

Diagnosing Instrument-Induced Bias: A Test for Control Variable Contamination in 2SLS Models¹

Asad Dossani

Jeffrey Dotson

Rob Schonlau

October 18, 2025

Abstract

In two-stage least squares (2SLS) models, most empirical researchers instrument only the key variable of interest but then include, as though exogenous, an assortment of potentially endogenous control variables. We show analytically and via simulation that if the instrument is correlated with the other controls this can create bias in the key 2SLS estimate, even when the instrument is strong. The bias in the key estimate can be substantial even with low-level correlations and is exacerbated by weak instruments. We develop a new diagnostic test to detect and quantify the size of bias in the key 2SLS coefficient of interest. Applying these tools to the diversification discount literature, we illustrate how the presence of instrument-control-variable-induced bias casts doubt on the validity of a set of published results.

Keywords: Endogeneity, 2SLS, Instrumental Variables, Corporate Finance, Diversification Discount

¹Dossani: Colorado State University, asad.dossani@colostate.edu. Dotson: The Ohio State University, dotson.83@osu.edu. Schonlau: Colorado State University, schonlau@colostate.edu. We are grateful for the comments and suggestions we received at the 2025 American Economic Association conference, the 2024 Western Economic Association International conference, the 2024 Society for Economic Measurement conference, and at seminars at Colorado State University, and Brigham Young University. We also thank Wei Jiang, Bharadwaj Kannan, Kyoo il Kim, Juan Rubio Ramirez, Garrett Schaller, Harry Turtle, Jeff Wooldridge, and Eric Zivot for their comments and suggestions. All errors are our own.

1 Introduction

In the presence of endogeneity it is challenging to identify the causal effect that a key variable of interest has on a specific outcome. A common empirical approach in finance and other related disciplines uses instrumental variables in two-stage least squares (2SLS) models to address endogeneity. Standard practice often leads researchers to include an assortment of control variables in addition to the key variable of interest on the right-hand side of the equation to mitigate the potential for omitted variable bias. The instrument(s) are typically well motivated in the various papers' discussions of the relevancy and exclusion conditions insofar as the instrument(s) relate specifically to the key endogenous variable of interest and the error term. However, in most papers, minimal consideration is given to the possibility that the control variables in the model might also be endogenous, and if also correlated with the instrument, a direct source of bias when estimating the key 2SLS coefficient of interest. The correlation between the instrument for the key variable of interest and the other control variables is observable and can be helpful in thinking about the robustness of the key 2SLS estimate. If the control variables are endogenous, then ignoring this correlation can have a significant effect on the researcher's ability to draw inference from the main 2SLS estimate. This paper calls attention to the widespread prevalence of this issue—which we will refer to as a contaminated controls problem—and suggests both a new diagnostic test and a formula for estimating the amount of potential bias in the main 2SLS coefficient of interest coming from the inclusion of non-instrumented contaminated control variables in the model.

The idea that endogenous control variables create problems for identification is not new. Various theoretical papers and econometric textbooks clearly indicate that there need to be at least as many excluded instruments as there are endogenous variables in order for the parameters in a system of equations to be identified.² But it is clear from a survey of even recent empirical work that

²For examples of several textbooks and papers that discuss this point see chapter 5 of Wooldridge (2002), Chapter 8 of Davidson and MacKinnon (2004), Murray (2006), or Section 3 of Roberts and Whited (2013). In related work, Chen and il Kim (2025) proposes a nonparametric estimator to identify treatment effects in the presence of endogenous controls, both in OLS and IV settings. Our paper differs in the proposed diagnostics and solutions for endogenous

there is ongoing disagreement in practice about how best to operationalize this point with many researchers including multiple control variables in 2SLS models with minimal discussion of their potential endogeneity and others simply dropping the control variables altogether. Indeed, of the hundreds of papers using 2SLS models in recent years that we surveyed in the *Journal of Finance*, *Journal of Financial Economics*, and the *Review of Financial Studies* the vast majority provide minimal or no discussion of the potential endogeneity of the control variables or simply assert that the controls are exogenous without any supporting analysis. Thus, while the ideal is clearly to have at least as many excluded instruments as there are endogenous variables, the challenge in finding even one good instrument for the key variable of interest is apparently leading empirical researchers to compromise with a narrow focus on that single variable and its instrument while ignoring the potential endogeneity of the other variables in the model. This approach of essentially ignoring the fact that the main 2SLS estimate may have significant bias from the inclusion of endogenous control variables may have become accepted, in part, due to the lack of a simple and direct way to estimate the potential size of the bias affecting the main 2SLS coefficient in the model. In this paper we propose a specific set of calculations that address this issue.

The discussion above highlights several questions that empirical researchers using 2SLS confront. For example, if the research focus is on one key variable of interest and there exists both a strong instrument for that specific variable as well as a set of potentially endogenous control variables that might also relate to the outcome of interest, is the researcher better off estimating the overall 2SLS model with or without the other control variables?³ What effect does the inclusion of the other endogenous control variables have on the 2SLS estimate for the key variable of interest given a strong instrument for that one variable that is itself not correlated with the error term? Is it possible to quantify the potential bias in the estimated marginal effect of interest coming from the inclusion of specific control variables? Is there information to be gained by estimating the key

control variables being included in 2SLS systems.

³Note that this question is not about whether to drop the control variables from the first stage alone. Rather this question is about whether to drop the potentially endogenous control variables from the overall system of equations.

2SLS coefficient both with and without the other control variables in the system and then comparing the results? If so, then what does the comparison reveal? Is there a statistical test that reveals whether the other control variables are endogenous and might be affecting the inference around the key variable of interest?⁴ Is there a cost to using multiple instruments for the key variable of interest in an overidentified system if some of the instruments are correlated with the other control variables? And, if more than one strong instrument is available, but different instruments lead to different inferences for the key variable of interest, how should one decide which instrument should be used? Given that literally hundreds of papers at top finance journals have used 2SLS methods in recent years combined with (1) the widespread lack of consideration of the potential endogeneity of the control variables in these papers, (2) the implicit disagreement in practice evidenced by the existence of many papers that either include or exclude the control variables from the analysis, (3) the lack of discussion or even of reporting of whether the instrument(s) on the key variable of interest are correlated with the other control variables, and (4) the common use of multiple instruments in overidentified 2SLS systems, there is obviously a need in the literature for a paper that discusses the exact tradeoffs involved in these decisions and provides clear practical advice for empirical researchers.

Our paper adds to the recent literature⁵ focused on the empirical methods used in finance by addressing these specific questions and makes several broad contributions. First, we draw attention to a common problem affecting inference with 2SLS that has been largely ignored in recent empirical work. Given the prevalence of this problem, with more than 600 papers using 2SLS models published in the *Journal of Finance*, the *Journal of Financial Economics*, and the *Review of Financial Studies* in the last 15 years alone, a discussion of the issues and consequences of the inclusion or exclusion of contaminated controls for inference with 2SLS seems important.⁶ In exploring

⁴Note that this question is not about whether the key variable of interest is endogenous. Rather it is about whether or not a control variable in the system might be creating bias in the 2SLS estimate for the key variable of interest.

⁵A partial list of other recent examples of papers focused on empirical methods in corporate finance include Petersen (2008), Thompson (2011), Erickson and Whited (2012), Roberts and Whited (2013), Gormley and Matsa (2014), Jiang (2017), Bazdresch et al. (2018), Grieser and Hadlock (2019), Berg et al. (2021), Huber (2023).

⁶The number of papers using 2SLS was calculated using textual analysis to identify papers published in the *Journal*

this issue we provide intuition from both analytical expressions for the bias related to endogenous control variables as well as simulation exercises. Second, we propose a new diagnostic test that will allow researchers to directly test whether the contaminated controls problem might exist in their data. Thus, unlike the exclusion condition, which is not directly testable, it is possible to ascertain whether the inclusion of specific endogenous controls might be affecting the key estimate in specific models. As part of this discussion we also provide a new formula for the maximum possible bias (MPB) in the main 2SLS coefficient of interest coming from each endogenous control variable.⁷ The combination of the new statistical test together with the new MPB calculation will not only help researchers better understand how robust their 2SLS inferences are for their main variable of interest but will also direct their attention to which specific control variables need further consideration. To our knowledge, we are the first to propose both the test and the MPB calculations when thinking about inference in a 2SLS system. Third, using simulations we address the question of whether the key 2SLS estimate is better estimated with or without the inclusion of the potentially endogeneous control variables. Although not contended in theory, this question is clearly contended in practice given many recent examples of papers that either include or drop the control variables, or present 2SLS results for models side-by-side with varying numbers of control variables. As part of this discussion we explore practical suggestions for what to do if the control variables are contaminated and show how using multiple instruments for the main variable of interest can lead to bias in the main 2SLS estimate if some of the instruments are also correlated with the other control variables.

of Finance, the *Journal of Financial Economics*, and the *Review of Financial Studies* between 2010 and 2024 that use terminology consistent with 2SLS models (e.g., “2SLS” or “two-stage least” or “IV regression” together with other words such as “instrument” and “stage” or “valid” etc.) Based on this approach approximately 14% of the total papers published in these journals in these years use 2SLS as part of their empirical approach. Of these papers, 10.5% have more than one instrument and 6.5% discuss overidentification tests related to the validity of their instruments. Using textual analysis to identify papers may in a few cases classify papers as using 2SLS when, in fact, their use of these word combinations may refer to another paper’s empirical approach or, alternatively, their use of these word combinations may be part of their explanation for how their empirical approach relates to traditional 2SLS models.

⁷In related work, Cinelli and Hazlett (2025) compute the maximum possible bias coming from omitted variables in an IV regression.

The test for contaminated control bias together with the MPB calculations introduced in this paper provide applied researchers with a new and detailed way to evaluate which among several instruments in an overidentified system provide the least biased 2SLS estimate for the key variable of interest, provide an explanation for why different instruments that are each “strong” may sometimes point to different 2SLS results, and highlight a potential issue with using an overidentified model—as is commonly done in practice if the researcher has more than one instrument—if one or more of the instruments is strongly correlated with the control variables. In addition to the analytical and simulation-based results, we also provide an empirical example of our diagnostic test and MPB calculations using a result from the diversification discount literature. For this example, we use 2SLS models with financial data from Compustat and report results that are similar to the results published in a paper in the *Journal of Finance* suggesting that firms with multiple divisions experience a valuation premium rather than the diversification discount commonly reported in this literature. We then show that this unexpected result can be explained by the contaminated control variables in the model and that using the diagnostic test and MPB calculations proposed in this paper would have clearly identified the issue. We use this example as the basis of a practical discussion for how researchers can explore their 2SLS results if they find possible evidence of contaminated control bias.

The paper is organized as follows. In Section 2 we describe the 2SLS estimator under ideal conditions and provide a detailed description of how the inclusion of endogenous control variables affects these estimates. In Section 3 we propose a test for contaminated controls and derive the relevant distribution for the test statistic. As part of this discussion, we illustrate the magnitude of the problem using simulated data. In Section 4 we suggest a way to calculate the maximum possible bias that can occur in the key 2SLS coefficient of interest due to the observed correlations among the instrument(s) and the other control variables. In Section 5 we run a series of simulations to validate the proposed test, and show the effect of contaminated controls on 2SLS estimates of the key variable of interest. In Section 6 we present an empirical example to illustrate the use of

our test and the MPB calculations. In Section 7 we conclude.

2 OLS and 2SLS Estimates

To facilitate the discussion of how contaminated control variables affect 2SLS estimation, it is helpful first to briefly review the equations involved. In this section, we start with a general regression model and show the form of the bias created when using OLS to estimate marginal effects in the presence of endogeneity. The setting we consider is general and could be motivated using omitted variables, measurement error, or simultaneity (e.g., see discussion in Section 4.1 Wooldridge (2002)). After showing the bias in an OLS setting, we show the form of the bias in a 2SLS setting and demonstrate how the bias is affected by the inclusion of endogenous controls.⁸

2.1 Bias in OLS Estimates

Suppose we are interested in explaining the effect that a particular explanatory variable x_1 has on the outcome of interest y , where x_2 is a control variable. Assume the data generating process for y is a function of x_1 , x_2 , and w as shown in Equation 1 with $E(w|x_1, x_2) \neq 0$.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + w \tag{1}$$

In a multivariate setting, the formula for the OLS estimate of β_1 measures the partial effect that x_1 has on y after netting out x_2 . As noted in the literature (e.g., see Davidson and MacKinnon (2004), Greene (2003)), the formula for the OLS estimate of β_1 can be written using a “double residual” expression by application of the Frisch-Waugh-Lovell (FWL) theorem.⁹ This is shown

⁸Throughout this paper, when we use the term bias, we are referring to the asymptotic or large sample bias computed using the probability limit of the estimator. We note that 2SLS estimators are known to be biased in finite samples but can be consistent in large samples (e.g., see Angrist and Krueger (2001) and chapter 5 of Wooldridge (2002)). Hence the focus in the literature on the large sample properties of 2SLS estimators.

⁹For a detailed discussion of the multivariate OLS coefficient formula and the Frisch-Waugh-Lovell theorem see Wooldridge (2003) pages 78-79, Davidson and MacKinnon (2004) Section 2.4, Greene (2003) page 27, and Lovell

in Equation 2 with x_1^* and y^* representing the residuals from the regressions of x_1 and y on the other explanatory variables from the model, respectively. x_1^* is the portion of x_1 uncorrelated with the other control variables and y^* is the portion of y uncorrelated with the other control variables (not including the key variable of interest x_1). Substituting the full model for y , from Equation 1, into the $\hat{\beta}_1$ formula highlights the factors that affect the bias in the OLS estimate as shown in the equations below.¹⁰ Angrist and Pischke (2009) describe this version of the formula as the “regression anatomy formula” and refer to it as an “important formula [that]...describe[s] the anatomy of a multivariate regression coefficient because it reveals much more than the matrix formula...” We will use this “double residual regression” notation for the coefficient formulas throughout the paper because this approach lends itself to intuitive analytical expressions for exactly how the inclusion of endogenous control variables in the model affects the bias in the key coefficient of interest.

$$\begin{aligned}
\text{plim } \hat{\beta}_{1,OLS} &= \frac{\text{cov}(x_1^*, y^*)}{\text{var}(x_1^*)} = \frac{\text{cov}(x_1^*, y)}{\text{var}(x_1^*)} \\
&= \frac{\beta_1 \text{cov}(x_1^*, x_1) + \beta_2 \text{cov}(x_1^*, x_2) + \text{cov}(x_1^*, w)}{\text{var}(x_1^*)} \\
&= \beta_1 + \underbrace{\frac{\text{cov}(x_1^*, w)}{\text{var}(x_1^*)}}_{\text{bias}} \tag{2}
\end{aligned}$$

The bias in the OLS estimate is a function of the variance of x_1^* and the covariance of the error term w with the portion of the key variable of interest that is uncorrelated with the other controls. The bias can be either positive or negative depending on the sign of the covariance between the variable of interest and the error term, and is increasing in the magnitude of the covariance.

(1963). Using asterisks to identify the residual is similar to the notation used by Greene but written without the matrix notation. See Greene page 27 for a matrix version of this formula. This approach is sometimes called “the double residual regression” (e.g., see Section 17.3 in Goldberger (1991) for example).

¹⁰In simplifying the $\hat{\beta}_1$ expression, by construction $\text{cov}(x_1^*, x_2) = 0$ and $\text{cov}(y^*, x_2) = 0$. This implies that $\text{cov}(x_1^*, y^*) = \text{cov}(x_1^*, y)$ and $\text{cov}(x_1^*, x_1) = \text{var}(x_1^*)$. See Filoso (2013) for additional discussion.

2.2 Bias in 2SLS Estimates With One Control and One Instrument

Empirical researchers often rely on instruments in a 2SLS framework to address endogeneity. In this section we focus on the analytical expression for the bias in a 2SLS model that has one control variable (x_2) in addition to the key variable of interest (x_1) and a single instrument (z) for the endogenous variable of interest. The underlying data generating process (DGP) in this case would be consistent with Equation 1. The first stage in the 2SLS model is a regression of the endogenous variable of interest x_1 on the instrument z and control variable x_2 . From this, we compute the fitted values \hat{x}_1 .

$$\begin{aligned}x_1 &= \gamma_0 + \gamma_1 z + \gamma_2 x_2 + e \\ \hat{x}_1 &= \hat{\gamma}_0 + \hat{\gamma}_1 z + \hat{\gamma}_2 x_2\end{aligned}\tag{3}$$

The second stage in the 2SLS model is a regression of y on \hat{x}_1 and x_2 .

$$y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + v\tag{4}$$

The theoretical literature is clear on the conditions required for the 2SLS estimate of β_1 to be consistent. These conditions are discussed in econometric textbooks (e.g., see Wooldridge (2002), Angrist and Pischke (2009)) as well as in various well-known papers (e.g., see Bound et al. (1995), Angrist and Krueger (2001), Murray (2006), and Roberts and Whited (2013)) and typically focus on the relevancy and exclusion conditions. The relevancy condition requires that the instrument (strongly) correlate with the endogenous variable after controlling for the effects of the other variables, i.e. $\hat{\gamma}_1 \neq 0$ in the first stage equation. The exclusion condition requires that the first stage regressors (the instrument and control variable) not be correlated with the error term, which in this example would mean that $\text{cov}(z, w) = 0$ and $\text{cov}(x_2, w) = 0$.¹¹ The exclusion

¹¹These conditions mean that z and x_2 are also not correlated with the second stage error term v .

condition is not directly testable and hence is motivated based on logic and theory.

Most of the empirical papers we surveyed in top finance journals tend to discuss the exclusion condition solely in terms of whether the instrument for the key variable of interest is correlated with the error term, and not whether the other control variables may also be correlated with the error. We will refer to this narrow focus as the “narrow exclusion restriction” (i.e. $\text{cov}(z, w) = 0$) to distinguish it from the complete set of exclusion conditions noted in the econometric textbooks which in our example would also require that $\text{cov}(x_2, w) = 0$. Distinguishing the narrow exclusion condition from the full set of exclusion conditions seems important given that standard practice in recent applied papers tends to discuss the endogeneity and instrument(s) only in terms of the key variable while the rest of the variables are carried along in the analysis simply as “controls” without careful consideration of either their potential endogeneity or their (observable) correlations with the instrument(s) applied to the main endogenous variable of interest. Indeed, of the hundreds of papers that use instrumental variables with 2SLS in the *Journal of Finance*, the *Journal of Financial Economics*, and the *Review of Financial Studies* in recent years, a large majority of them include various control variables with minimal or no discussion of the potential endogeneity of the control variables, and no discussion of whether the instruments are correlated with the other controls.

The potential bias in the 2SLS estimate of β_1 is of similar form as the OLS estimate in Equation 2 but with both \hat{x}_1 and \hat{x}_1^* used in place of x_1 and x_1^* . Consistent with the asterisks notation used above, \hat{x}_1^* represents the portion of \hat{x}_1 uncorrelated with the other controls.

$$\begin{aligned}
 \text{plim } \hat{\beta}_{1,2SLS} &= \frac{\text{cov}(\hat{x}_1^*, y)}{\text{var}(\hat{x}_1^*)} \\
 &= \frac{\beta_1 \text{cov}(\hat{x}_1^*, x_1) + \beta_2 \text{cov}(\hat{x}_1^*, x_2) + \text{cov}(\hat{x}_1^*, w)}{\text{var}(\hat{x}_1^*)} \\
 &= \beta_1 + \underbrace{\frac{\text{cov}(\hat{x}_1^*, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias}} \tag{5}
 \end{aligned}$$

Given the widespread inclusion of other control variables in 2SLS models in the literature without corresponding discussion of the control variables' potential correlation with the error term, a common implicit assumption in the literature must be that if an instrument z for the key variable of interest x_1 satisfies the narrow exclusion condition, i.e. if $\text{cov}(z, w) = 0$, then $\text{cov}(\hat{x}_1^*, w) = 0$. But this is often not true. Indeed, in the discussion below we show that even if the narrow exclusion condition is satisfied with $\text{cov}(z, w) = 0$, and the instrument for the key variable of interest is strong, the bias in the main 2SLS coefficient of interest can be large if the instrument(s) are also correlated with the other control variables.¹² In the above expressions, \hat{x}_1^* are the residuals from the regression of \hat{x}_1 on the control variables (i.e., x_2 in this example), and hence are orthogonal to whatever control variables are included in the model. The relation between \hat{x}_1 and \hat{x}_1^* is shown below for a regression model with a single control variable x_2 .

$$\begin{aligned}\hat{x}_1 &= \lambda_1 + \lambda_2 x_2 + \xi \\ &= \hat{\lambda}_1 + \hat{\lambda}_2 x_2 + \hat{x}_1^*\end{aligned}\tag{6}$$

To facilitate understanding for how the 2SLS bias in Equation 5 is directly affected by endogenous controls we rewrite \hat{x}_1^* as a function of the control variable x_2 and instrument z . To do this we set the two expressions for \hat{x}_1 from Equations 3 and 6 equal and solve for \hat{x}_1^* .

$$\begin{aligned}\hat{\gamma}_0 + \hat{\gamma}_1 z + \hat{\gamma}_2 x_2 &= \hat{\lambda}_1 + \hat{\lambda}_2 x_2 + \hat{x}_1^* \\ \hat{x}_1^* &= (\hat{\gamma}_0 - \hat{\lambda}_1) + \hat{\gamma}_1 z + (\hat{\gamma}_2 - \hat{\lambda}_2) x_2\end{aligned}\tag{7}$$

We now substitute Equation 7 into Equation 5 to show how the 2SLS bias is affected by en-

¹²Jiang (2017) questioned whether instrumental variable methods have helped empirical researchers identify true marginal effects in recent years given that in about 80% of the studies the instrumented estimates were larger, and often significantly larger, than the noninstrumented estimates. The contaminated control bias and the MPB calculations discussed in our paper could help explain that observation and will help future researchers be able to probe the robustness of their 2SLS results.

ogenous controls – even in the case that the $\text{cov}(z, w) = 0$.

$$\begin{aligned}
\text{plim } \hat{\beta}_{1,2SLS} &= \beta_1 + \frac{\text{cov}(\hat{x}_1^*, w)}{\text{var}(\hat{x}_1^*)} \\
&= \beta_1 + \underbrace{\hat{\gamma}_1 \frac{\text{cov}(z, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias related to the narrow exclusion condition}} + \underbrace{(\hat{\gamma}_2 - \hat{\lambda}_2) \frac{\text{cov}(x_2, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias related to endogenous controls}}
\end{aligned} \tag{8}$$

The bias in the 2SLS estimate is a function of several factors: Focusing on the narrow exclusion condition related term, the bias is increasing in the magnitude of the covariance of the instrument z with the error term w . Focusing on the relevancy condition, the size of the denominator in the bias expression is increasing in the $\text{var}(\hat{x}_1^*)$ which is increasing in the strength of the instrument and specifically increasing in the magnitude of $\hat{\gamma}_1$.¹³ Equation 8 shows mechanically how both the exclusion and relevancy conditions affect the bias in the key coefficient of interest with weak instruments causing the denominators in the second and last terms to be close to 0, and exclusion condition violations causing the numerators in the second and last terms to be far different from zero. The last term in Equation 8 shows how the bias in the main 2SLS coefficient of interest is directly related to the covariance of the control variables with the error and hence highlights the cost of including endogenous control variables in a 2SLS model. Like the bias that comes from violations of the narrow exclusion condition, the bias in the key coefficient of interest that comes from the inclusion of endogenous control variables is also exacerbated by weak instruments for the key variable of interest.

There are two situations where the bias in $\hat{\beta}_{1,2SLS}$ from the control variables will be zero. The first situation occurs if the control variables are exogenous and hence $\text{cov}(x_2, w) = 0$. This outcome is not testable for the same reason that the narrow exclusion condition is not testable: w is not observed. In contrast, the second situation is empirically testable and occurs when $\text{cov}(z, x_2) = 0$. When the instrument is not correlated with the control variable, $\hat{\gamma}_2 = \hat{\lambda}_2$ in Equation 8, causing the

¹³Proof of this result is in Appendix A.

last term in the bias expression to be zero.¹⁴ Researchers can check whether their key estimate is possibly affected by endogenous control variable bias by checking whether $(\hat{\gamma}_2 - \hat{\lambda}_2)$ is close to 0. The new diagnostic test we propose in this paper builds on this intuition. If this difference is close to zero then the bias from the contaminated controls is small. In Section 3 we discuss the details of how to use this difference as a diagnostic test for contaminated control bias.

2.3 Bias in 2SLS Estimates With Multiple Controls and Instruments

We now generalize the results to the case of multiple control variables and instruments. In the case of multiple instruments, we assume again that x_1 is the key variable of interest and is the only variable being instrumented in a first stage equation. Suppose we have a vector of J controls $\mathbf{x}_2 \equiv (x_{21}, \dots, x_{2J})'$ and K instruments $\mathbf{z} \equiv (z_1, \dots, z_K)'$. Let $\boldsymbol{\beta}_2 \equiv (\beta_{21}, \dots, \beta_{2J})'$. Generalizing Equation 1, the data generating process is given by:

$$y = \beta_0 + \beta_1 x_1 + \boldsymbol{\beta}_2' \mathbf{x}_2 + w \quad (9)$$

Let $\boldsymbol{\lambda}_2 \equiv (\lambda_{21}, \dots, \lambda_{2J})'$, $\boldsymbol{\gamma}_2 \equiv (\gamma_{21}, \dots, \gamma_{2J})'$, and $\boldsymbol{\gamma}_1 \equiv (\gamma_{11}, \dots, \gamma_{1K})'$. Generalizing Equations 3 and 4, the first and second stage estimates are given by:

$$\begin{aligned} x_1 &= \gamma_0 + \boldsymbol{\gamma}_1' \mathbf{z} + \boldsymbol{\gamma}_2' \mathbf{x}_2 + e \\ \hat{x}_1 &= \hat{\gamma}_0 + \hat{\boldsymbol{\gamma}}_1' \mathbf{z} + \hat{\boldsymbol{\gamma}}_2' \mathbf{x}_2 \\ y &= \beta_0 + \beta_1 \hat{x}_1 + \boldsymbol{\beta}_2' \mathbf{x}_2 + v \end{aligned} \quad (10)$$

Generalizing Equations 6 and 7, we solve for \hat{x}_1^* . Note that \hat{x}_1^* now partials out the effects of all

¹⁴Proof of this result is in Appendix A.

control variables \mathbf{x}_2 .

$$\begin{aligned}
\hat{x}_1 &= \lambda_1 + \lambda_2' \mathbf{x}_2 + \xi \\
&= \hat{\lambda}_1 + \hat{\lambda}_2' \mathbf{x}_2 + \hat{x}_1^* \\
\hat{\gamma}_0 + \hat{\gamma}_1' \mathbf{z} + \hat{\gamma}_2' \mathbf{x}_2 &= \hat{\lambda}_1 + \hat{\lambda}_2' \mathbf{x}_2 + \hat{x}_1^* \\
\hat{x}_1^* &= (\hat{\gamma}_0 - \hat{\lambda}_1) + \hat{\gamma}_1' \mathbf{z} + (\hat{\gamma}_2 - \hat{\lambda}_2)' \mathbf{x}_2
\end{aligned} \tag{11}$$

Generalizing Equation 8, the expression for the bias is given by:

$$\begin{aligned}
\hat{\beta}_{1,2SLS} &= \beta_1 + \underbrace{\frac{\text{cov}(\hat{x}_1^*, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias}} \\
&= \beta_1 + \underbrace{\hat{\gamma}_1' \frac{\text{cov}(\mathbf{z}, w)}{\text{var}(\hat{x}_1^*)}}_{\substack{\text{bias related} \\ \text{to the narrow} \\ \text{exclusion condition}}} + \underbrace{(\hat{\gamma}_2 - \hat{\lambda}_2)' \frac{\text{cov}(\mathbf{x}_2, w)}{\text{var}(\hat{x}_1^*)}}_{\substack{\text{bias related to} \\ \text{endogenous controls}}}
\end{aligned} \tag{12}$$

Thus the overall bias in the 2SLS β_1 estimate coming from endogenous controls can come from as many channels as there are control variables, with some channels potentially increasing and others potentially decreasing the overall bias. Being able to test whether bias in the key variable of interest might be coming from each of the control variables could be useful in understanding the model. Alternatively it may be useful for applied researchers to perform a single test of the net effect of all the control variables together. Either can be accomplished by testing whether the difference $(\hat{\gamma}_2 - \hat{\lambda}_2)$ is close to zero using a Wald test, with varying restrictions depending on the set of control variables to be tested. We note that this test is for a necessary condition for bias from contaminated controls and not for a sufficient condition; i.e, showing the difference is statistically different from zero signals that there may be contaminated control bias in the 2SLS estimate of interest whereas showing that the difference is not statistically different than zero indicates that there is negligible bias from the control variables even if they are also endogenous. Later in the

paper we provide an analytical expression for the maximum possible size of this bias.

2.4 Bias in 2SLS Estimates With No Control Variables

The implication from the above discussion is that even if an instrument is strongly correlated with the key endogenous variable of interest and even if the instrument itself is not directly correlated with the error term, the inclusion of other endogenous control variables in the system can cause the 2SLS estimate for the main variable of interest to be biased if the instrument is correlated with the endogenous controls. Given the widespread inclusion of potentially endogenous control variables in 2SLS specifications even in recent applied work in top finance and economics journals and the lack of focus on whether the instrument is correlated with the other controls, this point has not been fully appreciated in the empirical literature.

One natural reaction to the prior discussion might be to drop the potentially endogenous controls from the model. But this can lead to other problems. The tradeoff is that dropping the controls can lead to omitted variable bias but including them leads to contaminated control bias that is exacerbated by weak instruments. This issue is contended in practice with some researchers actively advocating the inclusion of as many controls as possible whereas others implicitly disagreeing with this logic by showing their results without controls. So the question we consider in this section is if the control variables are possibly endogenous, is it better to drop the controls from the 2SLS system?

To explore this issue we begin with the case of a single control variable and derive the analytical expression for the bias when the control variable is dropped. The bias in this case would be driven by the correlation between the fitted values (\hat{x}_1 – now estimated without controls) used in the second stage model and the error term which now includes the effects of the omitted controls. Because the second stage model is now also estimated without control variables, the $\hat{\beta}_1$ expression from Equation 8 would include \hat{x}_1 rather than \hat{x}_1^* . Unlike \hat{x}_1^* , which is orthogonal to x_2 , \hat{x}_1 can be correlated with x_2 , which is now part of the second stage error term. The expression for the bias in

this case would be given by:

$$\begin{aligned}
\text{plim } \hat{\beta}_{1,2SLS \text{ without } x_2} &= \frac{\text{cov}(\hat{x}_1, y)}{\text{var}(\hat{x}_1)} \\
&= \frac{\beta_1 \text{cov}(\hat{x}_1, x_1) + \beta_2 \text{cov}(\hat{x}_1, x_2) + \text{cov}(\hat{x}_1, w)}{\text{var}(\hat{x}_1)} \\
&= \beta_1 + \underbrace{\beta_2 \frac{\text{cov}(\hat{x}_1, x_2)}{\text{var}(\hat{x}_1)} + \frac{\text{cov}(\hat{x}_1, w)}{\text{var}(\hat{x}_1)}}_{\text{bias}} \tag{13}
\end{aligned}$$

If the narrow exclusion restriction is satisfied, the third term on the right hand side in Equation 13 is zero. The second term is zero only if x_2 is uncorrelated with the instrument z in which case \hat{x}_1 is uncorrelated with x_2 . As shown in Equation 14, comparing the bias expressions in Equations 8 and 13 shows that dropping the controls from the 2SLS system of equations does not guarantee in any way that the 2SLS estimate will be less biased without the endogenous controls than it is with the endogenous controls in the model. Indeed, without knowing the signs or sizes of β_2 , $\text{cov}(\hat{x}_1, w)$, $\text{cov}(z, w)$, and $\text{cov}(x_2, w)$, all of which are unobservable, it is impossible to know whether dropping the endogenous control variable(s) results in an increase or decrease in the overall bias.

$$\underbrace{\hat{\gamma}_1 \frac{\text{cov}(z, x_m)}{\text{var}(\hat{x}_1^*)} + (\hat{\gamma}_2 - \hat{\lambda}_2) \frac{\text{cov}(x_2, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias including control variable}} \text{ versus } \underbrace{\beta_2 \frac{\text{cov}(\hat{x}_1, x_2)}{\text{var}(\hat{x}_1)} + \frac{\text{cov}(\hat{x}_1, w)}{\text{var}(\hat{x}_1)}}_{\text{bias not including control variable}} \tag{14}$$

It is worth noting that if the narrow exclusion condition holds (or is almost satisfied) for the instrument on the key variable of interest, the instrument is strong, and $(\hat{\gamma}_2 - \hat{\lambda}_2)$ is close to zero then the overall bias is likely smaller in the 2SLS estimate with controls than in the estimate without controls. It is also worth noting that one cannot conclude that the 2SLS estimate for β_1 estimated with controls is biased based simply on whether the 2SLS estimate changes after dropping the controls from the system because the resulting change could be entirely attributable to omitted variable bias associated with the dropped variable(s) which were accounted for when the controls were included as part of the model but are not accounted for when estimating the model without

controls. In the case of multiple control variables and instruments, Equations 13 and 14 generalize to:

$$\begin{aligned}
\text{plim } \hat{\beta}_{1,2SLS \text{ without } x_2} &= \beta_1 + \underbrace{\beta_2' \frac{\text{cov}(\hat{x}_1, x_2)}{\text{var}(\hat{x}_1)} + \frac{\text{cov}(\hat{x}_1, w)}{\text{var}(\hat{x}_1)}}_{\text{bias}} \\
\underbrace{\hat{\gamma}_1' \frac{\text{cov}(z, w)}{\text{var}(\hat{x}_1^*)} + (\hat{\gamma}_2 - \hat{\lambda}_2)' \frac{\text{cov}(x_2, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias including control variables}} &\text{ versus } \underbrace{\beta_2' \frac{\text{cov}(\hat{x}_1, x_2)}{\text{var}(\hat{x}_1)} + \frac{\text{cov}(\hat{x}_1, w)}{\text{var}(\hat{x}_1)}}_{\text{bias not including control variables}} \quad (15)
\end{aligned}$$

It is possible for individual control variables to have opposite effects on the bias for the main 2SLS estimate. As in the single control variable case, if the narrow exclusion condition holds for the instrument(s), and $(\hat{\gamma}_2 - \hat{\lambda}_2)$ is close to zero, the bias is likely smaller in the 2SLS estimate with controls than in the estimate without controls.

3 Testing for Contaminated Controls

In this section, we propose a test for contaminated controls and derive the relevant test statistic. Our test is motivated by the analytical expressions for the bias derived in the previous section. The test is related to the concept of coefficient stability, i.e. the effect of the inclusion of instruments on first stage control variable coefficients in the first stage regression. In related work, Altonji et al. (2005) and Oster (2019) propose methods to estimate omitted variable bias based on coefficient movements after the inclusion of control variables in OLS regressions. Their method is based on the idea that if the estimate on the key variable of interest is not sensitive to the inclusion of control variables, the effect of omitted variable bias is likely to be limited. Our test differs, in part, in that it examines coefficient stability in the first stage regression of a 2SLS system. Our method applies to any 2SLS or IV model, as opposed to an OLS regression model. Specifically, coefficient stability in our test implies that contaminated control bias in the key 2SLS coefficient of interest is likely to be limited.

We begin with the case of one control variable and one instrument. We then proceed to the general case of multiple control variables. The test statistic focuses on the quantity $(\hat{\gamma}_2 - \hat{\lambda}_2)$ from Equation 8, or $(\hat{\gamma}_2 - \hat{\lambda}_2)$ from Equation 12 if there are multiple control variables, and tests whether this term is significantly different from zero. Under the null hypothesis, the expression is equal to zero and there is no contaminated control bias in the estimate for the key coefficient of interest. If the null hypothesis is rejected, the term is different from zero, suggesting contamination, and the possibility of bias coming from endogenous control variables.

3.1 One Control and One Instrument

We need the distribution of $(\hat{\gamma}_2 - \hat{\lambda}_2)$ to be able to determine whether the difference is statistically different from zero. To find the distribution we estimate the coefficients jointly using a modified version of the technique of seemingly unrelated regressions (SUR) of Zellner (1962). The SUR setup consists of a set of independent regression equations with correlated error terms. While the equations can be estimated independently using OLS, the original SUR method proposes estimating the regression equations jointly using feasible GLS to get more efficient parameter estimates. For our purposes, both $\hat{\gamma}_2$ and $\hat{\lambda}_2$ can be viewed as OLS parameter estimates of two regression equations with correlated error terms. The correlation structure between the error terms can be derived analytically, allowing us to obtain the joint distribution of $\hat{\gamma}_2$ and $\hat{\lambda}_2$. Thus, we employ the SUR setup, but estimate the regressions individually by OLS.

First, in the case of one instrument and one control variable, we note that the OLS estimate $\hat{\lambda}_2$ can be computed by regressing x_1 on x_2 and a constant, rather than by regressing \hat{x}_1 on x_2 and a constant. Let \hat{e} denote the residual from the first stage regression, as per Equation 3, so that $x_1 = \hat{x}_1 + \hat{e}$. The two approaches are numerically identical since \hat{e} is orthogonal to x_2 . To derive the correct distribution of $\hat{\lambda}_2$, we use x_1 rather than \hat{x}_1 as the dependent variable. Using \hat{x}_1 would

remove variation from the residuals and would result in standard errors that are too small.¹⁵

$$\text{plim } \hat{\lambda}_2 = \frac{\text{cov}(x_1, x_2)}{\text{var}(x_2)} = \frac{\text{cov}(\hat{x}_1 + \hat{e}, x_2)}{\text{var}(x_2)} = \frac{\text{cov}(\hat{x}_1, x_2)}{\text{var}(x_2)} \quad (16)$$

The SUR model stacks the observations of two regressions, and is set up using our example as follows:

$$\begin{bmatrix} x_1 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 & z & x_2 & 0 & 0 \\ 0 & 0 & 0 & 1 & x_2 \end{bmatrix} \begin{bmatrix} \gamma_0 \\ \gamma_1 \\ \gamma_2 \\ \lambda_1 \\ \lambda_2 \end{bmatrix} + \begin{bmatrix} e \\ \varepsilon \end{bmatrix} \quad (17)$$

Suppose N is the sample size. Let \mathbf{X} denote the matrix consisting of the observed data as per the model above, where the first N rows of \mathbf{X} correspond to the first regression, and the second N rows to the second regression. The covariance matrix of the OLS estimates Σ is given by:

$$\Sigma = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\Omega\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \quad (18)$$

Assuming that the errors in each regression are homoskedastic, letting I_N denote the identity matrix of order N , and noting that $\text{cov}(e, \varepsilon) = \text{var}(e)$:¹⁶

$$\Omega = \begin{bmatrix} \text{var}(e) I_N & \text{cov}(e, \varepsilon) I_N \\ \text{cov}(e, \varepsilon) I_N & \text{var}(\varepsilon) I_N \end{bmatrix} = \begin{bmatrix} \text{var}(e) I_N & \text{var}(e) I_N \\ \text{var}(e) I_N & \text{var}(\varepsilon) I_N \end{bmatrix} \quad (19)$$

Let $\hat{\Sigma}$ denote the estimated covariance matrix, computed as the sample analog of Σ . Using a

¹⁵Though the setting is different, the logic is analogous to the FWL decomposition used in Equation 2, whereby $\hat{\beta}_{1,OLS} = \frac{\text{cov}(x_1^*, y^*)}{\text{var}(x_1^*)} = \frac{\text{cov}(x_1^*, y)}{\text{var}(x_1^*)}$. Either of the two expressions recovers the OLS coefficient, but using y instead of y^* results in different residuals. See Chapter 2 in Davidson and MacKinnon (2004) for further discussion on the residuals in the context of the FWL theorem.

¹⁶ $\text{cov}(e, \varepsilon) = \text{cov}(e, x_1 - \lambda_1 - \lambda_2 x_2) = \text{cov}(e, x_1) = \text{cov}(e, \gamma_0 + \gamma_1 z_1 + \gamma_2 x_2 + e) = \text{var}(e)$

Wald test, under the null hypothesis that $\gamma_2 = \lambda_2$, the test statistic follows a chi-square distribution with one degree of freedom.¹⁷

$$\begin{aligned}
(\mathbf{R}\hat{\boldsymbol{\theta}})'(\mathbf{R}\hat{\boldsymbol{\Sigma}}\mathbf{R}')^{-1}(\mathbf{R}\hat{\boldsymbol{\theta}}) &\sim \chi^2(1) \\
\hat{\boldsymbol{\theta}} &\equiv (\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_2, \hat{\lambda}_1, \hat{\lambda}_2)' \\
\mathbf{R} &\equiv \begin{bmatrix} 0 & 0 & 1 & 0 & -1 \end{bmatrix}
\end{aligned} \tag{20}$$

3.2 Multiple Controls and Instruments

We now derive and present the test statistic for multiple controls and instruments. Generalizing Equation 17, the two stacked regressions are as follows:

$$\begin{bmatrix} x_1 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 & z' & x_2' & 0 & 0 \\ 0 & 0 & 0 & 1 & x_2' \end{bmatrix} \begin{bmatrix} \gamma_0 \\ \gamma_1 \\ \gamma_2 \\ \lambda_1 \\ \lambda_2 \end{bmatrix} + \begin{bmatrix} e \\ \varepsilon \end{bmatrix} \tag{21}$$

The structure for the covariance matrix is identical to the single control variable case. The joint distribution of $(\hat{\gamma}_2 - \hat{\lambda}_2)$ asymptotically follows a chi-square distribution with J degrees of freedom. This quantity is a joint test that each $\gamma_{2j} = \lambda_{2j}$, where $j = 1, \dots, J$. Let 0_J denote a column vector of zeros of length J , 0_{JK} denote a $(J \text{ by } K)$ matrix of zeros (K is the number of instruments), and I_J the identity matrix of order J . Generalizing Equation 20, the test statistic for

¹⁷ \mathbf{R} combined with $\hat{\boldsymbol{\theta}}$ computes the test statistic as the square of $(\gamma_2 - \lambda_2)$ divided by its variance.

all of the control variables considered together is given by:

$$\begin{aligned}
(\mathbf{R}\hat{\boldsymbol{\theta}})'(\mathbf{R}\hat{\boldsymbol{\Sigma}}\mathbf{R}')^{-1}(\mathbf{R}\hat{\boldsymbol{\theta}}) &\sim \chi^2(J) \\
\hat{\boldsymbol{\theta}} &\equiv (\hat{\gamma}_0, \hat{\gamma}'_1, \hat{\gamma}'_2, \hat{\lambda}_1, \hat{\lambda}'_2)' \\
\mathbf{R} &\equiv \begin{bmatrix} 0_J & 0_{JK} & I_J & 0_J & -I_J \end{bmatrix}
\end{aligned} \tag{22}$$

Rather than performing one joint test across all control variables, researchers can also test individual control variables within a multivariate setting or subsets of control variables using the following test statistics. Suppose we wish to test an individual control variable j . Let \mathbf{R}_j denote the j^{th} row of the matrix \mathbf{R} . The test statistic is given by:

$$(\mathbf{R}_j\hat{\boldsymbol{\theta}})'(\mathbf{R}_j\hat{\boldsymbol{\Sigma}}\mathbf{R}'_j)^{-1}(\mathbf{R}_j\hat{\boldsymbol{\theta}}) \sim \chi^2(1) \tag{23}$$

To test subsets of control variables, let k denote the number of control variables to be tested, and \mathbf{k} denote the row indices corresponding to the k control variables to be jointly tested. Let \mathbf{R}_k denote the \mathbf{R} matrix with the relevant k rows. The test statistic is given by:^{18, 19}

$$(\mathbf{R}_k\hat{\boldsymbol{\theta}})'(\mathbf{R}_k\hat{\boldsymbol{\Sigma}}\mathbf{R}'_k)^{-1}(\mathbf{R}_k\hat{\boldsymbol{\theta}}) \sim \chi^2(k) \tag{24}$$

¹⁸A note on computing $\hat{\boldsymbol{\Sigma}}$: The dimensions of \mathbf{X} are $[2N \text{ by } (2J + K + 2)]$. $\text{var}(e)$ and $\text{var}(\varepsilon)$ can be estimated by computing the variance of the residuals using the first half and the second half of the observations from the SUR regression, respectively. Due to the high dimension of $\boldsymbol{\Omega}$ ($2N$ by $2N$), it is computationally preferable to directly compute the matrix $\mathbf{X}'\boldsymbol{\Omega}\mathbf{X}$, which is $[(2J + K + 2) \text{ by } (2J + K + 2)]$. Let $\mathbf{X}_1 \equiv (1, \mathbf{z}', \mathbf{x}'_2)$ and $\mathbf{X}_2 \equiv (1, \mathbf{x}'_2)$, where \mathbf{X}_1 is $[N \text{ by } (J + K + 1)]$ and \mathbf{X}_2 is $[N \text{ by } (J + 1)]$. Then $\mathbf{X}'\boldsymbol{\Omega}\mathbf{X} = \begin{bmatrix} \text{var}(e) \mathbf{X}'_1 \mathbf{X}_1 & \text{var}(e) \mathbf{X}'_1 \mathbf{X}_2 \\ \text{var}(e) \mathbf{X}'_2 \mathbf{X}_1 & \text{var}(\varepsilon) \mathbf{X}'_2 \mathbf{X}_2 \end{bmatrix}$.

¹⁹Example code for calculating these statistics in Python and Stata is available upon request.

4 Maximum Possible Bias

4.1 Single Instrumented Variable

The diagnostic test derived in Section 3 will allow researchers to be able to test whether the correlations between the instruments and the non-instrumented control variables are large enough to potentially cause bias in the 2SLS estimate for the key variable of interest. This test is valid when a single endogenous variable is being instrumented. In this section we derive a formula to show how large the potential bias in the key coefficient of interest could be. Subtracting β_1 from both sides of Equation 8 provides an expression for bias in the 2SLS coefficient related to both the violation of the narrow exclusion condition and the presence of endogenous control variables. To derive an expression of the size of bias coming from the contaminated control variables, we first assume that the narrow exclusion condition holds, such that $\text{cov}(\mathbf{z}, \mathbf{w}) = 0$. This allows the bias to be written directly as a function of the endogenous controls.²⁰ Let σ_j , σ_w , σ_v , and σ_y denote the standard deviation of x_{2j} , w , v , and y , respectively, and let ρ_j denote the correlation between x_{2j} and v . Given the assumptions noted above:

$$\begin{aligned}\hat{\beta}_{1,2SLS} - \beta_1 &= (\hat{\gamma}_2 - \hat{\lambda}_2)' \frac{\text{cov}(\mathbf{x}_2, w)}{\text{var}(\hat{x}_1^*)} \\ &= (\hat{\gamma}_2 - \hat{\lambda}_2)' \frac{\text{cov}(\mathbf{x}_2, v)}{\text{var}(\hat{x}_1^*)} \\ &= \sum_{j=1}^J \frac{(\hat{\gamma}_{2j} - \hat{\lambda}_{2j}) \sigma_j \sigma_v \rho_j}{\text{var}(\hat{x}_1^*)}\end{aligned}\tag{25}$$

Next, we assume that $\sigma_y \geq \sigma_v$. This follows if the regressors explain some variation in the dependent variable, and rules out certain cases of severe correlation between the error term and the

²⁰The bias can be written equivalently using the error term v from the second stage regression or w from the DGP. To see this, let \hat{e} denote the residual from the first stage regression, so that $x_1 = \hat{x}_1 + \hat{e}$. Then the second stage regression error term $v = w + \beta_1 \hat{e}$. Since \hat{e} is orthogonal to \mathbf{x}_2 , we have that $\text{cov}(\mathbf{x}_2, v) = \text{cov}(\mathbf{x}_2, w)$.

regressors. This assumption allows the bias expression to be written as follows:

$$\hat{\beta}_{1,2SLS} - \beta_1 \leq \sum_{j=1}^J \frac{(\hat{\gamma}_{2j} - \hat{\lambda}_{2j})\sigma_j\sigma_y\rho_j}{\text{var}(\hat{x}_1^*)} \quad (26)$$

The next step is to bound the correlation ρ_j . The discussion around Equations 27 - 32 describe the steps involved in doing this. To bound the correlation, we first derive an expression relating σ_v^2 and its least squares estimate $\hat{\sigma}_v^2$. Let $\mathbf{x} \equiv (\hat{x}_1, \mathbf{x}'_2)'$ and $\beta \equiv (\beta_1, \beta'_2)'$. The second stage regression model and its variance decomposition are given by:

$$\begin{aligned} y &= \beta_0 + \beta' \mathbf{x} + v \\ \sigma_y^2 &= \beta' \text{var}(\mathbf{x}) \beta + 2\beta' \text{cov}(\mathbf{x}, v) + \sigma_v^2 \end{aligned} \quad (27)$$

Let \hat{v} , $\hat{\beta}_0$ and $\hat{\beta}$ denote the least squares estimates of v , β_0 and β , respectively. The least squares estimates of the second stage regression model and its variance decomposition are given by:²¹

$$\begin{aligned} y &= \hat{\beta}_0 + \hat{\beta}' \mathbf{x} + \hat{v} \\ \sigma_y^2 &= \hat{\beta}' \text{var}(\mathbf{x}) \hat{\beta} + \hat{\sigma}_v^2 \end{aligned} \quad (28)$$

Substituting $\hat{\beta} = \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y)$, we further expand Equation 28.

$$\begin{aligned} \sigma_y^2 &= [\text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y)]' \text{var}(\mathbf{x}) [\text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y)] + \hat{\sigma}_v^2 \\ \sigma_y^2 &= [\beta + \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, v)]' \text{var}(\mathbf{x}) [\beta + \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, v)] + \hat{\sigma}_v^2 \\ \sigma_y^2 &= \beta' \text{var}(\mathbf{x}) \beta + 2\beta' \text{cov}(\mathbf{x}, v) + \text{cov}(\mathbf{x}, v)' \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, v) + \hat{\sigma}_v^2 \end{aligned} \quad (29)$$

²¹While $\text{cov}(\mathbf{x}, v)$ may be nonzero, $\text{cov}(\mathbf{x}, \hat{v}) = 0$ by construction.

Next, we equate σ_y^2 using Equations 27 and 29 and solve for $\hat{\sigma}_v^2$ as a function of σ_v^2 .

$$\begin{aligned}\sigma_v^2 &= \hat{\sigma}_v^2 + \text{cov}(\mathbf{x}, v)' \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, v) \\ \hat{\sigma}_v^2 &= \sigma_v^2 - \text{cov}(\mathbf{x}, v)' \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, v)\end{aligned}\tag{30}$$

In defining the maximum possible bias, we assume that control variable j is correlated with the error term, while the other control variables are not which means that $\text{cov}(\hat{x}_1, v) = \hat{\gamma}_{2j} \text{cov}(x_{2j}, v)$. Let $\tilde{\gamma}_{2j}$ be a $[(J+1) \text{ by } 1]$ vector with $\hat{\gamma}_{2j}$ the first element, 1 the $(j+1)^{th}$ element, and 0 the remaining elements. Finally, let $\tilde{x}_j \equiv \sqrt{\tilde{\gamma}_{2j}' \text{var}(\mathbf{x})^{-1} \tilde{\gamma}_{2j}}$. Given this setup the second part of Equation 30 can be written as follows:

$$\text{cov}(\mathbf{x}, v)' \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, v) = (\rho_j \sigma_j \sigma_v \tilde{x}_j)^2\tag{31}$$

Substituting Equation 31 into Equation 30 and applying $\sigma_y \geq \sigma_v$, we bound the correlation ρ_j :

$$\begin{aligned}\hat{\sigma}_v^2 &= \sigma_v^2 [1 - (\rho_j \sigma_j \tilde{x}_j)^2] \\ \hat{\sigma}_v^2 &\leq \sigma_y^2 [1 - (\rho_j \sigma_j \tilde{x}_j)^2] \\ |\rho_j| &\leq (\sigma_y \sigma_j \tilde{x}_j)^{-1} \sqrt{\sigma_y^2 - \hat{\sigma}_v^2}\end{aligned}\tag{32}$$

Finally, we substitute Equation 32 into Equation t26 to bound the total bias and thereby define the maximum possible bias MPB_j as follows:

$$|\hat{\beta}_{1,2SLS} - \beta_1| \leq \left| \frac{(\hat{\gamma}_{2j} - \hat{\lambda}_{2j}) \sqrt{\sigma_y^2 - \hat{\sigma}_v^2}}{\text{var}(\hat{x}_1^*) \tilde{x}_j} \right| \equiv MPB_j\tag{33}$$

4.2 Multiple Instrumented Variables

Thus far in the paper, we have assumed the 2SLS system includes a single instrumented endogenous variable. In Appendix B, we relax this assumption and allow for multiple instrumented endogenous variables and derive expressions for the bias coming from contaminated controls and the maximum possible bias. Future research is needed to be able to define the appropriate test statistic in a setting with multiple first stages. The appropriate test statistic in this type of setting would be nonlinear and different from the test proposed in this paper for a single endogenous variable.

5 Simulation

In this section we use simulated data to assess the performance of our proposed test, and examine how the inclusion of contaminated controls affects the bias in the 2SLS coefficient for the key variable of interest. We start with a simple case and then consider a variety of extensions and robustness checks. We vary the magnitude of the correlations and the instrument strength in the simulations to provide better understanding for what matters in practice.

5.1 Baseline

We start the simulation exercise considering a setting similar to that observed in empirical work, as described in Equation 1, with one key endogenous variable of interest x_1 , one control variable x_2 , and an exogenous instrument z for x_1 . We model the endogenous error term w as the sum of an omitted variable effect plus an exogenous error term, such that $w = \beta_m x_m + u$, and $E(u|x_1, x_2, x_m) = 0$. x_m is an omitted variable in the observable model and β_m captures the partial effect of the omitted variable on y . The generated data is created using the following assumptions and calibrations: (1) the instrument z is constructed to be correlated with x_1 but not correlated with x_m or u and hence satisfies the narrow exclusion condition, (2) β_1 , β_2 , and β_m are set equal to one, (3) x_1 , x_2 , x_m , z ,

and u are jointly normally distributed with each variable set to have mean zero and variance of one, and (4) u is uncorrelated with any of the other variables. Using simulations built around these assumptions we consider various scenarios that differ in the level of correlation assumed between x_1 and x_2 as well as between z and x_2 . In each simulation scenario the underlying data is generated using the true data generating process (including x_m and β_m) described above and then the estimated model in the simulation attempts to recover a non-biased estimate of β_1 using a 2SLS model that does not include the omitted variable x_m and uses z as an instrument for the endogenous variable x_1 . In each scenario we run 100,000 simulations with a sample size of 10,000 in each case. The low-level correlations assumed to exist among the variables are similar to those seen in many real-world datasets and hence the simulation results reported in this section provide direct evidence for how contaminated control variables, with even low-level correlations, can create significant bias in the key 2SLS coefficient of interest.

Table 1 reports the simulation results for sixteen different scenarios where the estimated model in each case involves one control variable (x_2) but the scenarios differ in terms of the correlations that exist among the variables. The correlation information for each scenario is reported under the Correlation Scenarios columns in the table. Columns 1-3 report the fraction of times the null hypothesis is rejected for tests at the 10%, 5%, and 1% levels for each scenario based on 100,000 simulations, as per Equation 20. Theoretically, for a sufficiently large sample, when there is no correlation between x_2 and z , the fraction of times the null hypothesis is rejected should be equal to the significance level. The results from the simulation tell us how well the asymptotic distribution approximates the finite sample distribution for empirically relevant samples sizes and correlation structures. When there is a nonzero correlation between x_2 and z , the fraction of times the null hypothesis is rejected is the power of the test, with higher fractions indicating better performance. For comparison purposes, we report the bias in the β_1 estimate for 2SLS models that either include or exclude the control variable. In the last column of the table we report the maximum possible bias that could exist in the key coefficient of interest in each scenario coming from the contaminated

control variable per Equation 33. The bias and the maximum possible bias values are computed as the average across all simulations within each scenario.

Table 1: Simulation results: one instrument one control

	Correlation Scenarios				(1)	(2)	(3)	(4)	(5)	(6)
	$\rho_{x_2,z}$	$\rho_{x_1,z}$	ρ_{x_1,x_2}	ρ_{x_2,x_m}	$R_{0.10}$	$R_{0.05}$	$R_{0.01}$	bias	bias_{nc}	MPB
(1)	0.0	0.1	0.1	0.1	0.097	0.046	0.008	-0.004	-0.005	0.097
(2)	0.0	0.1	0.1	0.3	0.096	0.046	0.007	-0.003	-0.004	0.113
(3)	0.0	0.1	0.3	0.1	0.096	0.047	0.008	-0.003	-0.007	0.113
(4)	0.0	0.1	0.3	0.3	0.096	0.046	0.008	-0.002	-0.006	0.129
(5)	0.0	0.3	0.1	0.1	0.100	0.050	0.010	-0.000	-0.000	0.033
(6)	0.0	0.3	0.1	0.3	0.102	0.050	0.010	-0.000	-0.000	0.038
(7)	0.0	0.3	0.3	0.1	0.099	0.049	0.010	-0.000	-0.001	0.038
(8)	0.0	0.3	0.3	0.3	0.099	0.049	0.009	-0.000	-0.001	0.043
(9)	0.2	0.1	0.1	0.1	1.000	1.000	1.000	-0.258	2.016	2.989
(10)	0.2	0.1	0.1	0.3	1.000	1.000	1.000	-0.766	2.016	3.483
(11)	0.2	0.1	0.3	0.1	0.996	0.990	0.954	-0.552	2.014	7.328
(12)	0.2	0.1	0.3	0.3	0.996	0.990	0.954	-1.617	2.014	8.370
(13)	0.2	0.3	0.1	0.1	1.000	1.000	1.000	-0.072	0.667	0.861
(14)	0.2	0.3	0.1	0.3	1.000	1.000	1.000	-0.215	0.667	0.994
(15)	0.2	0.3	0.3	0.1	1.000	1.000	1.000	-0.085	0.666	1.160
(16)	0.2	0.3	0.3	0.3	1.000	1.000	1.000	-0.251	0.667	1.317

This table reports baseline simulation results for a 2SLS model with one instrument and one control variable using 16 different scenarios. In all models $\rho_{x_1,x_m} = 0.3$. The correlation information reported under the Correlation Scenarios section of the table describe the correlations that exist among the variables in the generated data in each of the 16 scenarios. The next three columns (columns 1 - 3) report the rejection rate at the 10%, 5%, and 1% levels, respectively. Columns 4 - 5 report the bias in the 2SLS estimate of β_1 with and without (nc) the inclusion of the control variable. The last column reports the maximum possible bias in the β_1 estimate.

The data generated for the simulations used in each of the scenarios reported in rows 1-8 of Table 1 have zero correlation between the control variable x_2 and the instrument, meaning there is no contaminated control bias in these scenarios per the discussion in Section 2.2. Rejection rates would therefore be expected to equal the significance level of the test in these rows, and the bias expected to equal zero with or without controls. As reported, the rejection rates from the simulated data in these rows are indeed close to significance levels. When the instrument is weak (rows 1-4), the rejection rates are modestly lower than significance levels, but always within half a

percentage point. In these scenarios, the bias is small and negative when the instrument is weak, though very close to zero. When the instrument is strong (rows 5-8), rejection rates are equal to their significance levels and the bias is zero regardless of whether the control variable is included. The simulated data results reported in rows 1-8 corroborate the intuition from the earlier analytical expressions showing that if the instrument used with the key variable of interest in the 2SLS system is uncorrelated with the other control variable in the model then there is no contaminated control bias in the key 2SLS coefficient estimate regardless of whether the control variable is correlated with the key variable of interest or whether the control variable is correlated with the omitted variable.

The data generated for the simulations used in each of the scenarios reported in rows 9-16 of Table 1 have nonzero correlations between the control variable and the instrument and hence the bias in the β_1 estimates is expected to be nonzero in these scenarios based on the analytical expressions for contaminated control bias described earlier in this paper (e.g., see analytical form of the bias in Equation 12). Given the known bias that exists in these scenarios, to be useful in applied settings our proposed test would ideally reject the null hypothesis for each of these scenarios. Higher rejection rates indicate greater power and better performance of the test. As reported in columns 1-3 in rows 9-16 the test rejects the null hypothesis almost 100% of the time showing that the proposed test for contaminated control bias correctly identifies the existence of the bias in these scenarios. In rows 11 and 12, a small fraction of non rejections occur when the instrument is weak and the control variable is strongly correlated with the variable of interest,

The bias reported in rows 9-16 of Table 1 varies depending on the scenario, but is substantial in all cases compared to the true β_1 value of 1. The bias is larger when the control variable is not included (column 5 versus column 4), when the instrument is weaker (rows 9-12), and when the correlation between the control variable and the key variable of interest is stronger (rows 11, 12, 15, and 16).²² As expected, the maximum possible bias in all 16 scenarios is generally much

²²The patterns regarding the bias in Table 1 are specific to this particular example, and should not be assumed to

larger than the actual bias, with the actual bias in the models that include the control variable being no more than a quarter of the theoretical maximum in the simulated data. This result is expected because the MPB formula is based on a series of assumptions that are intended to calculate the maximum possible bias and hence provides intuition about the worst-case scenario rather than an estimate of the actual bias.

Overall, the results in Table 1 indicate that the test statistic has an accurate rejection rate when the null hypothesis is true, and almost always rejects when the null hypothesis is false. The results also indicate, even when using relatively small correlations, as are often seen in the real-world data, that in the presence of contaminated controls with an instrument that is correlated with those control variables, the bias in the variable of interest can be substantial, regardless of whether the control variable is included in the regression.

Figures 1 and 2 plot the distribution of the test statistic with and without contamination, respectively. The plots visually describe the entire distribution, thus providing information beyond the rejection rates in the table. Without contamination, the distribution of the test statistic appears close to a chi-square distribution with one degree of freedom for all specifications, in line with the theory. With contamination, the test statistics have a bell-shaped distribution that varies in location and scale depending on the specification.²³ A weaker instrument (top four plots) reduces the mean of the test statistic. Increased correlation between the control variable and the key variable of interest (second and fourth row) modestly reduces both the mean and the variance of the test statistic.

Figures 3 and 4 plot the distribution of $\hat{\beta}_{1,2SLS}$ with and without contamination, respectively. Each figure plots the distribution with and without the inclusion of the control variable. The plots provide insight into how the distribution of $\hat{\beta}_{1,2SLS}$ is affected by the inclusion of potentially contaminated controls. Without contamination, the distribution of the β_1 estimates is narrower when

necessarily hold more generally. For example, the sign on the bias could be positive or negative depending on the situation.

²³The theoretical distribution of the test statistic is unknown under the alternative hypothesis.

the control variable is included, indicating a more efficient estimate. With contamination, the distribution is generally further away from the true value when the control is not included, indicating generally better performance when the control variable is included.

5.2 Extensions

We consider a series of extensions to the baseline simulation results. To save space, we omit the figures for all extensions considered, and only report the tables.

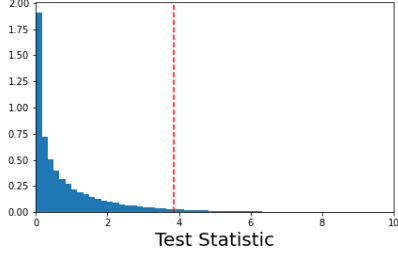
Table 2: Simulation results: one instrument one control, small sample

	Correlation Scenarios				(1)	(2)	(3)	(4)	(5)	(6)
	$\rho_{x_2,z}$	$\rho_{x_1,z}$	ρ_{x_1,x_2}	ρ_{x_2,x_m}	$R_{0.10}$	$R_{0.05}$	$R_{0.01}$	bias	bias _{nc}	MBP
(1)	0.0	0.1	0.1	0.1	0.048	0.014	0.001	-0.036	-0.089	0.379
(2)	0.0	0.1	0.1	0.3	0.049	0.015	0.001	-0.041	-0.067	0.443
(3)	0.0	0.1	0.3	0.1	0.052	0.016	0.001	-0.046	-0.097	0.427
(4)	0.0	0.1	0.3	0.3	0.051	0.015	0.000	0.012	-0.117	0.520
(5)	0.0	0.3	0.1	0.1	0.095	0.047	0.007	-0.003	-0.004	0.105
(6)	0.0	0.3	0.1	0.3	0.095	0.045	0.007	-0.003	-0.005	0.122
(7)	0.0	0.3	0.3	0.1	0.098	0.047	0.008	-0.003	-0.007	0.122
(8)	0.0	0.3	0.3	0.3	0.096	0.047	0.008	-0.002	-0.006	0.139
(9)	0.2	0.1	0.1	0.1	0.813	0.702	0.408	-0.843	2.271	5.801
(10)	0.2	0.1	0.1	0.3	0.813	0.702	0.411	-1.002	2.393	5.678
(11)	0.2	0.1	0.3	0.1	0.364	0.240	0.071	-1.000	1.660	35.032
(12)	0.2	0.1	0.3	0.3	0.364	0.240	0.071	-2.502	2.185	44.310
(13)	0.2	0.3	0.1	0.1	1.000	1.000	1.000	-0.075	0.670	0.872
(14)	0.2	0.3	0.1	0.3	1.000	1.000	1.000	-0.219	0.671	1.007
(15)	0.2	0.3	0.3	0.1	1.000	1.000	1.000	-0.090	0.667	1.179
(16)	0.2	0.3	0.3	0.3	1.000	1.000	1.000	-0.257	0.667	1.338

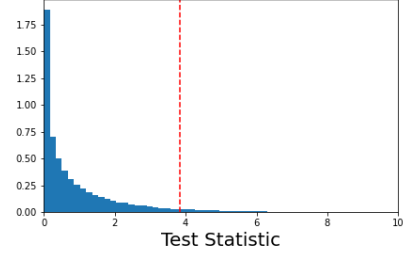
This table reports simulation results using a small sample size for a model with one instrument and one control variable using 16 different scenarios. In all models $\rho_{x_1,x_m} = 0.3$. The correlation information reported under the Correlation Scenarios section of the table describe the correlations that exist among the variables in the generated data in each of the 16 scenarios. The next 3 columns (columns 1 - 3) report the rejection rate at the 10%, 5%, and 1% levels, respectively. Columns 4 - 5 report the bias in the 2SLS estimate of β_1 with and without (nc) the inclusion of the control variable. The last column reports the maximum possible bias.

Figure 1: Test statistic: $\rho(x_2, z) = 0$ (no contamination)

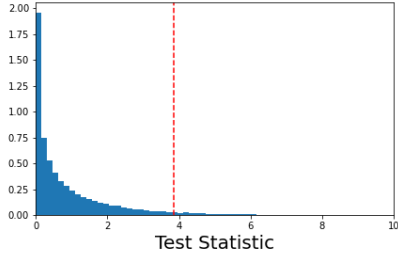
(a) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
weak instrument, less relevant control, low endogeneity



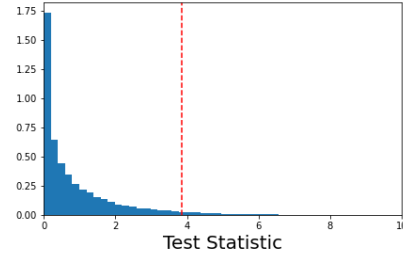
(b) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
weak instrument, less relevant control, high endogeneity



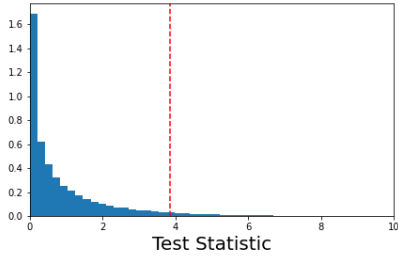
(c) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
weak instrument, more relevant control, low endogeneity



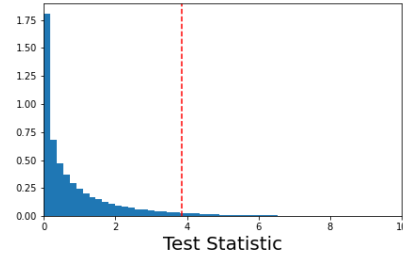
(d) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
weak instrument, more relevant control, high endogeneity



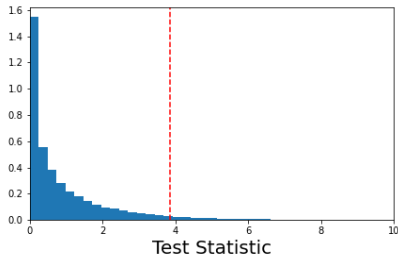
(e) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
strong instrument, less relevant control, low endogeneity



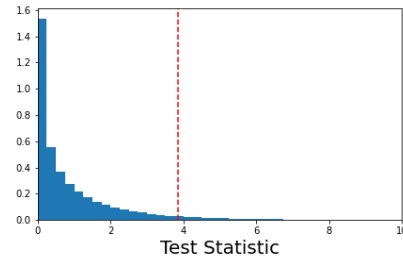
(f) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
strong instrument, less relevant control, high endogeneity



(g) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
strong instrument, more relevant control, low endogeneity



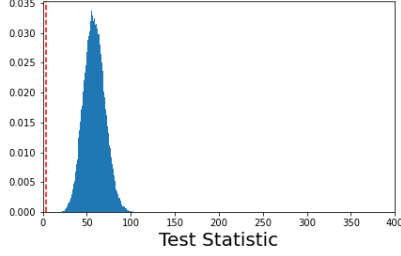
(h) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
strong instrument, more relevant control, high endogeneity



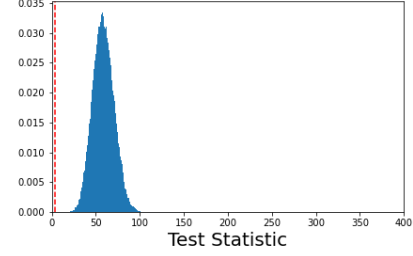
The figures plot the distribution of the test statistic when the instrument is uncorrelated with the other control variable. The vertical dashed red line is the critical value for a 5% significance level hypothesis test.

Figure 2: Test statistic: $\rho(x_2, z) = 0.2$ (contaminated control)

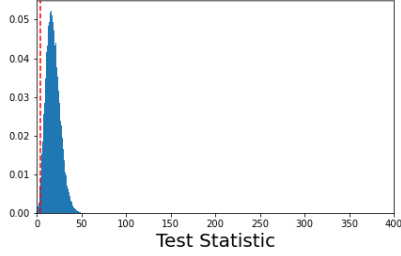
(a) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
weak instrument, less relevant control, low endogeneity



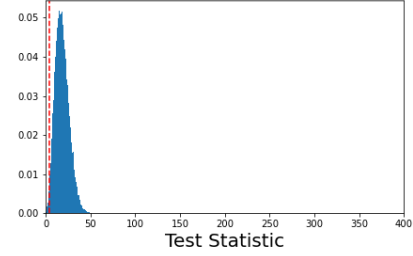
(b) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
weak instrument, less relevant control, high endogeneity



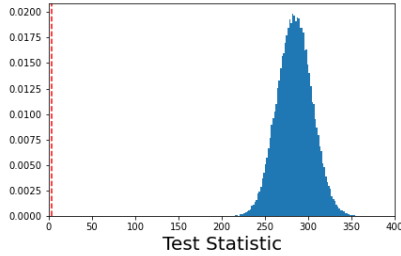
(c) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
weak instrument, more relevant control, low endogeneity



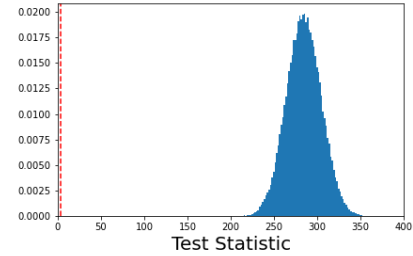
(d) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
weak instrument, more relevant control, high endogeneity



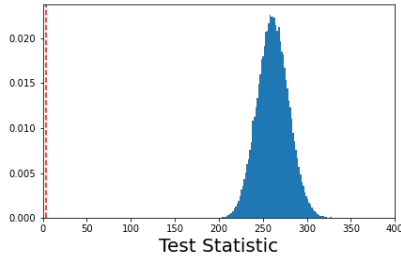
(e) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
strong instrument, less relevant control, low endogeneity



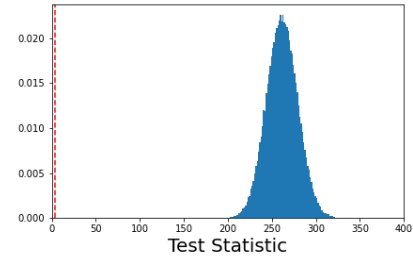
(f) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
strong instrument, less relevant control, high endogeneity



(g) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
strong instrument, more relevant control, low endogeneity



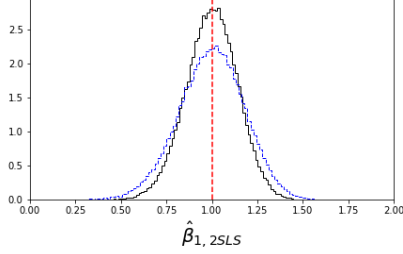
(h) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
strong instrument, more relevant control, high endogeneity



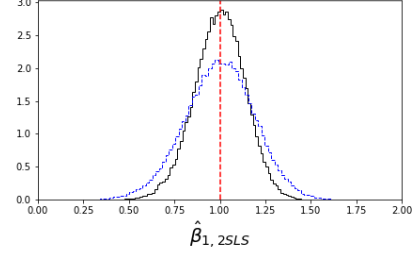
The figures plot the distribution of the test statistic when the instrument is correlated with the other control variable. The vertical dashed red line is the critical value for a 5% significance level hypothesis test.

Figure 3: 2SLS estimate: $\rho(x_2, z) = 0$ (no contamination)

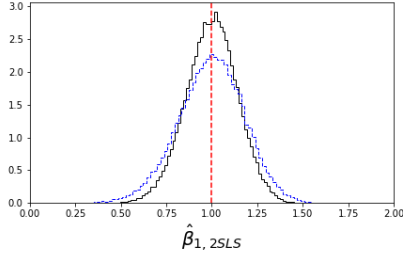
(a) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
weak instrument, less relevant control, low endogeneity



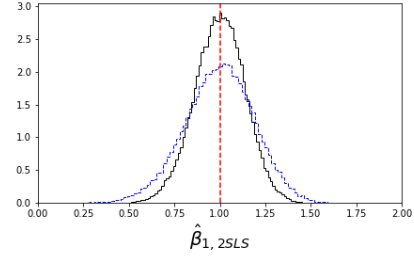
(b) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
weak instrument, less relevant control, high endogeneity



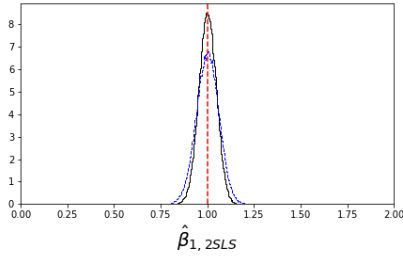
(c) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
weak instrument, more relevant control, low endogeneity



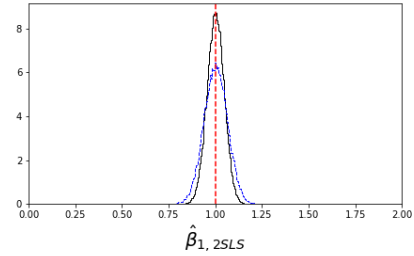
(d) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
weak instrument, more relevant control, high endogeneity



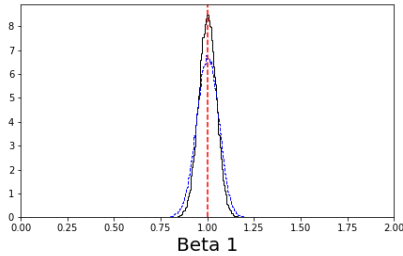
(e) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
strong instrument, less relevant control, low endogeneity



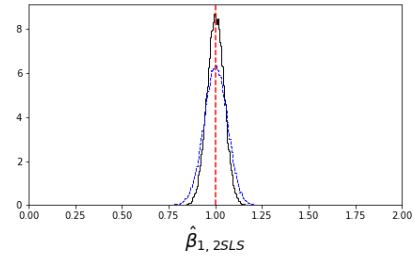
(f) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
strong instrument, less relevant control, high endogeneity



(g) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
strong instrument, more relevant control, low endogeneity



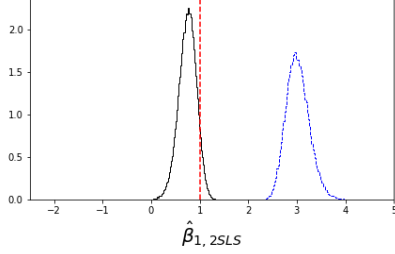
(h) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
strong instrument, more relevant control, high endogeneity



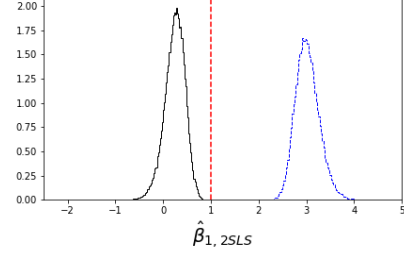
The figures plot the distribution of the 2SLS estimate of β_1 when the instrument is uncorrelated with the other control variable, with (solid black) and without (dashed blue) the inclusion of the control variable in the regression. The vertical dashed red line is the true value of β_1 .

Figure 4: 2SLS estimate: $\rho(x_2, z) = 0.2$ (contaminated control)

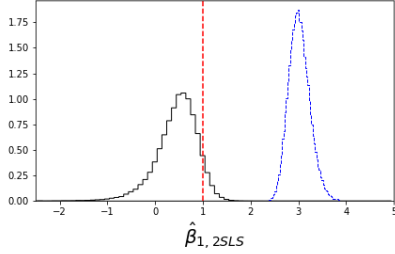
(a) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
weak instrument, less relevant control, low endogeneity



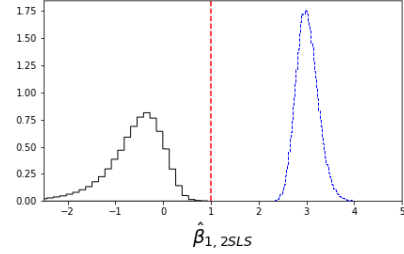
(b) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
weak instrument, less relevant control, high endogeneity



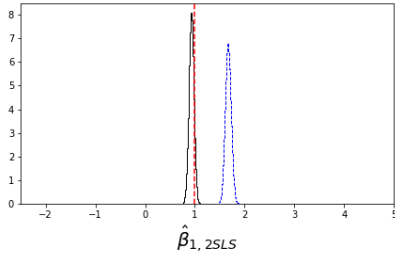
(c) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
weak instrument, more relevant control, low endogeneity



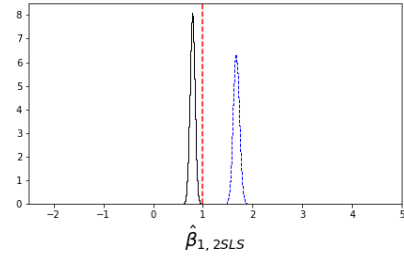
(d) $\rho(x_1, z) = 0.1, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
weak instrument, more relevant control, high endogeneity



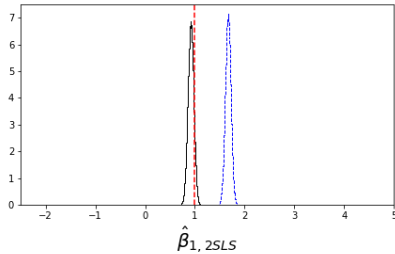
(e) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.1$
strong instrument, less relevant control, low endogeneity



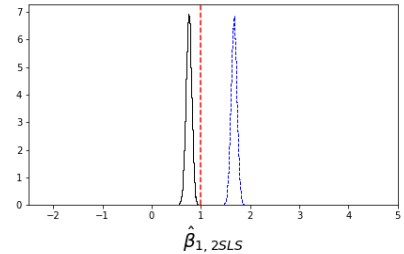
(f) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.1, \rho(x_2, x_m) = 0.3$
strong instrument, less relevant control, high endogeneity



(g) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.1$
strong instrument, more relevant control, low endogeneity



(h) $\rho(x_1, z) = 0.3, \rho(x_1, x_2) = 0.3, \rho(x_2, x_m) = 0.3$
strong instrument, more relevant control, high endogeneity



The figures plot the distribution of the 2SLS estimate of β_1 when the instrument is correlated with the control variable, with (solid black) and without (dashed blue) the inclusion of the control variable in the regression. The vertical dashed red line is the true value of β_1 .

5.2.1 Small Sample

First, we consider the impact of a smaller sample size. We repeat the earlier baseline simulations but this time using a sample size of 1,000. The results are reported in Table 2. The small sample impacts the performance of the test primarily when the instrument is weak. In rows 1-4, the rejection rates are well below their significance levels. In rows 9-12 the rejection rates are well below one, indicating lower power for the test. When the instrument is strong (rows 5-8, 13-16) there is little difference in the results compared with the baseline larger sample.

5.2.2 High Impact Omitted Variable

In this section, we set $\beta_m = 3$ when generating the data rather than $\beta_m = 1$. This change increases the impact of the omitted variable on y compared to the previous simulation, such that a greater fraction of y is now explained by the omitted variable. The results are reported in Table 3. The test statistic performs similarly as in the baseline case with a low impact omitted variable. As expected, the bias is larger across all specifications. In some cases, the bias is larger when the control variable is included relative to when the control variable is not included. The ratio of the actual bias to the maximum possible bias is larger than in the earlier simulation, and in some cases the actual bias is close to one-half the MPB calculation.²⁴

5.2.3 Additional Control Variables

In this section, we add an additional control variable but otherwise retain the same structure as the first simulation. For this section, let $\mathbf{x}_2 \equiv (x_{21}, x_{22})'$ be control variables and let $\beta_2 \equiv (\beta_{21}, \beta_{22})'$. In generating the data, we set $\beta_1, \beta_{21}, \beta_{22}$ and β_m equal to one, $x_1, x_{21}, x_{22}, x_m, z$, and u are jointly normally distributed, and each variable has mean zero and variance one. u is uncorrelated with any of the other variables. The sample size is 10,000.

²⁴In some cases (rows 13 and 14), the MPB is much larger than other cases, but this is consistent with the actual bias also being much larger in these cases.

Table 3: Simulation results: one instrument one control, high impact of omitted variable

	Correlation Scenarios				(1)	(2)	(3)	(4)	(5)	(6)
	$\rho_{x_2,z}$	$\rho_{x_1,z}$	ρ_{x_1,x_2}	ρ_{x_2,x_m}	$R_{0.10}$	$R_{0.05}$	$R_{0.01}$	bias	bias _{nc}	MBP
(1)	0.0	0.1	0.1	0.1	0.095	0.045	0.008	-0.009	-0.011	0.113
(2)	0.0	0.1	0.1	0.3	0.095	0.046	0.008	-0.007	-0.008	0.162
(3)	0.0	0.1	0.3	0.1	0.097	0.047	0.008	-0.008	-0.012	0.129
(4)	0.0	0.1	0.3	0.3	0.095	0.045	0.008	-0.007	-0.013	0.177
(5)	0.0	0.3	0.1	0.1	0.099	0.050	0.010	-0.001	-0.001	0.038
(6)	0.0	0.3	0.1	0.3	0.099	0.049	0.010	-0.001	-0.001	0.054
(7)	0.0	0.3	0.3	0.1	0.101	0.050	0.010	-0.000	-0.001	0.043
(8)	0.0	0.3	0.3	0.3	0.103	0.050	0.010	-0.000	-0.001	0.059
(9)	0.2	0.1	0.1	0.1	1.000	1.000	1.000	-0.775	2.009	3.482
(10)	0.2	0.1	0.1	0.3	1.000	1.000	1.000	-2.298	2.012	4.986
(11)	0.2	0.3	0.1	0.1	1.000	1.000	1.000	-0.215	0.667	0.994
(12)	0.2	0.3	0.1	0.3	1.000	1.000	1.000	-0.645	0.665	1.403
(13)	0.2	0.1	0.3	0.1	0.996	0.990	0.954	-1.661	2.008	8.371
(14)	0.2	0.1	0.3	0.3	0.995	0.989	0.954	-4.854	2.008	11.547
(15)	0.2	0.3	0.3	0.1	1.000	1.000	1.000	-0.252	0.666	1.317
(16)	0.2	0.3	0.3	0.3	1.000	1.000	1.000	-0.752	0.666	1.800

This table reports simulation results for a model with one instrument, one control variable, and a high impact omitted variable using 16 different scenarios. In all models $\rho_{x_1,x_m} = 0.3$. The correlation information reported under the Correlation Scenarios section of the table describe the correlations that exist among the variables in the generated data in each of the 16 scenarios. The next three columns (columns 1-3) report the rejection rate at the 10%, 5%, and 1% levels, respectively. Columns 4 - 5 report the bias in the 2SLS estimate of β_1 with and without (nc) the inclusion of the control variable. The last column reports the maximum possible bias.

For each scenario we run three tests. The first is a joint test of whether all the controls as a group are contaminated. The second and third are tests of whether each of the two controls are individually contaminated. As before, we report the fraction of times the null hypothesis is rejected, for tests at the 10%, 5%, and 1% levels, and we report the bias in the estimate of β_1 , when the control variables are included, and when they are excluded. Finally, we report the maximum possible bias coming from each control variable.

In the simulation, both controls are positively correlated with the omitted variable and hence are endogenous. In the scenarios where both controls are also uncorrelated with the instrument, there will be no contaminated control bias based on the analytical expressions for the bias discussed

earlier. In contrast, if either or both endogenous controls are also correlated with the instrument, then there will be contamination and nonzero bias. In the simulation, when the controls are correlated with the instrument, the first control variable is calibrated to have a high correlation with the instrument (high correlation control), and the second control is calibrated to have a low correlation with the instrument (low correlation control).

When testing each control individually, the test will show evidence of contamination if either that specific control variable is correlated with the instrument, or if that control variable is correlated with the other control variable that is correlated with the instrument. Thus, for individual control variable tests, there are two channels by which the test can suggest contamination. The first is via a direct correlation with the instrument, and the second is via correlation with another control variable that is itself correlated with the instrument.

Table 4 reports the results from the joint test. Table 5 reports the results for each control variable individually. In both tables, in rows 1-4, given the assumptions used to create the data there should be no contamination and the bias should be close to zero. Rejection rates should be close to or equal to their nominal levels. The results reported in the tables show the contaminated control test has higher power when the instrument is stronger. In rows 5-8, the low correlation (second) control is contaminated and hence we would expect the contaminated control test to reject the null. As reported, the rejection rates for both the joint and individual tests are equal to one for these rows. For the uncontaminated control variable, as expected the rejection rates are close to their nominal levels when it is not correlated with the contaminated control, and closer to one when it is correlated. In rows 9-12, the high-correlation (first) control is contaminated. The pattern of results in these rows is similar to the previous four rows, with improved power on individual tests of the uncontaminated control. In rows 13-16, both controls are contaminated. Rejection rates are close to or equal to one for all tests showing that the test is very good at identifying evidence of contamination in these situations.

Across all specifications, the bias is worse when the control variables are not included, regard-

Table 4: Simulation results: one instrument two controls, joint test

	Correlation Scenarios				(1)	(2)	(3)	(4)	(5)	(6)	(7)
	$\rho_{x_{21},z}$	$\rho_{x_{22},z}$	$\rho_{x_1,z}$	$\rho_{x_{21},x_{22}}$	$R_{0.10}$	$R_{0.05}$	$R_{0.01}$	bias	bias _{nc}	MPB ₁	MPB ₂
(1)	0.0	0.0	0.1	0.0	0.091	0.042	0.007	-0.002	-0.007	0.160	0.160
(2)	0.0	0.0	0.1	0.2	0.089	0.042	0.006	-0.002	-0.006	0.166	0.167
(3)	0.0	0.0	0.3	0.0	0.100	0.050	0.010	-0.000	-0.000	0.053	0.053
(4)	0.0	0.0	0.3	0.2	0.098	0.049	0.010	-0.000	-0.001	0.056	0.056
(5)	0.0	0.1	0.1	0.0	1.000	1.000	1.000	-0.257	1.003	0.200	2.501
(6)	0.0	0.1	0.1	0.2	1.000	1.000	1.000	-0.207	1.003	0.516	2.553
(7)	0.0	0.1	0.3	0.0	1.000	1.000	1.000	-0.072	0.333	0.057	0.711
(8)	0.0	0.1	0.3	0.2	1.000	1.000	1.000	-0.059	0.333	0.151	0.748
(9)	0.2	0.0	0.1	0.0	1.000	1.000	1.000	-0.690	2.013	6.636	0.265
(10)	0.2	0.0	0.1	0.2	1.000	1.000	1.000	-0.516	2.012	6.322	1.291
(11)	0.2	0.0	0.3	0.0	1.000	1.000	1.000	-0.154	0.667	1.504	0.060
(12)	0.2	0.0	0.3	0.2	1.000	1.000	1.000	-0.125	0.667	1.560	0.319
(13)	0.2	0.1	0.1	0.0	0.983	0.965	0.887	-1.621	3.025	10.378	5.271
(14)	0.2	0.1	0.1	0.2	0.999	0.998	0.987	-1.049	3.024	7.760	2.626
(15)	0.2	0.1	0.3	0.0	1.000	1.000	1.000	-0.251	1.000	1.626	0.825
(16)	0.2	0.1	0.3	0.2	1.000	1.000	1.000	-0.201	1.000	1.502	0.508

This table reports simulation results for a model with one instrument and two control variables using 16 different scenarios. In all models $\rho_{x_1,x_m} = 0.3$, $\rho_{x_1,x_{21}} = 0.2$ and $\rho_{x_1,x_{22}} = 0.3$. The correlation information reported under the Correlation Scenarios section of the table describe the correlations that exist among the variables in the generated data in each of the 16 scenarios. The next three columns (columns 1 - 3) report the rejection rate at the 10%, 5%, and 1% levels, respectively. Columns 4 - 5 report the bias in the 2SLS estimate of β_1 with and without (nc) the inclusion of the control variables. The last 2 columns report the maximum possible bias for the first and second control variables.

less of whether they are contaminated. The bias is largest when both controls are contaminated, followed by when the high correlation control is contaminated, followed by when the low correlation control is contaminated. As expected, the bias is zero when both control variables are not contaminated and a weaker instrument always results in greater bias. When the control variables are correlated, the bias is marginally smaller. Overall, the simulation results indicate that the contaminated control test has the correct size when there is no contamination, and adequate power to detect contamination when contaminated control bias exists. The results hold for a wide variety of specifications.

Table 5: Simulation results: one instrument two controls, individual tests

	Correlation Scenarios				(1)	(2)	(3)	(4)	(5)	(6)
	$\rho_{x_{21},z}$	$\rho_{x_{21},z}$	$\rho_{x_1,z}$	$\rho_{x_{21},x_{22}}$	$R_{1,0.10}$	$R_{1,0.05}$	$R_{1,0.01}$	$R_{2,0.10}$	$R_{2,0.05}$	$R_{2,0.01}$
(1)	0.0	0.0	0.1	0.0	0.097	0.046	0.008	0.096	0.046	0.008
(2)	0.0	0.0	0.1	0.2	0.094	0.045	0.008	0.096	0.046	0.008
(3)	0.0	0.0	0.3	0.0	0.098	0.049	0.010	0.100	0.050	0.010
(4)	0.0	0.0	0.3	0.2	0.099	0.049	0.010	0.099	0.049	0.010
(5)	0.0	0.1	0.1	0.0	0.094	0.044	0.007	1.000	1.000	1.000
(6)	0.0	0.1	0.1	0.2	0.646	0.513	0.257	1.000	1.000	1.000
(7)	0.0	0.1	0.3	0.0	0.099	0.050	0.010	1.000	1.000	1.000
(8)	0.0	0.1	0.3	0.2	0.659	0.536	0.296	1.000	1.000	1.000
(9)	0.2	0.0	0.1	0.0	1.000	1.000	1.000	0.088	0.039	0.005
(10)	0.2	0.0	0.1	0.2	1.000	1.000	1.000	0.993	0.983	0.915
(11)	0.2	0.0	0.3	0.0	1.000	1.000	1.000	0.098	0.049	0.010
(12)	0.2	0.0	0.3	0.2	1.000	1.000	1.000	0.994	0.986	0.943
(13)	0.2	0.1	0.1	0.0	0.995	0.989	0.954	0.995	0.989	0.947
(14)	0.2	0.1	0.1	0.2	1.000	1.000	0.997	1.000	0.999	0.992
(15)	0.2	0.1	0.3	0.0	1.000	1.000	1.000	1.000	1.000	1.000
(16)	0.2	0.1	0.3	0.2	1.000	1.000	1.000	1.000	1.000	1.000

This table reports simulation results for a model with one instrument and two control variables using 16 different scenarios. In all models $\rho_{x_1,x_m} = 0.3$, $\rho_{x_1,x_{21}} = 0.2$ and $\rho_{x_1,x_{22}} = 0.3$. The correlation information reported under the Correlation Scenarios section of the table describe the correlations that exist among the variables in the generated data in each of the 16 scenarios. The next three columns (columns 1 - 3) report the rejection rate of the individual test for the first control variable at the 10%, 5%, and 1% levels, respectively. Columns 4 - 6 report the rejection rate of the individual test for the second control variable at the 10%, 5%, and 1% levels, respectively.

5.2.4 Additional Instrument

We now consider an overidentified simulation with two instruments ($z \equiv (z_1, z_2)'$) for the key variable of interest and a single additional control variable. In the data generating process we set β_1 , β_2 , and β_m equal to one, and x_1 , x_2 , x_m , z_1 , z_2 and u are jointly normally distributed, and each variable has mean zero and variance one. u is uncorrelated with any of the other variables. The sample size is 10,000.

In the scenarios considered below, the first instrument z_1 is correlated with the endogenous control variable x_2 and hence is not strictly exogenous, while the second instrument z_2 is uncor-

related with the endogenous control variable. Both z_1 and z_2 are correlated with the key variable of interest x_1 with varying levels of correlation across scenarios. The results are reported in Table 6. In addition to reporting the rejection rates from the contaminated control tests, bias, and maximum possible bias that occur across the 16 different correlation scenarios, we also report the 5% rejection rate of the Sargan overidentification test for comparison. We include the Sargan test here because this test is commonly used in the literature with overidentified models to test the validity of the instruments. We would expect the Sargan test to help identify scenarios when one or more of the instruments is invalid. The null hypothesis for the Sargan test is that the instruments are exogenous and hence uncorrelated with the error term. The Sargan test is known to have less power when instruments are weak.

Since the first instrument is correlated with the control variable in the simulated data, we have contamination in all model specifications. As reported in Table 6, the rejection rates for the contaminated control test are equal to one whenever the first instrument is strong (rows 9-16). When the first instrument is weak and the second is strong (rows 5-8), the power of the test is substantially reduced. This is however of limited concern, as the actual bias (and the MPB) are both very close to zero in these cases. When both instruments are weak (rows 1-4), the test has good power when the control is less correlated with the key variable of interest (rows 1-2) and less power when the control is more correlated with the key variable of interest (rows 3-4). Again, our test has low power generally when the actual bias itself is economically small and hence unlikely, with or without being able to detect the bias, to change the inference on the key coefficient estimate. In most specifications, the Sargan test tends to have lower power when the correlation between the control variable and the omitted variable is low, and vice versa. The contaminated control test and the Sargan test capture different features of the data and are complementary in their usefulness. The Sargan test is only possible in overidentified models, whereas the contaminated control test can be calculated in both just identified and overidentified models.

To show what the just identified results would be using the same data, in Table 7 we report

Table 6: Simulation results: two instruments one control, overidentified model

	Correlation Scenarios				(1)	(2)	(3)	(4)	(5)	(6)	(7)
	ρ_{x_1,z_1}	ρ_{x_1,z_2}	ρ_{x_1,x_2}	ρ_{x_2,x_m}	$R_{0.10}$	$R_{0.05}$	$R_{0.01}$	bias	bias _{nc}	MPB	$R_{s,0.05}$
(1)	0.1	0.1	0.1	0.1	1.000	1.000	0.999	-0.094	1.000	1.123	0.341
(2)	0.1	0.1	0.1	0.3	1.000	1.000	0.999	-0.281	1.000	1.307	0.980
(3)	0.1	0.1	0.3	0.1	0.622	0.498	0.264	-0.041	1.001	0.590	0.316
(4)	0.1	0.1	0.3	0.3	0.618	0.494	0.260	-0.125	1.000	0.672	0.993
(5)	0.1	0.3	0.1	0.1	0.335	0.228	0.086	-0.005	0.091	0.065	0.316
(6)	0.1	0.3	0.1	0.3	0.333	0.227	0.085	-0.015	0.091	0.075	0.995
(7)	0.1	0.3	0.3	0.1	0.338	0.231	0.087	0.005	0.091	0.075	0.316
(8)	0.1	0.3	0.3	0.3	0.338	0.230	0.087	0.014	0.091	0.085	0.995
(9)	0.3	0.1	0.1	0.1	1.000	1.000	1.000	-0.068	0.636	0.818	0.554
(10)	0.3	0.1	0.1	0.3	1.000	1.000	1.000	-0.203	0.636	0.943	0.821
(11)	0.3	0.1	0.3	0.1	1.000	1.000	1.000	-0.076	0.637	1.065	0.513
(12)	0.3	0.1	0.3	0.3	1.000	1.000	1.000	-0.229	0.636	1.209	0.844
(13)	0.3	0.3	0.1	0.1	1.000	1.000	1.000	-0.034	0.333	0.417	0.354
(14)	0.3	0.3	0.1	0.3	1.000	1.000	1.000	-0.101	0.333	0.479	0.978
(15)	0.3	0.3	0.3	0.1	1.000	1.000	1.000	-0.031	0.333	0.448	0.343
(16)	0.3	0.3	0.3	0.3	1.000	1.000	1.000	-0.094	0.333	0.507	0.984

This table reports simulation results for two instruments and one control variable. In all models $\rho_{x_1,x_m} = 0.3$, $\rho_{x_2,z_1} = 0.2$, $\rho_{x_2,z_2} = 0$, and $\rho_{z_1,z_2} = 0.2$. The correlation information reported under the Correlation Scenarios section of the table describe the correlations that exist among the variables in the generated data in each of the 16 scenarios. The next three columns (columns 1 - 3) report the rejection rate at the 10%, 5%, and 1% levels, respectively. Columns 4 - 5 report the bias in the 2SLS estimate of β_1 with and without (nc) the inclusion of the control variable. Column 6 reports the maximum possible bias. The last column reports the rejection rate of the Sargan test at the 5% level.

the results from two justidentified models using each of the instruments separately. Rejection rates at the 5% levels, bias, and maximum possible bias are reported for each model. Our test performs well across all specifications with the just identified models. Rejection rates are close to the significance level when we use the second instrument that is not correlated with the endogenous control variable. Rejection rates are close to one when we use the first instrument that is correlated with the endogenous control variable and hence the test is effective in identifying the contaminated control bias in these settings.

Table 7: Simulation results: two instruments one control, just identified models

	Correlation Scenarios				(1)	(2)	(3)	(4)	(5)	(6)
	ρ_{x_1,z_1}	ρ_{x_1,z_2}	ρ_{x_1,x_2}	ρ_{x_2,x_m}	$R_{1,0.05}$	$R_{2,0.05}$	bias ₁	bias ₂	MPB ₁	MPB ₂
(1)	0.1	0.1	0.1	0.1	1.000	0.046	-0.258	-0.003	2.990	0.097
(2)	0.1	0.1	0.1	0.3	1.000	0.046	-0.766	-0.003	3.485	0.113
(3)	0.1	0.1	0.3	0.1	0.990	0.046	-0.552	-0.003	7.349	0.113
(4)	0.1	0.1	0.3	0.3	0.989	0.046	-1.612	-0.002	8.403	0.129
(5)	0.1	0.3	0.1	0.1	1.000	0.049	-0.258	-0.000	2.987	0.033
(6)	0.1	0.3	0.1	0.3	1.000	0.049	-0.766	-0.001	3.483	0.038
(7)	0.1	0.3	0.3	0.1	0.990	0.049	-0.554	-0.000	7.323	0.038
(8)	0.1	0.3	0.3	0.3	0.989	0.050	-1.619	-0.000	8.384	0.043
(9)	0.3	0.1	0.1	0.1	1.000	0.046	-0.072	-0.003	0.861	0.097
(10)	0.3	0.1	0.1	0.3	1.000	0.046	-0.215	-0.003	0.994	0.113
(11)	0.3	0.1	0.3	0.1	1.000	0.047	-0.084	-0.002	1.160	0.114
(12)	0.3	0.1	0.3	0.3	1.000	0.046	-0.250	-0.003	1.317	0.129
(13)	0.3	0.3	0.1	0.1	1.000	0.051	-0.072	-0.000	0.861	0.033
(14)	0.3	0.3	0.1	0.3	1.000	0.050	-0.215	-0.000	0.994	0.038
(15)	0.3	0.3	0.3	0.1	1.000	0.050	-0.084	-0.000	1.160	0.038
(16)	0.3	0.3	0.3	0.3	1.000	0.050	-0.251	-0.000	1.317	0.043

This table reports simulation results for two instruments and one control variable. In all models $\rho_{x_1,x_m} = 0.3$, $\rho_{x_2,z_1} = 0.2$, $\rho_{x_2,z_2} = 0$, and $\rho_{z_1,z_2} = 0.2$. The correlation information reported under the Correlation Scenarios section of the table describe the correlations that exist among the variables in the generated data in each of the 16 scenarios. The next two columns (columns 1 -2) report the rejection rate at the 5% levels for the two just identified models. Columns 3 - 4 report the bias in the 2SLS estimate of β_1 for the two just identified models. The subscript numbers 1 and 2 shown in the headers for columns 1 - 6 indicate whether z_1 or z_2 was used as the instrument for that result. The last two columns report the bias in the 2SLS estimate of β_1 for the two just identified models.

6 Empirical Example

The results in Section 5 were based on simulated data. The discussion in that section was important to show how the contaminated control test and MPB calculations performed as expected when the exact differences between the observed model and true data generating process were known. That exercise was also important to show empirically how the power of the test and the usefulness of the MPB calculation can be affected by weak instruments. In this section, we leave the simulated data aside and illustrate the use of the new contaminated control test together with the MPB calculations with a well-known empirical example based on a paper published in the *Journal of Finance* in 2002.

6.1 Illustration of How to Use the Contaminated Control Test and MPB Calculations

There is a large literature dating back several decades that explores the diversification discount of multidivisional firms with publications in both economics and finance journals. Across the years different studies have utilized different samples and econometric approaches to explore the diversification discount, and depending on the specific sample and approach used, have reported varying levels of a discount with many papers in this literature finding at least some evidence of a discount consistent with a multidivisional firm's overall market value being less than the sum of the imputed values of its individual segments if each segment had existed outside the conglomerate. Various explanations for the discount have been suggested including the idea that corporate diversification could be associated with inefficient investment and/or internal capital market policies (e.g., Shin and Stulz (1998); Rajan et al. (2000), Ozbas and Scharfstein (2009)), lower acquisition market reactions and/or lower valued target firms (e.g., Morck et al. (1990); Graham and Wolf (2002)), or agency and governance issues (e.g., Denis et al. (1997); Hoechle et al. (2012); Ellis et al. (2018); Andreou et al. (2019)).

For our purposes we are interested in a 2SLS result reported by Campa and Kedia (2002) suggesting a diversification premium rather than a discount. This paper called attention to the fact that the decision to diversify is endogenous and suggested several instrumental variables to account for the endogeneity. In their empirical approach they include the various instruments in a pre-first-stage probit model and then use the predicted probability from this model as a single generated instrument in the first stage equation of a 2SLS system. The dependent variable in the second stage is a measure of the excess value at the firm compared to the sum of the imputed values of the firm's segments. The dependent variable in the first stage in the 2SLS system is an indicator variable for whether the firm is diversified in that year ($D=1$). We use this setting to illustrate how the diagnostic test and MPB formula suggested in this paper can help researchers explore 2SLS

results in important ways.

To facilitate the discussion and to streamline the example, we make some simplifying assumptions for the empirical approach. The first change we make is to drop the probit model that was used in advance of the first stage equation in Campa and Kedia (2002). Rather than using the instruments in a pre-first-stage model to generate a single instrument, we include the instruments directly in the first stage to instrument D and then estimate a traditional 2SLS system of equations. This change allows us to model the effect of different instruments individually rather than altogether as part of a single generated instrument and recasts the three-equation approach using a simpler two-equation approach. The second change we make is to include each control variable once in the model rather than including the controls along with their respective lagged values. This change makes the example more parsimonious and eliminates control variables that are highly correlated.

Our simplified 2SLS system of equations is shown below. Following Campa and Kedia (2002) and Berger and Ofek (1995) we estimate a firm's excess value each year as the log of the ratio of the firm's total capital to the firm's imputed value.²⁵ D is the key variable of interest and is an indicator for the firm having more than one business division in a given year. The control variables x_2 are similar to the control variables used by Campa and Kedia (2002) and include the capital expenditure-to-sales ratio, the EBIT-to-sales ratio, book leverage, the log of assets, and an indicator for whether the firm is one of the largest 500 US firms by market value each year. Fixed effects are included to control for year, industry, and country of incorporation. The first and second stage equations are shown below.

²⁵The firm's excess value each year is measured as the log of the ratio of the firm's total capital to the firm's imputed value. A positive (negative) ratio suggests the firm is trading at a premium (discount) compared to what it would trade if its various segments existed as separate entities. To calculate the imputed value, a sales multiplier is calculated for each industry each year as the median total capital-to-sales ratio based only on US single segment firms in that industry. The sales multiplier is then used to find the imputed value of segments that are in the same industry by multiplying the segment sales by the sales multiplier. The firm's overall imputed value in a given year is the sum of the imputed segment values. Campa and Kedia (2002) use both a sales multiplier and an asset multiplier in their analysis. For the purpose of demonstrating the effect that contaminated controls may have on the 2SLS results we focus only on the sales-based calculation in our empirical example.

$$ExcessValue = \beta_0 + \beta_1 D + \beta_2' x_2 + w$$

$$D = \gamma_0 + \gamma_1' z + \gamma_2' x_2 + e \quad (34)$$

To create our sample we follow the approach described in Campa and Kedia (2002) but using more recent data from 1986 through 2022²⁶. Using this sample we model the firm's excess value as a function of D and the other control variables and obtain 2 empirical results that are similar to the 2002 paper. First, in untabulated results using the second stage equation above as a simple OLS model rather than as a two-stage model, consistent with Campa and Kedia (2002), we obtain a negative $\beta_{1,OLS}$ estimate for the effect of D on excess value suggesting a diversification discount. Second, as reported in Table 8 and consistent with Campa and Kedia (2002), the sign on the effect of D on excess value switches when using a 2SLS model suggesting a diversification premium rather than a discount. Thus using three of the same instruments and using similar control variables as in Campa and Kedia (2002), we obtain very similar OLS and 2SLS coefficients on the key variable of interest (D) as reported in that paper despite having a different sample and using a simpler model.²⁷

²⁶To create our sample we start with the full Compustat segment database and then generally follow the sample creation criteria described in Campa and Kedia (2002) and Berger and Ofek (1995) using data starting in 1986 and extending forward through 2022. This means we eliminate segments that do not report sales information, are missing a SIC code, or that are not identified as business segments. Following these earlier papers we also eliminate any firm-year if the overall sales are less than \$20 million, if the firm reports segments in the financial sector (SIC 6000-6999), if the sum of the segment sales is more than 1% different than the overall sales reported for the firm, or if the inputs to calculate total capital are missing. Total capital is calculated as the sum of Compustat's market value of equity, long-term debt, current portion of long-term debt, and preferred stock. The preferred stock is assumed to be 0 if missing. Following these earlier papers we eliminate any firm-year where the estimated excess firm value is above 1.386 or below -1.386. Control variables are winsorized at the 1% level. If a firm appears in the main Compustat file but not in the segments database, we assume it is a single segment firm.

²⁷Campa and Kedia (2002) suggest multiple firm-level, year-level, and industry-year level instruments that are all included in their pre-first-stage probit model. In our approach, we include year and industry controls in the first stage equation and hence do not include the various year and industry-level instruments used in their paper in our first stage model. Of the remaining instruments suggested by Campa and Kedia we include (exclude) only the subset of strong instruments that have F-statistics from justidentified first stage models above (below) 10. This approach leads to 3 strong instruments (PNDIV, PSDIV, and MAJOREX). In our model we include indicator variables for which country the firm is incorporated in rather than using whether the firm is incorporated outside the US as an instrument for D.

Given the various other papers in this literature that document a discount rather than a premium, it is worth exploring whether contaminated control bias could explain this unexpected result. It is important to note that Campa and Kedia's identification strategy, and indeed the empirical approach embraced by most finance papers that utilize 2SLS models, critically requires that the other control variables included in the 2SLS system either be exogenous variables or that they at least not be correlated with the specific instrument(s) being used in conjunction with the key endogenous variable of interest. We use the contaminated control test and MPB formula proposed in this paper to (1) show that the instruments used in the Campa and Kedia (2002) model are correlated with the control variables and hence the potential for contaminated control bias in the key coefficient of interest exists in this model, (2) estimate the maximum potential size of the bias coming from the contaminated controls using the MPB formula described above, and (3) illustrate how the above tools can help a researcher explore the robustness of their 2SLS results.

In column 1 of Table 8 we report the second stage results from Equation 34 using an overidentified 2SLS model where D (the dependent variable of the first stage) has been instrumented with the PNDIV, PSDIV, and MAJOREX instruments described in Campa and Kedia (2002).²⁸ Similar to Campa and Kedia, we find that the coefficient on D is positive after using a 2SLS approach. In untabulated tests, the p-values associated with the contaminated control test statistics for each of the control variables for the model reported in column 1 of Table 8 were each less than 1% indicating that in each case we reject the null hypothesis of no contaminated control bias for each of these variables. Given the strength of the instruments in our example (F-statistics reported in Table 8), and the simulation results in this paper, the contaminated control tests should have sufficient power in this setting to correctly identify whether contaminated control bias could be affecting

²⁸Following Campa and Kedia (2002), PNDIV is intended to capture the attractiveness of a firm being diversified and is defined as the "fraction of all firms in the industry which are conglomerates" that year. PSDIV provides similar information but is defined as the "fraction of sales by other firms in the industry accounted for by diversified firms" that year. Campa and Kedia (2002) argue that firms are more likely to diversify if they are more visible to investors due to a reduction in information asymmetries and that being listed on the NYSE, Nasdaq, or AMEX exchanges would lead to this visibility. MAJOREX is a indicator variable for whether a firm is on one of these exchanges in a given year.

Table 8: Diversification Discount Empirical Example

	(1) 3 IVs 1 End	(2) MPB in β_1	(3) Just Identified Models PNDIV	(4) PSDIV	(5) MAJOREX	(6) 3 IVs 2 End	(7) 2 IVs 1 End
D	0.062** (0.029)		-0.099*** (0.001)	-0.003 (0.963)	2.062*** (0.000)	-0.177*** (0.000)	-0.101*** (0.000)
Log(Assets)	0.043*** (0.000)	0.229	0.050*** (0.000)	0.046*** (0.000)	-0.047*** (0.000)	0.198*** (0.000)	0.050*** (0.000)
CAPX/Sales	0.770*** (0.000)	0.098	0.741*** (0.000)	0.758*** (0.000)	1.121*** (0.000)	0.621*** (0.000)	0.741*** (0.000)
EBIT/Sales	0.073*** (0.000)	0.053	0.080*** (0.000)	0.076*** (0.000)	-0.008 (0.628)	-0.062*** (0.000)	0.080*** (0.000)
Leverage	0.029*** (0.002)	0.031	0.034*** (0.000)	0.031*** (0.001)	-0.029 (0.130)	-0.211*** (0.000)	0.034*** (0.000)
SP500	0.198*** (0.000)	0.099	0.196*** (0.000)	0.197*** (0.000)	0.232*** (0.000)	-0.247*** (0.000)	0.195*** (0.000)
Constant	-0.780** (0.013)		-0.775*** (0.008)	-0.778** (0.011)	-0.847 (0.190)	-1.458*** (0.000)	-0.775*** (0.008)
Year FE	Yes		Yes	Yes	Yes	Yes	Yes
Industry FE	Yes		Yes	Yes	Yes	Yes	Yes
Country FE	Yes		Yes	Yes	Yes	Yes	Yes
Observations	108,782		108,782	108,782	108,782	108,782	108,782
1st Stage F	631.982		1,715.891	300.179	149.281	740.933	859.826
Sargan χ^2	341.053					20.341	2.641
Sargan p-value	<.001					<.001	0.104

The second stage dependent variable in Equation 34 is a measure of excess firm value. Column 1 reports the second stage results from an overidentified 2SLS model that instruments the endogenous variable D using all 3 instruments (PNDIV, PSDIV, and MAJOREX). Column 2 reports the maximum possible bias (MPB) that could exist in the coefficient on D in column 1 due to the contaminated control bias coming from each control variable. Columns 3 - 5 report the 2SLS results from just identified models with the instruments listed in the column headers. Column 6 reports the 2SLS results from an overidentified model that uses the 3 instruments to instrument both D and Log(Assets). Column 7 reports the 2SLS results from an overidentified model that uses 2 instruments (PNDIV, PSDIV) to instrument D. P-values are shown below the coefficients in parenthesis. Significance is shown at the 1%, 5% and 10% levels using ***, **, and * superscripts, respectively. Industry controls are defined using 2-digit SICs. Sargan (1958) and Basmann (1960) chi-squared overidentification test results are reported for the overidentified models.

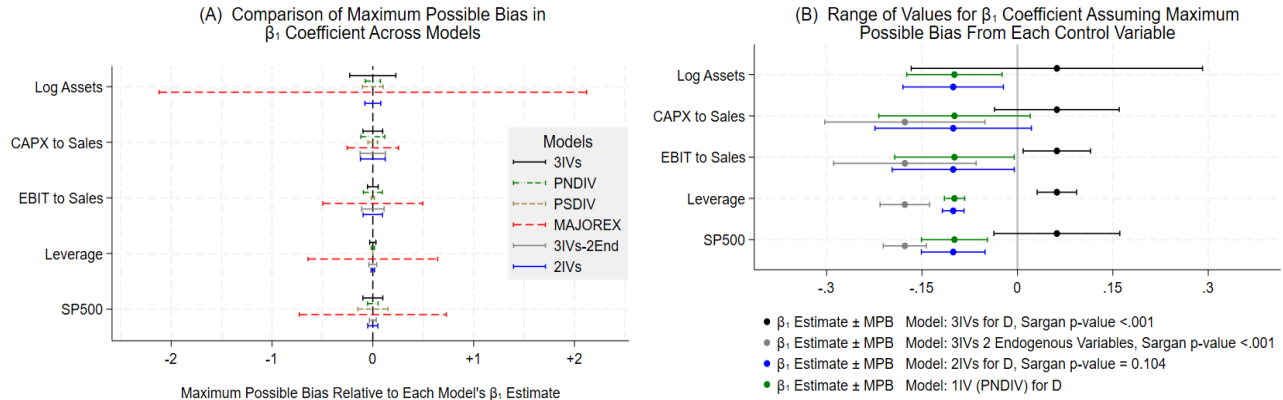
the β_1 estimate. The Sargan overidentification test p-value reported for this model in column 1 is below 1% indicating a problem with the validity of at least one of the instruments. The Sargan test result corroborates the intuition from the contaminated control tests which also suggest that

the 2SLS estimate for β_1 coefficient in column 1 is not reliable. Column 2 of Table 8 reports the maximum possible bias that could be affecting the $\beta_{1,2SLS}$ estimate in column 1 coming from each control variable. Comparing the size of the MPB for each control variable with the size of the $\beta_{1,2SLS}$ estimate helps clarify whether the bias could be large enough to change the sign on the key coefficient with a positive coefficient suggesting a diversification premium and a negative coefficient suggesting a discount. In this case the size of the MPB on several of the control variables is large enough to flip the sign on the key coefficient of interest. This is shown visually in Figure 5. In Panel A the size of the MPB from each variable is shown extending above and below 0 for the various models from Table 8. In Panel B the ranges of values that are possible in the true $\beta_{1,2SLS}$ estimate, after accounting for the MPB, are shown visually as the $\beta_{1,2SLS}$ estimate from each model \pm the MPB from each control variable. As shown, after accounting for the bias from the model in column 1 of Table 8, the true marginal effect of D on excess value could in fact be negative, zero, or positive – it is not possible to know the correct sign based on the information from the overidentified model in column 1.

So what should a researcher do in a situation like this when one or more of the control variables fails the contaminated control test and the MPB is large enough to change the sign on the key 2SLS estimate?

First, the contaminated control test derived in Section 3 identifies control variables that *may* be creating bias in the 2SLS coefficient of interest if the control variables are indeed also correlated with an omitted variable. In current practice, authors that use 2SLS already spend time discussing the (narrow) exclusion condition as it applies to the instrument(s) on the key variable of interest. If one or more of the control variables fails the contaminated control test then the researcher may still be able to use the 2SLS result if they can argue that the flagged control variable is likely to have little or no correlation with the error. In this sense, the diagnostic test developed in this paper helps researchers know which control variables need to be discussed in terms of their possible endogeneity. In our example, it would not be plausible to argue that the affected controls are

Figure 5: Maximum Possible Bias



Panel A plots the maximum possible bias coming from each of the control variables in the different models listed in the legend. In this figure if the horizontal bands are close to 0 (e.g., the narrow \pm MPB bands shown for Leverage in the 2IVs model) this means the maximum possible bias in the $\beta_{1,2SLS}$ estimate for that model coming from that specific control variable is close to 0. In contrast, if the horizontal bands are large, for example the ± 2.12 MPB bands shown in the just identified MAJOREX model for Log Assets, then the possible bias in the $\beta_{1,2SLS}$ estimate for that model could be on the order of ± 2.12 which is much larger than the size of the $\beta_{1,2SLS}$ estimate for that model. The 3IVs, PNDIV, PSDIV, MAJOREX, 3IVs-2End, and 2IVs models in the figure refer to the same models also described in header of columns 1, 3, 4, 5, 6, and 7 in Table 8, respectively. The circles plotted in the center of the bands in Panel B represent the $\beta_{1,2SLS}$ estimates for each of the respective models listed at the bottom of the panel. The bands show how different the true $\beta_{1,2SLS}$ estimate could be based on the MPB calculation for each control variable. The 3IVs, 3IVs 2 Endogenous, 2IVs, and 1IV models refer to the same models reported in columns 1, 6, 7, and 3 in Table 8, respectively.

exogenous given that assets, CAPX, earnings, leverage, and being part of the S&P500 are all plausibly related to firm valuation and are also likely co-determined with other factors not included in the model that are in the error term.

In exploring the 2SLS system, if the control variables are not plausibly exogenous then it is important to examine whether the MPB is large enough to change the sign on the main 2SLS coefficient. As noted above, the MPB is the maximum possible bias rather than the actual bias and hence represents the worst possible scenario. Based on the simulation results, it may be reasonable to assume in practice that the actual bias is likely less than half the MPB. This would suggest that a researcher may still be able to use the inference from the 2SLS result if the sign on the key variable

of interest remains unchanged after adjusting the coefficient $\pm 0.5 \times \text{MPB}$. On the other hand, there may be contaminated control bias coming from more than one control variable and together the bias may still be large enough to potentially create problems when drawing strong inference around the key variable of interest if the size of the MPB is relatively large for multiple control variables.²⁹

Another possible way forward would be to find a different instrument that is less correlated with the control variables. Or, if the 2SLS system is already overidentified, then researchers can either explore which of their instruments creates the least bias and use that instrument in a just identified model, or choose to use the various instruments to instrument not only the key endogenous variable of interest but also the control variable with the largest possible bias using two first stages in the same 2SLS system. The results reported in column 1 of Table 8 were based on a 2SLS model that was overidentified with 3 strong instruments. Having multiple instruments in this case, we can use both of the approaches suggested above to further explore our results by re-estimating the 2SLS models using 3 separate just identified models (columns 3-5 of Table 8) and by re-estimating an overidentified 2SLS model that has 2 endogenous variables and hence 2 first stage equations in one system (column 6 of Table 8). Given that the MPB was largest for the Log(Assets) control variable in the model reported in column 1, we use the same 3 instruments in column 6 to instrument both Log(Assets) and D in a new 2SLS model. As reported in column 6 of Table 8 after addressing the control variable with the largest possible bias, the $\beta_{1,2SLS}$ estimate turns negative and, as shown in Figure 5, the MPB values for this model are considerably smaller than the possible bias values related to the model in column 1 suggesting that the reason for the original positive coefficient on D was in fact due to bias.

To better understand the relation between D and excess value we also estimate 3 just identified models. As reported in columns 3 - 5 each of the 3 instruments is strong based on the F-statistic from the first stage but the estimated marginal effect of D on a firm's excess value ranges from a

²⁹The MPB calculation is done variable-by-variable assuming in each case that the other variables are not contaminated. This means that the MPB is not additive across variables.

-0.099 in column 3 to a positive 2.06 in column 5 suggesting that at least one of these instruments is invalid. The conclusion that at least one of the instruments is invalid is also supported by the small p-values from the Sargan overidentification tests reported for the overidentified models in columns 1 and 6 of Table 8. To find out which of the 3 instruments is likely invalid we compute the MPB for each of the control variables in each of the just identified models reported in columns 3 - 5 of Table 8. The MPB values for these models are shown side-by-side in Panel A of Figure 5. As shown in the figure, the MPB values associated with the control variables in the MAJOREX just identified model are much larger than the possible bias associated with the models that use the other instruments suggesting that MAJOREX is the problematic instrument. This conclusion is also supported by 3 other observations: (1) the marginal effect of D reported in column 5 of Table 8 is too large given the range of values in the dependent variable, (2) MAJOREX is the only instrument that suggests a diversification premium instead of a discount and this result contradicts many other studies, and (3) in column 7 of Table 8 we report the 2SLS results from a model that uses only PNDIV and PSDIV to instrument D and find that the Sargan overidentification test does not fail suggesting again that MAJOREX was the instrument creating validity issues in column 1. The logic and discussion above suggests the 2SLS models in columns 3, 6, and 7 of Table 8 would be better for finding the marginal effect of diversification on excess firm value compared to the other 2SLS models in that these models exhibit less contaminated control bias. Panel B of Figure 5 presents this same information visually showing that (1) the $\beta_{1,2SLS}$ estimate in each of these models is negative, and (2) the MPB values in these models are too small to flip the sign positive. From this analysis we conclude that the true marginal effect of D on excess value is likely negative, or possibly zero, but not positive suggesting that a diversification discount is possible but not a premium.

7 Conclusion

Identifying a causal relationship between variables is often difficult given the many unobservable factors that relate to most financial topics. In recent years, many researchers have used instrumental variables in 2SLS settings to deal with the endogeneity. A survey of the use of 2SLS in papers at the *Journal of Finance*, the *Journal of Financial Economics*, and the *Review of Financial Studies* indicates that literally hundreds of papers have utilized 2SLS as part of their analysis in recent years and that most of them provide minimal or no discussion of the potential endogeneity of the control variables included in the model and no consideration for whether their instruments may be correlated with those control variables. Indeed, recent standard practice tends to include a discussion of the relevancy and exclusion conditions for a given instrument insofar as these conditions relate specifically to the key variable of interest and then to include, as though exogenous, an assortment of other control variables that may themselves also be endogenous. Many of these papers simply assert or assume that the control variables are exogenous.

Yet, despite these assertions and the general lack of discussion around the potential endogeneity of the control variables, it is likely that most of these empirical settings and models have at least weakly endogenous control variables. Our paper shows both analytically and via simulation that ignoring the low-level correlations that can exist between the control variables and the error term can have a direct and strong effect on the researcher's ability to draw inference from the 2SLS results if the control variable(s) are also correlated with the instrument for the key variable of interest.

Along these lines, our paper provides guidance related to the following questions: First, what effect does the inclusion of potentially endogenous control variables have on the 2SLS estimate for the key variable of interest given a strong instrument for the key variable that itself is plausibly not correlated with the error term? Answer: Including endogenous control variables can generate large bias in the 2SLS estimate of interest even if the instrument for that key variable is strong

and is itself not directly correlated with the error term. The bias that comes from the inclusion of endogenous control variables only affects the 2SLS estimate on the key variable of interest if the control variables are both endogenous and correlated with the instrument for the key variable of interest. Contaminated control variable bias is exacerbated by weak instrument(s) for the key variable of interest. The contaminated control test and the MPB formula introduced in this paper can help researchers assess whether the size of the bias is likely large enough to affect the inference for the main variable of interest.

Second, is the bias in the 2SLS estimate for the key variable of interest made larger or smaller with or without including the other potentially endogenous control variables in the system of equations? Answer: It is not possible to say whether the bias in the 2SLS estimate will increase or decrease when dropping the endogenous control variables from the system. In some settings the bias increases whereas in others it decreases. However, based on the analytical form of the bias discussed in Section 2, if both (1) the narrow exclusion condition holds (or is almost satisfied) for the instrument on the key variable of interest and (2) the difference described in this paper, $(\hat{\gamma}_2 - \hat{\lambda}_2)$, is close to zero then the overall bias is likely smaller in the 2SLS estimate for the key variable of interest with the controls included than in the estimate without controls.

Third, what information can be inferred from estimating the 2SLS estimate both with and without the control variables and then comparing the estimates? Answer: Dropping an important variable (endogenous or not) from the system creates the potential for omitted variable bias in the key estimate if the dropped variable is correlated with the other variables and the instrument(s). Thus estimating the 2SLS estimate both with and without control variables is trading off potential bias from the inclusion of endogenous control variables that are correlated with the instrument against overall omitted variable bias. Hence, observing a large change in the 2SLS estimate when comparing the key result with and without the control variables need not indicate that the 2SLS estimate with controls is biased given that the change could be attributable to omitted variable bias created when dropping the control variables. However, observing little or no change in the 2SLS

estimate both with and without the control variables provides some corroborating evidence that the 2SLS estimate is unlikely to be largely affected by bias from the control variables.

And, fourth, is there a test that would reveal whether the 2SLS estimate for the key variable of interest could be affected by contaminated controls? Answer: Yes. In order for the inclusion of other control variables to affect the key 2SLS estimate, two conditions need to occur: (1) the control variable(s) must be endogenous, and (2) the control variable(s) must be correlated with the instrument for the key variable of interest. It is not possible to ascertain the first condition but the second is testable. We propose testing whether $(\hat{\gamma}_2 - \hat{\lambda}_2)$, as discussed in Section 3, is statistically different from zero when investigating this bias. Failing to reject the null hypothesis of this test statistic being equal to 0 supports the conclusion that the control variables are not creating material bias in the key variable of interest. This test has the advantage of being able to rule out the presence of contaminated control bias but is limited in that it rules out a necessary but not a sufficient condition for this type of bias which means that failing this test indicates that contaminated control bias may exist in the key coefficient of interest, not that it necessarily does.

As highlighted in Section 6 with the diversification discount example, this paper suggests an approach for researchers using 2SLS to first test for the possibility of contaminated control bias in specific control variables and then to calculate the maximum possible bias in the key coefficient coming from the flagged variables. Using the new contaminated control test together with the proposed MPB calculations will allow researchers in the future to examine whether their 2SLS results are robust or whether contaminated control bias may be affecting the results in a material way. As noted above, hundreds of recent papers in top finance journals have previously simply assumed that the control variables in their 2SLS systems are not biasing their key 2SLS result, but this assumption is unlikely true in many if not most cases. The tools and approaches proposed in this paper will allow researchers to examine this assumption in detail going forward.

The discussion in Section 6 also provides practical advice for how to explore 2SLS results in the event that one or more control variables are flagged by the contaminated control test. If specific

control variables fail the proposed test, and the MPB on those variables is relatively large compared to the key coefficient of interest, the researcher can either instrument those variables or explain why they are unlikely to be endogenous using arguments similar to the literature's current approach when motivating the narrow exclusion condition for the instrument on the key variable of interest. The proposed test and MPB calculations will help researchers understand which control variables may need additional discussion. One of the advantages of using this approach in an overidentified model is that the researcher can check whether specific instruments worsen the potential bias. Researchers often use the Sargan overidentification test to identify validity issues in 2SLS models, but this test is only possible in overidentified models and does not identify which of the instruments is invalid. The test and MPB calculations proposed in this paper provide an alternative way to identify potential problems with validity and have the additional benefits of not only allowing researchers to assess which of several instruments is likely causing the problems in the 2SLS system but also by being available for use with just identified in addition to overidentified models.

In summary, the diagnostic test derived in this paper and the related MPB calculations will help researchers know when they need to explore their 2SLS system of equations in more detail, know which control variables need specific consideration, be able to choose between specific instruments to minimize bias if the different instruments are suggesting different inferences, and be able to assess the robustness of their main 2SLS results and whether the inference is likely affected by contaminated control bias.

References

- Altonji, J. G., Elder, T. E., and Taber, C. R. (2005). Selection on observed and unobserved variables: Assessing the effectiveness of Catholic schools. *Journal of Political Economy*, 113(1):151–184.
- Andreou, P. C., Doukas, J. A., Koursaros, D., and Louca, C. (2019). Valuation effects of overconfident CEOs on corporate diversification and refocusing decisions. *Journal of Banking and Finance*, 100:182–204.
- Angrist, J. and Krueger, A. B. (2001). Instrumental variables and the search for identification: From supply and demand to natural experiments. *Journal of Economic Perspectives*, pages 69–85.
- Angrist, J. and Pischke, J.-S. (2009). *Mostly Harmless Econometrics*. Princeton University Press.
- Basmann, R. L. (1960). On finite sample distributions of generalized classical linear identifiability test statistics. *Journal of the American Statistical Association*, 55(292):650–659.
- Bazdresch, S., Kahn, R. J., and Whited, T. M. (2018). Estimating and testing dynamic corporate finance models. *The Review of Financial Studies*, 31(1):322–361.
- Berg, T., Reisinger, M., and Streitz, D. (2021). Spillover effects in empirical corporate finance. *Journal of Financial Economics*, 142(3):1109–1127.
- Berger, P. G. and Ofek, E. (1995). Diversification’s effect on firm value. *Journal of Financial Economics*, 37:39–65.
- 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.

- Campa, J. M. and Kedia, S. (2002). Explaining the diversification discount. *The Journal of Finance*, 57:1731–1762.
- Chen, K. and il Kim, K. (2025). Identification of partial effects with endogenous controls. *Working Paper*.
- Cinelli, C. and Hazlett, C. (2025). An omitted variable bias framework for sensitivity analysis of instrumental variables. *Biometrika*, page asaf004.
- Davidson, R. and MacKinnon, J. G. (2004). *Econometric Theory and Methods*. Oxford University Press.
- Denis, D. J., Denis, D. K., and Sarin, A. (1997). Agency problems, equity ownership, and corporate diversification. *The Journal of Finance*, 52:135.
- Ellis, J. A., Fee, C. E., and Thomas, S. (2018). Playing favorites? industry expert directors in diversified firms. *Journal of Financial and Quantitative Analysis*, 53:1679–1714.
- Erickson, T. and Whited, T. M. (2012). Treating measurement error in Tobin’s q . *The Review of Financial Studies*, 25(4):1286–1329.
- Filoso, V. (2013). Regression anatomy, revealed. *The Stata Journal*, 13(1):92–106.
- Goldberger, A. (1991). *A Course in Econometrics*. Harvard University Press.
- Gormley, T. A. and Matsa, D. A. (2014). Common errors: How to (and not to) control for unobserved heterogeneity. *The Review of Financial Studies*, 27(2):617–661.
- Graham, M. L. and Wolf, J. (2002). Does corporate diversification destroy value. *The Journal of Finance*.
- Greene, W. H. (2003). *Econometric Analysis*. Prentice Hall.

- Grieser, W. D. and Hadlock, C. J. (2019). Panel-data estimation in finance: Testable assumptions and parameter (in) consistency. *Journal of Financial and Quantitative Analysis*, 54(1):1–29.
- Hoechle, D., Schmid, M., Walter, I., and Yermack, D. (2012). How much of the diversification discount can be explained by poor corporate governance? *Journal of Financial Economics*, 103:41–60.
- Huber, K. (2023). Estimating general equilibrium spillovers of large-scale shocks. *The Review of Financial Studies*, 36(4):1548–1584.
- Jiang, W. (2017). Have instrumental variables brought us closer to the truth. *Review of Corporate Finance Studies*, 6(2):127–140.
- Lovell, M. H. (1963). Seasonal adjustment of economic time series and multiple regression analysis. *Journal of the American Statistical Association*, 58(304):993–1010.
- Morck, R., Shleifer, A., and Vishny, R. W. (1990). Do managerial objectives drive bad acquisitions? *The Journal of Finance*, 45:31.
- Murray, M. P. (2006). Avoiding invalid instruments and coping with weak instruments. *Journal of Economic Perspectives*, 20(4):111–132.
- Oster, E. (2019). Unobservable selection and coefficient stability: Theory and evidence. *Journal of Business & Economic Statistics*, 37(2):187–204.
- Ozbas, O. and Scharfstein, D. S. (2009). Evidence on the dark side of internal capital markets. *Review of Financial Studies*, 23:581–599.
- Petersen, M. A. (2008). Estimating standard errors in finance panel data sets: Comparing approaches. *The Review of Financial Studies*, 22(1):435–480.

- Rajan, R., Servaes, H., and Zingales, L. (2000). The cost of diversity: The diversification discount and inefficient investment. *The Journal of Finance*, 55:35–80.
- Roberts, M. R. and Whited, T. M. (2013). Endogeneity in empirical corporate finance. *Handbook of the Economics of Finance*, 2:493–572.
- Sargan, J. D. (1958). The estimation of economic relationships using instrumental variables. *Econometrica: Journal of the Econometric Society*, pages 393–415.
- Shin, H.-H. and Stulz, R. M. (1998). Are internal capital markets efficient? *The Quarterly Journal of Economics*, 113:531–552.
- Thompson, S. B. (2011). Simple formulas for standard errors that cluster by both firm and time. *Journal of Financial Economics*, 99(1):1–10.
- Wooldridge, J. (2003). *Introductory Econometrics: A Modern Approach*, 2Ed. South-Western/Thomson Learning.
- Wooldridge, J. M. (2002). *Econometric Analysis of Cross Section and Panel Data*. MIT Press.
- Zellner, A. (1962). An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias. *Journal of the American statistical Association*, pages 348–368.

A Bias in 2SLS Estimates - Additional Results

We show that $\hat{\gamma}_2 = \hat{\lambda}_2$ if $\text{cov}(z, x_2) = 0$.

$$\begin{aligned}
 \hat{\lambda}_2 &= \frac{\text{cov}(\hat{x}_1, x_2)}{\text{var}(x_2)} \\
 &= \frac{\text{cov}(\hat{\gamma}_0 + \hat{\gamma}_1 z + \hat{\gamma}_2 x_2, x_2)}{\text{var}(x_2)} \\
 &= \hat{\gamma}_1 \frac{\text{cov}(z, x_2)}{\text{var}(x_2)} + \hat{\gamma}_2 \\
 \hat{\gamma}_2 - \hat{\lambda}_2 &= -\hat{\gamma}_1 \frac{\text{cov}(z, x_2)}{\text{var}(x_2)}
 \end{aligned} \tag{35}$$

Let σ_z and σ_{x_2} denote the standard deviations of z and x_2 respectively, and ρ_{z, x_2} denote the correlation between z and x_2 . We show that $\text{var}(\hat{x}_1^*)$ is increasing in instrument strength $\hat{\gamma}_1^2$.

$$\begin{aligned}
 \text{var}(\hat{x}_1^*) &= \text{var}(\hat{\gamma}_0 - \hat{\lambda}_1 + \hat{\gamma}_1 z + (\hat{\gamma}_2 - \hat{\lambda}_2)x_2) \\
 &= \hat{\gamma}_1^2 \text{var}(z) + (\hat{\gamma}_2 - \hat{\lambda}_2)^2 \text{var}(x_2) + 2\hat{\gamma}_1 (\hat{\gamma}_2 - \hat{\lambda}_2) \text{cov}(z, x_2) \\
 &= \hat{\gamma}_1^2 \text{var}(z) + \left(-\hat{\gamma}_1 \frac{\text{cov}(z, x_2)}{\text{var}(x_2)}\right)^2 \text{var}(x_2) + 2\hat{\gamma}_1 \left(-\hat{\gamma}_1 \frac{\text{cov}(z, x_2)}{\text{var}(x_2)}\right) \text{cov}(z, x_2) \\
 &= \hat{\gamma}_1^2 \text{var}(z) - \hat{\gamma}_1^2 \frac{(\text{cov}(z, x_2))^2}{\text{var}(x_2)} \\
 &= \hat{\gamma}_1^2 \text{var}(z) - \hat{\gamma}_1^2 \frac{(\sigma_z \sigma_{x_2} \rho_{z, x_2})^2}{\text{var}(x_2)} \\
 &= \hat{\gamma}_1^2 \text{var}(z) (1 - \rho_{z, x_2}^2)
 \end{aligned} \tag{36}$$

B MPB with Multiple Instrumented Endogenous Variables

Let \mathbf{x}_3 denote an endogenous (G by 1) vector of control variables that are instrumented with a vector of instruments. $\mathbf{x}_3 \equiv (x_{31}, \dots, x_{3G})'$.³⁰ Note that $K \geq G + 1$. Suppose the DGP is given by:

$$y = \beta_0 + \beta_1 x_1 + \beta_2' \mathbf{x}_2 + \beta_3' \mathbf{x}_3 + w \quad (37)$$

The first stage regressions and estimates for x_1 are given by:

$$\begin{aligned} x_1 &= \gamma_0 + \gamma_1' z + \gamma_2' \mathbf{x}_2 + e_1 \\ \hat{x}_1 &= \hat{\gamma}_0 + \hat{\gamma}_1' z + \hat{\gamma}_2' \mathbf{x}_2 \end{aligned} \quad (38)$$

Let δ_0 be a (G by 1) vector, δ_1 be a (G by K) matrix, and δ_2 be a (G by J) matrix, of parameters. e_3 is a (G by 1) vector of errors. The first stage regressions and estimates for \mathbf{x}_3 are given by:

$$\begin{aligned} \mathbf{x}_3 &= \delta_0 + \delta_1 z + \delta_2 \mathbf{x}_2 + e_3 \\ \hat{\mathbf{x}}_3 &= \hat{\delta}_0 + \hat{\delta}_1 z + \hat{\delta}_2 \mathbf{x}_2 \end{aligned} \quad (39)$$

Let λ_3 be a (G by 1) vector of parameters. \hat{x}_1^* is the portion of \hat{x}_1 uncorrelated with \mathbf{x}_2 and $\hat{\mathbf{x}}_3$.

$$\hat{x}_1 = \lambda_1 + \lambda_2' \mathbf{x}_2 + \lambda_3' \hat{\mathbf{x}}_3 + \xi = \hat{\lambda}_1 + \hat{\lambda}_2' \mathbf{x}_2 + \hat{\lambda}_3' \hat{\mathbf{x}}_3 + \hat{x}_1^* \quad (40)$$

Equating equations 38 and 40, we solve for \hat{x}_1^* .

$$\begin{aligned} \hat{\gamma}_0 + \hat{\gamma}_1' z + \hat{\gamma}_2' \mathbf{x}_2 &= \hat{\lambda}_1 + \hat{\lambda}_2' \mathbf{x}_2 + \hat{\lambda}_3' \hat{\mathbf{x}}_3 + \hat{x}_1^* \\ \hat{x}_1^* &= (\hat{\gamma}_0 - \hat{\lambda}_1) + \hat{\gamma}_1' z + (\hat{\gamma}_2 - \hat{\lambda}_2)' \mathbf{x}_2 - \hat{\lambda}_3' \hat{\mathbf{x}}_3 \end{aligned} \quad (41)$$

³⁰The remaining variables have identical definitions as in previous sections.

Substituting using Equation 39, we solve for \hat{x}_1^* as a function of \mathbf{z} , \mathbf{x}_2 , and the parameters.

$$\begin{aligned}\hat{x}_1^* &= (\hat{\gamma}_0 - \hat{\lambda}_1) + \hat{\gamma}_1' \mathbf{z} + (\hat{\gamma}_2 - \hat{\lambda}_2)' \mathbf{x}_2 - \hat{\lambda}_3' (\hat{\delta}_0 + \hat{\delta}_1 \mathbf{z} + \hat{\delta}_2 \mathbf{x}_2) \\ &= (\hat{\gamma}_0 - \hat{\lambda}_1 - \hat{\delta}_0' \hat{\lambda}_3) + (\hat{\gamma}_1 - \hat{\delta}_1' \hat{\lambda}_3)' \mathbf{z} + (\hat{\gamma}_2 - \hat{\lambda}_2 - \hat{\delta}_2' \hat{\lambda}_3)' \mathbf{x}_2\end{aligned}\quad (42)$$

The second stage estimates and expression for the bias are given by:

$$\begin{aligned}\hat{\beta}_{1,2SLS} &= \beta_1 + \underbrace{\frac{\text{cov}(\hat{x}_1^*, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias}} \\ &= \beta_1 + \underbrace{(\hat{\gamma}_1 - \hat{\delta}_1' \hat{\lambda}_3)' \frac{\text{cov}(\mathbf{z}, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias related to the narrow exclusion condition}} + \underbrace{(\hat{\gamma}_2 - \hat{\lambda}_2 - \hat{\delta}_2' \hat{\lambda}_3)' \frac{\text{cov}(\mathbf{x}_2, w)}{\text{var}(\hat{x}_1^*)}}_{\text{bias related to endogenous controls}}\end{aligned}\quad (43)$$

A test for contaminated controls would focus on the quantity $(\hat{\gamma}_2 - \hat{\lambda}_2 - \hat{\delta}_2' \hat{\lambda}_3)$. Unlike the case of a single instrumented variable, this quantity is nonlinear in the parameters, and would require a different and more complex test. We leave this as an area of future research.

Last, we derive the maximum possible bias. Let $\hat{\delta}_{2j}$ denote the j^{th} column of $\hat{\delta}_2$. $\tilde{\gamma}_{2j}$ is now a $[(J + G + 1) \text{ by } 1]$ vector with $\hat{\gamma}_{2j}$ the first element, 1 the $(j + 1)^{th}$ element, $\hat{\delta}_{2j}$ the last G elements, and 0 the remaining elements. Generalizing equations 26 and 33, the maximum possible bias is given by:

$$\begin{aligned}\hat{\beta}_{1,2SLS} - \beta_1 &= (\hat{\gamma}_2 - \hat{\lambda}_2 - \hat{\delta}_2' \hat{\lambda}_3)' \frac{\text{cov}(\mathbf{x}_2, w)}{\text{var}(\hat{x}_1^*)} \\ \hat{\beta}_{1,2SLS} - \beta_1 &\leq \sum_{j=1}^J \frac{(\hat{\gamma}_{2j} - \hat{\lambda}_{2j} - \hat{\delta}_{2j}' \hat{\lambda}_3) \sigma_j \sigma_y \rho_j}{\text{var}(\hat{x}_1^*)} \\ MPB_j &\equiv \left| \frac{(\hat{\gamma}_{2j} - \hat{\lambda}_{2j} - \hat{\delta}_{2j}' \hat{\lambda}_3) \sqrt{\sigma_y^2 - \hat{\sigma}_v^2}}{\text{var}(\hat{x}_1^*) \tilde{x}_j} \right|\end{aligned}\quad (44)$$