



When is it justifiable to ignore explanatory variable endogeneity in a regression model?



Richard A. Ashley^a, Christopher F. Parmeter^{b,*}

^a Department of Economics, Virginia Polytechnic Institute and State University, United States

^b Department of Economics, University of Miami, United States

HIGHLIGHTS

- Sound empirical practice to question plausibility of exogeneity.
- Develop sensitivity analysis to check robustness of inference.
- Apply method to classic economic growth setting.
- Find many key hypotheses are robust to exogeneity violations.

ARTICLE INFO

Article history:

Received 1 July 2015

Received in revised form

13 September 2015

Accepted 23 September 2015

Available online 9 October 2015

JEL classification:

C2

C15

Keywords:

Robustness

Exogeneity

Instruments

ABSTRACT

The point of empirical work is commonly to test a very small number of crucial null hypotheses in a linear multiple regression setting. Endogeneity in one or more model explanatory variables is well known to invalidate such testing using OLS estimation. But attempting to identify credibly valid (and usefully strong) instruments for such variables is an enterprise which is arguably fraught and invariably subject to (often justified) criticism. As a modeling step prior to such an attempt at instrument identification, we propose a sensitivity analysis which quantifies the minimum degree of correlation between these possibly-endogenous explanatory variables and the model errors which is sufficient to overturn the rejection (or non-rejection) of a particular null hypothesis at, for example, the 5% level. An application to a classic model in the empirical growth literature illustrates the practical utility of the technique.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

No issue in econometrics has evoked as much literature (and angst) over the years as that of the likely endogeneity of explanatory variables in our regression models. The profession's main response – instrumental variables (IV) regression – has ameliorated this concern in some settings, albeit at the cost of a decrease in estimation precision. In large part, however, the use of IV regression has simply shifted the focus of attention to the exogeneity of the instruments, spawning a search for ‘clever’ instruments whose exogeneity can be argued—e.g., see Angrist and Krueger (1991) and Acemoglu et al. (2001). For evidence that instrument validity is a continuing concern, see Angrist and Pischke (2010), Keane (2010),

Leamer (2010), Murray (2006), Sims (2010) and Stock (2010). Even more recently, Bazzi and Clemens (2013) have strongly criticized the way IV is applied in growth regressions. In the present paper we suggest a different approach.

Applied economic analysis almost always culminates in the rejection of (or, occasionally, in the failure to reject) a very small number of crucial null hypotheses at some nominal level of significance, usually 5%. For any particular one of these hypothesis tests, this translates into a rejection p -value of less than or equal to 0.05. The endogeneity issue then becomes: is the reported rejection of the null hypothesis actually just an artifact of unaccounted for (or improperly accounted for) endogeneity in the supposedly exogenous explanatory variables?

But suppose that it was possible to determine that any likely degree of such endogeneity was insufficient to overturn our small set of key inferences? We could then base our analysis on OLS regression without having to expend resources (or credibility) on finding plausible instruments or on worrying about their validity.

* Corresponding author.

E-mail addresses: ashleyr@vt.edu (R.A. Ashley), cparmeter@bus.miami.edu (C.F. Parmeter).

Ashley and Parmeter (forthcoming) provides an empirically implementable algorithm for performing exactly this kind of sensitivity analysis with respect to the validity (exogeneity) of instruments used in linear GMM regression modeling. Where the algorithm finds that a key hypothesis test rejection is overturned by very small amounts of correlation between the instruments and the (unobserved) model errors, this inference is deemed ‘fragile’. In contrast, where it is found that quite substantial levels of instrument-error correlation – e.g., in excess of 0.50 in magnitude – are necessary in order to overturn this hypothesis test rejection, then this inference is deemed ‘robust’.¹ Clearly, inference with respect to some null hypotheses may be fragile whereas others are robust, even within the same regression model.

Here we observe that OLS regression is equivalent to letting regressors act as instruments for themselves and apply the Ashley/Parmeter algorithm to the underlying model estimated via OLS. As an illustrative example, in the next section we analyze the impact of explanatory variable endogeneity on the inferential conclusions obtained in Mankiw et al. (1992), a foundational paper in the economic growth literature an area that is routinely criticized for endogeneity.

2. A sensitivity analysis for exogeneity

2.1. Estimation/inference in the presence of unaddressed endogeneity

Consider the standard linear model, with the ‘structural equation’

$$Y_1 = Y_2\alpha + W_1\beta + \varepsilon, \tag{1}$$

where Y_2 is an $n \times m$ matrix of (potentially) endogenous variables, W_1 is an $n \times k$ matrix of variables whose exogeneity is not in question, α and β are $m \times 1$ and $k \times 1$ vectors of coefficients, respectively, and ε is the structural error. For the present purpose we do not assume the presence of additional (instrumental) variables to correct for the endogeneity of Y_2 .²

Accordingly, the moment conditions assumed here are:

$$\begin{aligned} E[Y_2'\varepsilon_i] &= 0 \\ E[W_1'\varepsilon_i] &= 0. \end{aligned} \tag{2}$$

The first of the two conditions in Eq. (2) incorporates the assumed exogeneity of the m potentially endogenous variables in Y_2 ; the second condition reflects the assumption that the remaining k variables are clearly exogenous. Thus, Y_2 and W_1 are defined in such a way that we need only concern ourselves with violations of exogeneity for the m variables in Y_2 .

Letting $\gamma = [\alpha'\beta']'$ and $X = [Y_2W_1]$, the structural equation (1) can be written more compactly as:

$$Y_1 = X\gamma + \varepsilon. \tag{3}$$

The OLS estimator of γ is thus:

$$\hat{\gamma}_{OLS} = (X'X)^{-1}X'Y_1. \tag{4}$$

¹ Where, as is common, the validity of multiple instruments is in question, the algorithm also provides a sensible indication as to which of the instruments are the source of any fragility found. R code implementing the algorithm is available from the authors.

² As noted above, see Ashley and Parmeter (forthcoming) for a related treatment explicitly allowing the use of (possibly flawed) instruments in 2SLS/GMM estimation; here the focus is on OLS estimation in the absence of credibly valid instruments.

This estimator is consistent, asymptotically efficient and asymptotically normal under the standard assumptions, including (at least asymptotic) exogeneity of **all** the regressors X .³

When some or all of the variables in X are *not* exogenous, then

$$E[X_i'\varepsilon_i] = n\Sigma_{X\varepsilon} \neq 0. \tag{5}$$

The factor n is introduced here so that $\Sigma_{X\varepsilon}$ can be interpreted as the population covariance vector between the structural error ε and the $g + k$ supposedly exogenous variables; $\Sigma_{X\varepsilon}$ can thus sensibly be referred to as “the exogeneity flaw covariance vector”.

For a given value of the exogeneity flaw covariance vector, then it is easy to show that the modified estimator of γ ,

$$\tilde{\gamma} = (X'X)^{-1}(X'Y_1 - n\Sigma_{X\varepsilon}) \tag{6}$$

is now consistent, asymptotically normal, and asymptotically efficient and that (conditional on the ‘flaw’ vector, $\Sigma_{X\varepsilon}$) $\tilde{\gamma}$ has asymptotic sampling distribution,

$$\sqrt{n}(\tilde{\gamma} - \gamma - E_{XX}^{-1}\Sigma_{X\varepsilon}) \sim N(0, \sigma_\varepsilon^2 E_{XX}^{-1}), \tag{7}$$

where $E_{XX}^{-1} = \text{plim } n^{-1}(X'X)^{-1}$. Thus, obtaining an asymptotically valid p -value at which any particular null hypothesis regarding γ could be rejected would be straightforward if $\Sigma_{X\varepsilon}$ were known.⁴

2.2. Quantifying the sensitivity of inference to endogeneity in Y_2

Suppose that a particular null hypothesis regarding γ can be rejected – using $\hat{\gamma}_{OLS}$ and its asymptotic sampling distribution, under the assumption that all of the variables in Y_2 are exogenous – at, say, the 5% level.⁵ This is equivalent to saying that the rejection p -value for this null hypothesis is less than 0.05 using the asymptotic sampling distribution of $\tilde{\gamma}$ given in Eq. (7) with the exogeneity flaw covariance vector ($\Sigma_{X\varepsilon}$) set equal to zero.

The key issue is how sensitive this 5% rejection of the null hypothesis is to values of $\Sigma_{X\varepsilon}$ which are non-zero, but “plausible”. It is straightforward to re-calculate this rejection p -value for alternative values of $\Sigma_{X\varepsilon}$, but difficult to have any intuition as to how large such a covariance is likely. In contrast, one might well have some intuition as to how large plausible values of the components of the concomitant *correlation* vector are likely to be. Thus, the crucial issue in a useful sensitivity analysis is to numerically characterize this rejection p -value as a function of the first m of these correlations, which are considered to be possibly non-zero.

Converting the covariance vector $\Sigma_{X\varepsilon}$ into the corresponding correlation vector merely involves dividing each of its components by the square root of the product of the variance of ε and the variance of the explanatory variable corresponding to this $\Sigma_{X\varepsilon}$ component. Since the columns of $X = [Y_2W_1]$ are observed, it is straightforward to consistently estimate the variance of the corresponding m explanatory variables (Y_2) considered to be possibly endogenous. The model errors, ε , in contrast, are not observed. But

³ These standard assumptions also include that of a homoscedastic and non-autocorrelated error term, a correct specification of the conditional mean of Y_1 , and full rank of the covariate matrix X .

⁴ Note that our sensitivity analysis remains squarely within the usual (asymptotic) inference framework for OLS/2SLS/GMM estimation and inference; see Du-four (2003) with regard to finite-sample alternatives.

⁵ The analysis would be essentially identical for a rejection at the 1% (or any other) level: the description in this section is made definite for the 5% level solely to enhance the clarity of the exposition. Similarly, the procedure described below can be readily modified to instead analyze the case where the null hypothesis is *not* rejected at the 5% level and the issue is whether this *failure* to reject is due to a flaw in the exogeneity of one of the g variables in question.

Download English Version:

<https://daneshyari.com/en/article/5058659>

Download Persian Version:

<https://daneshyari.com/article/5058659>

[Daneshyari.com](https://daneshyari.com)