_{tsls} is weakly causal, then $\gamma' \mathbb{E}[Z|X = x] = \gamma' \mathbb{L}[Z|X = x]$ for every $x \in \mathcal{X}$.

Proof of Corollary SA.1. If β_{tsls} is weakly causal, then Proposition 4 implies that $\omega_0(g,x) = 0$ for all g and x . Using the same argument as in the proof of Theorem 1 with the form of ω_0 given in (37), this implies that for every x ,

$$0 = \mathbb{E}[\tilde{Z}|X = x] \equiv \mathbb{E}[\dot{Z}|X = x] - \mathbb{L}[\dot{Z}|X = x] = \gamma' \mathbb{E}[Z|X = x] - \gamma' \mathbb{L}[Z|X = x].$$

The condition in Corollary SA.1 is a generalization of rich covariates from the case in which $i(Z, X) = Z$ is scalar, so that γ cancels out, to TSLS specifications with vectors of excluded variables.

SA.2.3 Monotonicity-correct first stages

Corollary SA.1 shows that rich covariates is necessary for β_{tsls} to be weakly causal, but unlike Theorem 1, it does not establish sufficiency. The reason is that with more general TSLS specifications one also has to consider the specification of the first stage relative to the maintained monotonicity condition. In this section, we derive the additional sufficient-and-necessary characterization of the missing piece.

We begin by stating a weaker form of the monotonicity condition (Assumption MON). We follow [Słoczyński \(2020, 2024\)](#) in calling this “weak” monotonicity.

Assumption WMON. (Weak monotonicity) For all $x \in \mathcal{X}$, and all $z, \bar{z} \in \mathcal{Z}$, either

$$\begin{aligned} & \mathbb{P}[T(\bar{z}) \geq T(z)|X = x] = 1 \\ \text{or} \quad & \mathbb{P}[T(z) \geq T(\bar{z})|X = x] = 1. \end{aligned}$$

Assumption WMON is weaker than Assumption MON monotonicity because it allows the direction of monotonicity to depend on x . For example, if $\mathcal{Z} = \{0, 1\}$ and $\mathcal{X} = \{0, 1\}$, then Assumption WMON allows for

$$\begin{aligned} & \mathbb{P}[T(1) \geq T(0)|X = 0] = 1 \\ \text{and} \quad & \mathbb{P}[T(0) \geq T(1)|X = 1] = 1. \end{aligned} \tag{38}$$

If, for example, $\mathcal{T} = \{0, 1\}$ is also binary, then group $G = (0, 1)$ would be compliers conditional on $X = 0$, but they would be defiers conditional on $X = 1$, and conversely for $G = (1, 0)$.

For any x , the order in which Assumption WMON holds between two instrument values can be determined by the conditional mean of T ,

$$p(z, x) \equiv \mathbb{E}[T|Z = z, X = x].$$

If $p(\bar{z}, x) \geq p(z, x)$ then $T(\bar{z}) \geq T(z)$ conditional on $X = x$, and conversely ([Imbens and Angrist, 1994](#); [Vytlacil, 2002](#)). We say that the first stage of the TSLS specification is monotonicity-correct if the first stage fitted values reproduce this ordering, in the sense of predicting higher values of treatment when the instrument is such that individuals choose higher values of treatment.

Definition MC. Let $\dot{t}(z, x) \equiv \gamma' i(z, x) + \lambda' x$ denote the population fitted values in the first stage regression for a realization with $Z = z$ and $X = x$. Suppose that (z, \bar{z}) are both in the support of Z , conditional on $X = x$. Then a TSLS first stage is *monotonicity-correct* for (z, \bar{z}) conditional on $X = x$, if

$$(p(\bar{z}, x) - p(z, x)) \times (\dot{t}(\bar{z}, x) - \dot{t}(z, x)) \geq 0.$$

Definition MC is easiest to appreciate in the case with $\mathcal{T} = \{0, 1\}$, $I = Z$, and $Z \in \{0, 1\}$, so that $\dot{t}(1, x) - \dot{t}(0, x) = \gamma$ is the scalar coefficient on Z in the first stage regression. If Assumption WMON holds with $T(1) \geq T(0)$ conditional on $X = x$, then $p(1, x) - p(0, x) \geq 0$. The TSLS first stage is monotonicity-correct conditional on $X = x$ if and only if $\gamma \geq 0$, so that the linear projection in the first stage reproduces the same sign as the (nonparametric) treatment propensity score.

Including interactions between covariates and instruments in the first stage can help ensure monotonicity-correctness. For example, suppose that X contains a binary component, $X_1 \in \{0, 1\}$, and that $I = [Z, ZX_1]'$ now has two components with first stage coefficient vector $[\gamma_1, \gamma_2]'$, so that for any realization of the other components x_{-1} of X ,

$$\begin{aligned} \dot{t}(1, x_1 = 0, x_{-1}) - \dot{t}(0, x_1 = 0, x_{-1}) &= \gamma_1 \\ \text{and } \dot{t}(1, x_1 = 1, x_{-1}) - \dot{t}(0, x_1 = 1, x_{-1}) &= \gamma_1 + \gamma_2. \end{aligned}$$

This first stage can still be monotonicity-correct conditional on all values of $X = (x_1, x_{-1})$, even if the direction of Assumption WMON is positive when $x_1 = 0$ and negative when $x_1 = 1$, as in (38). The requirement is that $\gamma_1 \geq 0$ and $\gamma_1 + \gamma_2 \leq 0$. Whether this requirement holds depends on the stochastic relationship between Z and the other components, X_{-1} .

The following proposition shows that the missing sufficient condition for Corollary SA.1 has to do with the monotonicity-correctness of its first stage specification. It also shows that if Assumption MON is weakened to Assumption WMON, then monotonicity-correctness characterizes the additional necessary condition beyond rich covariates that the TSLS specification must satisfy in order for β_{tsls} to be weakly causal.

Proposition SA.2. Suppose that Assumptions EXO and WMON are satisfied. Suppose that the TSLS specification for β_{tsls} has rich covariates in the sense of Corollary SA.1. If the TSLS specification is monotonicity-correct for every (z, \bar{z}) , conditional on every x , then β_{tsls} is weakly causal. Conversely, if β_{tsls} is weakly causal, then for every $x \in \mathcal{X}$ the TSLS first stage must be monotonicity-correct for at least one pair (z, \bar{z}) .

Proof of Proposition SA.2. Because rich covariates is assumed, we already know

that $\omega_0(g, x) = 0$ by Corollary SA.1. We also have an expression for $\omega_j(g, x)$ that is similar to (29) in the proof of Theorem 1:

$$\omega_j(g, x) = \mathbb{E}[\tilde{Z}^2]^{-1} \mathbb{C} \left[\dot{t}(Z, X), \mathbb{1}[Z \in \mathcal{Z}_j(g)] \middle| X = x \right] \mathbb{P}[G = g | X = x] \mathbb{P}[X = x], \quad (39)$$

where only differences are that $\dot{t}(Z, X)$ replaces Z , and that the denominator term, $\mathbb{E}[\tilde{Z}^2]$, is already clearly non-negative given the way we've redefined $\tilde{Z} \equiv \dot{Z} - \mathbb{L}[\dot{Z} | X] \equiv \gamma'(Z - \mathbb{L}[Z | X])$ in terms of the effective instrument. The sign of $\omega_j(g, x)$ is determined solely by the covariance term if $\mathbb{P}[G = g | X = x] > 0$, while $\omega_j(g, x) = 0$ for any other (g, x) pairs.

Fix any $x \in \mathcal{X}$. Assumption WMON allows the direction of monotonicity to vary with x , so first enumerate the support of Z as $\{z_0, z_1, \dots, z_K\}$ in order of the treatment propensity score, that is, such that $p(z_{k-1}, x) \leq p(z_k, x)$ for $k = 1, \dots, K$, with any ties being broken arbitrarily. This ordering depends on x , but x has been fixed on the outset, so we keep that implicit in the notation. Now we are going to re-parameterize the covariance term by the indices of Z , so that the second term is monotonic in the treatment propensity score. Do this by defining the bijective function $\zeta : \{0, 1, \dots, K\} \mapsto \{z_0, z_1, \dots, z_K\}$ that maps indices to support points of Z , then let $M \equiv \zeta^{-1}(Z)$, so that $Z = \zeta(M)$. Then we can write the covariance term as

$$\mathbb{C} \left[\dot{t}(Z, X), \mathbb{1}[Z \in \mathcal{Z}_j(g)] \middle| X = x \right] = \mathbb{C} \left[\dot{t}(\zeta(M), X), \mathbb{1}[\zeta(M) \in \mathcal{Z}_j(g)] \middle| X = x \right].$$

The same argument as in the proof of Theorem 1 now shows that the function $k \mapsto \mathbb{1}[\zeta(k) \in \mathcal{Z}_j(g)]$ must be weakly increasing under Assumption MON for any group g with $\mathbb{P}[G = g | X = x] > 0$, because $p(\zeta(k), x) \geq p(\zeta(k-1), x)$, by construction.

Now suppose that the TSLS specification is monotonicity-correct for every (z, \bar{z}) pair, conditional on any x . Then $k \mapsto \dot{t}(\zeta(k), x)$ is also a weakly increasing function of k : as k increases, $p(\zeta(k), x)$ increases, by construction, and hence so does $\dot{t}(\zeta(k), x)$, by hypothesis. It follows that the covariance term is also non-negative (e.g. [Thorisson, 1995](#), Section 2), so that $\omega_j(g, x)$ is non-negative as well. This statement holds for any j , as well as for any g , because $\omega_j(g, x) = 0$ if $\mathbb{P}[G = g | X = x] = 0$. It also holds for any x after defining the indices and ζ function as above. By Proposition 4, we conclude that β_{tsls} is weakly causal.

Conversely, suppose that β_{tsls} is weakly causal, so that $\omega_j(g, x) \geq 0$ for all j, g , and x . If x is such that $\dot{t}(z, x)$ is a constant function of z , then the TSLS specification is trivially monotonicity-correct for every (z, \bar{z}) given x . So, focus on any x for which $\dot{t}(z, x)$ is non-constant. Using this x , re-order Z by the treatment propensity score in the same fashion as above. Suppose $k \mapsto \dot{t}(\zeta(k), x)$ were weakly decreasing. Then the covariance term would be strictly negative, because $k \mapsto \dot{t}(\zeta(k), x)$ is non-constant, and

$k \mapsto \mathbb{1}[\zeta(k) \in \mathcal{Z}_j(g)]$ is weakly increasing and non-constant. This would contradict the hypothesis that $\omega_j(g, x) \geq 0$, so it must be that $k \mapsto \dot{t}(\zeta(k), x)$ is not weakly decreasing. As a consequence, there must exist a k such that $\dot{t}(\zeta(k), x) > \dot{t}(\zeta(k-1), x)$. Because $p(\zeta(k), x) \geq p(\zeta(k-1), x)$, this shows that the TSLS specification is monotonicity-correct for $(\zeta(k), \zeta(k-1))$ given x . *Q.E.D.*

Special cases of the sufficient condition in Proposition SA.2 appear in Angrist and Imbens (1995), Angrist and Pischke (2009), Kolesár (2013), and Słoczyński (2020, 2024). The saturate and weight specification in Angrist and Imbens (1995) and Angrist and Pischke (2009) has a first stage specification that is saturated in both the instruments and covariates, which is automatically monotonicity-correct for any instrument pair, conditional on any covariate value, because it has $\dot{t}(z, x) = p(z, x)$. Kolesár (2013) relaxes this to Definition MC, although stated somewhat differently, and provides a result like the sufficient direction of Proposition SA.2. See also Heckman and Vytlačil (2005, Section 4.3) and Heckman (2010, Section 3.4).

A special case of the necessary direction of Proposition SA.2 was shown by Słoczyński (2020, 2024) for the case when both Z and T are binary. Słoczyński (2020, 2024) observes that if one were to use the linear IV specification with $i(Z, X) = Z$ discussed in the main text, if X were saturated, and if Assumption WMON held but Assumption MON did not hold, then β_{tsls} would not be weakly causal, because some covariate groups would have negatively-weighted treatment effects. This case follows from Proposition SA.2 because there is only a single coefficient $\gamma = \dot{t}(1, x) - \dot{t}(0, x)$ on the excluded variables $i(Z, X)$. This single coefficient cannot have the same sign as $p(1, x) - p(0, x)$ for all x if this sign is different for some x , as would happen if Assumption WMON were satisfied, but Assumption MON were not.

Proposition SA.2 generalizes Słoczyński’s argument to include TSLS specifications with excluded variables that are combinations of multivalued instruments and covariates, and which have non-binary treatments. This opens up a gap between the sufficient and necessary conditions because it is possible, at least in principle, for the first stage specification to be monotonicity-incorrect for some instrument contrasts, as long as it is monotonicity-correct “on average” across all instrument contrasts. This type of fortuitous averaging seems difficult to defend, so for practical purposes we view the gap between sufficient and necessary in Proposition SA.2 as empirically irrelevant.

While Proposition SA.2 generalizes the setting considered by Słoczyński (2020, 2024), it doesn’t change the basic takeaway of his analysis. Even if rich covariates is satisfied, β_{tsls} is not weakly causal unless the first stage specification is sufficiently flexible to reproduce the assumed direction of the monotonicity condition for each covariate value. In the main text, we maintained Assumption MON, and considered the case where $i(Z, X) = Z$; in this setting, the TSLS specification is always monotonicity-

correct and the only consideration for weak causality is the rich covariates condition. If Assumption MON is weakened to Assumption WMON and/or a different specification of the excluded variables $i(Z, X)$ are used, then monotonicity-correctness needs to be considered in addition to the rich covariates condition.

SA.2.4 Weak monotonicity and (ordered) strong monotonicity

An implication of Proposition SA.2 is that if one is only willing to maintain Assumption WMON, then the full saturate and weight (SW) specification (Angrist and Pischke, 2009; Angrist and Imbens, 1995) must be used to ensure that β_{tsls} is weakly causal. This specification has X saturated and specifies the excluded variables $i(Z, X)$ to include interactions between each of the instrument and covariate values. Even if Z is binary, this results in a number of excluded variables equal to the number of covariates. This makes the SW specification vulnerable to many instruments bias. In our working paper (Blandhol et al., 2022), we investigated the extent of many instruments bias in a simulation and found that it was a serious problem for the TSLS estimator. In those simulations, it also remained a problem even when employing various forms of jackknife estimators (Angrist et al., 1999; Akerberg and Devereux, 2009; Kolesár, 2013). Słoczyński (2024) provides some evidence that many instruments bias in the SW specification can be reliably detected through the pre-test developed by Mikusheva and Sun (2022).

The large number of excluded variables in SW are created by interacting X and Z . If these interactions are removed, then the resulting TSLS specification will be monotonicity-correct and thus weakly causal under Assumption MON, but not necessarily under Assumption WMON. Our Assumption MON is actually a bit stronger than the usual statement of the monotonicity condition, such as the original statement in Imbens and Angrist (1994), because it requires Z to be ordered. In Blandhol et al. (2022), we called Assumption MON *ordered* strong monotonicity to reflect this additional requirement. If we drop this restriction, we get what Słoczyński (2024) calls strong monotonicity.

Assumption SMON. (Strong monotonicity) For all $z, \bar{z} \in \mathcal{Z}$, either

$$\begin{aligned} & \mathbb{P}[T(\bar{z}) \geq T(z)|X = x] = 1 \\ \text{or} \quad & \mathbb{P}[T(z) \geq T(\bar{z})|X = x] = 1 \quad \text{for all } x. \end{aligned}$$

Assumption SMON is the same as Assumption MON when Z is binary, but not for more general types of instruments.

Given the preceding discussion, one might think that Assumption SMON is sufficient to ensure that β_{tsls} is monotonicity-correct and therefore weakly causal in specifications

with no interactions, so that $i(Z, X) = Z$. Perhaps surprisingly, this turns out not to be true. The reason is that omitted interaction terms can bias the coefficients on $I = Z$ in a way that contradicts the sign of the propensity score.

For an example of this, suppose that $\mathcal{Z} = \{0, 1, 2\}$ and that X is binary, then specify $I \equiv [\mathbb{1}[Z = 1], \mathbb{1}[Z = 2]] \equiv [Z_1, Z_2]'$ as indicators. Then

$$\dot{t}(2, x) - \dot{t}(1, x) = \gamma_2 - \gamma_1,$$

where $\gamma \equiv [\gamma_1, \gamma_2]'$ is the vector of population coefficients on I for the first stage regression. Even if $p(2, x) - p(1, x) > 0$ for both values of x , it is still possible to have $\gamma_2 - \gamma_1 < 0$, so that the TSLS specification has a monotonicity-incorrect first stage.

To see the intuition, let $V \equiv T - p(Z, X)$ be the difference between T and its conditional mean, then enumerate:

$$\begin{aligned} T &= p(0, 0) + (p(0, 1) - p(0, 0))X + (p(1, 0) - p(0, 0))Z_1 + (p(2, 0) - p(0, 0))Z_2 \\ &\quad (p(1, 1) - p(0, 1))Z_1X + (p(2, 1) - p(0, 1))Z_2X \\ &\equiv X'\lambda^* + I'\gamma^* + W'\zeta + V, \end{aligned}$$

where $W \equiv [Z_1X, Z_2X]'$ and the coefficient vectors collect the appropriate values of $p(z, x)$. Letting $\tilde{I} \equiv I - \mathbb{E}[I|X]$, $\tilde{T} \equiv T - \mathbb{E}[T|X]$, and $\tilde{W} \equiv W - \mathbb{E}[W|X]$, then applying the Frisch-Waugh-Lovell Theorem,

$$\gamma = \mathbb{E}[\tilde{I}\tilde{I}']^{-1} \mathbb{E}[\tilde{I}\tilde{T}] = \mathbb{E}[\tilde{I}\tilde{I}']^{-1} \mathbb{E}[\tilde{I}(\tilde{I}'\gamma^* + \tilde{W}'\zeta + V)] = \gamma^* + \underbrace{\mathbb{E}[\tilde{I}\tilde{I}']^{-1} \mathbb{E}[\tilde{I}\tilde{W}']\zeta}_{\text{omitted variables bias}}.$$

If the bias term is zero, then $\gamma = \gamma^*$ and $\gamma_2 - \gamma_1 = p(2, 0) - p(1, 0) > 0$. However, the bias term is not zero in general.

As a numerical example, suppose that $\mathbb{P}[X = 1] = .5$, with

$$\mathbb{P}[Z = z|X = x] = \begin{cases} .5, & \text{if } z = 0 \\ .05 + .4x, & \text{if } z = 1 \\ .45 - .4x, & \text{if } z = 2 \end{cases}.$$

and set

$$p(z, 0) = \begin{cases} 0, & \text{if } z = 0 \\ .085, & \text{if } z = 1 \\ .170, & \text{if } z = 2 \end{cases} \quad \text{and} \quad p(z, 1) = \begin{cases} 0, & \text{if } z = 0 \\ .425, & \text{if } z = 1 \\ .510, & \text{if } z = 2 \end{cases}.$$

Then it can be shown through some tedious calculations that $\gamma = [.355, .24]'$, so that $\gamma_2 - \gamma_1 < 0$ even while $p(z, x)$ is strictly increasing in z for both values of x . Intuitively, when $Z = 1$ it is overwhelmingly likely that $X = 1$, and when $Z = 2$, it is overwhelmingly likely that $X = 0$. So γ_1 , the regression coefficient on Z_1 , is mostly determined by variation in the $X = 1$ group, while γ_2 , the regression coefficient on Z_2 , is mostly driven by variation in the $X = 0$ group. Yet the change in the conditional mean of T from $Z = 0$ to $Z = 1$ conditional on $X = 1$ is much larger than the change from $Z = 0$ to $Z = 2$ conditional on $X = 2$. As a consequence, γ_1 ends up being larger than γ_2 , violating monotonicity-correctness.

SA.3 Details on the simulation design

In this section, we discuss in detail how we constructed the DGP used in Section 5.4.

We set $X = (X_1, X_2)$ to be a two-dimensional vector of covariates, where X_1 takes many values and X_2 takes nine values. The support of X_1 , which we vary in the simulations, is determined by a Halton sequence on $[0, 1]$, while the support of X_2 is $0, 1/8, 2/8, \dots, 1$. The distribution of both X_1 and X_2 is taken to be uniform, with X_1 and X_2 independent.

We calibrate $\mathbb{E}[Z|X]$ to Card's data by setting Z to be the binary indicator for near four-year college, X_1 to be experience divided by 20, which is roughly the maximum in the data, and X_2 to be one of nine geographic regions. There are $9! = 362,880$ possible ways to map region to the numerical support of X_2 . For each one, we regress Z onto a fully interacted cubic polynomial between X_1 and X_2 , weighting each observation with $X = x$ by the inverse empirical probability that $X = x$. We select the region mapping that yields the regression with the smallest sum of squared residuals. The resulting specification of $\mathbb{E}[Z|X]$ can be written as

$$\mathbb{E}[Z|X = (x_1, x_2)] = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \end{bmatrix} \begin{pmatrix} 1.07 & -2.71 & 7.61 & -5.87 \\ -1.96 & 6.69 & -10.64 & 8.32 \\ 1.72 & -1.91 & -1.65 & -3.19 \\ 0.21 & -8.30 & 20.76 & -9.71 \end{pmatrix} \begin{bmatrix} 1 \\ x_2 \\ x_2^2 \\ x_2^3 \end{bmatrix}, \quad (40)$$

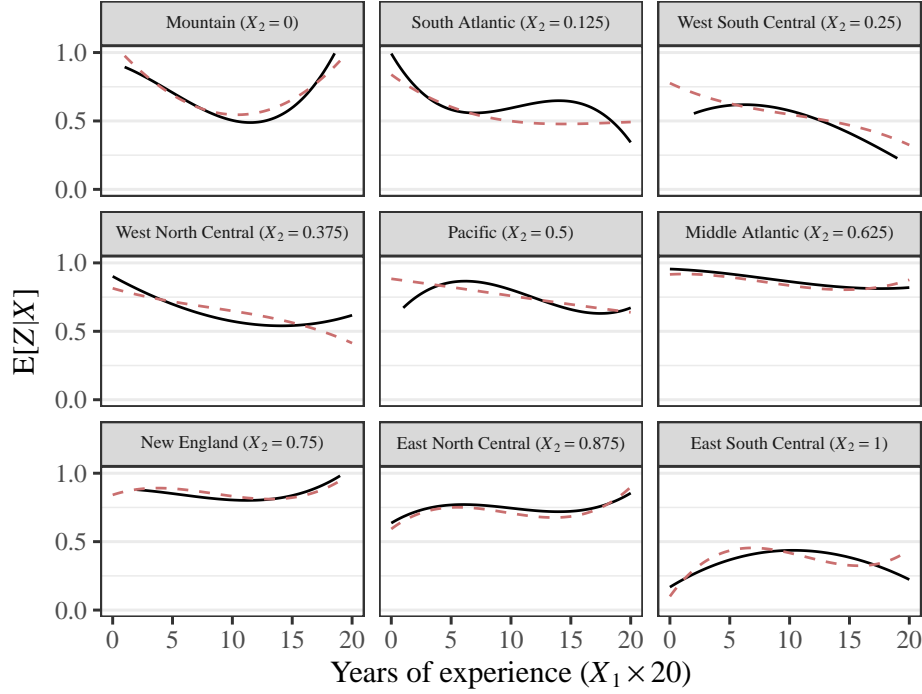
which is linear in 16 terms. Figure SA.1 plots (40) for $20 \times x_1$ against region-specific cubic regressions of the four-year college indicator onto years of experience.

We generate the binary treatment, T , by the threshold-crossing equation

$$T = \mathbb{1}[U \leq p(Z)], \quad (41)$$

where U is distributed uniformly over $[0, 1]$, independently of Z and X . Assumption MON is satisfied under (41) (Vytlacil, 2002). We take $p(0) = .42$ and $p(1) = .54$, which

Figure SA.1: Relationship between college-presence instrument and covariates



Notes: This figure plots the mean of the college-presence instrument used by Card (1995), conditional on region and years of experience. The solid black line is the line of best fit in the data, obtained by regressing $\mathbb{E}[Z|X]$ on a set of region-specific cubic polynomials in years of experience. The dashed red line is $\mathbb{E}[Z|X = x]$ for the DGP in our simulations.

matches the propensity score in Card's data when T is defined as 13 years or more of completed schooling (some college). The group indicator is determined directly from (41) as

$$G = \begin{cases} \text{AT} & \text{if } U \leq p(0) \\ \text{CP} & \text{if } U \in (p(0), p(1)] \\ \text{NT} & \text{if } U > p(1) \end{cases}.$$

To generate potential outcomes $Y(t)$, we let

$$\mathbb{E}[Y(t)|G = g, X = x] = \theta'_0 h(t|g, x)', \quad (42)$$

where $h(t|g, x)$ are basis functions that contain cubic terms in $x = (x_1, x_2)$ that vary freely with t and g . The coefficients on these basis functions, θ_0 , are found as solutions to the optimization problem described ahead. The dimension of h (and θ_0) is $96 = 2 \times 3 \times 16$ for two treatment arms, three groups, and sixteen cubic polynomial coefficients for each group. We generate $Y(t)$ by adding a normal error with mean zero and variance .2 to

(42). The variance of .2 is roughly equal to the sample variance of log wages in Card’s data.

The optimization problem we use to find θ_0 is set up to match some key estimates in Card’s data. To implement the problem, we utilize an observation from [Mogstad et al. \(2018\)](#) that many estimands can be written as weighted averages of θ_0 . We write the weights in these weighted averages as $w\{\text{estimand}\}$. The form of w can be complicated, so we do not provide explicit expressions here, but they depend on h and the joint distribution of (G, T, X, Z) , for which we use the distribution implied by the DGP through the above constructions when X_1 has 24 points of support. Having these linear-in- θ expressions is useful because it allows us to define the optimization problem as a convex quadratic program with linear constraints.

The objective of the optimization problem is to match a weighted average of treated outcomes for always-takers and average untreated outcomes for never-takers. Letting \bar{Y}_{tz} denote the sample average of Y among the subpopulation with $T = t$ and $Z = z$ in Card’s data, the objective we minimize is:

$$\Omega(\theta) \equiv (\bar{Y}_{10} - \theta'w\{\mathbb{E}[Y|T = 1, Z = 0]\})^2 + (\bar{Y}_{01} - \theta'w\{\mathbb{E}[Y|T = 0, Z = 1]\})^2.$$

The constraints involve the following estimates from Card’s data:

- $\bar{Y} \approx 6.26$ is the sample average of log wages.
- $\hat{\beta}_{\text{ols}} \approx .24$ is the OLS estimate of the coefficient on the some college indicator (defined as above) in a regression of log wage on some college, controlling for the covariates used by [Card \(1995, Table 3A, column \(5\)\)](#), which is the same specification we consider in Section 6.1.
- $\hat{\beta}_{\text{iv}} \approx .66$ is the corresponding IV estimate where the near college indicator is used to instrument for some college.
- $\hat{\beta}_{\text{rich}} \approx .43$ is the DDML-PLIV estimate, constructed using the same DDML estimator as in the simulations.
- $\hat{\beta}_{\text{late}} \approx .20$ is the DDML estimate of the unconditional LATE, constructed using the same algorithms as the DDML-PLIV estimate.

We use these estimators to impose the following linear constraints on θ :

$$\theta'w\{\mathbb{E}[Y]\} = \bar{Y} \quad (43)$$

$$\theta'w\{\beta_{\text{ols}}\} = \hat{\beta}_{\text{ols}} \quad (44)$$

$$\theta'w\{\beta_{\text{iv}}\} = \hat{\beta}_{\text{iv}} \quad (45)$$

$$\theta'w\{\beta_{\text{rich}}\} = \hat{\beta}_{\text{rich}} \quad (46)$$

$$\theta'w\{\beta_{\text{late}}\} = \hat{\beta}_{\text{late}}. \quad (47)$$

We additionally constrain θ so that the implied values of $\mathbb{E}[Y(t)|X = x]$ are linear (Assumption LIN) to match the special cases discussed in Propositions 1 and 2:

$$\theta'w\{\mathbb{E}[Y(t)|X = x]\} = \vartheta_0(t) + \vartheta_1(t)x_1 + \vartheta_2(t)x_2 \quad \text{for all } t \text{ and } x, \quad (48)$$

where $\vartheta_1(t), \vartheta_2(t), \vartheta_3(t)$ are additional variables of optimization.¹⁴ We also impose three additional constraints that restrict treatment effects:

$$\theta'w\{\mathbb{E}[Y(1) - Y(0)|G = g, X = x]\} \in [-2, 2] \quad \text{for all } g \text{ and } x. \quad (49)$$

$$\theta'w\{\mathbb{E}[Y(1) - Y(0)|G = \text{NT}]\} = 0 \quad (50)$$

$$\theta'w\{\mathbb{E}[Y(1) - Y(0)|G = \text{AT}, X = x]\} \geq 0 \quad \text{for all } x. \quad (51)$$

The overall optimization problem that we solve is then:

$$\theta_0 = \arg \min_{\theta, \vartheta} \Omega(\theta) \quad \text{s.t.} \quad (43)-(51). \quad (52)$$

The problem is a linearly-constrained convex quadratic program.

¹⁴In practice, we do this by restricting the nonlinear terms of $\theta'w\{\mathbb{E}[Y(t)|X = x]\}$ to be zero, so that the number of imposed constraints does not depend on the number of support points that X has.

SA.4 Additional figures and tables

Figure SA.2: Alternative weights for β_{iv} in the simulation DGP

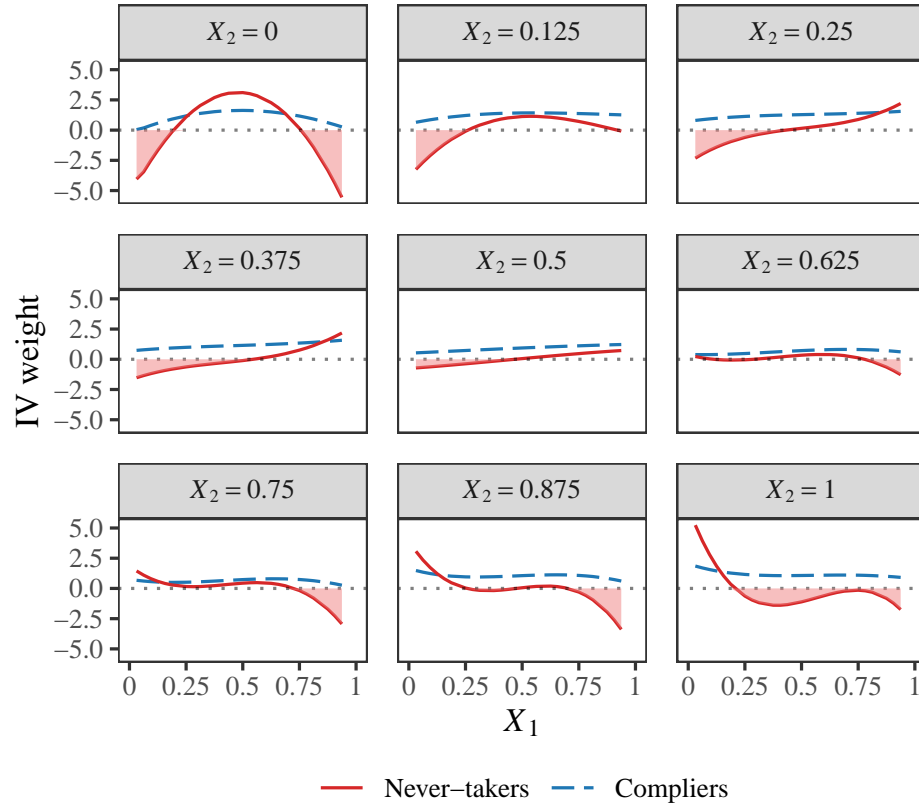


Figure SA.3: Population values of the estimands in the simulation DGP

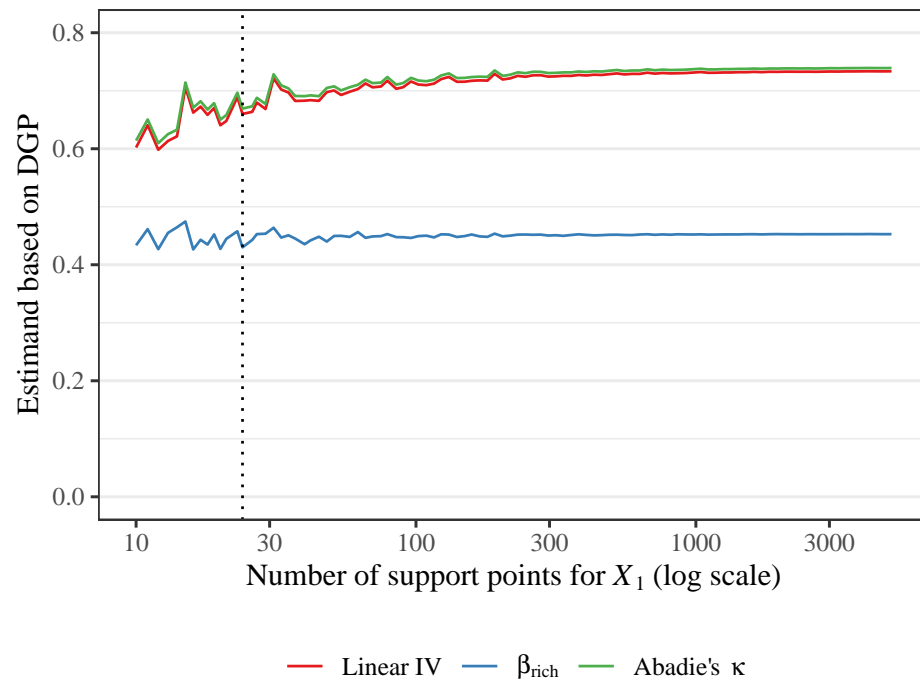


Table SA.1: Detailed simulation results

Estimator	Estimand	Mean (SD)	RMSE	10%	25%	Median	75%	90%	$p < .05$	Avg. CI
$N = 500, \mathcal{X}_1 = 24$										
Linear IV	0.660	0.719 (0.788)	0.840	-0.260	0.165	0.611	1.183	1.958	0.006	4.021
Correctly specified	0.430	0.484 (0.797)	0.799	-0.533	-0.055	0.396	0.955	1.639	0.018	4.355
Saturated	0.430	0.449 (1.251)	1.252	-1.077	-0.251	0.356	1.010	1.845	0.004	17.539
PLIV (DDML)	0.430	0.588 (0.783)	0.799	-0.435	0.050	0.514	1.065	1.734	0.016	3.987
Abadie's κ	0.669	0.726 (0.930)	0.976	-0.382	0.110	0.617	1.230	2.030	0.022	13.309
$N = 500, \mathcal{X}_1 = 100$										
Linear IV	0.709	0.845 (0.867)	0.955	-0.177	0.263	0.710	1.262	2.208	0.010	4.420
Correctly specified	0.445	0.544 (0.801)	0.807	-0.445	0.038	0.438	0.975	1.742	0.014	4.471
Saturated	0.445	0.546 (2.231)	2.234	-1.476	-0.442	0.424	1.280	2.866	0.000	—
PLIV (DDML)	0.445	0.683 (0.779)	0.815	-0.337	0.180	0.600	1.119	1.841	0.012	4.226
Abadie's κ	0.716	0.875 (0.931)	1.025	-0.248	0.305	0.734	1.348	2.206	0.016	13.780
$N = 3,000, \mathcal{X}_1 = 24$										
Linear IV	0.660	0.663 (0.249)	0.341	0.339	0.478	0.641	0.848	1.072	0.088	1.130
Correctly specified	0.430	0.434 (0.238)	0.238	0.098	0.262	0.409	0.601	0.803	0.048	1.085
Saturated	0.430	0.428 (0.245)	0.245	0.079	0.247	0.404	0.601	0.801	0.030	1.153
PLIV (DDML)	0.430	0.522 (0.242)	0.259	0.192	0.346	0.492	0.694	0.917	0.040	1.096
Abadie's κ	0.669	0.685 (0.290)	0.386	0.319	0.474	0.647	0.873	1.168	0.054	1.664
$N = 3,000, \mathcal{X}_1 = 1,000$										
Linear IV	0.731	0.748 (0.249)	0.386	0.399	0.553	0.749	0.922	1.118	0.098	1.170
Correctly specified	0.452	0.454 (0.232)	0.232	0.109	0.278	0.457	0.632	0.788	0.042	1.104
Saturated	0.452	0.451 (0.824)	0.824	-0.694	-0.046	0.381	0.985	1.604	0.000	6.163
PLIV (DDML)	0.452	0.546 (0.237)	0.255	0.203	0.366	0.548	0.723	0.900	0.034	1.117
Abadie's κ	0.737	0.768 (0.296)	0.433	0.359	0.554	0.755	0.968	1.234	0.086	1.707
$N = 10,000, \mathcal{X}_1 = 24$										
Linear IV	0.660	0.659 (0.133)	0.265	0.460	0.557	0.664	0.756	0.862	0.298	0.604
Correctly specified	0.430	0.427 (0.126)	0.126	0.235	0.330	0.430	0.529	0.620	0.038	0.577
Saturated	0.430	0.426 (0.127)	0.127	0.236	0.324	0.431	0.532	0.619	0.040	0.586
PLIV (DDML)	0.430	0.499 (0.128)	0.146	0.304	0.400	0.500	0.602	0.688	0.050	0.584
Abadie's κ	0.669	0.660 (0.138)	0.268	0.456	0.553	0.664	0.766	0.871	0.182	0.870
$N = 10,000, \mathcal{X}_1 = 3,000$										
Linear IV	0.733	0.737 (0.133)	0.314	0.555	0.621	0.735	0.837	0.928	0.422	0.613
Correctly specified	0.453	0.453 (0.127)	0.127	0.262	0.355	0.455	0.552	0.641	0.044	0.579
Saturated	0.453	0.440 (0.325)	0.326	-0.043	0.197	0.431	0.677	0.918	0.006	1.996
PLIV (DDML)	0.453	0.522 (0.127)	0.144	0.340	0.425	0.523	0.624	0.713	0.048	0.586
Abadie's κ	0.739	0.730 (0.145)	0.313	0.531	0.607	0.728	0.841	0.939	0.262	0.928

Notes: Simulations based on 500 replications. Confidence intervals are constructed using HC3 estimators. We do not report an average length for the saturated specification with $N = 500$ and $|\mathcal{X}_1| = 100$ because the standard errors are not defined in many replications.

Table SA.2: Simulation results using different learners separately

Estimator	Mean (SD)	RMSE	10%	25%	Median	75%	90%	$p < .05$	Avg. CI
$N = 500$									
Neural network (2 neurons)	0.564 (0.758)	0.770	-0.432	0.036	0.481	0.998	1.626	0.014	3.959
Neural network (10 neurons)	0.455 (0.796)	0.797	-0.623	-0.083	0.395	0.947	1.564	0.024	4.257
Gradient boosting (stumps)	0.560 (0.520)	0.536	-0.018	0.260	0.442	0.761	1.377	0.010	4.281
Gradient boosting (trees)	0.259 (0.382)	0.419	-0.155	0.065	0.228	0.453	0.734	0.036	3.353
Random forest (mtry = 3)	0.378 (0.581)	0.583	-0.363	-0.004	0.365	0.710	1.213	0.016	4.690
Random forest (mtry = 4)	0.433 (0.630)	0.630	-0.420	-0.024	0.401	0.858	1.380	0.010	4.351
Lasso	0.957 (1.179)	1.291	-0.284	0.294	0.863	1.539	2.646	0.002	14.575
$N = 3,000$									
Neural network (2 neurons)	0.512 (0.240)	0.254	0.186	0.337	0.486	0.686	0.896	0.034	1.096
Neural network (10 neurons)	0.437 (0.240)	0.240	0.105	0.264	0.416	0.613	0.812	0.050	1.087
Gradient boosting (stumps)	0.879 (0.245)	0.512	0.597	0.683	0.827	1.032	1.304	0.311	0.997
Gradient boosting (trees)	0.366 (0.130)	0.145	0.207	0.265	0.342	0.454	0.578	0.036	0.767
Random forest (mtry = 3)	0.324 (0.241)	0.263	0.045	0.170	0.293	0.444	0.700	0.136	1.239
Random forest (mtry = 4)	0.372 (0.183)	0.192	0.136	0.232	0.341	0.481	0.675	0.076	0.911
Lasso	0.714 (0.390)	0.483	0.255	0.404	0.629	0.946	1.359	0.002	2.391
$N = 10,000$									
Neural network (2 neurons)	0.489 (0.127)	0.140	0.302	0.392	0.488	0.589	0.680	0.044	0.584
Neural network (10 neurons)	0.427 (0.127)	0.127	0.235	0.327	0.430	0.533	0.618	0.040	0.577
Gradient boosting (stumps)	1.110 (0.165)	0.700	0.895	0.978	1.102	1.229	1.370	1.000	0.578
Gradient boosting (trees)	0.450 (0.078)	0.081	0.345	0.388	0.443	0.509	0.572	0.012	0.423
Random forest (mtry = 3)	0.318 (0.146)	0.184	0.149	0.211	0.290	0.407	0.541	0.306	0.527
Random forest (mtry = 4)	0.293 (0.092)	0.165	0.174	0.219	0.283	0.351	0.450	0.422	0.365
Lasso	0.428 (0.148)	0.148	0.250	0.312	0.396	0.534	0.665	0.038	0.776

Notes: Simulations based on 500 replications with $|\mathcal{X}_1| = 24$. If one views these algorithms as nonparametric, then their estimand is the value of β_{rich} reported in Table SA.1. The packages used to implement the learners are documented in footnote 9. The penalty parameter for the lasso is selected via K -fold cross validation with $K = 5$.

Table SA.3: Sensitivity to covariate specification in [Dube and Harish \(2020\)](#)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
IV estimate	1.011 (0.523) [0.011]	0.511 (0.231) [0.005]	0.681 (0.355) [0.029]	0.984 (0.519) [0.015]	1.220 (0.640) [0.013]	0.262 (0.170) [0.142]	1.190 (0.639) [0.014]	0.400 (0.211) [0.039]
Polity fixed effects		✓				✓		✓
Decade fixed effects			✓			✓		✓
Missing gender control				✓			✓	✓
Previous monarch controls					✓		✓	✓

Notes: Clustered standard errors are reported in parentheses. Brackets contain p -values for the clustered wild bootstrap procedure implemented by [Dube and Harish \(2020\)](#) with 1000 replications. Column (8) replicates Table 3, column (3) of [Dube and Harish \(2020\)](#). The sample size is 3,586.

Table SA.4: Detailed results from all applications

Application	(1) $\hat{\beta}_{ols}$	(2) $\hat{\beta}_{iv, \text{ no } X}$	(3) $\hat{\beta}_{iv}$	(4) $\hat{\beta}_{rich}$	(5) RESET p -value	(6) Included variables	(7) Sample size
<i>Panel A. Illustrative examples</i>							
Card (1995)	0.075 (0.004)	0.188 (0.026)	0.132 (0.054)	0.122 (0.053)	0.000	14	3,010
Numm and Wantchekon (2011)	-0.203 (0.033)	-0.190 (0.111)	-0.271 (0.088)	-0.071 (0.091)	0.000	99	16,679
Dube and Harish (2020)	0.115 (0.035)	1.011 (0.522)	0.400 (0.211)	0.318 (0.240)	0.000	66	3,586
<i>Panel B. IV survey</i>							
Alesina and Zhuravskaya (2011)	-1.984 (0.639)	-5.727 (1.289)	-3.646 (1.307)	-2.919 (1.115)	0.182	14	97
Autor et al. (2013)	-0.171 (0.028)	-0.666 (0.143)	-0.596 (0.099)	-0.547 (0.091)	0.000	15	1,444
Becker and Woesmann (2009)	0.099 (0.010)	0.422 (0.071)	0.189 (0.027)	0.186 (0.031)	0.012	12	452
Bloom et al. (2012)	1.669 (0.789)	2.708 (1.918)	3.071 (1.253)	2.152 (1.304)	0.000	160	422
Condra et al. (2018)	-0.016 (0.007)	-0.135 (0.128)	-0.092 (0.047)	-0.097 (0.067)	0.923	18	410
Dal Bo et al. (2009)	0.027 (0.006)	-0.015 (0.030)	0.083 (0.037)	0.058 (0.032)	0.000	141	5,502
Dinkelman (2011)	-0.001 (0.005)	0.025 (0.045)	0.095 (0.055)	0.118 (0.118)	0.004	22	1,816
Dippel (2014)	-0.295 (0.048)	-0.676 (0.326)	-0.443 (0.103)	-0.462 (0.235)	0.000	44	182
Gilchrist and Sands (2016)	0.619 (0.058)	0.939 (0.245)	0.843 (0.279)	0.828 (0.292)	0.401	213	2,064
Hornung (2014)	1.741 (0.287)	5.437 (4.180)	3.380 (1.137)	0.892 (0.773)	0.004	10	150

Table SA.5: Details on specifications used in all applications

Study	Specification	Sample size	Num. included variables
Alesina and Zhuravskaya (2011)	Table 7, Panel A, Column (2)	97	14
Autor et al. (2013)	Table 3, Panel I, Column (6)	1,444	15
Becker and Woessmann (2009)	Table 3, Column (2)	452	12
Bloom et al. (2012)	Table 2, Column (7)	422	160
Card (1995)	Table 3, Panel A, Column (5)	3,010	14
Condra et al. (2018)	Table 2, Panel A, Column (3)	410	18
Dal Bo et al. (2009)	Table 5, Panel B, Column (3)	5,501	141
Dinkelman (2011)	Table 4, Column (8)	1,816	22
Dippel (2014)	Table 5, Panel B, Column (6)	182	44
Dube and Harish (2020)	Table 3, Column (3)	3,586	66
Gilchrist and Sands (2016)	Table 4, Column (6)	2,064	213
Hornung (2014)	Table 4, Column (5)	150	10
Nunn and Wantchekon (2011)	Table 6, Column (2)	16,679	99