

Supplementary material for the paper:

Instrument approval by the Sargan test and its consequences for coefficient estimation

by Jan F. Kiviet and Sebastian Kripfganz

22 May 2021

Simulation design and invariance properties

The variables $\{u, x, z_1, z_2\}$ introduced in Section 2 of the paper are obtained as linear transformations of four mutually independent series ε_i , ξ_i , ζ_{i1} and ζ_{i2} (for $i = 1, \dots, n$), where all elements are $iid(0, 1)$ drawings from a distribution we chose to be normal. The four transformations are given by

$$u_i = \sigma_u \varepsilon_i \sim iid(0, \sigma_u^2), \quad (\text{S.1})$$

$$x_i = \sigma_x [(1 - \rho_{xu}^2)^{1/2} \xi_i + \rho_{xu} \varepsilon_i] \sim iid(0, \sigma_x^2), \quad (\text{S.2})$$

$$z_{ij} = \sigma_{z_j} (\rho_{z_j \zeta_j} \zeta_{ji} + \rho_{z_j \xi} \xi_i + \rho_{z_j u} \varepsilon_i) \sim iid(0, \sigma_{z_j}^2) \text{ for } j = 1, 2, \quad (\text{S.3})$$

where all ρ coefficients do not exceed 1 in absolute value; moreover,

$$\rho_{z_j \zeta_j}^2 + \rho_{z_j \xi}^2 + \rho_{z_j u}^2 = 1 \text{ for } j = 1, 2. \quad (\text{S.4})$$

Obviously, $\sigma_{xu} = \rho_{xu} \sigma_x \sigma_u$, $\sigma_{z_j u} = \rho_{z_j u} \sigma_{z_j} \sigma_u$ and $\sigma_{z_j x} = \sigma_{z_j} \sigma_x [\rho_{z_j \xi} (1 - \rho_{xu}^2)^{1/2} + \rho_{z_j u} \rho_{xu}]$, hence $\rho_{z_j x} = \rho_{z_j \xi} (1 - \rho_{xu}^2)^{1/2} + \rho_{z_j u} \rho_{xu}$, which yields

$$\rho_{z_j \xi} = (\rho_{z_j x} - \rho_{z_j u} \rho_{xu}) (1 - \rho_{xu}^2)^{-1/2}, \quad (\text{S.5})$$

for $|\rho_{xu}| < 1$. From (S.4) we also have

$$\rho_{z_j \zeta_j} = (1 - \rho_{z_j \xi}^2 - \rho_{z_j u}^2)^{1/2}. \quad (\text{S.6})$$

Hence, when values for $\sigma_u > 0$, $\sigma_x > 0$, $\sigma_{z_j} > 0$, $|\rho_{xu}| < 1$, $|\rho_{z_j x}| \leq 1$ and $|\rho_{z_j u}| \leq 1$ are chosen, we can generate series for u_i and x_i , and find from (S.5) and (S.6) matching values for $\rho_{z_j \xi}$ and $\rho_{z_j \zeta_j}$, which enable to generate the series z_{i1} and z_{i2} as well.

However, values for ρ_{xu} , $\rho_{z_j x}$ and $\rho_{z_j u}$ are only compatible if they yield $|\rho_{z_j \xi}| \leq 1$ and $|\rho_{z_j \zeta_j}| \leq 1$. This requires

$$(\rho_{z_j x} - \rho_{z_j u} \rho_{xu})^2 (1 - \rho_{xu}^2)^{-1} \leq 1 \quad (\text{S.7})$$

and $0 \leq 1 - \rho_{z_j\xi}^2 - \rho_{z_ju}^2 \leq 1$, which – substituting (S.5) – boils down to

$$0 \leq (1 - \rho_{xu}^2)(1 - \rho_{z_ju}^2) - (\rho_{z_jx} - \rho_{z_ju}\rho_{xu})^2 \leq 1 - \rho_{xu}^2. \quad (\text{S.8})$$

From the two inequalities expressed by (S.8), it is obvious that the one yielding $-\rho_{z_ju}^2(1 - \rho_{xu}^2) - (\rho_{z_jx} - \rho_{z_ju}\rho_{xu})^2 \leq 0$ will always be met, while satisfying the other one, being

$$(\rho_{z_jx} - \rho_{z_ju}\rho_{xu})^2 \leq (1 - \rho_{xu}^2)(1 - \rho_{z_ju}^2), \quad (\text{S.9})$$

would imply (S.7). Hence, by choosing values $|\rho_{xu}| < 1$, and for $j = 1, 2$ values $|\rho_{z_jx}| \leq 1$ and $|\rho_{z_ju}| \leq 1$ obeying (S.9), we can find admissible $\rho_{z_j\xi}$ and $\rho_{z_j\zeta_j}$ values to generate u_i , x_i and z_{ij} accordingly.

From each realization in the 100,001 simulation replications of the series u_i , x_i and z_{ij} ($i = 1, \dots, n$; $j = 1, 2$), we may first subtract the respective sample average, next generate $y_i = \beta x_i + u_i$, skipping the intercept. In that way, just regressing y on x , the slope of a model with one regressor x and an arbitrary intercept can be estimated by

$$\begin{aligned} \hat{\beta}_{OLS} &= x'y/x'x = \beta + x'u/x'x, \\ \hat{\beta}_{IV}^{(j)} &= z_j'y/z_j'x = \beta + z_j'u/z_j'x, \quad j = 1, 2, \\ \hat{\beta}_{TSLS} &= x'Z(Z'Z)^{-1}Z'y/x'Z(Z'Z)^{-1}Z'x = \beta + x'P_Zu/x'P_Zx, \end{aligned}$$

where x , y , u , z_1 and z_2 are column vectors stacking the n sample observations, $Z = (z_1, z_2)$ and $P_Z = Z(Z'Z)^{-1}Z'$. Note that these estimators may be inconsistent, namely OLS when $\rho_{xu} \neq 0$, IV^(j) when $\rho_{z_ju} \neq 0$, and TSLS when $\rho_{z_1u} \neq 0$ or/and $\rho_{z_2u} \neq 0$.

For the estimation errors we find, writing $x_i = \sigma_x \xi_i^*$ where $\xi_i^* = (1 - \rho_{xu}^2)^{1/2} \xi_i + \rho_{xu} \varepsilon_i$,

$$\hat{\beta}_{OLS} - \beta = (\sigma_u/\sigma_x) \Sigma_i(\xi_i^* \varepsilon_i) / \Sigma_i(\xi_i^{*2}), \quad (\text{S.10})$$

$$\hat{\beta}_{IV}^{(j)} - \beta = (\sigma_u/\sigma_x) \Sigma_i(\rho_{z_j\zeta} \zeta_{ij} \varepsilon_i + \rho_{z_j\xi} \xi_i \varepsilon_i + \rho_{z_ju} \varepsilon_i^2) / \Sigma_i(\rho_{z_j\zeta} \zeta_{ij} \xi_i^* + \rho_{z_j\xi} \xi_i \xi_i^* + \rho_{z_ju} \varepsilon_i \xi_i^*), \quad (\text{S.11})$$

$$\hat{\beta}_{TSLS} - \beta = x'P_Zu/x'P_Zx. \quad (\text{S.12})$$

It is directly seen that all estimation errors are invariant regarding β , are a multiple of σ_u/σ_x , and that P_Z is invariant with respect to the scale of the vectors z_1 and z_2 . Hence, without loss of generality, we may choose in the simulations $\beta = 0$, $\sigma_x = 1$, and $\sigma_{z_1} = \sigma_{z_2} = 1$. Then the dispersion of all estimators can be regulated by just varying σ_u . However, relative differences between dispersions will be invariant with respect to σ_u . So, by just choosing $\sigma_u = 1$ all relevant information will be obtained through choosing compatible values for the remaining design parameters: n , ρ_{xu} , ρ_{z_jx} and ρ_{z_ju} , where the latter two determine $\rho_{z_j\xi}$ and $\rho_{z_j\zeta_j}$. Due to the symmetry of the distribution of all

variables, changing the sign of any of the correlations, while keeping their absolute value fixed, has mostly simple (anti-)symmetric effects just on the sign of the estimation errors. Therefore, by just investigating nonnegative values for ρ_{xu} , ρ_{z_jx} and ρ_{z_ju} we will already get a rather complete picture.

Since the TSLS residuals equal $\hat{u}_{TSLS} = y - \hat{\beta}_{TSLS}x = u - (x'P_Zu/x'P_Zx)x$, the Sargan test statistic can be expressed as

$$\begin{aligned} S &= n \cdot \hat{u}'_{TSLS}P_Z\hat{u}_{TSLS}/\hat{u}'_{TSLS}\hat{u}_{TSLS} \\ &= n \frac{u'P_Zu(x'P_Zx)^2 - (x'P_Zu)^2x'P_Zx}{u'u(x'P_Zx)^2 - 2u'x(x'P_Zu)x'P_Zx + x'x(x'P_Zu)^2}. \end{aligned} \quad (\text{S.13})$$

It is obvious that S is invariant with respect to β and to all scale factors, because all individual terms in both the numerator and denominator are multiples of $\sigma_u^2\sigma_x^4$.

It is well known that the Sargan test is equivalent to literally testing over-identification exclusion restrictions. In the present design, this amounts to estimating the model $y_i = \beta x_i + \delta_j z_{ij} + u_i$, where j is either 1 or 2, while using both instruments, and then testing the significance of δ_j . Under the null hypothesis we have $\delta_j = 0$, thus

$$\begin{pmatrix} \hat{\beta} \\ \hat{\delta}_j \end{pmatrix} = \begin{pmatrix} z'_1x & z'_1z_j \\ z'_2x & z'_2z_j \end{pmatrix}^{-1} \begin{pmatrix} z'_1y \\ z'_2y \end{pmatrix} = \begin{pmatrix} \beta \\ 0 \end{pmatrix} + \begin{pmatrix} z'_1x & z'_1z_j \\ z'_2x & z'_2z_j \end{pmatrix}^{-1} \begin{pmatrix} z'_1u \\ z'_2u \end{pmatrix},$$

so that

$$\text{plim } \hat{\delta}_j = \text{plim } \frac{(z'_1x)(z'_2u) - (z'_2x)(z'_1u)}{(z'_1x)(z'_2z_j) - (z'_2x)(z'_1z_j)},$$

which has numerator

$$\text{plim}[(z'_1x/n)(z'_2u/n) - (z'_2x/n)(z'_1u/n)] = \sigma_u\sigma_x(\rho_{z_1x}\rho_{z_2u} - \rho_{z_2x}\rho_{z_1u}).$$

So, $\text{plim } \hat{\delta}_j = 0$ when

$$\rho_{z_1u}/\rho_{z_1x} = \rho_{z_2u}/\rho_{z_2x}. \quad (\text{S.14})$$

This explains that when (S.14) holds, the Sargan test will asymptotically reject the exclusion restriction with probability equal to the chosen significance level, even when ρ_{z_1u} and ρ_{z_2u} are far away from zero.

Further simulation results for a smaller and a larger sample

Figure S.1 provides results for $n = 50$ and Figure S.2 for $n = 2500$. In each figure the four rows of panels correspond to the four examined different cases of (in)validity of the instruments z_1 and z_2 .

Figure S.1: Simulation results for $n = 50$, $\sigma_u/\sigma_x = 1$, and the chosen set of correlation values

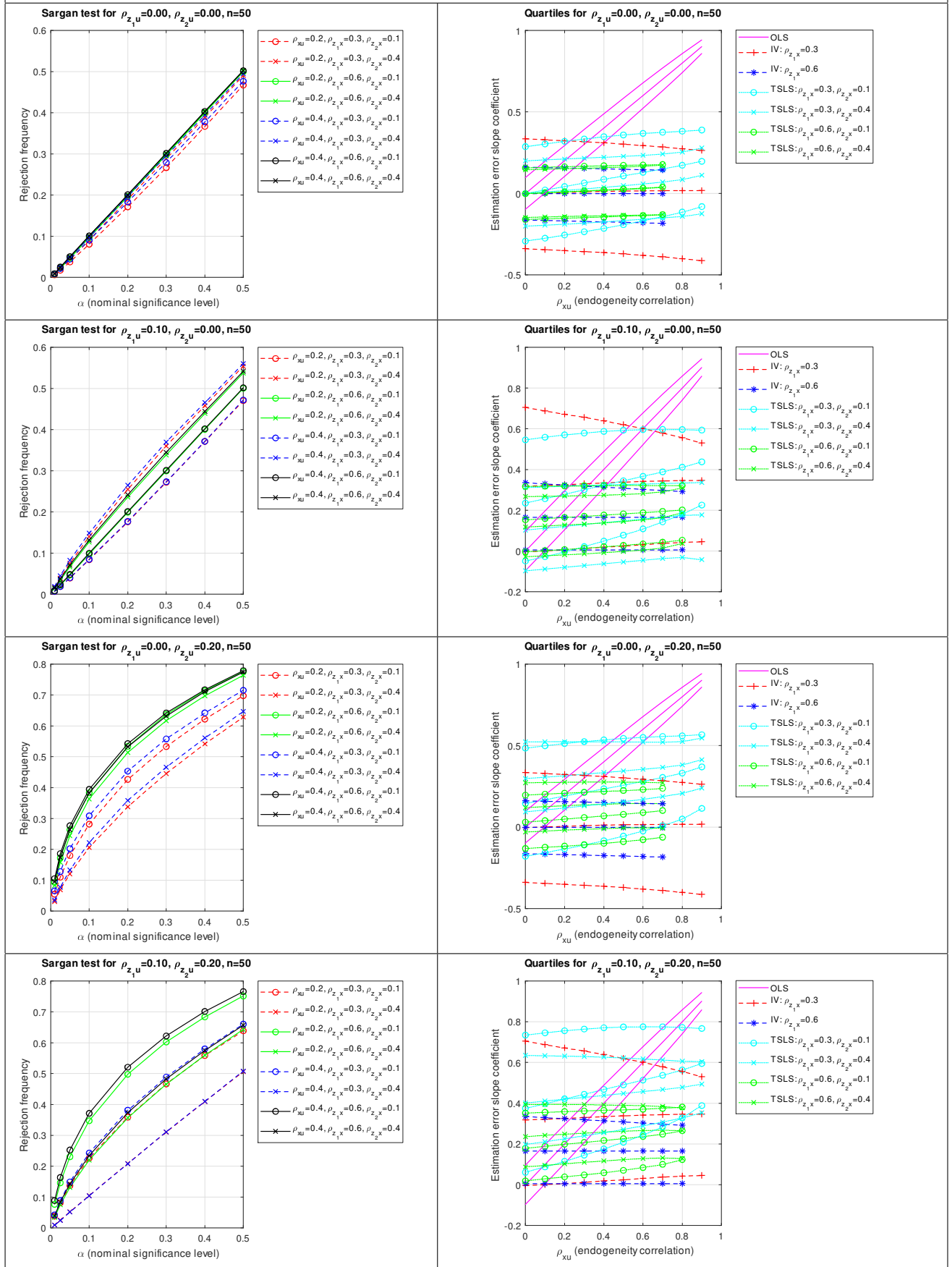


Figure S.2: Simulation results for $n = 2500$, $\sigma_u/\sigma_x = 1$, and the chosen set of correlation values

