

An Omitted Variable Bias Framework for Sensitivity Analysis of Instrumental Variables

BY CARLOS CINELLI

Department of Statistics, University of Washington, Seattle.

cinelli@uw.edu

5

AND CHAD HAZLETT

Department of Statistics and Data Science, University of California, Los Angeles.

chazlett@ucla.edu

ABSTRACT

We develop an omitted variable bias framework for sensitivity analysis of instrumental variable (IV) estimates that naturally handles multiple side-effects (violations of the exclusion restriction assumption) and confounders (violations of the ignorability of the instrument assumption) of the instrument, exploits expert knowledge to bound sensitivity parameters, and can be easily implemented with standard software. Specifically, we introduce sensitivity statistics for routine reporting, such as (extreme) *robustness values* for IV estimates, describing the minimum strength that omitted variables need to have to change the conclusions of an IV study. Next we provide visual displays that fully characterize the sensitivity of point estimates and confidence intervals to violations of the standard IV assumptions. Finally, we offer formal bounds on the worst possible bias under the assumption that the maximum explanatory power of omitted variables is no stronger than a multiple of the explanatory power of observed variables. Conveniently, many pivotal conclusions regarding the sensitivity of the IV estimate (e.g. tests against the null hypothesis of zero causal effect) can be reached simply through separate sensitivity analyses of the effect of the instrument on the treatment (the *first stage*) and the effect of the instrument on the outcome (the *reduced form*). We apply our methods in a running example that uses proximity to college as an instrumental variable to estimate the returns to schooling.

10

15

20

25

Some key words: Instrumental Variables; Omitted Variable Bias; Sensitivity Analysis; Robustness Values.

1. INTRODUCTION

30

Unobserved confounding often complicates efforts to make causal claims from observational data (e.g. [Pearl, 2009](#); [Imbens and Rubin, 2015](#)). Instrumental variable (IV) regression offers a powerful and widely used tool to address unobserved confounding, by exploiting exogenous sources of variation of the treatment (e.g. [Angrist et al., 1996](#); [Angrist and Pischke, 2009](#)). IV methods are “a central part of the econometrics canon since the first half of the twentieth century” ([Imbens, 2014](#), p.324), and, beyond economics, are now prominent tools in the arsenal of

35

investigators seeking to make causal claims across the social sciences, epidemiology, medicine, genetics, and other fields (see e.g. [Hernán and Robins, 2006](#); [Burgess and Thompson, 2015](#)).

Yet, IV methods carry their own set of demanding assumptions. Principally, conditionally on
 40 certain observed covariates, an instrumental variable must not be confounded with the outcome,
 and it should influence the outcome only by affecting uptake of the treatment. These assumptions
 can be violated by omitted confounders of the instrument-outcome association, and by omitted
 side-effects of the instrument that influence the outcome via paths other than through the treat-
 ment.¹ Although in certain cases the IV assumptions may entail testable implications ([Pearl,](#)
 45 [1995](#); [Gunsilius, 2020](#); [Kédagni and Mourifié, 2020](#)), they are often unverifiable and must be de-
 fended by appealing to domain knowledge. Whether a given IV study identifies the causal effect
 of interest, then, turns on debates as to whether these assumptions hold.

Particularly in recent years, economists and other scholars have adopted a more skeptical
 posture towards IV methods, emphasizing the importance of both defending the credibility of
 50 these assumptions as well as assessing the consequences of their failures (e.g., [Deaton, 2009](#);
[Heckman and Urzua, 2010](#)). Extensive reviews of many widely-used instrumental variables have
 catalogued several plausible violations of the exclusion restriction for such instruments (e.g.
[Gallen, 2020](#); [Mellon, 2020](#)). More worrisome, if the IV assumptions fail to hold, it is well
 known that the bias of the IV estimate may be *worse* than the original confounding bias of the
 55 simple regression estimate that the IV was supposed to address ([Bound et al., 1995](#)). Therefore,
 researchers are also advised to perform *sensitivity analyses* to assess the degree of violation of
 the IV assumptions that would be required to alter the conclusions of an IV study.²

While a number of sensitivity analyses for IV have been proposed ([DiPrete and Gangl, 2004](#);
[Small, 2007](#); [Small and Rosenbaum, 2008](#); [Conley et al., 2012](#); [Wang et al., 2018](#); [Masten and](#)
 60 [Poirier, 2021](#)), they have rarely been employed in practice.³ We suggest several reasons for this
 slow uptake. First, the traditional approach for the sensitivity of IV has focused on parameteriz-
 ing violations of the IV assumptions with a single parameter summarizing the overall bias in the
 association of the instrument with the outcome. While this parameterization may be well-suited
 when the bias is only due to the direct effect of the instrument on the outcome (not through
 65 the treatment), it is not as straightforward to use when reasoning about multiple side-effects or
 confounders of the instrument, in which case that sensitivity parameter is a complicated com-
 posite of many sources of bias (see Supplementary Material for a comparison of our proposal
 with the traditional approach to the sensitivity of IV). Second, while users of IV methods are
 instructed to routinely report quantities to diagnose certain inferential problems such as “weak
 70 instruments” (e.g. the F-statistic, [Stock and Yogo, 2002](#)), we lack sensitivity statistics that can

¹ In the recent IV literature, the first assumption is usually called *exogeneity*, *ignorability*, or *unconfoundedness* of the instrument, whereas the second assumption is called the *exclusion restriction* ([Angrist and Pischke, 2009](#); [Imbens and Rubin, 2015](#)). In earlier econometric works, these two assumptions were often combined into one, also labeled the exclusion restriction ([Imbens, 2014](#)).

² In this paper, we focus exclusively on sensitivity to unobserved variables that violate the IV assumptions of ignorability and exclusion. In doing so we are concerned with identification challenges, which are separate from inferential issues posed by weak instruments (see e.g. [Nelson and Startz, 1990](#); [Staiger and Stock, 1994](#); [Kleibergen, 2002](#); [Moreira, 2003, 2009](#); [Andrews et al., 2019](#)), or robustness to different choices of control variables, different estimators, outliers or effect heterogeneity (see, e.g. [Blundell et al., 2001](#); [Belzil and Hansen, 2002](#); [Jaeger and Parys, 2009](#) for examples applied to returns to education). Furthermore, we note that conventional statistical concerns regarding weak instruments depend on sample size. By contrast, our sensitivity analysis examines whether such instruments would be fragile in the face of omitted variables, a property which is unaffected by sample size. Consequently, an instrument deemed “strong” by conventional statistical tests may prove extremely fragile to identification violations. See Remark 4.

³ In economics, only 1 out of 27 papers using IV published in the *American Economic Review* in 2020 performed formal sensitivity analysis (see Supplementary Materials for a description of the data collection procedure). In political science, this number was 1 out of 12 papers across the top three general interest journals (*American Political Science Review*, *American Journal of Political Science*, and *Journal of Politics*) for 2019 ([Cinelli and Hazlett, 2020](#)). In Sociology, it was zero out of 34, in the *American Journal of Sociology* and the *American Sociological Review* from 2004 to 2022 ([Felton and Stewart, 2022](#)).

quickly communicate how robust an IV study is to violations in the form of omitted confounders or side-effects of the instrument. Finally, it is often difficult to connect the formal results of a sensitivity analysis to a cogent argument about what types of biases can be ruled out by expert knowledge.

In this paper, we develop an omitted variable bias (OVB) framework for assessing the sensitivity of IV estimates that aims to address these challenges. Building on the Anderson-Rubin approach to IV estimation (Anderson and Rubin, 1949) and on recent developments of OVB for ordinary least squares (OLS) (Cinelli and Hazlett, 2020), we develop a simple suite of sensitivity analysis tools for IV that: (i) naturally handles violations due to multiple side-effects and confounders, possibly acting non-linearly; (ii) is well suited for routine reporting; and (iii) exploits expert knowledge to bound sensitivity parameters.⁴ Specifically, we introduce two main sensitivity statistics for IV estimates: (i) the *robustness value* (RV) describes the minimum strength of association (in terms of partial R^2) that omitted variables (side-effects or confounders) need to have, both with the instrument and with the outcome, in order to change the conclusions of the study; and (ii) the *extreme robustness value*, which describes the minimal strength of association that omitted variables need to have with the *instrument alone* in order to be problematic. Routine reporting of these quantities provides a quick and simple way to improve the transparency and facilitate the assessment of the credibility of IV studies. Next, we offer intuitive graphical tools for investigators to assess how postulated confounding of any degree would alter the IV hypothesis tests, as well as lower or upper limits of confidence intervals. These tools can be supplemented with formal bounds on the worst possible bias that side-effects or confounders could cause, under the assumption that the maximum explanatory power of these omitted variables is no stronger than a chosen multiple of the explanatory power of one or more observed variables.

A final contribution of this paper is the proposal of a novel *bias-adjusted critical value* that accounts for a postulated degree of omitted variable bias. Notably, this correction on the critical value does not depend on the observed data, and can be computed by simply postulating a hypothetical partial R^2 of the omitted variables with the dependent and independent variables of the OLS regression. Applied researchers can thus quickly and easily perform sensitivity analysis by simply substituting traditional thresholds with bias-adjusted thresholds, when testing a particular null hypothesis, or when constructing confidence intervals. All proofs and details can be found in the Supplementary Materials. Open-source software for R implements the methods discussed in this paper: <https://github.com/carloscinelli/iv.sensemakr>.

2. RUNNING EXAMPLE

We begin by introducing the running example and reviewing the required background on IV.

2.1. Ordinary least squares and the OVB problem

Many observational studies have established a positive and large association between educational achievement and earnings using regression analysis. Here we consider the work of Card (1993), which employed a sample of $n = 3,010$ individuals from the National Longitudinal Survey of Young Men (NLSYM).

⁴ Here we focus on the one treatment and one instrument (just-identified) case. We do so for two reasons. First, thoroughly considering how identification assumptions may be violated is complicated enough with one instrument (Angrist and Pischke, 2009). Second, most applied IV work uses this approach. Reviewing papers in the *American Economic Review* and 15 other journals of the *American Economic Association*, Young (2022) finds that 80% of IV regressions used a single instrument. Even in multiple-instrument studies, it is not uncommon for researchers to also report and give special focus to the analysis of their best single instrument (Angrist and Pischke, 2009), or to combine multiple instruments into a single instrument.

Considering the following multiple linear regression $Y = \hat{\tau}_{OLS,res}D + \mathbf{X}\hat{\beta}_{OLS,res} + \hat{\varepsilon}_{OLS,res}$, where Y denotes *Earnings* and measures the log transformed hourly wages of the individual D denotes *Education* and consists of an integer-valued variable indicating the completed years of education of the individual, and the matrix \mathbf{X} comprises race, experience, and a set of regional factors, Card concluded that each additional year of schooling was associated with approximately 7.5% higher wages.

Educational achievement, however, is not randomly assigned; perhaps individuals who obtain more education have higher wages for other reasons, such as family background, or higher levels of some other unobserved characteristic such as *Ability* or *Motivation*. If data on these variables were available, then further adjustment for such variables would capture the causal effect of educational attainment on schooling, as in $Y = \hat{\tau}_{OLS}D + \mathbf{X}\hat{\beta}_{OLS} + \mathbf{U}\hat{\gamma}_{OLS} + \hat{\varepsilon}_{OLS}$, where \mathbf{U} is a set of variables that, along with \mathbf{X} , eliminates confounding concerns.⁵ Unfortunately, such detailed information on individuals is not available, and researchers may not agree on which variables \mathbf{U} are needed. Regression estimates that adjust for only a partial list of characteristics (such as \mathbf{X}) may suffer from OVB, likely overestimating the true returns to schooling.

2.2. Instrumental variables as a solution to the OVB problem

Instrumental variable methods offer an alternative route to estimate the causal effect of schooling on earnings without having data on the unobserved variables \mathbf{U} . The key for such methods to work is to find a new variable (the *instrument*) that changes the incentives to educational achievement, but is associated with earnings only through its effect on education. To that end, Card (1993) proposed exploiting the role of geographic differences in college accessibility. In particular, consider the variable *Proximity*, encoding an indicator of whether the individual grew up in an area with a nearby accredited 4-year college, which we denote by Z . Students who grow up far from the nearest college may face higher educational costs, discouraging them from pursuing higher level studies. Next, and most importantly, Card (1993) argues that, conditional on the set of observed variables \mathbf{X} (available on the NLSYM), whether one lives near a college is not itself confounded with earnings, nor does proximity to college affect earnings apart from its effect on years of education. If we believe such assumptions hold it is possible to recover a valid estimate of the (weighted average of local) average treatment effect(s) of *Education* on *Earnings* by simply taking the ratio of two OLS coefficients⁶, one measuring the effect of *Proximity* on *Earnings*, and another measuring the effect of *Proximity* on *Education*, as in the two OLS models

$$\text{First Stage: } D = \hat{\theta}_{res}Z + \mathbf{X}\hat{\psi}_{res} + \hat{\varepsilon}_{d,res}, \quad (1)$$

$$\text{Reduced Form: } Y = \hat{\lambda}_{res}Z + \mathbf{X}\hat{\beta}_{res} + \hat{\varepsilon}_{y,res}. \quad (2)$$

Throughout the paper we refer to these equations as the *first stage* (Equation 1) and the *reduced form* (Equation 2), as these are now common usage (Angrist and Pischke, 2009; Imbens and Rubin, 2015; Andrews et al., 2019). The coefficient for *Proximity* (Z) on the first-stage regression reveals that those who grew up near a college indeed have higher educa-

⁵ I.e., the set $\{\mathbf{X}, \mathbf{U}\}$ is sufficient to render the treatment assignment ignorable. In graphical terms, the set would satisfy the backdoor criterion (see, e.g. Pearl, 2009; Angrist and Pischke, 2009). Beyond ignorability, if the treatment effect is heterogeneous, this may affect the causal interpretation of $\hat{\tau}_{OLS}$ (e.g. Angrist and Pischke, 2009).

⁶ Conditions that allow a causal interpretation of the traditional IV estimand (also known as the 2SLS *estimand*) are extensively discussed elsewhere and will not be reviewed here, see Angrist et al. (1996); Angrist and Pischke (2009); Imbens (2014); Swanson et al. (2018); Słoczyński (2020) and Blandhol et al. (2022). In particular, Blandhol et al. (2022) provides conditions for a “weakly causal” interpretation of the traditional IV estimand. Here we start from the premise that the researcher has already decided she is interested in the results of Equations (4)–(6). We note the bulk of current applied work using instrumental variables takes this form, and non-parametric estimation is still rare in practice (Blandhol et al., 2022, p.11). It is nevertheless possible to extend our tools to nonparametric settings leveraging recent results in Chernozhukov et al. (2022). We leave this to future work.

tional attainment, having completed an additional 0.32 years of education, on average. Likewise, the coefficient for *Proximity* (Z) on the reduced-form regression suggests that those who grew up near a college have 4.2% higher earnings. The IV estimate is then given by the ratio, $\hat{\tau}_{\text{res}} := \hat{\lambda}_{\text{res}}/\hat{\theta}_{\text{res}} \approx 0.042/0.319 \approx 0.132$. The value of $\hat{\tau}_{\text{res}} \approx 0.132$ suggests that, contrary to the OLS estimate of 7.5%, and perhaps surprisingly, each additional year of schooling instead raises wages by much more—13.2%.

The ratio $\hat{\lambda}_{\text{res}}/\hat{\theta}_{\text{res}}$ is sometimes called the *indirect least squares* (ILS) estimator. Inference in the ILS framework is usually performed using the delta-method. A closely related approach is denoted by *two-stage least squares* (2SLS), in which one saves the predictions of the first-stage regression, and then regress the outcome on these fitted values. By the Frisch-Waugh-Lovell (FWL) theorem (Frisch and Waugh, 1933; Lovell, 1963) one can readily show that 2SLS and ILS are numerically identical.

2.3. Anderson-Rubin regression and Fieller's theorem

The methods of ILS and 2SLS may prove unreliable when the first-stage coefficient is too close to zero relative to the sampling variability of its estimator. This is known as the “weak instrument” problem.⁷ The Anderson-Rubin (AR) regression (Anderson and Rubin, 1949) provides one approach to constructing confidence intervals with correct coverage, regardless of the strength of the first stage. Additionally, it also yields the uniformly most powerful unbiased test under this setup (Moreira, 2009). The approach starts by creating the random variable $Y_{\tau_0} := Y - \tau_0 D$ in which we subtract from Y a putative causal effect of D , namely, τ_0 . If Z is a valid instrument, under the null hypothesis $H_0 : \tau = \tau_0$, we should not see an association between Y_{τ_0} and Z , conditional on \mathbf{X} . In other words, if we run the OLS model

$$\text{Anderson-Rubin: } Y_{\tau_0} = \hat{\phi}_{\tau_0, \text{res}} Z + \mathbf{X} \hat{\beta}_{\tau_0, \text{res}} + \hat{\varepsilon}_{\tau_0, \text{res}}, \quad (3)$$

we should find that $\hat{\phi}_{\tau_0, \text{res}}$ is equal to zero, but for sampling variation. To test the null hypothesis $H_0 : \phi_{\tau_0, \text{res}} = 0$ in the Anderson-Rubin regression is thus equivalent to test the null hypothesis $H_0 : \tau = \tau_0$. The $1 - \alpha$ confidence interval is constructed by collecting all values τ_0 such that the null hypothesis $H_0 : \phi_{\tau_0, \text{res}} = 0$ is not rejected at the chosen significance level α . This approach is numerically identical to Fieller's theorem (Fieller, 1954). Finally, it is convenient to define the point estimate $\hat{\tau}_{\text{AR}, \text{res}}$ as the value τ_0 which makes $\hat{\phi}_{\tau_0, \text{res}}$ exactly equal to zero. By the FWL theorem, we can easily show that $\hat{\tau}_{\text{AR}, \text{res}}$ is numerically identical to 2SLS and ILS.

2.4. The IV estimate may suffer from OVB

The previous IV estimate relies on the assumption that, conditional on \mathbf{X} , *Proximity* and *Earnings* are unconfounded, and the effect of *Proximity* on *Earnings* must go entirely through *Education*. As is often the case, neither assumption is easy to defend. First, the same factors that might confound the relationship between *Education* and *Earnings* could similarly confound the relationship of *Proximity* and *Earnings* (e.g. family wealth or connections). Second, as argued in Card (1993), the presence of a college nearby may be associated with high school quality, which in turn also affects earnings. Finally, other geographic confounders can make some localities likely to both have colleges nearby and lead to higher earnings. These are only coarsely conditioned on by the observed regional indicators, and residual biases may still remain.

⁷ See Andrews et al. (2019) for an extensive review of inference with weak instruments. See also Kleibergen (2002) Moreira (2003) and Moreira (2009). An intuitive visual comparison between the delta-method and Fieller's approach is given by Hirschberg and Lye (2010, 2017).

Therefore, instead of adjusting for \mathbf{X} only, as in the previous regressions, we should have adjusted for both the observed covariates \mathbf{X} and unobserved covariates \mathbf{W} as in

$$\textbf{First Stage: } D = \hat{\theta}Z + \mathbf{X}\hat{\psi} + \mathbf{W}\hat{\delta} + \hat{\varepsilon}_d, \quad (4)$$

$$\textbf{Reduced Form: } Y = \hat{\lambda}Z + \mathbf{X}\hat{\beta} + \mathbf{W}\hat{\gamma} + \hat{\varepsilon}_y, \quad (5)$$

$$\textbf{Anderson-Rubin: } Y_{\tau_0} = \hat{\phi}_{\tau_0}Z + \mathbf{X}\hat{\beta}_{\tau_0} + \mathbf{W}\hat{\gamma}_{\tau_0} + \hat{\varepsilon}_{\tau_0}, \quad (6)$$

where \mathbf{W} stands for all unobserved factors necessary to make *Proximity* a valid instrument for the effect of *Education* on *Earnings* (e.g., *Family Wealth*, *High School Quality*, *Place of Residence*).⁸ Our task is thus to characterize how the IV point estimates and confidence intervals, as given by the OLS regressions (4)-(6), would have changed due to the inclusion of omitted variables \mathbf{W} .

3. EXTENSIONS TO THE OVB FRAMEWORK FOR OLS

As we have seen, IV estimates are based on OLS estimates. Therefore, we should be able to leverage sensitivity analysis tools for OLS to examine the sensitivity of IV. To that end, this section first extends and refines several results for the sensitivity analysis of arbitrary OLS estimates. These results are not only useful on their own right, but, importantly, they will later be applied to the development of a suite sensitivity analysis tools for IV in Section 4.

Throughout the paper, we impose the following regularity condition.

Assumption 1 (Full Rank). The matrices of independent variables in (4)-(6) have full rank.

This ensures all relevant quantities discussed below are finite. Finally, for concreteness, in this section we discuss the OVB framework in the context of the reduced-form regression. Readers should keep in mind, however, that all results presented here hold for *arbitrary OLS estimates*—including, but not limited to, the first stage and the Anderson-Rubin regression. The logical implications of the sensitivity of these auxiliary regressions for the sensitivity of IV itself are deferred to Section 4.

3.1. Preliminaries

We start by briefly establishing key ideas, formulae, and notation from prior work (Cinelli and Hazlett, 2020). Consider the regression coefficient estimate $\hat{\lambda}$ and the classical (i.e., homoskedastic) standard error estimate $\widehat{\text{se}}(\hat{\lambda})$ of Equation (5), namely, the OLS regression of the outcome Y on the instrument Z , adjusting for a set of observed covariates \mathbf{X} and (for now) a single *unobserved* covariate W (we generalize to multivariate W below). Here Y , Z and W are $(n \times 1)$ vectors, \mathbf{X} is an $(n \times p)$ matrix (including a constant), with n observations; $\hat{\lambda}$, $\hat{\beta}$ and $\hat{\gamma}$ are the OLS coefficient estimates and $\hat{\varepsilon}_y$ the corresponding residuals. As W is unobserved, the investigator instead estimates the *restricted* model of Equation (2) where $\hat{\lambda}_{\text{res}}$ and $\hat{\beta}_{\text{res}}$ are the coefficients of the restricted OLS adjusting for Z and \mathbf{X} alone, and $\hat{\varepsilon}_{y,\text{res}}$ the corresponding residuals. The OVB framework seeks to answer the following question: how do the estimates from the restricted OLS model compare with the estimates from the full OLS model?

Let $R_{Y \sim W|Z, \mathbf{X}}^2$ denote the (sample) partial R^2 of W with Y , after controlling for Z and \mathbf{X} , and let $R_{Z \sim W|\mathbf{X}}^2$ denote the partial R^2 of W with Z after adjusting for \mathbf{X} . It is also useful to define Cohen's partial f^2 , e.g., $f_{Z \sim W|\mathbf{X}}^2 := \frac{R_{Z \sim W|\mathbf{X}}^2}{1 - R_{Z \sim W|\mathbf{X}}^2}$ which will appear frequently through-

⁸ See Supplementary Material for canonical causal diagrams illustrating settings in which $\{\mathbf{X}, \mathbf{W}\}$ renders Z a valid instrument for the effect of D and Y . Equivalent assumptions can be articulated in the potential outcomes framework (Angrist et al., 1996; Pearl, 2009; Swanson et al., 2018).

out derivations. Given the point estimate and (estimated) standard error of the restricted model actually run, $\hat{\lambda}_{\text{res}}$ and $\widehat{\text{se}}(\hat{\lambda}_{\text{res}})$, these two R^2 values are sufficient to recover $\hat{\lambda}$ and $\widehat{\text{se}}(\hat{\lambda})$. 225

THEOREM 1 (OVb IN THE PARTIAL R^2 PARAMETERIZATION). *Under Assumption 1, the absolute difference between the restricted and full OLS estimates is given by,*

$$|\hat{\lambda}_{\text{res}} - \hat{\lambda}| = \underbrace{\sqrt{R_{Y \sim W|Z, \mathbf{X}}^2 f_{Z \sim W|\mathbf{X}}^2}}_{BF} \times \frac{sd(Y^{\perp \mathbf{X}, Z})}{sd(Z^{\perp \mathbf{X}})}; \quad (7)$$

moreover, the (classical) standard error of the full OLS estimate is given by 230

$$\widehat{\text{se}}(\hat{\lambda}) = \underbrace{\sqrt{\frac{1 - R_{Y \sim W|Z, \mathbf{X}}^2}{1 - R_{Z \sim W|\mathbf{X}}^2}}}_{SEF} \times \frac{sd(Y^{\perp \mathbf{X}, Z})}{sd(Z^{\perp \mathbf{X}})} \times \sqrt{\frac{1}{\text{df} - 1}}, \quad (8)$$

where $sd(Y^{\perp \mathbf{X}, Z})$ is the (sample) residual standard deviation of Y after removing the part linearly explained by $\{\mathbf{X}, Z\}$, $sd(Z^{\perp \mathbf{X}})$ is the (sample) residual standard deviation of Z after removing the part linearly explained by \mathbf{X} , and $\text{df} = n - p - 1$ is the residual degrees of freedom from the restricted model (2). To aid interpretation, we call the term BF in (7) the “bias factor” of W , and the term SEF in (8) the “standard error factor” of W . 235

For simplicity of exposition, throughout the text we usually refer to a single omitted variable W . These results, however, can be used for performing sensitivity analyses considering multiple omitted variables $\mathbf{W} = [W_1, W_2, \dots, W_l]$, and thus also non-linearities and functional form misspecification of observed variables. In such cases, barring an adjustment in the degrees of freedom, the equations are conservative, and reveal the maximum bias a multivariate \mathbf{W} with such pair of partial R^2 values could cause (Cinelli and Hazlett, 2020, Sec. 4.5). 240

Note Theorem 1 is stated in terms of sample estimates. All results presented in this paper are of this type: they are exact algebraic results of how traditional OLS coefficients and standard error estimates change due to the inclusion of omitted variables. Conditions under which traditional OLS estimates yield valid inferences are well-known and thus omitted. 245

3.2. Bias-adjusted critical values

We now introduce a novel correction to traditional critical values that researchers can use to account for omitted variable bias. Let $t_{\alpha, \text{df} - 1}^* > 0$ denote the (absolute value of) the critical value for a standard t-test with significance level α and $\text{df} - 1$ degrees of freedom. Now let $\text{LL}_{1-\alpha}(\lambda)$ be the lower limit and $\text{UL}_{1-\alpha}(\lambda)$ be the upper limit of a $1 - \alpha$ confidence interval for λ in the full model, i.e., 250

$$\text{LL}_{1-\alpha}(\lambda) := \hat{\lambda} - t_{\alpha, \text{df} - 1}^* \times \widehat{\text{se}}(\hat{\lambda}), \quad \text{UL}_{1-\alpha}(\lambda) := \hat{\lambda} + t_{\alpha, \text{df} - 1}^* \times \widehat{\text{se}}(\hat{\lambda}). \quad (9)$$

Considering the worst-case direction of the bias that further reduces the lower limit (or increases the upper limit) in (9), Equations (7) and (8) of Theorem 1 imply that both quantities can be written as a function of the restricted estimates and a new multiplier. 255

THEOREM 2 (BIAS ADJUSTED CRITICAL VALUE). *Under Assumption 1, for given $\mathbf{R}^2 = \{R_{Y \sim W|Z, \mathbf{X}}^2, R_{Z \sim W|\mathbf{X}}^2\}$, and α , consider the direction of the bias that reduces $\text{LL}_{1-\alpha}(\lambda)$. Then*

$$\text{LL}_{1-\alpha}(\lambda) = \hat{\lambda}_{\text{res}} - t_{\alpha, \text{df} - 1}^{\dagger, \mathbf{R}^2} \times \widehat{\text{se}}(\hat{\lambda}_{\text{res}}). \quad (10)$$

Conversely, considering the direction of the bias that increases $UL_{1-\alpha}(\lambda)$, we have

$$UL_{1-\alpha}(\lambda) = \hat{\lambda}_{res} + t_{\alpha, df-1, R^2}^\dagger \times \widehat{se}(\hat{\lambda}_{res}). \quad (11)$$

Here $t_{\alpha, df-1, R^2}^\dagger$ denotes the bias-adjusted critical value

$$t_{\alpha, df-1, R^2}^\dagger := SEF \sqrt{df/(df-1)} \times t_{\alpha, df-1}^* + BF \sqrt{df}, \quad (12)$$

where BF and SEF are the bias and standard error factors of Theorem 1.

As the subscript $R^2 = \{R_{Y \sim W|Z, X}^2, R_{Z \sim W|X}^2\}$ conveys, $t_{\alpha, df-1, R^2}^\dagger$ depends on both sensitivity parameters. Notably, this correction does not depend on the observed data, but for the degrees of freedom. In other words, the bias correction is a function of the strength of unobserved confounding and the sample size alone. This allows one to quickly assess the robustness of reported findings to omitted variables of any postulated strength R^2 , by simply comparing the reported t-statistic with the desired adjusted critical value, even without access to the original data.

Example 1. It is instructive to consider the case in which the omitted variable W has equal strength of association with Y and Z , i.e., $R_{Y \sim W|Z, X}^2 = R_{Z \sim W|X}^2 = R^2$. We then have that $SEF = 1$ and $BF = R^2/\sqrt{1-R^2}$ resulting in a very simple correction formula,

$$t_{\alpha, df-1, R^2}^\dagger \approx t_{\alpha, df-1}^* + \frac{R^2}{\sqrt{1-R^2}} \sqrt{df}, \quad (13)$$

where we employ the approximation $\sqrt{df/(df-1)} \approx 1$. Table 1 shows the adjusted critical values for this case, considering different strengths of the omitted variable and various sample sizes.

R^2	Degrees of Freedom (sample size)				
	100	1,000	10,000	100,000	1,000,000
0.00	1.98	1.96	1.96	1.96	1.96
0.01	2.08	2.28	2.97	5.14	12.01
0.02	2.19	2.60	3.98	8.35	22.16
0.03	2.29	2.92	5.01	11.59	32.42
0.04	2.39	3.25	6.04	14.87	42.78
0.05	2.50	3.58	7.09	18.18	53.26

Table 1: Bias-adjusted critical values, $t_{\alpha, df-1, R^2, R^2}^\dagger$, for different strengths of the omitted variable W (with $R_{Y \sim W|Z, X}^2 = R_{Z \sim W|X}^2 = R^2$) and various sample sizes; $\alpha = 5\%$.

Tests using these new critical values account both for sampling uncertainty and residual biases with the postulated strength. Note $t_{\alpha, df-1, R^2}^\dagger$ increases the larger the sample size. This behaviour is simply a consequence of the well-known, but often overlooked fact that in large samples any signal will eventually be detected, even if it is spurious. Thus, as the sample size grows, a higher threshold is needed in order to protect inferences against systematic biases.

We note Table 1 picks $R_{Y \sim W|Z, X}^2 = R_{Z \sim W|X}^2 = R^2$ for illustrative purposes only. Researchers can construct bias-adjusted critical values for any arbitrary pair of R^2 values—see, e.g., Supplementary Material for 2×2 tables of $t_{\alpha, df-1, R^2}^\dagger$ where both $R_{Y \sim W|Z, X}^2$ and $R_{Z \sim W|X}^2$ are varied simultaneously.

Remark 1. Sensitivity analysis cannot reveal the strength of confounding present, only the strength of confounding required to alter a research conclusion. For instance, Table 1 reveals that in a study with 1 million observations, one needs a t-value of at least 12 in order to guarantee that the results are robust to latent variables that explain 1% of the residual variation both of the dependent and independent variables. The table also tells us that any study with a t-value less than 12 is vulnerable to such biases. The table *does not* tell us whether latent variables with such strength do exist in any particular study—this needs to be adjudicated using expert knowledge. Note, however, that knowing what one needs to know is useful, and represents an improvement over conventional analysis, which assumes $R^2 = 0$. See Section 6 for additional discussion.

3.3. Compatible inferences given bounds on partial R^2

Given hypothetical values for $R_{Y \sim W|Z, \mathbf{X}}^2$ and $R_{Z \sim W|\mathbf{X}}^2$, the previous results allow us to determine exactly how the inclusion of W with such strength would change inference regarding the parameter of interest. Often, however, the analyst does not know the exact strength of omitted variables, and wishes to investigate the *worst* possible inferences that could be induced by a W with bounded strength, for instance, $R_{Y \sim W|Z, \mathbf{X}}^2 \leq R_{Y \sim W|Z, \mathbf{X}}^{2 \max}$ and $R_{Z \sim W|\mathbf{X}}^2 \leq R_{Z \sim W|\mathbf{X}}^{2 \max}$. Writing $t_{\alpha, \text{df}-1, \mathbf{R}^2}^\dagger$ as a function of the sensitivity parameters $R_{Y \sim W|Z, \mathbf{X}}^2$ and $R_{Z \sim W|\mathbf{X}}^2$, we then solve the maximization problem,

$$\max_{R_{Y \sim W|Z, \mathbf{X}}^2, R_{Z \sim W|\mathbf{X}}^2} t_{\alpha, \text{df}-1, \mathbf{R}^2}^\dagger \quad \text{s.t.} \quad R_{Y \sim W|Z, \mathbf{X}}^2 \leq R_{Y \sim W|Z, \mathbf{X}}^{2 \max}, \quad R_{Z \sim W|\mathbf{X}}^2 \leq R_{Z \sim W|\mathbf{X}}^{2 \max}. \quad (14)$$

Denoting the solution to the optimization problem in expression (14) as $t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max}$, we obtain the maximum bias-adjusted critical value.

THEOREM 3 (MAXIMUM BIAS-ADJUSTED CRITICAL VALUE). Fix α , $R_{Y \sim W|Z, \mathbf{X}}^{2 \max}$ and $R_{Z \sim W|\mathbf{X}}^{2 \max} < 1$ in the optimization problem (14). Then,

$$t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max} = t_{\alpha, \text{df}-1, \mathbf{R}^{2*}}^\dagger,$$

with $\mathbf{R}^{2*} = \{R_{Y \sim W|Z, \mathbf{X}}^{2 \max}, R_{Z \sim W|\mathbf{X}}^{2 \max}\}$ if $R_{Z \sim W|\mathbf{X}}^{2 \max} \geq f_{\alpha, \text{df}-1}^{*2} f_{Y \sim W|Z, \mathbf{X}}^{2 \max}$, and $\mathbf{R}^{2*} = \{R_{Z \sim W|\mathbf{X}}^{2 \max} / (f_{\alpha, \text{df}-1}^{*2} + R_{Z \sim W|\mathbf{X}}^{2 \max}), R_{Z \sim W|\mathbf{X}}^{2 \max}\}$ otherwise, where here we define $f_{\alpha, \text{df}-1}^* := t_{\alpha, \text{df}-1}^* / \sqrt{\text{df}-1}$.

Once in possession of $t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max}$, the most extreme possible lower and upper limits of confidence intervals after adjusting for W are then given by

$$\text{LL}_{1-\alpha, \mathbf{R}^2}^{\max}(\lambda) = \hat{\lambda}_{\text{res}} - t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max} \times \widehat{\text{se}}(\hat{\lambda}_{\text{res}}), \quad \text{UL}_{1-\alpha, \mathbf{R}^2}^{\max} = \hat{\lambda}_{\text{res}} + t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max} \times \widehat{\text{se}}(\hat{\lambda}_{\text{res}}).$$

The interval composed of such limits,

$$\text{CI}_{1-\alpha, \mathbf{R}^2}^{\max}(\lambda) := \left[\text{LL}_{1-\alpha, \mathbf{R}^2}^{\max}(\lambda), \quad \text{UL}_{1-\alpha, \mathbf{R}^2}^{\max}(\lambda) \right], \quad (15)$$

retrieves the union of all confidence intervals for λ that are compatible with an omitted variable with such strengths.

3.4. Sensitivity statistics for routine reporting

Widespread adoption of sensitivity analysis benefits from simple and interpretable statistics that quickly convey the overall robustness of an estimate. To that end, Cinelli and Hazlett (2020) proposed two sensitivity statistics for routine reporting: (i) the partial R^2 of Z with Y , $R_{Y \sim Z|\mathbf{X}}^2$;

and, (ii) the *robustness value* (RV). Here we generalize the notion of a partial R^2 as a measure of robustness to extreme scenarios, by introducing the *extreme robustness value* (XRV), for which the partial R^2 is a special case. We also recast these sensitivity statistics as a solution to an “inverse” question regarding the interval $\text{CI}_{1-\alpha, R^2}^{\max}(\lambda)$. This framework facilitates extending these metrics to other contexts, in particular to the IV setting in Section 4.

The extreme robustness value

Our first inverse question is: what is the *bare minimum* strength of association of the omitted variable W with Z that could bring its estimated coefficient to a region where it is no longer statistically different than zero (or another threshold of interest)? To answer this question, we can see $\text{CI}_{1-\alpha, R^2}^{\max}(\lambda)$ as a function of the bound $R_{Z \sim W | \mathbf{X}}^2$ alone, obtained from maximizing the adjusted critical value in expression (14) where: (i) the parameter $R_{Y \sim W | Z, \mathbf{X}}^2$ is left completely unconstrained (i.e., $R_{Y \sim W | Z, \mathbf{X}}^2 \leq 1$); and, (ii) the parameter $R_{Z \sim W | \mathbf{X}}^2$ is bounded by XRV (i.e., $R_{Z \sim W | \mathbf{X}}^2 \leq \text{XRV}$). The *Extreme Robustness Value* $\text{XRV}_{q^*, \alpha}(\lambda)$ is defined as the greatest lower bound XRV such that the null hypothesis that a change of $(100 \times q^*)\%$ of the original estimate, $H_0 : \lambda = (1 - q^*)\hat{\lambda}_{\text{res}}$, is not rejected at the α level,

$$\text{XRV}_{q^*, \alpha}(\lambda) := \inf \left\{ \text{XRV}; (1 - q^*)\hat{\lambda}_{\text{res}} \in \text{CI}_{1-\alpha, 1, \text{XRV}}^{\max}(\lambda) \right\}. \quad (16)$$

The solution to this problem gives the following result.

THEOREM 4 (EXTREME ROBUSTNESS VALUE—OLS). *Under Assumption 1, for given q^* and α , the extreme robustness value equals*

$$\text{XRV}_{q^*, \alpha}(\lambda) = \begin{cases} 0, & \text{if } f_{q^*}(\lambda) \leq f_{\alpha, \text{df}-1}^*, \\ \frac{f_{q^*}^2(\lambda) - f_{\alpha, \text{df}-1}^{*2}}{1 + f_{q^*}^2(\lambda)}, & \text{otherwise,} \end{cases}$$

where $f_{q^*}(\lambda) := q^*|f_{Y \sim Z | \mathbf{X}}|$, and $f_{\alpha, \text{df}-1}^* := t_{\alpha, \text{df}-1}^* / \sqrt{\text{df}-1}$.

Remark 2. Beyond its procedural interpretation, $\text{XRV}_{q^*, \alpha}(\lambda)$ can also be interpreted as an “adjusted partial R^2 ” of Z with Y . To see why, consider the case of the minimal strength to bring the point estimate ($\alpha = 1$) to exactly zero ($q^* = 1$). We then have that $f_{\alpha=1, \text{df}-1}^* = 0$ and $f_{q^*=1}^2(\lambda) = f_{Y \sim Z | \mathbf{X}}^2$, resulting in $\text{XRV}_{q^*=1, \alpha=1}(\lambda) = \frac{f_{Y \sim Z | \mathbf{X}}^2}{1 + f_{Y \sim Z | \mathbf{X}}^2} = R_{Y \sim Z | \mathbf{X}}^2$. For the general case, we simply perform two adjustments that dampens the “raw” partial R^2 of Z with Y . First we adjust it by the proportion of reduction deemed to be problematic q^* through $f_{q^*} = q^*|f_{Y \sim Z | \mathbf{X}}|$; next, we subtract the threshold for which statistical significance is lost.

The robustness value

An alternative measure of robustness of the OLS estimate is to consider the minimal strength of association that the omitted variable needs to have, *both* with Z and Y , so that a $1 - \alpha$ confidence interval for λ will include a change of $(100 \times q^*)\%$ of the current restricted estimate. Write $\text{CI}_{1-\alpha, R^2}^{\max}(\lambda)$ as a function of both bounds varying simultaneously, $\text{CI}_{1-\alpha, \text{RV}, \text{RV}}^{\max}(\lambda)$, by maximizing the adjusted critical value with bounds given by $R_{Y \sim W | Z, \mathbf{X}}^2 \leq \text{RV}$ and $R_{Z \sim W | \mathbf{X}}^2 \leq \text{RV}$. The *Robustness Value* $\text{RV}_{q^*, \alpha}(\lambda)$ for not rejecting the null hypothesis that $H_0 : \lambda = (1 - q^*)\hat{\lambda}_{\text{res}}$,

at the significance level α , is defined as

$$\text{RV}_{q^*,\alpha}(\lambda) := \inf \left\{ \text{RV}; (1 - q^*)\hat{\lambda}_{\text{res}} \in \text{CI}_{1-\alpha, \text{RV}, \text{RV}}^{\max}(\lambda) \right\}. \quad (17)$$

The RV of OLS estimates has then the following characterization.

360

THEOREM 5 (ROBUSTNESS VALUE—OLS). *Under Assumption 1, for given q^* and α , the robustness value equals*

$$\text{RV}_{q^*,\alpha}(\lambda) = \begin{cases} 0, & \text{if } f_{q^*}(\lambda) \leq f_{\alpha, \text{df}-1}^*, \\ \frac{1}{2} \left(\sqrt{f_{q^*,\alpha}^4(\lambda) + 4f_{q^*,\alpha}^2(\lambda)} - f_{q^*,\alpha}^2(\lambda) \right), & \text{if } f_{\alpha, \text{df}-1}^* < f_{q^*}(\lambda) < f_{\alpha, \text{df}-1}^{*-1}, \\ \text{XRV}_{q^*,\alpha}(\lambda), & \text{otherwise.} \end{cases}$$

where $f_{q^*,\alpha}(\lambda) := q^*|f_{Y \sim Z|X}| - f_{\alpha, \text{df}-1}^*$, and $f_{\alpha, \text{df}-1}^* := t_{\alpha, \text{df}-1}^*/\sqrt{\text{df}-1}$.

The first case occurs when the confidence interval already includes $(1 - q^*)\hat{\lambda}_{\text{res}}$ or the mere change of one degree of freedom achieves this. In the second case, both associations of W reach the bound.⁹ The last case is an interior point solution—when the constraint on the partial R^2 with the outcome is not binding, the RV reduces to the XRV.

365

3.5. Bounding the plausible strength of omitted variables

One final result is required before turning to the sensitivity of IV estimates. Let X_j be a specific covariate of the set \mathbf{X} , and define

370

$$k_Z := \frac{R_{Z \sim W|X_{-j}}^2}{R_{Z \sim X_j|X_{-j}}^2}, \quad k_Y := \frac{R_{Y \sim W|Z, X_{-j}}^2}{R_{Y \sim X_j|Z, X_{-j}}^2}, \quad (18)$$

where X_{-j} represents the vector of covariates \mathbf{X} excluding X_j . These new parameters, k_Z and k_Y , stand for how much “stronger” W is relatively to the observed covariate X_j in terms of residual variation explained of Z and Y . Our goal in this section is to re-express (or bound) the sensitivity parameters $R_{Z \sim W|X}^2$ and $R_{Y \sim W|Z, X}^2$ in terms of the relative strength parameters k_Z and k_Y . Cinelli and Hazlett (2020) derived bounds considering the part of W not linearly explained by \mathbf{X} . These are particularly useful when contemplating X_j and W both *confounders* of Z (violations of the ignorability of the instrument). In the IV setting, however, W and X_j may be *side-effects* of Z , instead of causes of Z . In such cases, it may be more natural to reason about the orthogonality of \mathbf{X} and W after conditioning on Z . Therefore, here we additionally provide bounds under the condition $R_{W \sim X_j|Z, X_{-j}}^2 = 0$.

375

380

THEOREM 6 (RELATIVE BOUNDS ON THE STRENGTH OF W). *Under Assumption 1, for fixed k_Z and k_Y as defined in (18), if $R_{W \sim X_j|Z, X_{-j}}^2 = 0$ then*

$$R_{Z \sim W|X}^2 \leq \eta f_{Z \sim X_j|X_{-j}}^2, \quad R_{Y \sim W|Z, X}^2 = k_Y f_{Y \sim X_j|Z, X_{-j}}^2, \quad (19)$$

385

$$\text{where, } \eta = \left(\frac{\sqrt{k_Z} + |R_{Z \sim X_j|X_{-j}}^3|}{\sqrt{1 - k_Z R_{Z \sim X_j|X_{-j}}^4}} \right).$$

⁹ When the f statistic is very large, it may be numerically convenient to use the equivalent expression $2/(1 + \sqrt{1 + 4/f_{q^*,\alpha}^2(\lambda)})$ which avoids catastrophic cancellations.

These results allow investigators to leverage knowledge of *relative importance* of variables (Kruskal and Majors, 1989) when making plausibility judgments regarding sensitivity parameters, by setting $R_{Y \sim W|Z, X}^{2 \max} = k_Y f_{Y \sim X_j|Z, X_{-j}}^2$, $R_{Z \sim W|X}^{2 \max} = \eta f_{Z \sim X_j|X_{-j}}^2$ in $CI_{1-\alpha, R^2}^{\max}(\lambda)$.

4. AN OVB FRAMEWORK FOR THE SENSITIVITY OF IV

We are now ready to develop a suite of sensitivity analysis tools for instrumental variable regression. In this section, we first show how separate sensitivity analysis of the reduced form and first stage is sufficient to draw many valuable conclusions regarding the sensitivity of IV. We then construct a complete OVB framework for the sensitivity analysis of the IV estimate itself within the Anderson-Rubin approach.

4.1. What can be learned from the sensitivity analysis of the reduced form and first stage?

The critical examination of the first stage and the reduced form plays an important role for supporting the causal story behind a particular instrumental variable (Angrist and Krueger, 2001; Angrist and Pischke, 2009; Imbens, 2014). While investigating these separate regressions, all sensitivity analysis results discussed in the previous section can be readily deployed. Fortunately, such sensitivity analyses also answer many pivotal questions regarding the IV estimate itself. First, if the investigator is interested in assessing the strength of confounders or side-effects needed to bring the IV point estimate to zero, or to not reject the null hypothesis of zero effect, the results of the sensitivity analysis of the reduced form is all that is needed. Second, the sensitivity of the first stage (to confounding that could change its sign) reveals whether the IV estimate could be arbitrarily large in either direction.¹⁰ We now elaborate on these claims.

What the reduced form and first stage reveal about the IV point estimate

The IV estimators under consideration here equal to the ratio of the reduced-form and the first-stage regression coefficients, $\hat{\tau} := \hat{\lambda}/\hat{\theta}$. This simple algebraic fact leads to two immediate and practically important conclusions regarding the sensitivity of $\hat{\tau}$ from the sensitivity of $\hat{\lambda}$ and $\hat{\theta}$ alone. First, residual biases can bring the IV point estimate to zero *if and only if* they can bring the reduced-form point estimate to zero. Therefore, if sensitivity analysis of the reduced form reveals that omitted variables are not strong enough to explain away $\hat{\lambda}$, then they also cannot explain away the IV point estimate $\hat{\tau}$. Or, more worrisome, if analysis reveals that it takes weak confounding or side-effects to explain away $\hat{\lambda}$, the same holds for the IV estimate $\hat{\tau}$. Second, if we cannot rule out confounders or side-effects able to *change the sign* of the first stage, we cannot rule out that the IV point estimate $\hat{\tau}$ could be *arbitrarily large* in either direction. This can be immediately seen by letting $\hat{\theta}$ approach zero on either side of the limit. Thus, whenever we are interested in biases as large *or larger* than a certain amount, the robustness of the first stage to the zero null puts an upper bound on the robustness of the IV point estimate.

What the reduced form and first stage reveal about IV hypothesis tests

The Anderson-Rubin test for the null hypothesis $H_0 : \tau = \tau_0$ is based on the test of $H_0 : \phi_{\tau_0} = 0$. By the FWL theorem, the point estimate and (estimated) standard error for $\hat{\phi}_{\tau_0}$ can be expressed in terms of the first-stage and reduced-form estimates, namely, $\hat{\phi}_{\tau_0} = \hat{\lambda} - \tau_0 \hat{\theta}$ and, $\widehat{\text{se}}(\hat{\phi}_{\tau_0}) = \sqrt{\widehat{\text{var}}(\hat{\lambda}) + \tau_0^2 \widehat{\text{var}}(\hat{\theta}) - 2\tau_0 \widehat{\text{cov}}(\hat{\lambda}, \hat{\theta})}$. Testing $H_0 : \phi_{\tau_0} = 0$ requires comparing the t-value for $\hat{\phi}_{\tau_0}$ with a critical threshold $t_{\alpha, \text{df} - 1}^*$, and the null hypothesis is not rejected if

¹⁰ In the context of randomization inference, similar observations have been noted by Rosenbaum (1996); Imbens and Rosenbaum (2005); Small and Rosenbaum (2008); Keele et al. (2017) and Rosenbaum (2017).

$|t_{\hat{\phi}_{\tau_0}}| \leq t_{\alpha, \text{df}-1}^*$. Squaring and rearranging terms we obtain the quadratic inequality,

$$\underbrace{(\hat{\theta}^2 - \widehat{\text{var}}(\hat{\theta})t_{\alpha, \text{df}-1}^{*2})}_{a} \tau_0^2 + \underbrace{2(\widehat{\text{cov}}(\hat{\lambda}, \hat{\theta})t_{\alpha, \text{df}-1}^{*2} - \hat{\lambda}\hat{\theta})}_{b} \tau_0 + \underbrace{(\hat{\lambda}^2 - \widehat{\text{var}}(\hat{\lambda})t_{\alpha, \text{df}-1}^{*2})}_{c} \leq 0. \quad (20)$$

When considering the null hypothesis $H_0 : \tau_0 = 0$, only the term c remains, and c is less or equal to zero if and only if one cannot reject the null hypothesis $H_0 : \lambda = 0$ in the reduced-form regression. Also note that arbitrarily large values for τ_0 will satisfy the inequality in Equation (20) if, and only if, $a < 0$, meaning that we cannot reject the null hypothesis $H_0 : \theta = 0$ in the first-stage regression. Within the Anderson-Rubin framework, we thus reach analogous conclusions regarding hypothesis testing as those regarding the point estimate: (i) when interest lies in the zero null hypothesis, the sensitivity of the reduced form is exactly the sensitivity of the IV—no other analyses are needed; and, (ii) if one is interested in biases of a certain amount, or larger, then the sensitivity of the first stage to the zero null hypothesis needs also to be assessed.

4.2. A complete set of sensitivity tools for IV

We now build a complete set of sensitivity tools for IV within the Anderson-Rubin approach.

Testing a specific null hypothesis

A sensitivity analysis for the null hypothesis $H_0 : \tau = \tau_0$, for any arbitrary value τ_0 can be performed as follows.

Algorithm 1. Sensitivity analysis for IV given a specific null hypothesis.

- (1) Set $H_0 : \tau = \tau_0$, α , and $\mathbf{R}^2 = \{R_{Z \sim W|X}^2, R_{Y_{\tau_0} \sim W|Z, X}^2\}$;
- (2) Construct $Y_{\tau_0} = Y - \tau_0 D$;
- (3) Fit the OLS model $Y_{\tau_0} = \hat{\phi}_{\text{res}, \tau_0} Z + \mathbf{X} \hat{\beta}_{\text{res}, \tau_0} + \hat{\varepsilon}_{\tau_0, \text{res}}$;
- (4) Compare the t-value for testing $H_0 : \phi_{\text{res}, \tau_0} = 0$ against the critical value $t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max}$;
- (5) Compute $\text{XRV}_{q^*=1, \alpha}(\phi_{\tau_0})$ and $\text{RV}_{q^*=1, \alpha}(\phi_{\tau_0})$;
- (6) Report the results of (4) and (5).

The procedure above tells us how omitted variables no worse than $\mathbf{R}^2 = \{R_{Z \sim W|X}^2, R_{Y_{\tau_0} \sim W|Z, X}^2\}$ would alter inferences regarding the null $H_0 : \tau = \tau_0$, as well as the minimal strength of \mathbf{R}^2 required to not reject the null $H_0 : \tau = \tau_0$, as given by the RV or XRV. Note the bounds on \mathbf{R}^2 can be chosen to reflect the assumption that the omitted variables are no stronger than certain observed covariates, as per Section 3.5.

Compatible inferences for IV given bounds on partial R^2

More broadly, analysts can recover the set of inferences compatible with plausibility judgments on the maximum strength of W . For a critical threshold $t_{\alpha, \text{df}-1}^*$, the confidence interval for τ in the Anderson-Rubin framework is given by $\text{CI}_{1-\alpha}(\tau) = \{\tau_0; t_{\phi_{\tau_0}}^2 \leq t_{\alpha, \text{df}-1}^{*2}\}$. Thus, consider bounds on sensitivity parameters $R_{Y_{\tau_0} \sim W|Z, X}^2 \leq R_{Y_0 \sim W|Z, X}^{\max}$ (which should be judged to hold *regardless* of the value of τ_0) and $R_{Z \sim W|X}^2 \leq R_{Z \sim W|X}^{\max}$. Let $t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max}$ denote the maximum bias-adjusted critical value under the posited bounds on the strength of W . The set of compatible inferences for the IV estimate $\text{CI}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau)$ is then defined as

$$\text{CI}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau) := \left\{ \tau_0; t_{\hat{\phi}_{\text{res}, \tau_0}}^2 \leq \left(t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max} \right)^2 \right\}. \quad (21)$$

This interval can be found analytically using the same inequality as in Equation (20), but now with the parameters of the restricted regression actually run, and $t_{\alpha, \text{df}-1}^*$ replaced by $t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max}$.
 460 Note that users can easily obtain $\text{CI}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau)$ with any software that computes Anderson-Rubin or Fieller's confidence intervals by simply providing the modified critical threshold $t_{\alpha, \text{df}-1, \mathbf{R}^2}^{\dagger \max}$. Armed with the notion of a set of compatible inferences for IV, $\text{CI}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau)$, we are now able to formally define and derive (extreme) robustness values for instrumental variable estimates.

Extreme robustness values for IV

465 The extreme robustness value $\text{XRV}_{q^*, \alpha}(\tau)$ for the IV estimate is defined as the minimum strength of association of omitted variables with the instrument so that we cannot reject a reduction of $(100 \times q^*)\%$ of the original IV estimate; that is,

$$\text{XRV}_{q^*, \alpha}(\tau) := \inf \{ \text{XRV}; (1 - q^*)\hat{\tau}_{\text{res}} \in \text{CI}_{1-\alpha, 1, \text{XRV}}^{\max}(\tau) \}. \quad (22)$$

The $\text{XRV}_{q^*, \alpha}(\tau)$ computes the minimal strength of W required to not reject a particular null hypothesis of interest. However, we might be interested, instead, in asking about the minimal strength of omitted variables to not reject a specific value *or worse*. When confidence intervals are connected, such as the case of standard OLS, the two notions coincide. But in the Anderson-Rubin case, confidence intervals for the IV estimate can sometimes consist of disjoint intervals. Therefore, let the upper and lower limits of $\text{CI}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau)$ be $\text{LL}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau)$ and $\text{UL}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau)$ respectively. The extreme robustness value $\text{XRV}_{\geq q^*, \alpha}(\tau)$ for the IV estimate is defined as the minimum strength of association that confounders or side-effects need to have with the instrument so that we cannot reject a change of $(100 \times q^*)\%$ *or worse* of the original IV estimate,

$$\text{XRV}_{\geq q^*, \alpha}(\tau) := \inf \{ \text{XRV}; (1 - q^*)\hat{\tau}_{\text{res}} \in [\text{LL}_{1-\alpha, 1, \text{XRV}}^{\max}(\tau), \text{UL}_{1-\alpha, 1, \text{XRV}}^{\max}(\tau)] \}. \quad (23)$$

Both quantities can be obtained via the Anderson-Rubin and first-stage regressions as follows.

480 **THEOREM 7 (EXTREME ROBUSTNESS VALUE—IV).** *Under Assumption 1, for given q^* and α , the extreme robustness values for IV are given by*

$$\text{XRV}_{q^*, \alpha}(\tau) = \text{XRV}_{1, \alpha}(\phi_{\tau^*}), \quad \text{and}, \quad (24)$$

$$\text{XRV}_{\geq q^*, \alpha}(\tau) = \min\{\text{XRV}_{1, \alpha}(\phi_{\tau^*}), \text{XRV}_{1, \alpha}(\theta)\}, \quad (25)$$

where $\tau^* = (1 - q^*)\hat{\tau}_{\text{res}}$.

485 **Remark 3.** Theorem 7 corroborates the discussion of Section 4.1. The robustness of IV estimates against biases as large or larger than a certain amount is bounded by the robustness of the first stage assessed at the zero null. Moreover, for the special case of the null hypothesis of zero effect, $H_0 : \tau = 0$, we obtain $\text{XRV}_{\geq 1, \alpha}(\tau) = \min\{\text{XRV}_{1, \alpha}(\lambda), \text{XRV}_{1, \alpha}(\theta)\}$, that is, the XRV of the IV estimate, against biases that bring it to zero or worse, is equal to the minimum of the XRV of the first stage and the reduced form, both evaluated at the zero null ($q^* = 1$).
 490

Remark 4. Note that the XRV of the first stage $\text{XRV}_{1, \alpha}(\theta)$ can be arbitrarily different from traditional metrics of instrument strength. For a simple numerical example, consider $\text{df} = 100,000$ and suppose the first stage F statistic is $F = t^2 \approx 100$, which could be considered a strong instrument for statistical inference purposes. In this case, we still have $\text{XRV}_{1, \alpha}(\theta) \approx 0.001$.

Robustness values for IV

495 The definitions of the robustness value for IV follow the same logic discussed above, but now considering both bounds on $\text{CI}_{1-\alpha, \mathbf{R}^2}^{\max}(\tau)$ varying simultaneously. That is, the RV for not

rejecting a bias of exactly q^* is defined as

$$RV_{q^*,\alpha}(\tau) := \inf \{RV; (1 - q^*)\hat{\tau}_{\text{res}} \in \text{CI}_{1-\alpha, \text{RV}, \text{RV}}^{\max}(\tau)\}, \quad (26)$$

and the RV for not rejecting the null of a reduction of $(100 \times q^*)\%$ or worse is defined as,

$$RV_{\geq q^*,\alpha}(\tau) := \inf \{RV; (1 - q^*)\hat{\tau}_{\text{res}} \in [\text{LL}_{1-\alpha, \text{RV}, \text{RV}}^{\max}(\tau), \text{UL}_{1-\alpha, \text{RV}, \text{RV}}^{\max}(\tau)]\}. \quad (27)$$

We then have analogous results for robustness values, and similar discussion applies.

THEOREM 8 (ROBUSTNESS VALUE—IV). *Under Assumption 1, for given q^* and α , the robustness values for IV are given by*

$$RV_{q^*,\alpha}(\tau) = RV_{1,\alpha}(\phi_{\tau^*}), \quad \text{and}, \quad (28)$$

$$RV_{\geq q^*,\alpha}(\tau) = \min\{RV_{1,\alpha}(\phi_{\tau^*}), \quad RV_{1,\alpha}(\theta)\}, \quad (29)$$

where $\tau^* = (1 - q^*)\hat{\tau}_{\text{res}}$.

Bounds on the strength of omitted variables

When testing a specific null hypothesis $H_0 : \tau = \tau_0$ in the AR regression, we have k_Z as in Section 3.5, and instead of k_Y we now have $k_{Y_{\tau_0}} := R_{Y_{\tau_0} \sim W|Z, \mathbf{X}_{-j}}^2 / R_{Y_{\tau_0} \sim X_j|Z, \mathbf{X}_{-j}}^2$. The plausibility judgment one is making here is thus under the null $H_0 : \tau = \tau_0$. Since the judgment is made under a specific null, the bounds will be different when testing different hypotheses. Therefore, it is useful to compute bounds under a slightly more *conservative* assumption. We can posit that the omitted variables are no stronger than (a multiple of) the *maximum* explanatory power of an observed covariate, regardless of the value of τ_0 , i.e., $k_{Y_{\tau_0}}^{\max} := \frac{\max_{\tau_0} R_{Y_{\tau_0} \sim W|Z, \mathbf{X}_{-j}}^2}{\max_{\tau_0} R_{Y_{\tau_0} \sim X_j|Z, \mathbf{X}_{-j}}^2}$.

This has the useful property of providing a unique bound for any null hypothesis, and can be used to place bounds on the sensitivity contours of the lower and upper limit of the AR confidence intervals, as we show next.

5. USING THE OVB FRAMEWORK FOR THE SENSITIVITY OF IV

In this section we return to our running example and show how these tools can be deployed to assess the robustness of those findings to violations of the IV assumptions. Throughout, we focus the discussion on violations of the ignorability of the instrument due to confounders, as this is the main threat of the study under investigation. Readers should keep in mind, however, that mathematically all analyses performed here can be equally interpreted as assessing violations of the exclusion restriction (or both).

5.1. Minimal sensitivity reporting

Model	Param.	Estimate	LL _{1-α}	UL _{1-α}	t-value	XRV _{$\geq q^*, \alpha$}	RV _{$\geq q^*, \alpha$}
Inst. Variable (AR)	τ	0.132	0.025	0.285	2.33	0.05%	0.67%
First Stage	θ	0.320	0.148	0.492	3.64	0.31%	3.02%
Reduced Form	λ	0.042	0.007	0.078	2.33	0.05%	0.67%

Bound AR (1x SMSA): $R_{Y \sim W|Z, \mathbf{X}}^2 = 2\%$, $R_{W \sim Z|\mathbf{X}}^2 = 0.6\%$, $t_{\alpha, \text{df} - 1, R^2}^{\dagger \max} = 2.55$.

Note: df = 2994, $q^* = 1$, $\alpha = 0.05$.

Table 2: Minimal sensitivity reporting.

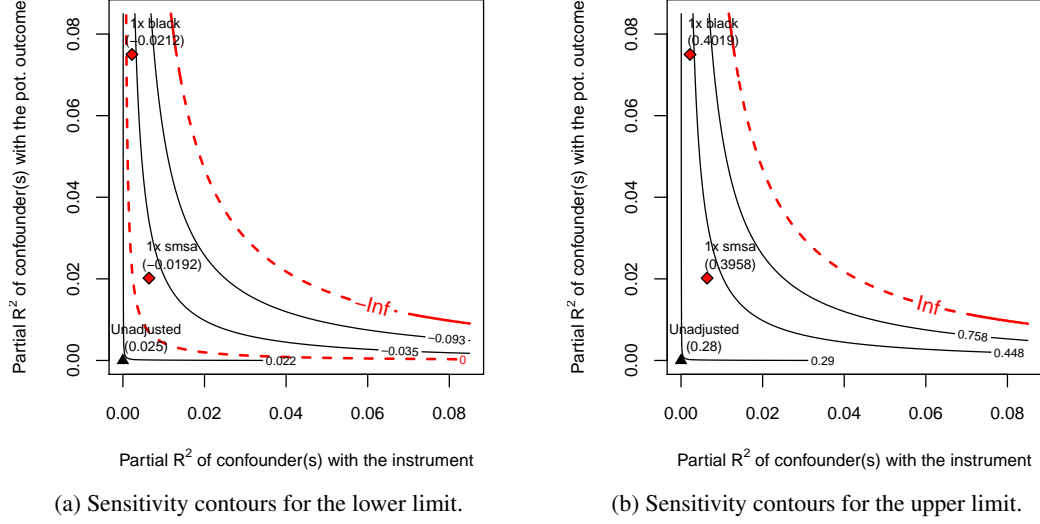


Figure 1: Sensitivity contours of the lower and upper limits of AR confidence interval.

Table 2 shows our proposed minimal sensitivity reporting for IV estimates. It starts by replicating the usual statistics, such as the point estimate (0.132), as well as the lower and upper limits of the Anderson-Rubin confidence interval [0.025, 0.285], and the t-value against the null hypothesis of zero effect (2.33). Next, we propose researchers report the extreme robustness value (XRV_{≥q*,α} = 0.05%) and the robustness value (RV_{≥q*,α} = 0.67%) required to bring the lower limit of the confidence interval to or beyond zero (or another meaningful threshold), at the 5% significance level. We also show these same statistics for the first stage and reduced form. As derived in Theorems 7 and 8, the (extreme) robustness value of the IV estimate required to bring the lower limit of the confidence interval to zero or below is the *minimum* of the (extreme) robustness value of the reduced form and the (extreme) robustness value of the first stage evaluated at the zero null. In our running example, the reduced form is more fragile, thus the sensitivity of the IV hinges critically on the sensitivity of the reduced form.¹¹

The RV reveals that confounders explaining 0.67% of the residual variation both of *proximity* and of (log) *Earnings* are already sufficient to make the IV estimate statistically insignificant. Further, the XRV shows that, if we are not willing to impose constraints on the partial R² of confounders with the outcome, they need only explain 0.05% of the residual variation of the instrument to be problematic. To aid users in making plausibility judgments, the note of the table provides bounds on the maximum strength of unobserved confounding if it were as strong as SMSA (an indicator variable for whether the individual lived in a metropolitan region) along with the bias-adjusted critical value for a confounder with such strength, $t_{\alpha, df-1, R^2}^{\dagger \max} = 2.55$. Since the observed t-value (2.33) is less than the adjusted critical threshold of 2.55, this immediately reveals that confounding as strong as SMSA (e.g. residual geographic confounding) is already sufficiently strong to be problematic.

5.2. Sensitivity contours plots

It will often be valuable to assess the sensitivity of the IV against hypothesis *other than zero*. To that end, investigators may wish to examine sensitivity contour plots showing the whole range

¹¹ See Supplementary Materials for a more detailed analysis of the robustness of the reduced form and first stage.

of adjusted lower and upper limits of the AR confidence interval against various strengths of the omitted variables W . These contours are shown in Figure 1. Here the horizontal axis indicates the bounds on $R^2_{Z \sim W|X}$ and the vertical axis indicates the bounds on $R^2_{Y_{\tau_0} \sim W|Z, X}$. Under a constant treatment effects model, $R^2_{Y_{\tau_0} \sim W|Z, X}$ has a simple interpretation—it stands for how much residual variation confounders explain of the untreated potential outcome. For simplicity, of exposition, we adopt this interpretation here. The contour lines show the worst lower (or upper) limit of the $CI^{\max}_{1-\alpha, R^2}(\tau)$, with omitted variables bounded by such strength. Red dashed lines shows a critical contour line of interest (such as zero) as well as the boundary beyond confidence intervals become unbounded. The red diamonds places bounds on strength of W as strong as *Black* (an indicator for race) and, again, *SMSA*, as per Section 4.2. As the plot reveals, both confounding as strong as *SMSA*, or as strong as *black*, could lead to an interval for the target parameter of $CI^{\max}_{1-\alpha, R^2}(\tau) = [-0.02, 0.40]$, which includes not only implausibly high values (40%), but also negative values (-2%), and is thus too wide for any meaningful conclusions. Since it is not very difficult to imagine residual confounders as strong or stronger than those (e.g., parental income, finer grained geographic location, etc), these results call into question the strength of evidence provided by this IV study.

6. DISCUSSION

Sensitivity analysis tools, such as those introduced in this paper, provide logical deductions aimed at: (i) revealing the consequences of varying degrees of violation of identifying assumptions (e.g., via bias-adjusted critical values), and (ii) determining the minimal degree of violation of those assumptions necessary to overturn certain conclusions (e.g., via robustness values). This shifts the scientific debate from arguing whether, say, latent confounders of an instrumental variable have exactly zero strength—an indefensible claim in most settings—to a more realistic discussion about whether we can confidently rule out strengths that are shown to be problematic.

The results of sensitivity analyses are not always self-evident and can often be surprising. They may reveal that certain studies are highly sensitive to plausible perturbations of identifying assumptions, while others remain robust despite such perturbations. Even when results fall in between these two extremes, sensitivity analyses still represent an improvement over simply assuming away the problem. They clarify what one needs to know, by transparently revealing how vulnerable the results are to violations of the exclusion and independence restrictions. This provides policymakers a better understanding of what remains unknown about an estimated effect, and offers researchers a roadmap for improving their analyses in future inquiries.

It is important to emphasize that plausibility judgments on the maximum strength of latent variables inevitably depend on expert knowledge and can thus vary substantially across scientific disciplines, fields of study, and the quality of the research design. For that reason, we do not propose any universal thresholds for the sensitivity statistics we propose here. For instance, in an observational study without randomization nor a rich set of measured confounders, it would be hard to rule out latent confounders that explain, say, 1% of the residual variation of the instrument. This indeed seems to be the case in our running example (Card, 1993), where residual geographic confounders could plausibly attain such strength. In other scientific contexts, however, a value of 1% may in fact be large. For example, in a Mendelian randomization study where the main concern is pleiotropy, it may be defensible to argue against genetic variants explaining 1% of the variation of a latent complex pleiotropic trait (Cinelli et al., 2022).

Finally, in this paper we focused on the sensitivity of the traditional IV estimate, consisting of the ratio of two OLS regression coefficients. We chose to do so because this reflects current

practices for IV analysis and encompasses the vast majority of applied work. These tools can thus be immediately put to use to improve the robustness of current research, without requiring any additional assumptions, beyond those that already justified the traditional IV analysis in the first place. Recent papers, however, have usefully questioned the causal interpretation of the traditional IV estimand, as it relies on strong parametric assumptions (Słoczyński, 2020; Blandhol et al., 2022). Extending the sensitivity tools we present here to the nonparametric case is possible by leveraging recent results in Chernozhukov et al. (2022), and offers an interesting direction for future work.

ACKNOWLEDGEMENT

Cinelli's research was supported in part by the Royalty Research Fund at the University of Washington, and by the National Science Foundation under Grant No. MMS-2417955.

BIBLIOGRAPHY

- Anderson, T. W. and Rubin, H. (1949). Estimation of the parameters of a single equation in a complete system of stochastic equations. *The Annals of Mathematical Statistics*, 20(1):46–63.
- Andrews, I., Stock, J. H., and Sun, L. (2019). Weak instruments in instrumental variables regression: Theory and practice. *Annual Review of Economics*, 11:727–753.
- Angrist, J. D., Imbens, G. W., and Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *Journal of the American statistical Association*, 91(434):444–455.
- Angrist, J. D. and Krueger, A. B. (2001). Instrumental variables and the search for identification: From supply and demand to natural experiments. *Journal of Economic perspectives*, 15(4):69–85.
- Angrist, J. D. and Pischke, J.-S. (2009). *Mostly harmless econometrics: An empiricist's companion*. Princeton university press.
- Belzil, C. and Hansen, J. (2002). Unobserved ability and the return to schooling. *Econometrica*, 70(5):2075–2091.
- Blandhol, C., Bonney, J., Mogstad, M., and Torgovitsky, A. (2022). When is TSLS actually late? Technical report, National Bureau of Economic Research.
- Blundell, R., Dearden, L., and Sianesi, B. (2001). Estimating the returns to education: Models, methods and results.
- Bound, J., Jaeger, D. A., and Baker, R. M. (1995). Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak. *Journal of the American statistical association*, 90(430):443–450.
- Burgess, S. and Thompson, S. G. (2015). *Mendelian randomization: methods for using genetic variants in causal estimation*. CRC Press.
- Card, D. (1993). Using geographic variation in college proximity to estimate the return to schooling. Technical report, National Bureau of Economic Research.
- Chernozhukov, V., Cinelli, C., Newey, W., Sharma, A., and Syrgkanis, V. (2022). Long story short: Omitted variable bias in causal machine learning. Technical report, National Bureau of Economic Research.
- Cinelli, C. and Hazlett, C. (2020). Making sense of sensitivity: Extending omitted variable bias. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- Cinelli, C., LaPierre, N., Hill, B. L., Sankararaman, S., and Eskin, E. (2022). Robust mendelian randomization in the presence of residual population stratification, batch effects and horizontal pleiotropy. *Nature communications*, 13(1):1–13.
- Conley, T. G., Hansen, C. B., and Rossi, P. E. (2012). Plausibly exogenous. *Review of Economics and Statistics*, 94(1):260–272.
- Deaton, A. S. (2009). Instruments of development: Randomization in the tropics, and the search for the elusive keys to economic development. Technical report, National bureau of economic research.
- DiPrete, T. A. and Gangl, M. (2004). Assessing bias in the estimation of causal effects: Rosenbaum bounds on matching estimators and instrumental variables estimation with imperfect instruments. *Sociological methodology*, 34(1):271–310.
- Felton, C. and Stewart, B. M. (2022). Handle with care: A sociologist's guide to causal inference with instrumental variables.
- Fieller, E. C. (1954). Some problems in interval estimation. *Journal of the Royal Statistical Society: Series B (Methodological)*, 16(2):175–185.
- Frisch, R. and Waugh, F. V. (1933). Partial time regressions as compared with individual trends. *Econometrica: Journal of the Econometric Society*, pages 387–401.

- Gallen, T. (2020). Broken instruments. *Available at SSRN*.
- Gunsilius, F. (2020). Non-testability of instrument validity under continuous treatments. *Biometrika*.
- Heckman, J. J. and Urzua, S. (2010). Comparing IV with structural models: What simple IV can and cannot identify. *Journal of Econometrics*, 156(1):27–37.
- Hernán, M. A. and Robins, J. M. (2006). Instruments for causal inference: an epidemiologist’s dream? *Epidemiology*, pages 360–372. 655
- Hirschberg, J. and Lye, J. (2010). A geometric comparison of the delta and fieller confidence intervals. *The American Statistician*, 64(3):234–241.
- Hirschberg, J. and Lye, J. (2017). Inverting the indirect—the ellipse and the boomerang: Visualizing the confidence intervals of the structural coefficient from two-stage least squares. *Journal of Econometrics*, 199(2):173–183. 660
- Imbens, G. (2014). Instrumental variables: An econometrician’s perspective. Technical report, National Bureau of Economic Research.
- Imbens, G. W. and Rosenbaum, P. R. (2005). Robust, accurate confidence intervals with a weak instrument: quarter of birth and education. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 168(1):109–126.
- Imbens, G. W. and Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press. 665
- Jaeger, D. A. and Parys, J. (2009). On the sensitivity of return to schooling estimates to estimation methods, model specification, and influential outliers if identification is weak.
- Kédagni, D. and Mourifié, I. (2020). Generalized instrumental inequalities: testing the instrumental variable independence assumption. *Biometrika*. 670
- Keele, L., Small, D., and Grieve, R. (2017). Randomization-based instrumental variables methods for binary outcomes with an application to the ‘improve’ trial. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 180(2):569–586.
- Kleibergen, F. (2002). Pivotal statistics for testing structural parameters in instrumental variables regression. *Econometrica*, 70(5):1781–1803. 675
- Kruskal, W. and Majors, R. (1989). Concepts of relative importance in recent scientific literature. *The American Statistician*, 43(1):2–6.
- Lovell, M. C. (1963). Seasonal adjustment of economic time series and multiple regression analysis. *Journal of the American Statistical Association*, 58(304):993–1010.
- Masten, M. A. and Poirier, A. (2021). Salvaging falsified instrumental variable models. *Econometrica*, 89(3):1449–1469. 680
- Mellon, J. (2020). Rain, rain, go away: 137 potential exclusion-restriction violations for studies using weather as an instrumental variable. *Available at SSRN*.
- Moreira, M. J. (2003). A conditional likelihood ratio test for structural models. *Econometrica*, 71(4):1027–1048.
- Moreira, M. J. (2009). Tests with correct size when instruments can be arbitrarily weak. *Journal of Econometrics*, 152(2):131–140. 685
- Nelson, C. and Startz, R. (1990). Some further results on the exact small sample properties of the instrumental variable estimator.
- Pearl, J. (1995). On the testability of causal models with latent and instrumental variables. In *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, pages 435–443. Morgan Kaufmann Publishers Inc. 690
- Pearl, J. (2009). *Causality*. Cambridge university press.
- Rosenbaum, P. R. (1996). Identification of causal effects using instrumental variables: Comment. *Journal of the American Statistical Association*, 91(434):465–468.
- Rosenbaum, P. R. (2017). *Observation and experiment: an introduction to causal inference*. Harvard University Press. 695
- Słoczyński, T. (2020). When should we (not) interpret linear iv estimands as late? *arXiv preprint arXiv:2011.06695*.
- Small, D. S. (2007). Sensitivity analysis for instrumental variables regression with overidentifying restrictions. *Journal of the American Statistical Association*, 102(479):1049–1058.
- Small, D. S. and Rosenbaum, P. R. (2008). War and wages: the strength of instrumental variables and their sensitivity to unobserved biases. *Journal of the American Statistical Association*, 103(483):924–933. 700
- Staiger, D. O. and Stock, J. H. (1994). Instrumental variables regression with weak instruments.
- Stock, J. H. and Yogo, M. (2002). Testing for weak instruments in linear iv regression.
- Swanson, S. A., Hernán, M. A., Miller, M., Robins, J. M., and Richardson, T. S. (2018). Partial identification of the average treatment effect using instrumental variables: review of methods for binary instruments, treatments, and outcomes. *Journal of the American Statistical Association*, 113(522):933–947. 705
- Wang, X., Jiang, Y., Zhang, N. R., and Small, D. S. (2018). Sensitivity analysis and power for instrumental variable studies. *Biometrics*.
- Young, A. (2022). Consistency without inference: Instrumental variables in practical application. *European Economic Review*, page 104112.