

# working paper

## 2210

### PML vs minimum $\chi^2$ : the comeback

Dante Amengual  
Gabriele Fiorentini  
Enrique Sentana

October 2022

## PML vs minimum $\chi^2$ : the comeback

### Abstract

Arellano (1989a) showed that valid equality restrictions on covariance matrices could result in efficiency losses for Gaussian PMLEs in simultaneous equations models. We revisit his two-equation example using finite normal mixtures PMLEs instead, which are also consistent for mean and variance parameters regardless of the true distribution of the shocks. Because such mixtures provide good approximations to many distributions, we relate the asymptotic variance of our estimators to the relevant semiparametric efficiency bound. Our Monte Carlo results indicate that they systematically dominate MD, and that the version that imposes the valid covariance restriction is more efficient than the unrestricted one.

JEL Codes: C30, C36.

Keywords: Covariance restrictions, distributional misspecification, efficiency bound, finite normal mixtures, partial adaptivity.

Dante Amengual  
CEMFI  
amengual@cemfi.es

Gabriele Fiorentini  
Università di Firenze  
gabriele.fiorentini@unifi.it

Enrique Sentana  
CEMFI  
sentana@cemfi.es

## **Acknowledgement**

We would like to thank Manuel Arellano, as well as participants at CEMFI Econometrics Workshop, the UNIA Workshop on Empirical Microeconomics and Applied Econometrics, and the Bologna-Waseda Time Series Workshop for useful comments and suggestions. Of course, the usual caveat applies. The first and third authors gratefully acknowledge financial support from the Spanish Ministry of Science and Innovation through grant PID2021-128963NB-I00, while the second one is thankful to MIUR through the PRIN project "High-dimensional time series for structural macroeconomic analysis in times of pandemic.

# 1 Introduction

Maximum likelihood and minimum chi-square methods have been competing for the estimator throne for a long time. At the turn of the 19th century, Legendre (1805) and Gauss (1809) put forward least squares estimation as a Gaussian-based alternative to Laplace's (1774) least absolute deviation method, which relied on his eponymous distribution. Almost a century later, Pearson proposed not only the method of moments (see Pearson (1894)), but also the chi-square criterion in the context of matching theoretical and empirical frequencies (see Pearson (1900)). In turn, the development of maximum likelihood estimation (MLE) by Fisher (1922, 1925) was one of the most important achievements in 20th century statistics. Under standard regularity conditions, MLE asymptotically achieves the Cramér-Rao lower bound (see Cramér (1946) and Rao (1945)), which makes it at least as good as any minimum  $\chi^2$  estimator. In addition, it achieves second-order efficiency after a bias correction (see Rao (1961)). Moreover, the imposition of valid equality restrictions on the parameters systematically leads to efficiency gains (see Rothenberg (1973)).

However, not everybody was convinced (see Neyman and Scott (1948) on the incidental parameter problem, as well as the inconsistent MLE examples in Basu (1955), Kraft and Le Cam (1956) and Bahadur (1958)), and minimum  $\chi^2$  methods remained popular. In fact, Berkson (1980) argued that ML often was just a special case of minimum  $\chi^2$ , and not necessarily the best one. Soon afterwards, White (1982), building on earlier work by Huber (1967), and Gouriéroux, Monfort and Trognon (1984) studied the properties of Pseudo MLEs, characterising their consistency and general inefficiency. Arellano (1989a) put another nail on the ML coffin by showing that valid equality restrictions could result in efficiency losses for Gaussian PMLEs. Arguably, the wooden stake to the heart was driven by Newey and Steigerwald (1997), who described the inconsistency of non-Gaussian PMLE procedures under distributional misspecification. Since then, graduate students with non-Bayesian teachers learn the normal distribution only, and Gaussian PMLE is just an example of Hansen's (1982) GMM. In this paper, though, we argue that non-Gaussian PMLE, like a B-movie vampire, deserves a second life (or death).

We do so by revisiting the two-equation textbook example in Arellano (1989a),<sup>1</sup> except that instead of basing PMLE on the Gaussian distribution, as he did, we use discrete mixtures of normals. The reason is twofold. First, Fiorentini and Sentana (2022b) show that under standard regularity conditions such estimators are consistent for the conditional mean and variance

---

<sup>1</sup>Surprisingly, Arellano (1989a), which should be mentioned in all graduate econometric textbooks, has received very few citations: Pollock (1988), Islam (1993), Monés and Ventura (1996), Calzolari, Fiorentini and Sentana (2004), and Sentana (2005), plus a handful of self-citations, and two more which really meant to cite Arellano (1989b).

parameters regardless of the true distributions of the shocks to the model and the number of mixture components, thereby nesting the results for Gaussian PMLE in Gouriéroux, Monfort and Trognon (1984) while simultaneously avoiding the concerns raised by Newey and Steigerwald (1997). Second, finite normal mixtures with a sufficiently large number of components can provide good approximations to many distributions (see Nguyen et al (2020)), so it is reasonable to expect that PMLEs based on them can get close to achieving the semiparametric (SP) efficiency bound, and therefore exploit the potential adaptivity of some of the parameters when it exists, at least asymptotically.<sup>2</sup>

The rest of the paper is organised as follows. Section 2 introduces the example in Arellano (1989a) and summarises his main results. Then, section 3 derives the relevant semiparametric efficiency bounds, which we use to benchmark the different estimators. Next, section 4 contains the results of our extensive Monte Carlo experiments while section 5 concludes. Proofs and auxiliary results are relegated to the appendices.

## 2 The example

Consider the following textbook example:

$$y_1 = \gamma + \alpha y_2 + \beta z_1 + u_1, \quad (1)$$

$$y_2 = \mu_0 + \mu_1 z_1 + \mu_2 z_2 + u_2, \quad (2)$$

with

$$\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \Big| z_1, z_2 \sim D \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} \right].$$

As is well known, the unrestricted Gaussian PMLE of  $\alpha$  and  $\beta$  coincides with the IV estimator that uses a constant,  $z_1$  and  $z_2$  as instruments in the first equation. In turn, the restricted Gaussian PMLE that imposes  $\sigma_{12} = 0$  coincides with the OLS estimator of the first equation.

When the joint conditional distribution of  $u_1$  and  $u_2$  is Gaussian, OLS is at least as efficient as IV, which justifies the Durbin-Wu-Hausman test.<sup>3</sup> But Arellano's (1989a) counterintuitive result says that when the true conditional distribution is not Gaussian, IV may be more efficient than OLS for  $\alpha$  and  $\beta$  even though  $\sigma_{12} = 0$ . Specifically, he showed that IV will beat OLS if and only if

$$\mu_{22} \geq 1 + \rho_{y_2 z_2, z_1}^{-2}, \quad (3)$$

---

<sup>2</sup>See Fiorentini and Sentana (2022a) for a related discussion in the context of structural VARs.

<sup>3</sup>Wu (1973) compared OLS with IV in linear single equation models to assess regressor exogeneity unaware that Durbin (1954) had already suggested this. Hausman (1978) provided a procedure with far wider applicability.

where

$$\mu_{22} = E \left( \frac{u_1^2 u_2^2}{\sigma_1^2 \sigma_2^2} \middle| z_1, z_2 \right)$$

is the co-kurtosis coefficient between the two structural shocks and  $\rho_{y_2 z_2, z_1}$  is the correlation coefficient between  $y_2$  and  $z_2$  after partialling out the effect of  $z_1$ . Intuitively,  $\mu_{22}$  affects the correct sandwich version of the asymptotic covariance matrix of the OLS estimators of the slope parameters.

Appendix A contains detailed expressions for the asymptotic variances of the OLS and IV estimators of  $\alpha$  and  $\beta$ . We have used those expressions to create Figure 1, which displays in  $(\rho_{y_2 z_2, z_1}, \mu_{22})$  space (minus one plus) the ratio of the asymptotic variances of the OLS and IV estimators of  $\alpha$  for positive values of  $\rho_{y_2 z_2, z_1}$ .<sup>4</sup> We do so for the special case in which the  $R^2$  of equation (2) coincides with  $\rho_{y_2 z_2, z_1}^2$ , which allows this parameter to vary freely from 0 to 1.<sup>5</sup> As expected, OLS is more/less efficient than IV to the left/right of the boundary line (3).

This figure also shows the locus of  $(\rho_{y_2 z_2, z_1}, \mu_{22})$  combinations for which the IV estimator of  $\alpha$  reaches its maximum asymptotic efficiency relative to the corresponding OLS estimator, which is given by the curve

$$\rho_{y_2 z_2, z_1}^2 = \frac{\mu_{22}}{2(\mu_{22} - 1)}.$$

Further increases in  $\rho_{y_2 z_2, z_1}$  for a given  $\mu_{22}$  result in decreases in relative efficiency, with OLS and IV becoming indistinguishable as  $\rho_{y_2 z_2, z_1} \rightarrow 1$ , in which case  $z_2$  becomes a perfect instrument for  $y_2$ .

In this context, Arellano's (1989a) proposed solution is to replace Gaussian PMLE by Minimum Distance estimators (MD), a special case of minimum chi-square methods popularised in econometrics by Malinvaud (1970). The rationale is as follows. Let  $\boldsymbol{\theta} = (\gamma, \alpha, \beta, \mu_0, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2)'$  denote the vector of structural parameters. Given that the reduced form of model (2) is

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \middle| z_1, z_2 \sim D[\boldsymbol{\mu}(z_1, z_2; \boldsymbol{\theta}), \boldsymbol{\Omega}(z_1, z_2; \boldsymbol{\theta})] \quad (4)$$

$$\boldsymbol{\mu}(z_1, z_2; \boldsymbol{\theta}) = \begin{bmatrix} (\gamma + \alpha\mu_0) + (\beta + \alpha\mu_1)z_1 + \alpha\mu_2 z_2 \\ \mu_0 + \mu_1 z_1 + \mu_2 z_2 \end{bmatrix} \quad (5)$$

$$\boldsymbol{\Omega}(z_1, z_2; \boldsymbol{\theta}) = \begin{pmatrix} \sigma_1^2 + \alpha^2 \sigma_2^2 + 2\sigma_{12}\alpha & \alpha\sigma_2^2 + \sigma_{12} \\ \alpha\sigma_2^2 + \sigma_{12} & \sigma_2^2 \end{pmatrix}, \quad (6)$$

which is exactly identified, the unrestricted MD coincides with IV, which is Indirect Least Squares. Then, Arellano (1989a) shows that imposing the restriction  $\sigma_{12} = 0$  leads to an overidentified optimal MD (weakly) more efficient than both IV and OLS for  $\alpha$  and  $\beta$ .

<sup>4</sup>The plot would be the mirror image of Figure 1 for negative values.

<sup>5</sup>As we shall see in Proposition 1 below, though, this special case is such that, asymptotically, the difference in efficiency between the IV and OLS estimators affects  $\alpha$  exclusively.

This optimal MD requires an asymptotic covariance of the reduced form parameter estimators which recognises that the third- and fourth-order multivariate cumulants of  $u_1$  and  $u_2$  are not usually 0 when they are jointly non-normally distributed.

Appendix A also contains detailed expressions for the asymptotic variances of the optimal MD estimators of  $\alpha$  and  $\beta$ . We have used those expressions to create Figure 2, which depicts in  $(\rho_{y_2 z_2, z_1}, \mu_{22})$  space (minus one plus) the ratio of the asymptotic variance of the restricted optimal MD of  $\alpha$  to the asymptotic variance of either the OLS estimator (to the left of (3)) or the IV one (to its right) in the same set up as Figure 1. As can be seen, the efficiency gains are relatively small over the displayed range, and they vanish when either the partial correlation goes to 0 or 1 or the co-kurtosis term goes to 0.<sup>6</sup>

Although the main focus of the analysis in Arellano (1989a) was  $\alpha$  and  $\beta$ , it is of some interest to study the asymptotic efficiency of the optimal MD estimators of the remaining structural model parameters relative to their OLS and IV counterparts. Given that the number of different bivariate cumulants of orders three and four is 4 and 5, respectively, we focus on the special case in which the joint distribution of the (standardised) structural shocks conditional on the instruments is spherical, or  $s(\mathbf{0}, \mathbf{I}_N, \boldsymbol{\eta})$  for short, where  $\boldsymbol{\eta}$  is the possibly infinite vector of shape parameters. More formally,

**Assumption 1**

$$\left. \frac{u_1}{\sigma_1}, \frac{u_2}{\sigma_2} \right| z_1, z_2; \boldsymbol{\theta}, \boldsymbol{\eta} \sim i.i.d. \ s(\mathbf{0}, \mathbf{I}_N, \boldsymbol{\eta}) \quad (7)$$

To simplify the expressions further, we are going to follow appendix B in Fiorentini and Sentana (2019) and reparametrise the unrestricted covariance matrix of the structural residuals as

$$\begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} = \sigma^2 \begin{pmatrix} 1 & 0 \\ \psi_{12} & 1 \end{pmatrix} \begin{pmatrix} e^\omega & 0 \\ 0 & e^{-\omega} \end{pmatrix} \begin{pmatrix} 1 & \psi_{12} \\ 0 & 1 \end{pmatrix}, \quad (8)$$

where  $\psi_{12}$  is the coefficient in the least squares projection of  $u_2$  on  $u_1$ , and  $\sigma^2$  and  $\omega$  the geometric mean of their variances and the natural log of the ratio of the standard deviations of these shocks, respectively, under the maintained assumption that they are uncorrelated.<sup>7</sup> Let  $\boldsymbol{\theta}^\dagger = (\gamma, \alpha, \beta, \mu_0, \mu_1, \mu_2, \omega, \sigma^2)'$  denote the vector of structural parameters implied by (8).

We can then show under standard regularity conditions that:

**Proposition 1** *Let  $(\tau_1, \tau_2)$  and  $(\sigma_{z_1}^2, \sigma_{z_2}^2, \sigma_{z_1 z_2})$  denote the means, variances and covariance of  $z_1$  and  $z_2$ . If Assumption 1 holds, then:*

(a) *The difference between the asymptotic covariance matrices of the OLS and MD estimators*

---

<sup>6</sup> Again, Proposition 1 below implies that the differences in asymptotic variances between the MD, IV and OLS estimators affect  $\alpha$  exclusively in the special case in which the (squared) partial correlation of  $y_2$  and  $z_2$  given  $z_1$  coincides with the  $R^2$  in the regression of  $y_2$  on  $z_1$  and  $z_2$

<sup>7</sup> More generally,  $\sigma^2 = \sqrt{\sigma_1^2 \sigma_2^2 - \sigma_{12}^2}$  and  $\omega = \ln \left( \sigma_1 / \sqrt{\sigma_2^2 - \sigma_{12}^2 / \sigma_1^2} \right)$ .

of  $\theta^\dagger$ ,  $\hat{\theta}_{LS}^\dagger$  and  $\hat{\theta}_{MD}^\dagger$ , respectively, is positive semidefinite of rank 1 at most, with a basis for its image given by

$$\{ -[\mu_0 + (\tau_2 - \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \tau_1) \mu_2], 1, -\mu_1 + \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \mu_2, \mathbf{0}_{1 \times 5} \}, \quad (9)$$

and a basis for its kernel by

$$[1, \mu_0 + (\tau_2 - \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \tau_1) \mu_2, 0, \mathbf{0}_{1 \times 5}]', \quad (10)$$

$$[\mu_1 + \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \mu_2, 0, \mu_0 + (\tau_2 - \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \tau_1) \mu_2, \mathbf{0}_{1 \times 5}]' \quad (11)$$

and

$$(\mathbf{0}_{5 \times 3}, \mathbf{I}_5)'. \quad (12)$$

(b) The difference between the asymptotic covariance matrices of the IV and MD estimators of  $\theta^\dagger$ ,  $\tilde{\theta}_{IV}^\dagger$  and  $\hat{\theta}_{MD}^\dagger$ , respectively, is positive semidefinite of rank 1 at most, with the same basis for image and kernel.

(c) The difference between the asymptotic covariance matrices of the OLS and IV estimators of  $\theta^\dagger$ ,  $\hat{\theta}_{LS}^\dagger$  and  $\tilde{\theta}_{IV}^\dagger$ , respectively, is positive/negative semidefinite of rank 1 depending of condition (3), with exactly the same basis for image and kernel.

This proposition considerably sharpens the results in Arellano (1989a) for the special case of spherically symmetric disturbances by showing that the asymptotic efficiency gains concentrate in a single linear combination of the parameters of the first equation  $\gamma$ ,  $\alpha$  and  $\beta$  given by (9). In contrast, any other linear combination of the parameters orthogonal to this one does not generate any efficiency gains. Specifically, the parameters of the second equation as well as the residual variances are estimated just as efficiently by the three procedures.

The predictable reaction of an ML believer to Figures 1 and 2 would be to argue that condition (3) requires the combination of a very good instrument (a high  $\rho_{y_2 z_2, z_1}^2$ ) with a substantial amount of non-normality (a large  $\mu_{22}$ ), in which case the Gaussian assumption would be very inappropriate. For example, a joint Student  $t$  distribution for  $u_1$  and  $u_2$  cannot satisfy this condition when the number of degrees of freedom is six or more, and the requirement becomes increasingly difficult for poor instruments.

A naïve ML solution would be to assume that  $u_1$  and  $u_2$  follow a bivariate Student  $t$  distribution and estimate the model parameters together with the degrees of freedom, which should dominate MD. In this respect, we have used the expressions in Appendix A to create Figures 3a and 3b, which display in  $(\rho_{y_2 z_2, z_1}, \mu_{22})$  space (minus one plus) the ratio of the asymptotic variances of the  $t$ -based MLE of  $\alpha$  and  $\beta$  that impose  $\sigma_{12} = 0$  to the asymptotic variances of the corresponding restricted optimal MD. As can be seen, these figures confirm that ML does indeed dominate MD in this case.

The problem with this naïve approach is that if the assumed joint distribution is incorrect, the resulting PMLEs may be inconsistent, as forcefully argued by Newey and Steigerwald (1997).

However, this does not mean that all parameters will be inconsistently estimated. Specifically, Proposition 3 in Fiorentini and Sentana (2019) implies that the unrestricted  $t$ -based PMLEs of  $\alpha$  and  $\beta$  are always consistent irrespective of the true distribution. Similarly, their Proposition 1 implies that the restricted  $t$ -based PMLEs of  $\alpha$  and  $\beta$  will remain consistent when (7) is true even though it does not coincide with the distribution assumed for estimation purposes. Besides, it may be possible to obtain two-step consistent estimators in closed-form along the lines of Fiorentini and Sentana (2019).

More importantly, Fiorentini and Sentana (2022b) show that all parameters will always be consistently estimated if one assumes for estimation purposes that  $u_1$  and  $u_2$  follow a finite mixture of bivariate normals regardless of the true distribution of those innovations and the number of components of the mixture. Thus, the consistency of the Gaussian PMLE is just a special case.

The ability of finite Gaussian mixtures to approximate many other distributions mentioned in the introduction means that we can relate these finite mixture PMLEs to SP estimators which simply exploit the independence of the shocks and the conditioning variables without making any parametric assumptions. For that reason, in the next section we take SP estimators as our benchmark to study:

1. the efficiency of the OLS, IV, MD and correct ML estimators relative to the SP ones,
2. the relative efficiency of restricted and unrestricted versions of these SP estimators, and
3. the relative efficiency of finite mixture-based PMLEs relative to the SP estimators

in the context of model (2).

### 3 Semiparametric estimation and efficiency bounds

The optimal instruments theory of Chamberlain (1987) implies that Arellano's (1989a) MD estimator achieves the SP efficiency bound which exploits the correct specification of the conditional mean and variance functions for  $y_1$  and  $y_2$  in the reduced form model (2) when the joint third- and fourth-order cumulants of  $u_1$  and  $u_2$  conditional on  $z_1$  and  $z_2$  are constant. However, if this last maintained assumption is true, one can in principle obtain an even more efficient MD estimator of the model parameters after augmenting it with equations for the third- and fourth-order cumulants of the reduced-form residuals under the assumption that the joint cumulants of  $u_1$  and  $u_2$  conditional on  $z_1$  and  $z_2$  are constant up to the eighth-order.

In fact, the results in Bickel et al (1993) allow us to obtain the SP efficiency bound that exploits that the joint distribution of  $u_1$  and  $u_2$  is independent of  $z_1$  and  $z_2$ . Moreover, we can

also consider a restricted version of this SP bound under the maintained assumption that (7) holds, as in Hodgson and Vorkink (2003), which will be bigger in the usual positive semidefinite sense. Henceforth, we shall refer to this bound and its associated estimator as SS.

An interesting question in this context is the possibility that some but not all of the parameters of model (2) can be partially adaptively estimated, in the sense that their SP estimators are as asymptotically efficient as the infeasible ML estimators which exploit the information of the true distribution of the shocks, including the values of their shape parameters. The following proposition provides a precise answer to this question under sphericity for the restricted estimators that impose  $\sigma_{12} = 0$ :

**Proposition 2** *If Assumption 1 holds, then:*

- (a) *The difference between the asymptotic covariance matrices of the restricted SS and infeasible ML estimators of  $\theta^\dagger$ ,  $\hat{\theta}_{SS}^\dagger$  and  $\hat{\theta}_{ML}^\dagger(\bar{\eta})$ , respectively, is positive semidefinite of rank 1 at most, with the basis for its image given by  $(\mathbf{0}_{1 \times 7}, 1)'$ , and a basis for its kernel by  $(\mathbf{I}_7, \mathbf{0}_{7 \times 1})'$ .*
- (b) *The difference between the asymptotic covariance matrices of the restricted SP and infeasible ML estimators of  $\theta^\dagger$ ,  $\hat{\theta}_{SP}^\dagger$  and  $\hat{\theta}_{ML}^\dagger(\bar{\eta})$ , respectively, is positive semidefinite of rank 5 at most, with the basis for its image given by  $(1, \mathbf{0}_{1 \times 7})'$ ,  $(0, -1, \mu_1 + \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \mu_2, \mathbf{0}_{1 \times 5})'$ ,  $(\mathbf{0}_{1 \times 3}, 1, \mathbf{0}_{1 \times 4})'$  and  $(\mathbf{0}_{1 \times 6}, \mathbf{I}_2)'$ , and a basis for its kernel by  $(0, \mu_1 + \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \mu_2, 1, \mathbf{0}_{1 \times 5})'$  and  $(\mathbf{0}_{2 \times 4}, \mathbf{I}_2, \mathbf{0}_{2 \times 2})'$ .*
- (c) *The difference between the asymptotic covariance matrices of the MD and SP estimators of  $\theta^\dagger$ ,  $\hat{\theta}_{MD}^\dagger$  and  $\hat{\theta}_{SP}^\dagger$ , respectively, is positive semidefinite of rank 4 at most, with the basis for its image given by  $[-(\mu_0 + \mu_1 \tau_1 + \mu_2 \tau_2), 1, 0, \mathbf{0}_{1 \times 5}]'$ ,  $(-\tau_1, 0, 1, \mathbf{0}_{1 \times 5})'$ ,  $(\mathbf{0}_{1 \times 3}, -\tau_1, 1, 0, \mathbf{0}_{1 \times 2})'$  and  $(\mathbf{0}_{1 \times 3}, -\tau_2, 0, 1, \mathbf{0}_{1 \times 2})'$ , and a basis for its kernel by  $(\mathbf{0}_{2 \times 6}, \mathbf{I}_2)'$ ,  $(1, \mu_0 + \mu_1 \tau_1 + \mu_2 \tau_2, \tau_1, \mathbf{0}_{1 \times 5})'$  and  $(\mathbf{0}_{1 \times 3}, 1, \tau_1, \tau_2, \mathbf{0}_{1 \times 2})'$ .*

The first part of the proposition implies that all the structural model parameters except the overall residual scale  $\sigma^2$  can be (partially) adaptively estimated by the SS estimator, as expected from Proposition 12 in Fiorentini and Sentana (2021).

More interestingly, the second part of the proposition implies that in addition to  $\mu_1$  and  $\mu_2$ , the coefficient of the linear projection of  $y_1$  onto a constant and  $z_1$ , which is given by

$$\beta + (\mu_1 + \sigma_{z_1 z_2} \sigma_{z_1}^{-2} \mu_2) \alpha,$$

will be adaptively estimated by the SP estimator. In this respect, a very important by-product of this proposition is that the model parameters that can be partially adaptively estimated often continue to be consistently estimated under distributional misspecification of the innovations, as shown by Fiorentini and Sentana (2019, 2021) in the context of multivariate location-scale models such as (2). We will revisit this issue in the Monte Carlo section.

Finally, the last part of the proposition says that the variances of the structural residuals, as well as  $E(y_1) = \gamma + E(y_2)\alpha + E(z_2)\beta$  and  $E(y_2) = \mu_0 + E(z_1)\mu_1 + E(z_1)\mu_2$ , are asymptotically equally efficiently estimated by the MD and SP estimators, and that the efficiency gains are

concentrated in  $\alpha - E(y_2)\gamma$ ,  $\beta - E(z_1)\gamma$ ,  $\mu_1 - E(z_1)\mu_0$  and  $\mu_2 - E(z_2)\mu_0$ .

Motivated by the seemingly counterintuitive result in Arellano (1989a), it is also of interest to analyze the effects of imposing the valid restriction  $\sigma_{12} = 0$  on these different estimators:

**Proposition 3** *If Assumption 1 holds, then:*

- (a) *The difference between the asymptotic covariance matrices of the unrestricted and restricted ML estimators of  $\theta^\dagger$ ,  $\tilde{\theta}_{ML}^\dagger$  and  $\hat{\theta}_{ML}^\dagger$ , respectively, is positive semidefinite of rank 1 at most, with the basis for its image given by (9), and a basis for its kernel by (10), (11) and (12).*
- (b) *The difference between the asymptotic covariance matrices of the unrestricted and restricted SS estimators of  $\theta^\dagger$ ,  $\tilde{\theta}_{SS}^\dagger$  and  $\hat{\theta}_{SS}^\dagger$ , respectively, is positive semidefinite of rank 1 at most, with the basis for its image given by (9), and a basis for its kernel by (10), (11) and (12).*
- (c) *The difference between the asymptotic covariance matrices of the unrestricted and restricted SP estimators of  $\theta^\dagger$ ,  $\tilde{\theta}_{SP}^\dagger$  and  $\hat{\theta}_{SP}^\dagger$ , respectively, is positive semidefinite of rank 1 at most, with the basis for its image given by (9), and a basis for its kernel by (10), (11) and (12).*

Therefore, the imposition of the valid covariance restriction  $\sigma_{12} = 0$  always leads to (weak) efficiency gains for exactly the same linear combination of the parameters of the first structural equation for which MD leads to an efficiency gain relative to both OLS and IV.

Finally, we study the extent to which PMLEs based on finite mixtures of normals with an increasing number of components could constitute the basis for a proper sieves-type SP procedure, as we argued in the introduction. To do so, we have conducted a simple exercise in which we look at the case the shocks to model (2) follow a bivariate Student  $t$  with 0 means, unit standard deviations, no correlation and 5 degrees of freedom but whose parameters are estimated by finite scale mixture-based log-likelihood functions with  $K = 2, 3$  and 4 components. For comparison purposes, we consider four different benchmarks: (i) the MLE based on the correctly specified log-likelihood function that fixes the number of degrees of freedom to 5, (ii) the SS estimator, (iii) the OLS estimator, and (iv) the optimal MD estimator.

We compute the expected value of the Hessian and variance of the score of the scale mixture-based PMLEs using the expressions in Fiorentini and Sentana (2021) evaluated at the true values of the mean and variance parameters in  $\theta$  and the pseudo true values of the shape parameters, which we numerically obtain from samples of 40 million simulated observations.

The results, which we report in Table 1, show that the scale mixture-based PMLEs of all the model parameters except the overall residual scale  $\sigma^2$  quickly approach the asymptotic efficiency of the infeasible MLE based on the Student  $t$  with 5 degrees of freedom despite the fact that finite Gaussian mixtures are thin tailed. In fact, although panel (a) in Figure 3 of Gallant and Tauchen (1999) clearly illustrates that a more complex misspecified model does not necessarily lead to more efficient estimators because one is not simply adding new elements to the score, but also changing the pseudo true values of the shape parameters at which one evaluates the

original components of the score, we find that the efficiency improvements occur monotonically.

In contrast, the asymptotic variances of the scale mixture-based PMLEs of  $\sigma^2$  coincides with the asymptotic variances of the OLS estimators irrespective of the number of components, which reflects (i) the block diagonality of the different asymptotic covariance matrices in Proposition 12 of Fiorentini and Sentana (2021) because the determinant of (8) is precisely  $\sigma^4$ , and (ii) the fact that the ML estimators of the mean in a scale mixture of  $K$  gammas is numerically the same regardless of  $K$ , as explained in Fiorentini and Sentana (2022b).

## 4 Monte Carlo analysis

In previous sections, we have derived several asymptotic results regarding the relative efficiency of the LS, IV and MD estimators, as well as the finite mixture-based PMLEs and the feasible and infeasible MLEs. In this section, in contrast, we make use of an extensive Monte Carlo simulation exercise to asses the small sample behaviour of all those estimators.

### 4.1 Design

We consider three different parameter configurations:

- a.  $\mu_{22} = 3$  and  $\rho_{y_2 z_2, z_1} = (\mu_{22} - 1)^{-\frac{1}{2}} = 1/\sqrt{2} \simeq 0.71$ , which is such that the IV and OLS estimators of  $\alpha$  and  $\beta$  have the same asymptotic efficiency (see the solid line in Figure 1);
- b.  $\mu_{22} = 3$  and  $\rho_{y_2 z_2, z_1} = 2^{-\frac{1}{2}} \sqrt{\mu_{22}/(\mu_{22} - 1)} = \sqrt{3}/2 \simeq 0.87$ , which is such that the ratio of the asymptotic variance of the IV and OLS estimators of  $\alpha$  is maximum for the chosen value of  $\mu_{22}$  (see the dotted line in Figure 1); and
- c.  $\mu_{22} = 7/3$  and  $\rho_{y_2 z_2, z_1} = (\mu_{22} - 1)^{-\frac{1}{2}} = \sqrt{3}/2 \simeq 0.87$ , which is another case of equal efficiency of IV and OLS, but with lower co-kurtosis.<sup>8</sup>

As for the distribution of the structural shocks, we consider four non-Gaussian possibilities in which  $(u_1, u_2)$  follow a:

1. Student  $t$  distribution with  $\nu = 5$  or  $\nu = 5.5$  degrees of freedom corresponding to  $\mu_{22} = 3$  and  $\mu_{22} = 7/3$ , respectively;
2. scale mixture of two normals in which the higher variance component has probability  $\lambda = .05$  and the ratio of the variances is either  $\varkappa = 0.094$  or  $\varkappa = 0.122$  corresponding to  $\mu_{22} = 3$  and  $\mu_{22} = 7/3$ , respectively;

---

<sup>8</sup>We do not consider the case in which  $\mu_{22} = 7/3$  and  $\rho_{y_2 z_2, z_1} = .5\sqrt{\mu_{22}/(\mu_{22} - 1)}$  because the maximum efficiency of IV relative to OLS for  $\alpha$  is just 1.02 in that case.

3. asymmetric Student  $t$  distribution with negative tail dependence  $\mathbf{b} = (-1, -1)'$  but degrees of freedom  $\nu = 9.65$  or  $\nu = 10.38$ , respectively;
4. location-scale mixture of two normals in which the higher variance component has probability  $\lambda = .05$ ,  $\mu_{22}$  is as in 1., and the marginal skewness of  $u_1$  and  $u_2$  is as in 3., which is achieved with

$$\boldsymbol{\delta} = \begin{pmatrix} -1.01 \\ -1.06 \end{pmatrix} \text{ or } \boldsymbol{\delta} = \begin{pmatrix} -1.16 \\ -1.24 \end{pmatrix} \text{ and } \boldsymbol{\Sigma}_L = \begin{pmatrix} 0.32 & 0 \\ 0 & 0.32 \end{pmatrix} \text{ or } \begin{pmatrix} 0.38 & 0 \\ 0 & 0.38 \end{pmatrix},$$

respectively (see Appendix D for further details on this parametrisation).

For illustrative purposes, we display the joint densities and contours for standardised versions of these distributions in comparison to the bivariate spherical Gaussian distribution in Figures 4 and 5 for the spherically symmetric and general cases, respectively.

In all simulated samples the exogenous variables  $\mathbf{z} = (z_1, z_2)'$  are generated according to a bivariate Student  $t$  distribution with 8 degrees of freedom with mean vector  $\boldsymbol{\tau} = (1, 1)'$  and an identity variance covariance matrix.<sup>9</sup>

Next, for each choice of the partial correlation  $\rho_{y_2 z_2 \cdot z_1}$  mentioned above, we choose

$$R_2^2 = \frac{2\rho_{y_2 z_2 \cdot z_1}}{1 + \rho_{y_2 z_2 \cdot z_1}} \text{ and } \rho_{y_2 z_1} = \rho_{y_2 z_2} = \sqrt{\frac{R_2^2 - \rho_{y_2 z_2 \cdot z_1}^2}{1 - \rho_{y_2 z_2 \cdot z_1}}},$$

so that the two slope coefficients of the second equations coincide. If we fix the variance of both  $y_1$  and  $y_2$  to 1 without loss of generality, these restrictions implicitly determine the variance of the error term of the second equation as  $\sigma_2^2 = 1 - R_2^2$ . We also impose the same balancing restriction on the slopes of the first equation by choosing

$$\alpha = \beta = \sqrt{\frac{(1 + \rho_{y_2 z_1})R_1^2}{2}}.$$

Then, we fix  $R_1^2$  to 0.5, which implies  $\sigma_1^2 = 1/2$ , an arbitrary choice that simply scales the asymptotic variances of all the different estimators of  $\alpha$  and  $\beta$  by the same amount.<sup>10</sup> Finally, we choose the values of the intercepts  $\gamma$  and  $\mu_0$  so that  $E(y_1) = E(y_2) = 1$  (see Appendix C for further details).

---

<sup>9</sup>Notice that the choice of  $\sigma_{z_1 z_2} = 0$  considerably simplifies some of the eigenvectors in Propositions 1, 2 and 3. For example the linear combination that according to Proposition 2.b can be adaptively estimated by the SP estimator and consistently estimated by a distributionally misspecified ML estimator becomes  $\beta + \mu_1 \alpha$ .

<sup>10</sup>In design *a.*, we then have  $R_2^2 = 2/3$ ,  $\sigma_2^2 = 1/3$ ,  $\gamma = 0.20$ ,  $\alpha = \beta = 0.40$ ,  $\mu_0 = 0.16$ , and  $\mu_1 = \mu_2 = 0.58$ . In turn, in designs *b* and *c.*,  $R_2^2 = 6/7$ ,  $\sigma_2^2 = 1/7$ ,  $\gamma = 0.22$ ,  $\alpha = \beta = 0.39$ ,  $\mu_0 = -0.31$ , and  $\mu_1 = \mu_2 = 0.66$ .

## 4.2 Simulation results

We simulate 10,000 samples of length  $N = 250$  and  $N = 1,000$  for each of the above designs. For each simulated sample, we compute the IV, LS and MD estimators, together with unrestricted and restricted versions of PMLE estimators that use either a discrete mixture of two normals –UPML(mn), RPML(mn)– or a Student  $t$  distribution –UPML(t) and RPML(t).

We display the finite sample results by means of the box-plots in Figures 6 to 11, which concentrate on  $\alpha$  and  $\beta$ , the two parameters of interest.

Figures 6 to 8 show the Monte Carlo results for 250 observations for cases  $a.$ ,  $b.$  and  $c.$ , respectively, while Figures 9 to 11 contain the results for 1,000 observations in the same order.

Our findings indicate that OLS is better in finite samples than what the asymptotic theory suggests because the sample co-kurtosis coefficient is downward biased for  $\mu_{22}$ . In fact, the asymptotic efficiency of the IV estimator of  $\alpha$  relative to LS can only be observed in panels b and d of Figure 10 when the sample length is large and the distribution of the shocks is either a spherical or a general finite mixture of normals, which is when there seems to be a lower small sample bias for  $\mu_{22}$ .

They also confirm that optimal MD dominates both OLS and IV in finite samples, but the need to estimate third- and fourth-order multivariate cumulants to compute the optimal weighting matrix handicaps it somewhat (see Altonji and Segal (1996) for analogous results in the context of optimal GMM estimators when the shocks are fat tailed)

Our results also indicate that non-Gaussian PML based on a restrictive parametric distribution like the Student  $t$  works well when the true distribution is spherical, but it generates inconsistencies otherwise. Notice though that the RPML(t) estimator seems to be consistent for  $\beta + \mu_1\alpha$  despite being inconsistent for both  $\alpha$  and  $\beta$ , which is in line with our theoretical discussion following Proposition 2.

More importantly, we find that non-Gaussian PMLEs based on a flexible distribution like finite mixture of normals works well in practice regardless of the true distribution, systematically dominating MD. In addition, the version that imposes the valid covariance restriction  $\sigma_{12} = 0$  is systematically more efficient than the the unrestricted one.

## 5 Directions for further research

Although we have seen that our proposed finite mixture-based PMLEs get close to achieving the SP efficiency bound, an obvious extension of our Monte Carlo experiments would be to consider standard two-step SP estimators that starting from a consistent estimator such as OLS carry out one BHHH iteration using the efficient SP score estimated non-parametrically. The

curse of dimensionality in estimating multivariate densities, though, might reduce the theoretical advantages of such methods in finite samples.

Another worthwhile exercise would be to extend the analysis in this paper to the general simultaneous equation model with  $n$  endogenous variables and  $k$  instrumental ones considered by Arellano (1989a). Aside from involving more complex analytical expressions than in the bivariate example we have considered, the main practical complication would be that the number of free parameters of a standardised multivariate mixture increases with the square of the cross-sectional dimension, as we explain in Appendix D.

## References

- Altonji, J.G and Segal, L.M. (1996): “Small-sample bias in GMM estimation of covariance structures”, *Journal of Business and Economic Statistics*, 14, 353–366
- Amengual, D., Fiorentini, G. and Sentana, E. (2022): “Tests for random coefficient variation in vector autoregressive models”, in J.J. Dolado, L. Gambetti and C. Matthes (eds.) *Essays in honour of Fabio Canova, Advances in Econometrics* 44B, 1–35, Emerald.
- Arellano (1989a): “On the efficient estimation of simultaneous equations with covariance restrictions”, *Journal of Econometrics* 42, 247–265.
- Arellano (1989b): “An efficient GLS estimator of triangular models with covariance restrictions”, *Journal of Econometrics* 42, 267–273.
- Bahadur, R.R. (1958): “Examples of inconsistency of maximum likelihood estimates”, *Sankhya* 20, 207–210.
- Basu, D. (1955): “An inconsistency of the method of maximum likelihood”, *Annals of Mathematical Statistics* 26, 144–145.
- Berkson, J. (1980): “Minimum chi-square, not maximum likelihood!”, *Annals of Statistics* 8, 457–487.
- Bickel, P.J., Klaassen, C.A.J., Ritov, Y., and Wellner, J.A. (1993): *Efficient and adaptive estimation for semiparametric models*. Johns Hopkins.
- Calzolari, G., Fiorentini, G. and Sentana, E. (2004): “Constrained indirect estimation”, *Review of Economic Studies* 71, 945–973.
- Chamberlain, G. (1987): “Asymptotic efficiency in estimation with conditional moment restrictions”, *Journal of Econometrics* 34, 305–334.
- Cramér, H. (1946): *Mathematical methods of statistics*, Princeton University Press.
- Durbin, J. (1954): “Errors in variables”, *Review International Statistical Institute* 22, 23–32.
- Fisher, R. A. (1922): “On the mathematical foundations of theoretical statistics”, *Philosophical Transactions of the Royal Society London Series A* 222, 309–368.
- Fisher, R. A. (1925): “Theory of statistical estimation”, *Proceedings of the Cambridge Philosophical Society* 22, 700–725.
- Fiorentini, G. and Sentana, E. (2019): “Consistent non-Gaussian pseudo maximum likelihood estimators”, *Journal of Econometrics* 213, 321–358.
- Fiorentini, G. and Sentana, E. (2021): “Specification tests for non-Gaussian maximum likelihood estimators”, *Quantitative Economics* 12, 683–742.
- Fiorentini, G. and Sentana, E. (2022a): “Discrete mixtures of normals pseudo maximum likelihood estimators of structural vector autoregressions”, forthcoming in the *Journal of Econo-*

metrics.

Fiorentini, G. and Sentana, E. (2022b): “Consistent estimation with finite mixtures”, mimeo, CEMFI.

Gallant, A.R. and Tauchen, G. (1999): “The relative efficiency of method of moments estimators”, *Journal of Econometrics* 92, 149–172.

Gauss, C.F. (1809): *Theoria motus corporum coelestium*, Perthes.

Gouriéroux, C., Monfort A. and Trognon, A. (1984): “Pseudo maximum likelihood methods: theory”, *Econometrica* 52, 681–700.

Hansen, L.P. (1982): “Large sample properties of generalized method of moments estimators”, *Econometrica* 50, 1029–1054.

Hausman, J. (1978): “Specification tests in econometrics”, *Econometrica* 46, 1273–1291.

Hodgson and Vorkink (2003): “Efficient estimation of conditional asset pricing models”, *Journal of Business and Economic Statistics*, 21, 269–283.

Huber, P. J. (1967): “The behavior of maximum likelihood estimates under nonstandard conditions”, in *Proceedings of the V Berkeley Symposium in Mathematical Statistics and Probability* 1, 221–233, University of California Press.

Islam, N. (1993): “Estimation of dynamic models from panel data”, unpublished Ph.D. dissertation, Harvard University.

Kraft, C.H., and Le Cam, L.M. (1956): “A remark on the roots of the maximum likelihood equation”, *Annals of Mathematical Statistics* 27, 1174–1177.

Laplace, P-S. (1774): “Mémoire sur la probabilité des causes par les évènements”, *Mémoires de l’Academie Royale des Sciences Présentés par Divers Savan* 6, 621–656.

Legendre, A-M. (1805): *Nouvelles méthodes pour la détermination des orbites des comètes*, F. Didot.

Magnus, J.R. and Neudecker, H. (2019): *Matrix differential calculus with applications in statistics and econometrics*, 3rd edition, Wiley.

Malinvaud, E. (1970): *Statistical methods in econometrics*, 2nd edition, North Holland.

Mardia, K.V. (1970): “Measures of multivariate skewness and kurtosis with applications”, *Biometrika* 57, 519–530.

Monés, M.A. and Ventura, E. (1996): “Saving decisions and fiscal incentives”, *Applied Economics* 28, 1105–1117

Newey, W.K. and Steigerwald, D.G. (1997): “Asymptotic bias for quasi-maximum-likelihood estimators in conditional heteroskedasticity models”, *Econometrica* 65, 587–99.

Neyman, J., and Scott, E.L. (1948): “Consistent estimation from partially consistent obser-

uations”, *Econometrica* 16, 1–32.

Nguyen, T.T., Nguyen, H.D., Chamroukhi, F. and McLachlan, G.J. (2020): “Approximation by finite mixtures of continuous density functions that vanish at infinity”, *Cogent Mathematics & Statistics* 7, 1750861.

Pearson, K. (1894): “Contributions to the mathematical theory of evolution”, *Philosophical Transactions of the Royal Society London Series A* 185, 71–110.

Pearson, K. (1900): “On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling”, *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science Series 5* 50, 157–175.

Pollock, D.S.G. (1988): “The estimation of linear stochastic models with covariance restrictions”, *Econometric Theory* 4, 403–427.

Rao, C.R. (1945): “Information and the accuracy attainable in the estimation of statistical parameters”, *Bulletin of the Calcutta Mathematical Society* 37, 81–89.

Rao, C.R. (1961): “Asymptotic efficiency and limiting information”, *Proceedings of the IV Berkeley Symposium on Mathematical Statistics and Probability* 1, 531–546, University of California Press.

Rothenberg, T.J. (1973): *Efficient estimation with a priori information*, Cowles Foundation Monograph 23, Yale.

Sentana, E. (2005): “Least squares predictions and mean-variance analysis”, *Journal of Financial Econometrics* 3, 56–78.

White, H. (1982): “Maximum likelihood estimation of misspecified models”, *Econometrica* 50, 1–25.

Wu, D-M. (1973): “Alternative tests of independence between stochastic regressors and disturbances”, *Econometrica*. 41, 733–750.

## Appendices

### A Asymptotic covariance matrices

#### A.1 Instrumental Variables (IV)

Let  $\mathbf{v}_i = (v_{1i}, v_{2i})'$  denote the reduced form innovations

$$\mathbf{v}_i = \mathbf{y}_i - \mathbf{C}\mathbf{z}_i = \mathbf{B}^{-1}\mathbf{u}_i,$$

where  $\mathbf{y}_i = (y_{1i}, y_{2i})'$  and  $\mathbf{z}_i = (z_{1i}, z_{2i})'$ , so that  $E(\mathbf{v}_i|\mathbf{z}_i) = \mathbf{0}$  and  $V(\mathbf{v}_i|\mathbf{z}_i) = \mathbf{B}^{-1}\mathbf{\Sigma}\mathbf{B}'^{-1} = \mathbf{\Omega}$ , with

$$\mathbf{B}'^{-1} = \begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix}.$$

In this context, the unrestricted Gaussian PMLE of  $\alpha$  and  $\beta$  coincides with the IV estimator that uses a constant,  $z_1$  and  $z_2$  as instruments in the first equation. To consider both equations at once, let  $\boldsymbol{\vartheta} = (\boldsymbol{\theta}', \sigma_{12})'$  and

$$\mathbf{Z}_{di}^U(\boldsymbol{\vartheta}) = [\mathbf{Z}_{li}^U(\boldsymbol{\vartheta}), \mathbf{Z}_{si}^U(\boldsymbol{\vartheta})], \quad (\text{A1})$$

where

$$\begin{aligned} \mathbf{Z}_{li}^U(\boldsymbol{\vartheta}) &= \frac{\partial \boldsymbol{\mu}'_i(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} \mathbf{\Omega}^{-\frac{1}{2}'}(\boldsymbol{\vartheta}), \\ \mathbf{Z}_{si}^U(\boldsymbol{\vartheta}) &= \frac{1}{2} \frac{\partial \text{vec}'[\mathbf{\Omega}(\boldsymbol{\vartheta})]}{\partial \boldsymbol{\vartheta}} \left[ \mathbf{\Omega}^{-\frac{1}{2}'}(\boldsymbol{\vartheta}) \otimes \mathbf{\Omega}^{-\frac{1}{2}'}(\boldsymbol{\vartheta}) \right], \\ \frac{\partial \boldsymbol{\mu}'_i(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} &= \begin{pmatrix} 1 & 0 \\ \mu_0 + \mu_1 z_{1i} + \mu_2 z_{2i} & 0 \\ z_{2i} & 0 \\ \alpha & 1 \\ \alpha z_{1i} & z_{1i} \\ \alpha z_{2i} & z_{2i} \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}, \quad \frac{\partial \text{vec}'[\mathbf{\Omega}(\boldsymbol{\vartheta})]}{\partial \boldsymbol{\vartheta}} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 2(\alpha\sigma_2^2 + \sigma_{12}) & \sigma_2^2 & \sigma_2^2 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ \alpha^2 & \alpha & \alpha & 1 \\ 2\alpha & 1 & 1 & 0 \end{pmatrix}, \end{aligned}$$

and

$$\mathbf{\Omega}^{-\frac{1}{2}}(\boldsymbol{\vartheta}) = \begin{pmatrix} \frac{1}{\sqrt{\sigma_1^2 + \alpha^2 \sigma_2^2 + 2\alpha\sigma_{12}}} & 0 \\ -\frac{\alpha\sigma_2^2 + \sigma_{12}}{\sigma_1^2 + \alpha^2 \sigma_2^2 + 2\alpha\sigma_{12}} \sqrt{\frac{\sigma_1^2 \sigma_2^2 - \sigma_{12}^2}{\sigma_1^2 + \alpha^2 \sigma_2^2 + 2\alpha\sigma_{12}}} & \sqrt{\frac{\sigma_1^2 + \alpha^2 \sigma_2^2 + 2\alpha\sigma_{12}}{\sigma_1^2 \sigma_2^2 - \sigma_{12}^2}} \end{pmatrix}$$

is the inverse of the (lower) Cholesky decomposition of  $\mathbf{\Omega}$ .

We can then exploit Proposition C2 in Supplementary Appendix C of Fiorentini and Sentana (2021) to obtain

$$AVar(\sqrt{n}\tilde{\boldsymbol{\vartheta}}_{IV}) = [\mathcal{A}_{U,\boldsymbol{\vartheta}\boldsymbol{\vartheta}}(\boldsymbol{\vartheta})]^{-1} \mathcal{B}_{U,\boldsymbol{\vartheta}\boldsymbol{\vartheta}}(\boldsymbol{\vartheta}, \boldsymbol{\varrho}) [\mathcal{A}_{U,\boldsymbol{\vartheta}\boldsymbol{\vartheta}}(\boldsymbol{\vartheta})]^{-1}, \quad (\text{A2})$$

where

$$\mathcal{A}_{U,\vartheta\vartheta}(\vartheta) = E[\mathbf{Z}_{di}^U(\vartheta)\mathcal{K}(\mathbf{0})\mathbf{Z}_{di}^{U'}(\vartheta)] \quad \text{and} \quad \mathcal{B}_{U,\vartheta\vartheta}(\vartheta, \varrho) = E[\mathbf{Z}_{di}^U(\vartheta)\mathcal{K}^\mathbf{v}(\vartheta, \varrho)\mathbf{Z}_{di}^{U'}(\vartheta)],$$

with

$$\mathcal{K}^\mathbf{v}(\vartheta, \varrho) = V[\mathbf{e}_{di}(\vartheta, \mathbf{0})] = \begin{bmatrix} \mathbf{I}_2 & \Phi^\mathbf{v}(\vartheta, \varrho) \\ \Phi^{\mathbf{v}'}(\vartheta, \varrho) & \Upsilon^\mathbf{v}(\vartheta, \varrho) \end{bmatrix}, \quad (\text{A3})$$

$\Phi^\mathbf{v}(\vartheta_0, \varrho_0) = E[\mathbf{v}_i^* \text{vec}'(\mathbf{v}_i^* \mathbf{v}_i^{*'})]$ ,  $\Upsilon^\mathbf{v}(\vartheta_0, \varrho_0) = E[\text{vec}(\mathbf{v}_i^* \mathbf{v}_i^{*'} - \mathbf{I}_2) \text{vec}'(\mathbf{v}_i^* \mathbf{v}_i^{*'} - \mathbf{I}_2)]$  and  $\mathbf{v}_i^* = \Omega^{-1/2} \mathbf{v}_i$ , so that  $\Phi^\mathbf{v}(\mathbf{0}) = \mathbf{0}$  and  $\Upsilon^\mathbf{v}(\mathbf{0}) = (\mathbf{I}_4 + \mathbf{K}_{22})$  if we use  $\varrho = \mathbf{0}$  to denote normality and  $\mathbf{K}_{mn}$  for the commutation matrix of orders  $m$  and  $n$  (see e.g. Magnus and Neudecker (2019)).

Given that the assumption of constant conditional higher-order cumulants applies to the structural model, though, we need to relate the higher-order moments of the reduced form residuals to those of the structural ones. Defining

$$\mathbf{F}(\theta) = \mathbf{L}_2[\mathbf{B}^{-1}(\theta) \otimes \mathbf{B}^{-1}(\theta)]\mathbf{D}_2 = \begin{bmatrix} 1 & 2\alpha & \alpha^2 \\ 0 & 1 & \alpha \\ 0 & 0 & 1 \end{bmatrix},$$

where  $\mathbf{L}_2$  and  $\mathbf{D}_2$  are the elimination and duplication matrices of order 2, respectively (see Magnus and Neudecker (2019)), we will have that

$$E[\mathbf{v}_i \text{vec}'(\mathbf{v}_i \mathbf{v}_i')] = -\mathbf{B}^{-1}(\theta) \Sigma(\theta)^{\frac{1}{2}} \Phi^\mathbf{u}(\varrho) [\Sigma(\theta)^{\frac{1}{2}'} \otimes \Sigma(\theta)^{\frac{1}{2}'}] \mathbf{F}'(\theta)$$

and

$$E[\text{vec}(\mathbf{v}_i \mathbf{v}_i' - \mathbf{I}_2) \text{vec}'(\mathbf{v}_i \mathbf{v}_i' - \mathbf{I}_2)] = \mathbf{F}(\theta) [\Sigma(\theta)^{\frac{1}{2}} \otimes \Sigma(\theta)^{\frac{1}{2}}] \Upsilon^\mathbf{u}(\varrho) [\Sigma(\theta)^{\frac{1}{2}'} \otimes \Sigma(\theta)^{\frac{1}{2}'}] \mathbf{F}'(\theta),$$

where  $\Phi^\mathbf{u}(\varrho_0) = E[\mathbf{u}_i^* \text{vec}'(\mathbf{u}_i^* \mathbf{u}_i^{*'})]$ ,  $\Upsilon^\mathbf{u}(\varrho_0) = E[\text{vec}(\mathbf{u}_i^* \mathbf{u}_i^{*'} - \mathbf{I}_2) \text{vec}'(\mathbf{u}_i^* \mathbf{u}_i^{*'} - \mathbf{I}_2)]$  and  $\mathbf{u}_i^* = \Sigma^{-1/2} \mathbf{u}_i$ .

After some tedious calculations, it is straightforward to prove that

$$AVar(\sqrt{n}\tilde{\alpha}_{IV}) = \frac{\sigma_1^2 \sigma_{z_1}^2}{\mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}$$

and

$$AVar(\sqrt{n}\tilde{\beta}_{IV}) = \frac{\sigma_1^2 (\mu_1^2 \sigma_{z_1}^2 + \mu_2^2 \sigma_{z_2}^2 + 2\mu_1 \mu_2 \sigma_{z_1 z_2})}{\mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}.$$

For our purposes, it is convenient to rewrite these expressions as

$$AVar(\sqrt{n}\tilde{\alpha}_{IV}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)}{(1 - R_2^2)\rho_{y_2 z_2, z_1}^2}$$

and

$$AVar(\sqrt{n}\tilde{\beta}_{IV}) = \frac{R_2^2(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)}{(1 - R_2^2)\rho_{y_2 z_2, z_1}^2},$$

where  $R_1^2$  and  $R_2^2$  are the population coefficients of determination in equations (1) and (2), respectively, and  $\rho_{y_2 z_2 \cdot z_1}$  the correlation coefficient between  $y_2$  and  $z_2$  after partialling out the effect of  $z_1$ .

## A.2 Ordinary Least Squares (LS)

As mentioned in Section 2, the restricted Gaussian PMLE that imposes  $\sigma_{12} = 0$  coincides with the OLS estimator of the first equation. To consider both equations at once, let

$$\mathbf{Z}_{di}^R(\boldsymbol{\theta}) = (\mathbf{I}_8, \mathbf{0}_{8 \times 1}) \mathbf{Z}_{di}^U(\boldsymbol{\vartheta}, 0). \quad (\text{A4})$$

Then, analogous calculations to the ones in the previous subsection imply that

$$AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{LS}) = \mathcal{A}_{R, \boldsymbol{\theta}\boldsymbol{\theta}}^{-1}(\boldsymbol{\theta}) \mathcal{B}_{R, \boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{\theta}, \boldsymbol{\varrho}) \mathcal{A}_{R, \boldsymbol{\theta}\boldsymbol{\theta}}^{-1}(\boldsymbol{\theta}), \quad (\text{A5})$$

where

$$\mathcal{A}_{R, \boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{\theta}) = E [\mathbf{Z}_{di}^R(\boldsymbol{\theta}) \mathcal{K}^v(\mathbf{0}) \mathbf{Z}_{di}^{R'}(\boldsymbol{\theta})] \quad \text{and} \quad \mathcal{B}_{R, \boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{\theta}, \boldsymbol{\varrho}) = E [\mathbf{Z}_{di}^R(\boldsymbol{\theta}) \mathcal{K}^v(\boldsymbol{\theta}, \boldsymbol{\varrho}) \mathbf{Z}_{di}^{R'}(\boldsymbol{\theta})],$$

After some straightforward calculations, it is easy to show that

$$AVar(\sqrt{n}\hat{\alpha}_{LS}) = \frac{\sigma_1^2 \sigma_{z_1} (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) \mu_2^2}{[\mu_2^2 \sigma_{z_1 z_2}^2 - \sigma_{z_1}^2 (\sigma_2^2 + \mu_2^2 \sigma_{z_2}^2)]^2} + \frac{\sigma_1^2 \sigma_2^2 \sigma_{z_1}^4 \mu_{22}}{[\mu_2^2 \sigma_{z_1 z_2}^2 - \sigma_{z_1}^2 (\sigma_2^2 + \mu_2^2 \sigma_{z_2}^2)]^2}$$

and

$$\begin{aligned} AVar(\sqrt{n}\hat{\beta}_{LS}) &= \frac{\sigma_1^2 \{ \sigma_{z_2}^2 \mu_2^2 [\sigma_{z_1}^4 \mu_1^2 + 2\sigma_{z_1} (\sigma_2^2 + \sigma_{z_1 z_2} \mu_1 \mu_2) - \sigma_{z_1 z_2}^2 \mu_2^2] + \sigma_{z_1}^2 \sigma_{z_2}^4 \mu_2^4 \}}{[\sigma_2^2 \sigma_{z_1}^2 + \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]^2} \\ &+ \frac{\sigma_1^2 (\sigma_2^2 + \mu_1 \mu_2 \sigma_{z_1 z_2}) [\sigma_2^2 \sigma_{z_1}^2 - \sigma_{z_1 z_2} (\sigma_{z_1}^2 \mu_1 \mu_2 + 2\sigma_{z_1 z_2} \mu_2^2)]}{[\sigma_2^2 \sigma_{z_1}^2 + \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]^2} \\ &+ \frac{\sigma_1^2 \sigma_2^2 (\sigma_{z_1}^2 \mu_1 + \sigma_{z_1 z_2} \mu_2)^2 \mu_{22}}{[\sigma_2^2 \sigma_{z_1}^2 + \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]^2}. \end{aligned}$$

Again, it is convenient to rewrite these expressions as

$$AVar(\sqrt{n}\hat{\alpha}_{LS}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2 \cdot z_1}^2) [\mu_{22}(1 - \rho_{y_2 z_2 \cdot z_1}^2) + \rho_{y_2 z_2 \cdot z_1}^2]}{1 - R_2^2}$$

and

$$AVar(\sqrt{n}\hat{\beta}_{LS}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2 \cdot z_1}^2) [1 + (\mu_{22} - 1)(R_2^2 - \rho_{y_2 z_2 \cdot z_1}^2)]}{1 - R_2^2},$$

### A.3 Optimum Minimum Distance (MD)

Let  $\mathbf{c} = \text{vec}(\mathbf{C})$  and  $\boldsymbol{\omega} = \text{vech}(\boldsymbol{\Omega})$  denote the parameters of the unrestricted reduced form model. From equations (5)-(6), we will have that

$$\begin{aligned} c_{10} &= \gamma + \alpha\mu_0, & c_{20} &= \mu_0, & \omega_{11} &= \sigma_{11} + \alpha^2\sigma_{22} + 2\alpha\sigma_{12}, \\ c_{11} &= \beta + \alpha\mu_1, & c_{21} &= \mu_1, & \omega_{12} &= \alpha\sigma_{22} + \sigma_{12}, \\ c_{12} &= \alpha\mu_2, & c_{22} &= \mu_2, & \omega_{22} &= \sigma_{22}. \end{aligned},$$

Let  $\tilde{\boldsymbol{\phi}}_{LS} = (\tilde{c}_{10}, \tilde{c}_{11}, \tilde{c}_{12}, \tilde{c}_{20}, \tilde{c}_{21}, \tilde{c}_{22}, \tilde{\omega}_{11}, \tilde{\omega}_{12}, \tilde{\omega}_{22})'$  denote their unrestricted Gaussian PML estimators, which coincide with equation by equation OLS. To obtain the asymptotic distributions of these estimators, we need the first derivatives of the conditional mean vector and covariance matrix with respect to the unrestricted reduced form parameters, which are given by

$$\frac{\partial \mathbf{C}\mathbf{z}_i}{\partial \mathbf{c}'} = \mathbf{z}_i' \otimes \mathbf{I}_2 \quad \text{and} \quad \frac{\partial \text{vec}[\boldsymbol{\Omega}(\boldsymbol{\theta})]}{\partial \boldsymbol{\omega}'} = \mathbf{D}_2.$$

In this notation, the contribution to the Gaussian log-likelihood scores for  $\mathbf{c}$  and  $\boldsymbol{\omega}$  corresponding to observation  $i$  will be given by

$$\mathbf{s}_{\mathbf{c}i}(\mathbf{c}, \boldsymbol{\omega}) = \mathbf{z}_i \otimes \boldsymbol{\Omega}^{-1}(\boldsymbol{\theta}) \mathbf{v}_i(\mathbf{c})$$

and

$$\mathbf{s}_{\boldsymbol{\omega}i}(\mathbf{c}, \boldsymbol{\omega}) = \frac{1}{2} \mathbf{D}_2' \text{vec}[\boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \mathbf{v}_i(\mathbf{c}) \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} - \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}].$$

Consequently, the outer product of the scores will be

$$\begin{aligned} \mathbf{s}_{\mathbf{c}i}(\mathbf{c}, \boldsymbol{\omega}) \mathbf{s}_{\mathbf{c}i}'(\mathbf{c}, \boldsymbol{\omega}) &= \mathbf{z}_i \mathbf{z}_i' \otimes \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \mathbf{v}_i(\mathbf{c}) \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}, \\ \mathbf{s}_{\boldsymbol{\omega}i}(\mathbf{c}, \boldsymbol{\omega}) \mathbf{s}_{\boldsymbol{\omega}i}'(\mathbf{c}, \boldsymbol{\omega}) &= \frac{1}{2} \mathbf{D}_2' \text{vec}[\boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \mathbf{v}_i(\mathbf{c}) \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} - \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}] [\mathbf{z}_i' \otimes \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}] \end{aligned}$$

and

$$\begin{aligned} \mathbf{s}_{\boldsymbol{\omega}i}(\mathbf{c}, \boldsymbol{\omega}) \mathbf{s}_{\boldsymbol{\omega}i}'(\mathbf{c}, \boldsymbol{\omega}) &= \frac{1}{4} \mathbf{D}_2' \text{vec}[\boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \mathbf{v}_i(\mathbf{c}) \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} - \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}] \\ &\quad \times \text{vec}'[\boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \mathbf{v}_i(\mathbf{c}) \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} - \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}] \mathbf{D}_2. \end{aligned}$$

Similarly, we can easily adapt the expressions in Amengual, Fiorentini and Sentana (2022) to write the contribution of observation  $i$  to the Hessian matrix  $\mathbf{h}_{\mathbf{c}, \boldsymbol{\omega}i}(\mathbf{c}, \boldsymbol{\omega})$  as

$$= - \left\{ \begin{array}{cc} (\mathbf{z}_i \mathbf{z}_i' \otimes \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}) & [\mathbf{z}_i \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}^{-1}(\boldsymbol{\theta}) \otimes \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}] \mathbf{D}_2 \\ \mathbf{D}_2' [\boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \mathbf{v}_i(\mathbf{c}) \mathbf{z}_i' \otimes \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}] & \mathbf{D}_2' \{ \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \otimes [\boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} \mathbf{v}_i(\mathbf{c}) \mathbf{v}_i'(\mathbf{c}) \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1} - \frac{1}{2} \boldsymbol{\Omega}(\boldsymbol{\theta})^{-1}] \} \mathbf{D}_2 \end{array} \right\}.$$

Thus, we have all the ingredients to compute  $AVar(\sqrt{n}\tilde{\boldsymbol{\phi}}_{LS})$  using the standard sandwich formula in White (1982) and Gouriéroux, Monfort and Trognon (1984).

On this basis, we can show that the asymptotic variance of Malinvaud's (1970) optimum MD

estimator will be given by

$$AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{MD}) = \left\{ \frac{\partial \boldsymbol{\phi}'(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \left[ AVar(\sqrt{n}\tilde{\boldsymbol{\phi}}_{LS}) \right]^{-1} \frac{\partial \boldsymbol{\phi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'} \right\}^{-1}, \quad (\text{A6})$$

where

$$\boldsymbol{\phi}(\boldsymbol{\theta}) = \begin{pmatrix} c_{10} - \gamma - \alpha\mu_0 \\ c_{11} - \beta - \alpha\mu_1 \\ c_{12} - \alpha\mu_2 \\ c_{20} - \mu_0 \\ c_{21} - \mu_1 \\ c_{22} - \mu_2 \\ \omega_{11} - \sigma_{11} - \alpha^2\sigma_{22} \\ \omega_{12} - \alpha\sigma_{22} \\ \omega_{22} - \sigma_{22} \end{pmatrix}.$$

Specifically, we obtain that

$$AVar(\sqrt{n}\hat{\alpha}_{MD}) = \frac{\sigma_1^2 \sigma_{z_1}^2 \mu_{22}}{\sigma_{z_1}^2 \sigma_2^2 + (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) \mu_2^2 \mu_{22}}$$

and

$$AVar(\sqrt{n}\hat{\beta}_{MD}) = \frac{\sigma_1^2 [\sigma_2^2 + (\sigma_{z_1}^2 \mu_1^2 + \sigma_{z_2}^2 \mu_2^2 + 2\sigma_{z_1 z_2} \mu_1 \mu_2) \mu_{22}]}{\sigma_{z_1}^2 \sigma_2^2 + (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) \mu_2^2 \mu_{22}},$$

which, rewritten in terms of the population coefficients of determination, become

$$AVar(\sqrt{n}\hat{\alpha}_{MD}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2) \mu_{22}}{(1 - R_2^2)[1 + \rho_{y_2 z_2, z_1}^2 (\mu_{22} - 1)]}$$

and

$$AVar(\sqrt{n}\hat{\beta}_{MD}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)[1 + R_2^2(\mu_{22} - 1)]}{(1 - R_2^2)[1 + \rho_{y_2 z_2, z_1}^2 (\mu_{22} - 1)]}.$$

#### A.4 Maximum likelihood with spherical innovations

Let  $\exp[c(\boldsymbol{\eta}) + g(\varsigma_t, \boldsymbol{\eta})]$  denote the assumed conditional density of  $\mathbf{v}_t^*$  given  $\mathbf{z}_t$ , where  $c(\boldsymbol{\eta})$  corresponds to the constant of integration,  $g(\varsigma_t, \boldsymbol{\eta})$  to its kernel and  $\varsigma_i = \mathbf{v}_i^{*'} \mathbf{v}_i^*$ . Invoking Proposition C1 in Supplementary Appendix C of Fiorentini and Sentana (2021), we can obtain the asymptotic variance of the ML estimator that imposes  $\sigma_{12} = 0$  when the true distribution of the shocks is spherically symmetric by

$$AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{ML}) = \mathcal{I}_R^{-1}(\boldsymbol{\theta}, \boldsymbol{\eta}), \text{ where } \mathcal{I}_R(\boldsymbol{\theta}, \boldsymbol{\eta}) = E[\mathbf{Z}_i^R(\boldsymbol{\theta}) \mathcal{M}(\boldsymbol{\eta}) \mathbf{Z}_i^{R'}(\boldsymbol{\theta})], \quad (\text{A7})$$

$$\mathbf{Z}_i^R(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{Z}_{di}^R(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_q \end{pmatrix}, \quad \mathcal{M}(\boldsymbol{\eta}) = \begin{pmatrix} \mathcal{M}_{ll}(\boldsymbol{\eta}) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathcal{M}_{ss}(\boldsymbol{\eta}) & \mathcal{M}_{sr}(\boldsymbol{\eta}) \\ \mathbf{0} & \mathcal{M}'_{sr}(\boldsymbol{\eta}) & \mathcal{M}_{rr}(\boldsymbol{\eta}) \end{pmatrix},$$

$$\begin{aligned}
\mathcal{M}_{ll}(\boldsymbol{\eta}) &= \mathbf{M}_{ll} \mathbf{I}_N, \\
\mathcal{M}_{ss}(\boldsymbol{\eta}) &= \mathbf{M}_{ss} (\mathbf{I}_{N^2} + \mathbf{K}_{NN}) + [\mathbf{M}_{ss} - 1] \text{vec}(\mathbf{I}_N) \text{vec}'(\mathbf{I}_N), \\
\mathcal{M}_{sr}(\boldsymbol{\eta}) &= \text{vec}(\mathbf{I}_N) \mathbf{M}_{sr},
\end{aligned}$$

with

$$\begin{aligned}
\mathbf{M}_{ll} &= E \left[ \delta^2(\varsigma_i, \boldsymbol{\eta}) \frac{\varsigma_i}{N} \right], \\
\mathbf{M}_{ss} &= 1 + \frac{N}{N+2} E \left[ \frac{2\partial\delta(\varsigma_i, \boldsymbol{\eta})}{\partial\varsigma} \left( \frac{\varsigma_i}{N} \right)^2 \right], \\
\mathbf{M}_{sr} &= -E \left[ \frac{\varsigma_i}{N} \frac{\partial\delta(\varsigma_i, \boldsymbol{\eta})}{\partial\boldsymbol{\eta}'} \right],
\end{aligned}$$

and  $\delta(\varsigma_i, \boldsymbol{\eta}) = -2\partial g(\varsigma_i, \boldsymbol{\eta})/\partial\varsigma$ .

Similarly, we can compute the asymptotic variance of the unrestricted ML estimator which also estimates  $\sigma_{12}$  as

$$AVar(\sqrt{n}\tilde{\boldsymbol{\theta}}_{ML}) = \mathcal{I}_U^{-1}(\boldsymbol{\theta}, \boldsymbol{\eta}), \quad \text{where } \mathcal{I}_U(\boldsymbol{\theta}, \boldsymbol{\eta}) = E[\mathbf{Z}_i^U(\boldsymbol{\theta}) \mathcal{M}(\boldsymbol{\eta}) \mathbf{Z}_i^{U'}(\boldsymbol{\theta})], \quad (\text{A8})$$

with

$$\mathbf{Z}_i^U(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{Z}_{di}^U(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_q \end{pmatrix}.$$

As a consequence,

$$AVar(\sqrt{n}\hat{\alpha}_{ML}) = \frac{\sigma_1^2 \sigma_{z_1}^2}{\mathbf{M}_{ss} \sigma_2^2 \sigma_{z_1}^2 + \mathbf{M}_{ll} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}$$

and

$$AVar(\sqrt{n}\hat{\beta}_{ML}) = \frac{\sigma_1^2 [\mathbf{M}_{ss} \sigma_2^2 + \mathbf{M}_{ll} (\mu_1^2 \sigma_{z_1}^2 + \mu_2^2 \sigma_{z_2}^2 + 2\mu_1 \mu_2 \sigma_{z_1 z_2})]}{\mathbf{M}_{ll} [\mathbf{M}_{ss} \sigma_2^2 \sigma_{z_1}^2 + \mathbf{M}_{ll} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]}.$$

Analogous calculations but using  $\mathbf{Z}_t^U(\boldsymbol{\theta})$  in place of  $\mathbf{Z}_t^R(\boldsymbol{\theta})$  for the unrestricted ML estimator yield

$$AVar(\sqrt{n}\tilde{\alpha}_{ML}) = \frac{\sigma_1^2 \sigma_{z_1}^2}{\mathbf{M}_{ll} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}$$

and

$$AVar(\sqrt{n}\tilde{\beta}_{ML}) = \frac{\sigma_1^2 (\mu_1^2 \sigma_{z_1} + \mu_2^2 \sigma_{z_2}^2 + 2\mu_1 \mu_2 \sigma_{z_1 z_2})}{\mathbf{M}_{ll} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}.$$

Once again, we can write these expressions as

$$AVar(\sqrt{n}\hat{\alpha}_{ML}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)}{(1 - R_2^2)[(1 - \rho_{y_2 z_2, z_1}^2) \mathbf{M}_{ss} + \rho_{y_2 z_2, z_1}^2 \mathbf{M}_{ll}]}$$

and

$$AVar(\sqrt{n}\hat{\beta}_{ML}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)[R_2^2 \mathbf{M}_{ll} + (1 - R_2^2) \mathbf{M}_{ss}]}{(1 - R_2^2) \mathbf{M}_{ll} [(1 - \rho_{y_2 z_2, z_1}^2) \mathbf{M}_{ss} + \rho_{y_2 z_2, z_1}^2 \mathbf{M}_{ll}]},$$

for the restricted estimator, and as

$$AVar(\sqrt{n}\tilde{\alpha}_{ML}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)}{(1 - R_2^2)\rho_{y_2 z_2, z_1}^2 M_{ll}}$$

and

$$AVar(\sqrt{n}\tilde{\beta}_{ML}) = \frac{R_2^2(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)}{(1 - R_2^2)\rho_{y_2 z_2, z_1}^2 M_{ll}},$$

for the unrestricted one.

## A.5 Spherically symmetric semiparametric estimator (SS)

From Proposition C3 in Supplementary Appendix C of Fiorentini and Sentana (2021), the spherically symmetric SP efficiency bound is given by

$$\hat{\mathcal{S}}_j(\boldsymbol{\theta}) = \mathcal{I}_{j, \boldsymbol{\theta}\boldsymbol{\theta}}(\phi_0) - \mathbf{W}_s^j(\boldsymbol{\theta})\mathbf{W}_s^{j'}(\boldsymbol{\theta}) \cdot \left\{ [2M_{ss} - 1] - \frac{2}{4\kappa + 2} \right\}$$

where  $\kappa_0$  denotes the population coefficient of multivariate excess kurtosis (see Mardia (1970) for details),

$$\mathbf{W}_s^j(\boldsymbol{\theta}) = \mathbf{Z}_d^j(\boldsymbol{\theta})[\mathbf{0}', \text{vec}'(\mathbf{I}_2)]' \quad \text{for } j = R, U,$$

and

$$\mathcal{I}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{\theta}, \boldsymbol{\rho}) = E \left[ \mathbf{Z}_{dt}(\boldsymbol{\theta}) \mathcal{M}_{dd}(\boldsymbol{\theta}, \boldsymbol{\rho}) \mathbf{Z}_{dt}'(\boldsymbol{\theta}) \right].$$

Under suitable regularity conditions, we have that

$$AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{SS}) = [\hat{\mathcal{S}}_R(\boldsymbol{\theta})]^{-1} \tag{A9}$$

and

$$AVar(\sqrt{n}\tilde{\boldsymbol{\theta}}_{SS}) = [\hat{\mathcal{S}}_U(\boldsymbol{\theta})]^{-1}. \tag{A10}$$

Tedious but otherwise straightforward calculations show that for the restricted estimator that imposes  $\sigma_{12} = 0$  we obtain

$$AVar(\sqrt{n}\hat{\alpha}_{SS}) = AVar(\sqrt{n}\hat{\alpha}_{ML}) \quad \text{and} \quad AVar(\sqrt{n}\hat{\beta}_{SS}) = AVar(\sqrt{n}\hat{\beta}_{ML}),$$

while for the unrestricted one we get

$$AVar(\sqrt{n}\tilde{\alpha}_{SS}) = AVar(\sqrt{n}\tilde{\alpha}_{ML}) \quad \text{and} \quad AVar(\sqrt{n}\tilde{\beta}_{SS}) = AVar(\sqrt{n}\tilde{\beta}_{ML}).$$

## A.6 Maximum likelihood with general innovations

Let  $f(\mathbf{v}_t^*, \boldsymbol{\varrho})$  denote the joint density of  $\mathbf{v}_t^*$  conditional on  $\mathbf{z}_t$ , where  $\boldsymbol{\varrho}$  are some  $q$  additional parameters that determine the shape of the distribution. If we use Proposition D3 in Supple-

mentary Appendix D of Fiorentini and Sentana (2021), we can obtain the asymptotic variance of the ML estimator that imposes  $\sigma_{12} = 0$  by computing

$$AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{ML}) = \mathcal{I}_{GR}^{-1}(\boldsymbol{\theta}, \boldsymbol{\varrho}), \text{ where } \mathcal{I}_{GR}(\boldsymbol{\theta}, \boldsymbol{\varrho}) = E[\mathbf{Z}_i^{GR}(\boldsymbol{\theta})\mathcal{M}(\boldsymbol{\varrho})\mathbf{Z}_i^{GR'}(\boldsymbol{\theta})],$$

where

$$\begin{aligned} \mathbf{Z}_{di}^{GR}(\boldsymbol{\theta}) &= [\mathbf{Z}_{li}^R(\boldsymbol{\theta}), \mathbf{Z}_{si}^{GR}(\boldsymbol{\theta})], \\ \mathbf{Z}_{si}^{GR}(\boldsymbol{\theta}) &= \frac{\partial vec'[\boldsymbol{\Omega}^{\frac{1}{2}}(\boldsymbol{\theta})]}{\partial \boldsymbol{\theta}} \left[ \mathbf{I}_2 \otimes \boldsymbol{\Omega}^{-\frac{1}{2}}(\boldsymbol{\theta}) \right], \\ \frac{\partial vec'[\boldsymbol{\Omega}^{\frac{1}{2}}(\boldsymbol{\theta})]}{\partial \boldsymbol{\theta}} &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ \frac{\alpha\sigma_2^2}{\sqrt{\sigma_1^2+\alpha\sigma_2^2}} & \frac{\sigma_1^2\sigma_2^2}{(\sigma_1^2+\alpha\sigma_2^2)^{3/2}} & 0 & -\frac{\alpha}{\sigma_1^2} \left( \frac{\sigma_1^2\sigma_2^2}{\sigma_1^2+\alpha\sigma_2^2} \right)^{3/2} \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2\sqrt{\sigma_1^2+\alpha\sigma_2^2}} & -\frac{\alpha\sigma_2^2}{2(\sigma_1^2+\alpha\sigma_2^2)^{3/2}} & 0 & \frac{\alpha^2\sigma_2^2}{2(\sigma_1^2+\alpha^2\sigma_1^2\sigma_2^2)} \sqrt{\frac{\sigma_1^2\sigma_2^2}{\sigma_1^2+\alpha\sigma_2^2}} \\ \frac{\alpha^2}{2\sqrt{\sigma_1^2+\alpha\sigma_2^2}} & \frac{2\alpha\sigma_1^2+\alpha^3\sigma_2^2}{2(\sigma_1^2+\alpha\sigma_2^2)^{3/2}} & 0 & \frac{\sigma_1^2}{2(\sigma_1^2\sigma_2^2+\alpha^2\sigma_2^2)} \sqrt{\frac{\sigma_1^2\sigma_2^2}{\sigma_1^2+\alpha\sigma_2^2}} \end{bmatrix}, \end{aligned} \tag{A11}$$

and

$$\mathcal{M}(\boldsymbol{\varrho}) = \begin{bmatrix} \mathcal{M}_{ll}(\boldsymbol{\varrho}) & \mathcal{M}_{ls}(\boldsymbol{\varrho}) & \mathcal{M}_{lr}(\boldsymbol{\varrho}) \\ \mathcal{M}'_{ls}(\boldsymbol{\varrho}) & \mathcal{M}_{ss}(\boldsymbol{\varrho}) & \mathcal{M}_{sr}(\boldsymbol{\varrho}) \\ \mathcal{M}'_{lr}(\boldsymbol{\varrho}) & \mathcal{M}'_{sr}(\boldsymbol{\varrho}) & \mathcal{M}_{rr}(\boldsymbol{\varrho}) \end{bmatrix},$$

with

$$\begin{aligned} \mathcal{M}_{ll}(\boldsymbol{\varrho}) &= E \left[ \partial^2 \ln f(\mathbf{v}_i^*; \boldsymbol{\varrho}) / \partial \mathbf{v}^* \partial \mathbf{v}^{*'} | \boldsymbol{\varrho} \right], \\ \mathcal{M}_{ls}(\boldsymbol{\varrho}) &= E \left[ \partial^2 \ln f(\mathbf{v}_i^*; \boldsymbol{\varrho}) / \partial \mathbf{v}^* \partial \mathbf{v}^{*'} \cdot (\mathbf{v}^{*'} \otimes \mathbf{I}_2) | \boldsymbol{\varrho} \right], \\ \mathcal{M}_{ss}(\boldsymbol{\varrho}) &= E \left[ (\mathbf{v}_i^* \otimes \mathbf{I}_2) \cdot \partial^2 \ln f(\mathbf{v}_i^*; \boldsymbol{\varrho}) / \partial \mathbf{v}_i^* \partial \mathbf{v}_i^{*'} \cdot (\mathbf{v}_i^{*'} \otimes \mathbf{I}_2) | \boldsymbol{\varrho} \right] - \mathbf{K}_{22}, \\ \mathcal{M}_{lr}(\boldsymbol{\varrho}) &= -E \left[ \partial^2 \ln f(\mathbf{v}_i^*; \boldsymbol{\varrho}) / \partial \mathbf{v}^* \partial \boldsymbol{\varrho}' | \boldsymbol{\varrho} \right], \\ \mathcal{M}_{sr}(\boldsymbol{\varrho}) &= -E \left[ (\mathbf{v}_i^* \otimes \mathbf{I}_2) \partial^2 \ln f(\mathbf{v}_i^*; \boldsymbol{\varrho}) / \partial \mathbf{v}^* \partial \boldsymbol{\varrho}' | \boldsymbol{\varrho} \right], \end{aligned}$$

and

$$\mathcal{M}_{rr}(\boldsymbol{\varrho}) = -E \left[ \partial^2 \ln f(\mathbf{v}_i^*; \boldsymbol{\varrho}) / \partial \boldsymbol{\varrho} \partial \boldsymbol{\varrho}' | \boldsymbol{\varrho} \right].$$

Analogously, we can obtain  $AVar(\sqrt{n}\tilde{\boldsymbol{\theta}}_{ML}) = \mathcal{I}_U^{-1}(\boldsymbol{\theta}, \boldsymbol{\varrho})$  by exploiting the expressions for the derivatives of the unrestricted model that we obtained when we discussed the IV estimators.

## A.7 Semiparametric estimator (SP)

We can make use of Proposition D3 in Supplementary Appendix D of Fiorentini and Sentana (2021), which indicates that the SP efficiency bound for  $j = R, U$  will be given by

$$\ddot{S}_j(\phi) = \mathcal{I}_{\theta\theta}(\theta, \varrho) - \mathbf{Z}_d^{Gj}(\theta) [\mathcal{M}_{dd}(\varrho) - \mathcal{K}(0)\mathcal{K}^{\mathbf{v}+}(\varrho)\mathcal{K}(0)] \mathbf{Z}_d^{Gj'}(\theta), \quad (\text{A12})$$

where  $+$  denotes the Moore-Penrose inverse, with

$$\mathcal{M}_{dd}(\varrho) = \begin{pmatrix} \mathcal{M}_{ll}(\varrho) & \mathcal{M}_{ls}(\varrho) \\ \mathcal{M}'_{ls}(\varrho) & \mathcal{M}_{ss}(\varrho) \end{pmatrix}$$

and the matrix of third and fourth order central moments  $\mathcal{K}^{\mathbf{v}}(\varrho)$  in (A3). Then, under suitable regularity conditions, we will have that

$$AVar(\sqrt{n}\hat{\theta}_{SP}) = [\ddot{S}_R(\theta)]^{-1} \quad (\text{A13})$$

and

$$AVar(\sqrt{n}\tilde{\theta}_{SP}) = [\ddot{S}_U(\theta)]^{-1}. \quad (\text{A14})$$

The expression for  $\mathcal{K}^{\mathbf{v}}(\rho)$  simplifies considerably in the spherically symmetric case because

$$E(\mathbf{v}_i^* \mathbf{v}_i^{*'} \otimes \mathbf{v}_i^*) = \mathbf{0}, \quad (\text{A15})$$

$$E(\mathbf{v}_i^* \mathbf{v}_i^{*'} \otimes \mathbf{v}_i^* \mathbf{v}_i^{*'}) = E[vec(\mathbf{v}_i^* \mathbf{v}_i^{*'}) vec'(\mathbf{v}_i^* \mathbf{v}_i^{*'})] = (\kappa_0 + 1)[(\mathbf{I}_4 + \mathbf{K}_{22}) + vec(\mathbf{I}_2) vec'(\mathbf{I}_2)]. \quad (\text{A16})$$

As a result, after some tedious calculations we obtain that for the estimator that imposes the restriction  $\sigma_{12} = 0$ ,

$$AVar(\sqrt{n}\hat{\alpha}_{SP}) = \frac{\sigma_1^2 \sigma_{z_1}^2 (1 + \kappa)}{\sigma_2^2 \sigma_{z_1}^2 + M_{ll}(1 + \kappa) \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}$$

and

$$AVar(\sqrt{n}\hat{\beta}_{SP}) = \frac{\sigma_1^2 [\sigma_2^2 + (1 + \kappa) M_{ll} (\mu_1^2 \sigma_{z_1}^2 + \mu_2^2 \sigma_{z_2}^2 + 2\mu_1 \mu_2 \sigma_{z_1 z_2})]}{M_{ll} [\sigma_2^2 \sigma_{z_1}^2 + (1 + \kappa) \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]},$$

while for the unrestricted one,

$$AVar(\sqrt{n}\tilde{\alpha}_{SP}) = \frac{\sigma_1^2 \sigma_{z_1}^2}{M_{ll} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}$$

and

$$AVar(\sqrt{n}\tilde{\beta}_{SP}) = \frac{\sigma_1^2 (\mu_1^2 \sigma_{z_1}^2 + \mu_2^2 \sigma_{z_2}^2 + 2\mu_1 \mu_2 \sigma_{z_1 z_2})}{M_{ll} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}.$$

Once again, we can rewrite these expressions as

$$AVar(\sqrt{n}\hat{\alpha}_{SP}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)(1 + \kappa)}{(1 - R_2^2)[(1 - \rho_{y_2 z_2, z_1}^2) + \rho_{y_2 z_2, z_1}^2 M_{ll}(1 + \kappa)]}$$

and

$$AVar(\sqrt{n}\hat{\beta}_{SP}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)(1 + \kappa)[1 - R_2^2 + R_2^2 M_{ll}(1 + \kappa)]}{(1 - R_2^2)[(1 - \rho_{y_2 z_2, z_1}^2) + \rho_{y_2 z_2, z_1}^2 M_{ll}(1 + \kappa)]M_{ll}(1 + \kappa)},$$

in the restricted case, and as

$$AVar(\sqrt{n}\tilde{\alpha}_{SP}) = \frac{(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)}{M_{ll}(1 - R_2^2)\rho_{y_2 z_2, z_1}^2}$$

and

$$AVar(\sqrt{n}\tilde{\beta}_{SP}) = \frac{R_2^2(1 - R_1^2)(1 - \rho_{y_2 z_2, z_1}^2)}{M_{ll}(1 - R_2^2)\rho_{y_2 z_2, z_1}^2}$$

when  $\sigma_{12}$  is also estimated.

## A.8 Reparametrisations

The results in the previous subsections can be used to derive the asymptotic distribution of alternative parametrisations. In the case of an estimator  $\hat{\boldsymbol{\theta}}$  that imposes  $\sigma_{12} = 0$ , the asymptotic covariance of the reparametrisation in (8) is simply

$$AVar(\sqrt{n}\hat{\boldsymbol{\theta}}^\dagger) = \mathbf{J}_{\boldsymbol{\theta}^\dagger \boldsymbol{\theta}} AVar(\sqrt{n}\hat{\boldsymbol{\theta}}) \mathbf{J}_{\boldsymbol{\theta}^\dagger \boldsymbol{\theta}}',$$

where

$$\mathbf{J}_{\boldsymbol{\theta}^\dagger \boldsymbol{\theta}} = \frac{\partial \boldsymbol{\theta}^\dagger}{\partial \boldsymbol{\theta}'} = \begin{bmatrix} \mathbf{I}_6 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{1}{2\sigma_1^2} & -\frac{1}{2\sigma_2^2} \\ \mathbf{0} & \frac{\sigma_2}{2\sigma_1} & \frac{\sigma_1}{2\sigma_2} \end{bmatrix}. \quad (\text{A17})$$

In contrast, for an unconstrained estimator  $\tilde{\boldsymbol{\vartheta}}$  that also estimates  $\sigma_{12}$ , we would have that  $\boldsymbol{\vartheta}^\dagger = (\boldsymbol{\vartheta}', \psi_{12})'$  so that

$$AVar(\sqrt{n}\tilde{\boldsymbol{\vartheta}}^\dagger) = \mathbf{J}_{\boldsymbol{\vartheta}^\dagger \boldsymbol{\vartheta}} AVar(\sqrt{n}\tilde{\boldsymbol{\vartheta}}) \mathbf{J}_{\boldsymbol{\vartheta}^\dagger \boldsymbol{\vartheta}}'$$

with

$$\mathbf{J}_{\boldsymbol{\vartheta}^\dagger \boldsymbol{\vartheta}} = \frac{\partial \boldsymbol{\vartheta}^\dagger}{\partial \boldsymbol{\vartheta}'} = \begin{bmatrix} \mathbf{I}_6 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\sigma_1^2 \sigma_2^2 - 2\sigma_{12}^2}{2\sigma_1^2(\sigma_1^2 \sigma_2^2 - \sigma_{12}^2)} & -\frac{\sigma_1^2}{2(\sigma_1^2 \sigma_2^2 - \sigma_{12}^2)} & \frac{\sigma_{12}}{\sigma_1^2 \sigma_2^2 - \sigma_{12}^2} \\ \mathbf{0} & \frac{\sigma_2^2}{2\sqrt{\sigma_1^2 \sigma_2^2 - \sigma_{12}^2}} & \frac{\sigma_1^2}{2\sqrt{\sigma_1^2 \sigma_2^2 - \sigma_{12}^2}} & -\frac{\sigma_{12}}{2\sqrt{\sigma_1^2 \sigma_2^2 - \sigma_{12}^2}} \\ \mathbf{0} & -\frac{\sigma_{12}}{\sigma_1^4} & 0 & \frac{1}{\sigma_1^2} \end{bmatrix}. \quad (\text{A18})$$

## B Proofs of Propositions

### Proof of Proposition 1

Computing in Mathematica the spectral decomposition of  $AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{LS}^\dagger) - AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{MD}^\dagger)$  using the expressions (A5) and (A6), as well as (A17) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$ , we find that it has

only one eigenvalue different from zero, namely,

$$\begin{aligned} & \frac{(\mu_{22} - 1)^2 \mu_2^2 \sigma_1^2 \sigma_2^2 \{ \sigma_{z_1 z_2}^2 (1 + \tau_1^2) \mu_2^2 + [1 + \mu_1^2 + (\mu_0 + \tau_2 \mu_2)^2] \sigma_{z_1}^4 \} (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}{[\mu_{22} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) - \sigma_2^2 \sigma_{z_1}^2] [\sigma_{z_1}^2 \sigma_{z_2}^2 - \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]^2} \\ & - \frac{2(\mu_{22} - 1)^2 \mu_2^2 \sigma_1^2 \sigma_2^2 \sigma_{z_1 z_2} \mu_2 [\tau_1 (\mu_0 + \tau_2 \mu_2) - \mu_1] \sigma_{z_1} (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)}{[\mu_{22} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) - \sigma_2^2 \sigma_{z_1}^2] [\sigma_{z_1}^2 \sigma_{z_2}^2 - \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]^2}, \end{aligned}$$

which is non-negative, with

$$\left( \frac{\sigma_{z_1} (1 + \mu_2 \tau_2) - \sigma_{z_1 z_2} \mu_2 \tau_1}{\sigma_{z_1} \mu_1 + \sigma_{z_1 z_2} \mu_2}, -\frac{\sigma_{z_1}}{\sigma_{z_1}^2 \mu_1 + \sigma_{z_1 z_2} \mu_2}, 1, \mathbf{0}_{1 \times 5} \right)' \quad (\text{B19})$$

as associated eigenvector.

Analogously, after computing the spectral decomposition of  $AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{IV}^\dagger) - AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{MD}^\dagger)$  using the expressions (A2) and (A6), as well as (A17) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$ , we find that it has only one eigenvalue different from zero, namely,

$$\frac{\sigma_1^2 \sigma_2^2 \{ \sigma_{z_1 z_2}^2 (1 + \tau_1^2) \mu_2^2 - 2\sigma_{z_1 z_2} \mu_2 [\tau_1 (\mu_0 + \tau_2 \mu_2) - \mu_1] \sigma_{z_1}^2 + [1 + \mu_1^2 + (\mu_0 + \tau_2 \mu_2)^2] \sigma_{z_1}^4 \}}{[\mu_{22} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) - \sigma_2^2 \sigma_{z_1}^2] \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)},$$

which is non-negative, with (B19) being again its associated eigenvector.

Finally, doing the same for  $AVar(\sqrt{n}\tilde{\boldsymbol{\theta}}_{IV}^\dagger) - AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{MD}^\dagger)$  by combining (A2) and (A5), as well as (A17) and (A18) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$  and from  $\boldsymbol{\vartheta}$  to  $\boldsymbol{\vartheta}^\dagger$ , respectively, we find that it has only one eigenvalue different from zero, namely,

$$\begin{aligned} & \frac{(\mu_{22} - 1)^2 \mu_2^2 \sigma_{z_1}^2 \sigma_{z_2}^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) \{ \sigma_{z_1 z_2}^2 (1 + \tau_1^2) \mu_2^2 - 2\sigma_{z_1 z_2} \mu_2 [(\mu_0 + \mu_2 \tau_2) \tau_1 - \mu_1] \sigma_{z_1}^2 \}}{[\mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) - \sigma_1^2 \sigma_{z_1}^2]^2 [\mu_{22} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) - \sigma_2^2 \sigma_{z_1}^2]} \\ & + \frac{(\mu_{22} - 1)^2 \mu_2^2 \sigma_{z_1}^2 \sigma_{z_2}^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) \{ [1 + \mu_1^2 + \sigma_{z_1}^2 (\mu_0 + \mu_2 \tau_2)^2] \}}{[\mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) - \sigma_1^2 \sigma_{z_1}^2]^2 [\mu_{22} \mu_2^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) - \sigma_2^2 \sigma_{z_1}^2]}, \end{aligned}$$

which can be positive or negative depending on  $\mu_{22}$ , and with the same eigenvector.  $\square$

## Proof of Proposition 2

Computing in Mathematica the spectral decomposition of  $AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{SS}^\dagger) - AVar[\sqrt{n}\hat{\boldsymbol{\theta}}_{ML}^\dagger(\bar{\boldsymbol{\eta}})]$  using the expression in (A9), as well as (A17) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$ , and exploiting the fact that

$$AVar[\sqrt{n}\hat{\boldsymbol{\theta}}_{ML}^\dagger(\bar{\boldsymbol{\eta}})] = [\mathcal{I}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{\theta}, \boldsymbol{\eta})]^{-1},$$

where  $\mathcal{I}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{\theta}, \boldsymbol{\eta})$  denotes the block of the information matrix of the mean and variance parameters, we find that it has only one eigenvalue different from zero, with associated eigenvector

$$(\mathbf{0}_{1 \times 7}, 1)'.$$

Similarly, we find that the spectral decomposition of  $AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{SP}^\dagger) - AVar[\sqrt{n}\hat{\boldsymbol{\theta}}_{ML}^\dagger(\bar{\boldsymbol{\eta}})]$ , which we obtain using also (A13), as well as (A17) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$ , has five eigenvalues

different from zero. By looking at the orthogonal basis for its null space, which is given by

$$(0, \sigma_{z_1}^2 \mu_1 + \sigma_{z_1 z_2} \mu_2, \sigma_{z_1}^2, \mathbf{0}_{1 \times 5})'$$

and

$$(\mathbf{0}_{2 \times 4}, \mathbf{I}_2, \mathbf{0}_{2 \times 2})',$$

we can immediately see that the parameters that are estimated adaptively are  $\mu_1$ ,  $\mu_2$ , and the linear combination of  $\alpha$  and  $\beta$  indicated by the first eigenvector. In turn, the orthogonal basis for its image is given by

$$(1, \mathbf{0}_{1 \times 7})',$$

$$(0, -\sigma_{z_1}^2, \mu_1 \sigma_{z_1}^2 + \sigma_{z_1 z_2} \mu_2, \mathbf{0}_{1 \times 5})',$$

$$(\mathbf{0}_{1 \times 3}, 1, \mathbf{0}_{1 \times 4})'$$

and

$$(\mathbf{0}_{1 \times 6}, \mathbf{I}_2)'. \quad \square$$

Finally, using an entirely analogous procedure with (A13) and (A6), together with (A17) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$ , we find that the spectral decomposition of  $AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{MD}^\dagger) - AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{SP}^\dagger)$  has four eigenvalues different from zero, and its null space is defined by

$$\left(0, \mu_1 + \frac{\sigma_{z_1 z_2}}{\sigma_{z_1}^2} \mu_2, 1, \mathbf{0}_{1 \times 5}\right)',$$

$$(\mathbf{0}_{1 \times 4}, 1, \mathbf{0}_{1 \times 3})' \text{ and } (\mathbf{0}_{1 \times 5}, 1, \mathbf{0}_{1 \times 2})'. \quad \square$$

### Proof of Proposition 3

Computing in Mathematica the spectral decomposition of  $AVar(\sqrt{n}\tilde{\boldsymbol{\theta}}_{ML}^\dagger) - AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{ML}^\dagger)$  using the expressions (A7) and (A8), as well as (A17) and (A18) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$  and from  $\boldsymbol{\vartheta}$  to  $\boldsymbol{\vartheta}^\dagger$ , respectively, we find that it has only one eigenvalue different from zero, namely,

$$\frac{\sigma_1^2 \sigma_2^2 M_{ss} \{ \sigma_{z_1 z_2}^2 \mu_2^2 (1 + \tau_1^2) - 2 \sigma_{z_1 z_2} \sigma_{z_1}^2 \mu_2 [(\mu_0 + \mu_2 \tau_2) \tau_1 - \mu_1] + \sigma_{Z_1}^4 [1 + \mu_1^2 + (\mu_0 + \mu_2 \tau_2)^2] \}}{\mu_2^2 M_{ll} (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) [M_{ss} \sigma_2^2 \sigma_{z_1}^2 + \mu_2^2 M_{ll} (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]}$$

with associated eigenvector (B19).

Using (A9) and (A10), as well as (A17) and (A18) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$  and from  $\boldsymbol{\vartheta}$  to  $\boldsymbol{\vartheta}^\dagger$ , respectively, we find that the same turns out to be true for  $AVar(\sqrt{n}\tilde{\boldsymbol{\theta}}_{SS}^\dagger) - AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{SS}^\dagger)$ .

Finally, if we do the same for  $AVar(\sqrt{n}\tilde{\boldsymbol{\theta}}_{SP}^\dagger) - AVar(\sqrt{n}\hat{\boldsymbol{\theta}}_{SP}^\dagger)$  using (A13) and (A14), as well as (A17) and (A18) to go from  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^\dagger$  and from  $\boldsymbol{\vartheta}$  to  $\boldsymbol{\vartheta}^\dagger$ , respectively, we also find that it

has only one eigenvalue different from zero, namely

$$\frac{\sigma_1^2 \sigma_2^2 \{ \sigma_{z_1 z_2}^2 \mu_2^2 (1 + \tau_1^2) - 2 \sigma_{z_1 z_2} \sigma_{z_1}^2 \mu_2 [(\mu_0 + \mu_2 \tau_2) \tau_1 - \mu_1] + \sigma_{Z_1}^4 [1 + \mu_1^2 + (\mu_0 + \mu_2 \tau_2)^2] \}}{\mu_{2\text{M}}^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2) [\sigma_2^2 \sigma_{z_1}^2 + (1 + \kappa) \mu_{2\text{M}}^2 (\sigma_{z_1}^2 \sigma_{z_2}^2 - \sigma_{z_1 z_2}^2)]},$$

and that its image is given by the same eigenvector as in the previous cases.  $\square$

## C Simplifying the DGP

### C.1 Standardised variables

We start by assuming that:

$$\begin{pmatrix} y_1 \\ y_2 \\ z_1 \\ z_2 \end{pmatrix} \sim \left[ \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho_{y_1 y_2} & \rho_{y_1 z_1} & \rho_{y_1 z_2} \\ \rho_{y_1 y_2} & 1 & \rho_{y_2 z_1} & \rho_{y_2 z_2} \\ \rho_{y_1 z_1} & \rho_{y_2 z_1} & 1 & \rho_{z_1 z_2} \\ \rho_{y_1 z_2} & \rho_{y_2 z_2} & \rho_{z_1 z_2} & 1 \end{pmatrix} \right],$$

where the correlation matrix is positive definite.

In this notation, the coefficients of the least squares projection of  $y_1$  onto  $y_2$  and  $z_1$  are

$$\begin{aligned} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} &= \begin{pmatrix} 1 & \rho_{y_2 z_1} \\ \rho_{y_2 z_1} & 1 \end{pmatrix}^{-1} \begin{pmatrix} \rho_{y_1 y_2} \\ \rho_{y_1 z_1} \end{pmatrix} \\ &= \frac{1}{1 - \rho_{y_2 z_1}^2} \begin{pmatrix} \rho_{y_1 y_2} - \rho_{y_1 z_1} \rho_{y_2 z_1} \\ \rho_{y_1 z_1} - \rho_{y_1 y_2} \rho_{y_2 z_1} \end{pmatrix}, \end{aligned}$$

the corresponding projection errors

$$u_1 = y_1 - \alpha y_2 - \beta z_1 = y_1 - \frac{\rho_{y_1 y_2} - \rho_{y_1 z_1} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} y_2 - \frac{\rho_{y_1 z_1} - \rho_{y_1 y_2} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} z_1$$

and the residual variance

$$\begin{aligned} V(u_1) &= 1 - \begin{pmatrix} \rho_{y_1 y_2} & \rho_{y_1 z_1} \end{pmatrix} \begin{pmatrix} 1 & \rho_{y_2 z_1} \\ \rho_{y_2 z_1} & 1 \end{pmatrix}^{-1} \begin{pmatrix} \rho_{y_1 y_2} \\ \rho_{y_1 z_1} \end{pmatrix} \\ &= 1 - \frac{\rho_{y_1 y_2}^2 + \rho_{y_1 z_1}^2 - 2 \rho_{y_2 z_1} \rho_{y_1 y_2} \rho_{y_1 z_1}}{1 - \rho_{y_2 z_1}^2}, \end{aligned}$$

so that the  $R^2$  becomes

$$R_1^2 = \frac{\rho_{y_1 y_2}^2 + \rho_{y_1 z_1}^2 - 2 \rho_{y_2 z_1} \rho_{y_1 y_2} \rho_{y_1 z_1}}{1 - \rho_{y_2 z_1}^2}.$$

In turn, the coefficients of the least squares projection of  $y_2$  onto  $z_1$  and  $z_2$  are

$$\begin{aligned} \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} &= \begin{pmatrix} 1 & \rho_{z_1 z_2} \\ \rho_{z_1 z_2} & 1 \end{pmatrix}^{-1} \begin{pmatrix} \rho_{y_2 z_1} \\ \rho_{y_2 z_2} \end{pmatrix} \\ &= \frac{1}{1 - \rho_{z_1 z_2}^2} \begin{pmatrix} \rho_{y_2 z_1} - \rho_{y_2 z_2} \rho_{z_1 z_2} \\ \rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2} \end{pmatrix}, \end{aligned}$$

the corresponding projection errors

$$u_2 = y_2 - \mu_1 z_1 - \mu_2 z_2 = y_2 - \frac{\rho_{y_2 z_1} - \rho_{y_2 z_2} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} z_1 - \frac{\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} z_2$$

and the residual variance

$$\begin{aligned} V(u_2) &= 1 - \begin{pmatrix} \rho_{y_2 z_1} & \rho_{y_2 z_2} \end{pmatrix} \begin{pmatrix} 1 & \rho_{z_1 z_2} \\ \rho_{z_1 z_2} & 1 \end{pmatrix}^{-1} \begin{pmatrix} \rho_{y_2 z_1} \\ \rho_{y_2 z_2} \end{pmatrix} \\ &= 1 - \frac{\rho_{y_2 z_1}^2 + \rho_{y_2 z_2}^2 - 2\rho_{z_1 z_2} \rho_{y_2 z_1} \rho_{y_2 z_2}}{1 - \rho_{z_1 z_2}^2}, \end{aligned}$$

so that the  $R^2$  becomes

$$R_2^2 = \frac{\rho_{y_2 z_1}^2 + \rho_{y_2 z_2}^2 - 2\rho_{z_1 z_2} \rho_{y_2 z_1} \rho_{y_2 z_2}}{1 - \rho_{z_1 z_2}^2}.$$

Finally, the covariance between the previous projection errors is

$$\begin{aligned} &E[(y_1 - \alpha y_2 - \beta z_1)(y_2 - \mu_1 z_1 - \mu_2 z_2)] \\ &= E \left[ \left( y_1 - \frac{\rho_{y_1 y_2} - \rho_{y_1 z_1} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} y_2 - \frac{\rho_{y_1 z_1} - \rho_{y_1 y_2} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} z_1 \right) \right. \\ &\quad \left. \left( y_2 - \frac{\rho_{y_2 z_1} - \rho_{y_2 z_2} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} z_1 - \frac{\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} z_2 \right) \right] \\ &= \rho_{y_1 y_2} - \frac{\rho_{y_2 z_1} - \rho_{y_2 z_2} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} \rho_{y_1 z_1} - \frac{\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} \rho_{y_1 z_2} \\ &\quad - \frac{\rho_{y_1 y_2} - \rho_{y_1 z_1} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} + \frac{\rho_{y_1 y_2} - \rho_{y_1 z_1} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} \frac{\rho_{y_2 z_1} - \rho_{y_2 z_2} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} \rho_{y_2 z_1} \\ &\quad + \frac{\rho_{y_1 y_2} - \rho_{y_1 z_1} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} \frac{\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} \rho_{y_2 z_2} - \frac{\rho_{y_1 z_1} - \rho_{y_1 y_2} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} \rho_{y_2 z_1} \\ &\quad + \frac{\rho_{y_1 z_1} - \rho_{y_1 y_2} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} \frac{\rho_{y_2 z_1} - \rho_{y_2 z_2} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} + \frac{\rho_{y_1 z_1} - \rho_{y_1 y_2} \rho_{y_2 z_1}}{1 - \rho_{y_2 z_1}^2} \frac{\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} \rho_{z_1 z_2} \\ &= \frac{(\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2})}{1 - \rho_{z_1 z_2}^2} \frac{[\rho_{y_1 y_2} (\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}) + \rho_{y_1 z_1} (\rho_{z_1 z_2} - \rho_{y_2 z_2} \rho_{y_2 z_1}) - \rho_{y_1 z_2} (1 - \rho_{y_2 z_1}^2)]}{1 - \rho_{y_2 z_1}^2} \end{aligned}$$

Therefore, for  $y_2$  to be exogenous in the first equation, we need either

$$\mu_2 = \frac{\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}}{1 - \rho_{z_1 z_2}^2} = 0,$$

which seems very restrictive, or

$$\rho_{y_1 z_2} = \frac{\rho_{y_1 y_2} (\rho_{y_2 z_2} - \rho_{y_2 z_1} \rho_{z_1 z_2}) + \rho_{y_1 z_1} (\rho_{z_1 z_2} - \rho_{y_2 z_1} \rho_{y_2 z_2})}{1 - \rho_{y_2 z_1}^2} = \frac{1 - \rho_{z_1 z_2}^2}{1 - \rho_{y_2 z_1}^2} \mu_2 \rho_{y_1 y_2} + \delta \rho_{y_1 z_1}, \quad (\text{C20})$$

where  $\delta$  is the coefficient of  $z_1$  in the least squares projection of  $z_2$  onto  $y_2$  and  $z_1$ , whose coefficients are given by

$$\begin{pmatrix} \gamma \\ \delta \end{pmatrix} = \begin{pmatrix} 1 & \rho_{y_2 z_1} \\ \rho_{y_2 z_1} & 1 \end{pmatrix}^{-1} \begin{pmatrix} \rho_{y_2 z_2} \\ \rho_{z_1 z_2} \end{pmatrix} = \frac{1}{1 - \rho_{y_2 z_1}^2} \begin{pmatrix} \rho_{y_2 z_2} - \rho_{z_1 z_2} \rho_{y_2 z_1} \\ \rho_{z_1 z_2} - \rho_{y_2 z_2} \rho_{y_2 z_1} \end{pmatrix}.$$

Therefore, if we assume  $\mu_2 \neq 0$ , then we need to choose  $\rho_{y_1 z_2}$  so that (C20) holds.

## C.2 Original variables

Let us now consider the least squares projection of  $y_2^o$  onto a constant,  $z_1^o$  and  $z_2^o$ , which is given by

$$y_2^o = \mu_0^o + \mu_1^o z_1^o + \mu_2^o z_2^o + u_1^o.$$

We can then individually centre and standardise each of the variables involved as follows

$$y_2 = \frac{y_2^o - \mu_0^o - \mu_1^o E(z_1^o) - \mu_2^o E(z_2^o)}{\sqrt{\mu_1^{o2} V(z_1^o) + \mu_2^{o2} V(z_2^o) + 2\mu_1^o \mu_2^o \text{Cov}(z_1^o, z_2^o) + V(u_1^o)}},$$

$$z_1 = \frac{z_1^o - E(z_1^o)}{\sqrt{V(z_1^o)}}, \quad \text{and} \quad z_2 = \frac{z_2^o - E(z_2^o)}{\sqrt{V(z_2^o)}},$$

which leads to the following transformed equation

$$y_2 = \mu_1 z_1 + \mu_2 z_2 + u_1,$$

where

$$\mu_1 = \mu_1^o \sqrt{\frac{V(z_1^o)}{\mu_1^{o2} V(z_1^o) + \mu_2^{o2} V(z_2^o) + 2\mu_1^o \mu_2^o \text{Cov}(z_1^o, z_2^o) + V(u_1^o)}},$$

$$\mu_2 = \mu_2^o \sqrt{\frac{V(z_2^o)}{\mu_1^{o2} V(z_1^o) + \mu_2^{o2} V(z_2^o) + 2\mu_1^o \mu_2^o \text{Cov}(z_1^o, z_2^o) + V(u_1^o)}},$$

and

$$V(u_1) = \frac{V(u_1^o)}{\mu_1^{o2} V(z_1^o) + \mu_2^{o2} V(z_2^o) + 2\mu_1^o \mu_2^o \text{Cov}(z_1^o, z_2^o) + V(u_1^o)} = 1 - R_2^2$$

The coefficients  $\mu_1$  and  $\mu_2$  are sometimes called the standardised regression coefficients, as they explain the ceteris paribus change in  $y_2^o$  (measured in standard deviation units) resulting from a unit standard deviation change in  $z_1^o$  or  $z_2^o$ .

Thus, once we standardise the three variables involved, the crucial ingredients of the first equation are the coefficient of determination  $R_2^2$ , the correlation between the regressors  $\rho_{z_1 z_2}$  and the partial correlations between  $y_2$  and each of the regressors, which are given by

$$\rho_{y_2 z_1 \cdot z_2} = \frac{E[(y_2 - \rho_{y_2 z_2} z_2)(z_1 - \rho_{z_1 z_2} z_2)]}{\sqrt{V(y_2 - \rho_{y_2 z_2} z_2)V(z_1 - \rho_{z_1 z_2} z_2)}} = \frac{\rho_{y_2 z_1} - \rho_{z_1 z_2} \rho_{y_2 z_2}}{\sqrt{(1 - \rho_{y_2 z_2}^2)(1 - \rho_{z_1 z_2}^2)}} = \mu_1 \sqrt{\frac{1 - \rho_{z_1 z_2}^2}{1 - \rho_{y_2 z_2}^2}},$$

$$\rho_{y_2 z_2 \cdot z_1} = \frac{E[(y_2 - \rho_{y_2 z_1} z_1)(z_2 - \rho_{z_1 z_2} z_1)]}{\sqrt{V(y_2 - \rho_{y_2 z_1} z_1)V(z_2 - \rho_{z_1 z_2} z_1)}} = \frac{\rho_{y_2 z_2} - \rho_{z_1 z_2} \rho_{y_2 z_1}}{\sqrt{(1 - d^2)(1 - \rho_{z_1 z_2}^2)}} = \mu_2 \sqrt{\frac{1 - \rho_{z_1 z_2}^2}{1 - \rho_{y_2 z_1}^2}}.$$

In fact, there are only three underlying parameters that determine these four quantities:

$\rho_{y_2 z_1}$ ,  $\rho_{y_2 z_2}$  and  $\rho_{z_1 z_2}$  because

$$\begin{aligned}\rho_{y_2 z_1 \cdot z_2}^2 &= \frac{R_2^2 - \rho_{y_2 z_2}^2}{1 - \rho_{y_2 z_2}^2}, \\ \rho_{y_2 z_2 \cdot z_1}^2 &= \frac{R_2^2 - \rho_{y_2 z_1}^2}{1 - \rho_{y_2 z_1}^2},\end{aligned}$$

or alternatively

$$\begin{aligned}\rho_{y_2 z_2}^2 &= \frac{R_2^2 - \rho_{y_2 z_1 \cdot z_2}^2}{1 - \rho_{y_2 z_1 \cdot z_2}^2}, \\ \rho_{y_2 z_1}^2 &= \frac{R_2^2 - \rho_{y_2 z_2 \cdot z_1}^2}{1 - \rho_{y_2 z_2 \cdot z_1}^2}.\end{aligned}$$

Thus, we can either select  $\rho_{y_2 z_1}$ ,  $\rho_{y_2 z_2}$  and  $\rho_{z_1 z_2}$ , or we can select  $R_2^2$ ,  $\rho_{y_2 z_1 \cdot z_2}^2$  and  $\rho_{y_2 z_2 \cdot z_1}^2$ .

## D On multivariate discrete mixture of normals

Consider the following mixture of two multivariate normals

$$\mathbf{u}_t^* \sim \begin{cases} N(\boldsymbol{\nu}_1, \boldsymbol{\Gamma}_1) & \text{with probability } \lambda, \\ N(\boldsymbol{\nu}_2, \boldsymbol{\Gamma}_2) & \text{with probability } 1 - \lambda. \end{cases} \quad (\text{D21})$$

Let  $s_t$  denote a Bernoulli variable which takes the value 1 with probability  $\lambda$  and 0 with probability  $1 - \lambda$ . As is well known, the unconditional mean vector and covariance matrix of the observed variables are:

$$E(\mathbf{u}_t^*) = \boldsymbol{\pi} = E[E(\mathbf{u}_t^* | s_t)] = \lambda \boldsymbol{\nu}_1 + (1 - \lambda) \boldsymbol{\nu}_2,$$

$$V(\mathbf{u}_t^*) = \boldsymbol{\Psi} = V[E(\mathbf{u}_t^* | s_t)] + E[V(\mathbf{u}_t^* | s_t)] = \lambda(1 - \lambda) \boldsymbol{\delta} \boldsymbol{\delta}' + \lambda \boldsymbol{\Sigma}_1 + (1 - \lambda) \boldsymbol{\Sigma}_2,$$

where  $\boldsymbol{\delta} = \boldsymbol{\nu}_1 - \boldsymbol{\nu}_2$ .

Therefore, the random vector  $\mathbf{u}_t^*$  will be standardised if and only if

$$\lambda \boldsymbol{\nu}_1 + (1 - \lambda) \boldsymbol{\nu}_2 = \mathbf{0}$$

and

$$\lambda(1 - \lambda)(\boldsymbol{\nu}_1 - \boldsymbol{\nu}_2)(\boldsymbol{\nu}_1 - \boldsymbol{\nu}_2)' + \lambda \boldsymbol{\Sigma}_1 + (1 - \lambda) \boldsymbol{\Sigma}_2 = \lambda(1 - \lambda) \boldsymbol{\delta} \boldsymbol{\delta}' + \lambda \boldsymbol{\Sigma}_1 + (1 - \lambda) \boldsymbol{\Sigma}_2 = \mathbf{I}.$$

To see how we can achieve these two conditions, let us initially assume that  $\boldsymbol{\nu}_1 = \boldsymbol{\nu}_2 = \mathbf{0}$ , so that  $\boldsymbol{\delta} = \mathbf{0}$ . Let  $\boldsymbol{\Gamma}_{1L} \boldsymbol{\Gamma}_{1L}'$  and  $\boldsymbol{\Gamma}_{2L} \boldsymbol{\Gamma}_{2L}'$  denote the Cholesky decompositions of the covariance matrices of the two components. Then, we can write

$$\lambda \boldsymbol{\Gamma}_1 + (1 - \lambda) \boldsymbol{\Gamma}_2 = \boldsymbol{\Gamma}_{1L} [\lambda \mathbf{I}_2 + (1 - \lambda) \boldsymbol{\Gamma}_{1L}^{-1} \boldsymbol{\Gamma}_{2L} \boldsymbol{\Gamma}_{2L}' \boldsymbol{\Gamma}_{1L}^{-1}] \boldsymbol{\Gamma}_{1L}' = \boldsymbol{\Gamma}_{1L} [\lambda \mathbf{I}_2 + (1 - \lambda) \boldsymbol{\aleph}_L \boldsymbol{\aleph}_L'] \boldsymbol{\Gamma}_{1L}',$$

where

$$\mathbf{\Gamma}_{1L}^{-1}\mathbf{\Gamma}_{2L} = \mathbf{\aleph}_L = \begin{bmatrix} \varkappa_{11} & 0 \\ \varkappa_{21} & \varkappa_{22} \end{bmatrix}$$

is a lower triangular matrix. Thus, it is not difficult to see that choosing

$$\mathbf{\Gamma}_1 = [\lambda\mathbf{I}_2 + (1 - \lambda)\mathbf{\aleph}_L\mathbf{\aleph}_L']^{-1} \quad \text{and} \quad \mathbf{\Gamma}_2 = \mathbf{\Gamma}_{1L}\mathbf{\aleph}_L\mathbf{\aleph}_L'\mathbf{\Gamma}_{1L}'$$

or, equivalently,

$$\mathbf{\Gamma}_{1L} = [\lambda\mathbf{I}_2 + (1 - \lambda)\mathbf{\aleph}_L\mathbf{\aleph}_L']^{-1/2} \quad \text{and} \quad \mathbf{\Gamma}_{2L} = \mathbf{\Gamma}_{1L}\mathbf{\aleph}_L$$

we can indeed obtain a standardised vector  $\mathbf{u}_t^*$ .

Now consider the case  $\boldsymbol{\delta} \neq \mathbf{0}$ , and let  $\mathbf{\Upsilon} = \lambda(1 - \lambda)\boldsymbol{\delta}\boldsymbol{\delta}' + \mathbf{I}_2$ . Then, it is easy to see that

$$\boldsymbol{\nu}_1^* = \mathbf{\Upsilon}^{-\frac{1}{2}}\boldsymbol{\nu}_1, \quad \boldsymbol{\nu}_2^* = \mathbf{\Upsilon}^{-\frac{1}{2}}\boldsymbol{\nu}_2, \quad \mathbf{\Gamma}_1^* = \mathbf{\Upsilon}_1^{-\frac{1}{2}}\mathbf{\Gamma}_1\mathbf{\Upsilon}'^{-\frac{1}{2}}, \quad \text{and} \quad \mathbf{\Gamma}_2^* = \mathbf{\Upsilon}^{-\frac{1}{2}}\mathbf{\Gamma}_2\mathbf{\Upsilon}'^{-\frac{1}{2}}$$

continue to generate a standardised vector. As a result, the vector of shape parameters of a standardised finite mixture distribution reduces to  $\boldsymbol{\varrho} = (\delta_1, \delta_2, \varkappa_{11}, \varkappa_{21}, \varkappa_{22}, \lambda)'$  in the bivariate case.

## Tables

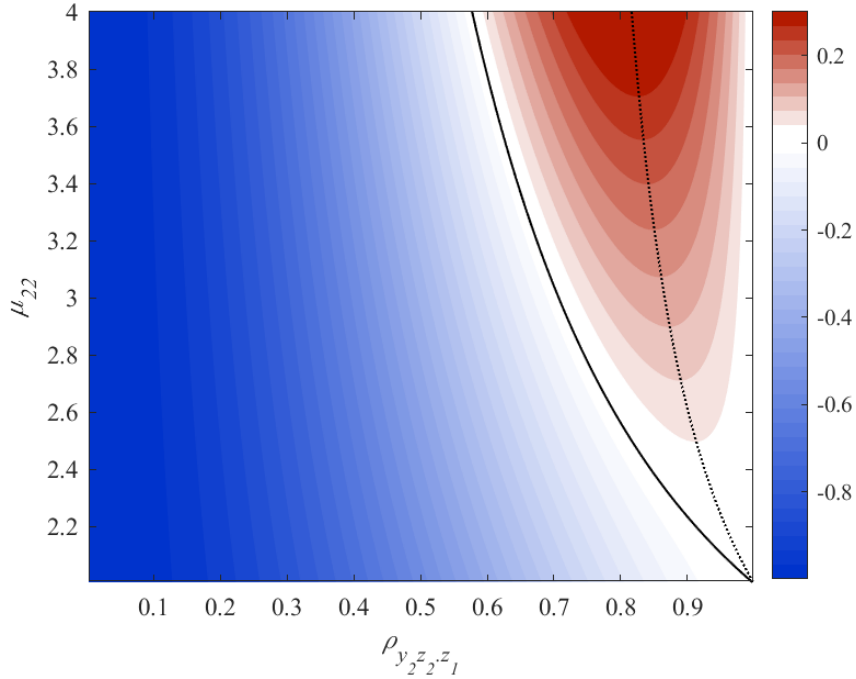
Table 1: Asymptotic variances of alternative estimators

Parameter	OLS	PML-SMN <sub>K</sub>			SS	ML	MD
		<i>K</i> = 2	<i>K</i> = 3	<i>K</i> = 4			
Mean parameters of equation 1a							
<i>γ</i>	1.268	0.931	0.905	0.902	0.901	0.901	1.201
<i>α</i>	1.500	0.782	0.731	0.725	0.723	0.723	1.125
<i>β</i>	1.000	0.656	0.631	0.627	0.627	0.627	0.875
Mean parameters of equation 1b							
<i>μ</i> <sub>0</sub>	1.000	0.792	0.775	0.772	0.771	0.771	1.000
<i>μ</i> <sub>1</sub>	0.333	0.264	0.258	0.257	0.257	0.257	0.333
<i>μ</i> <sub>2</sub>	0.333	0.264	0.258	0.257	0.257	0.257	0.333
(Reparametrised) variance parameters of structural innovations							
<i>ω</i>	3.000	1.493	1.313	1.290	1.286	1.286	3.000
<i>σ</i> <sup>2</sup>	0.833	0.833	0.833	0.833	0.833	0.300	0.833

*Notes:* DGP for structural innovations: bivariate Student  $t$  with 0 means, unit standard deviations, no correlation and 5 degrees of freedom. Parameter values:  $\gamma = 0.204$ ,  $\alpha = \beta = 0.398$ ,  $\mu_0 = 0.155$ ,  $\mu_1 = \mu_2 = 0.577$ ,  $\sigma_1^2 = 1/2$ ,  $\sigma_2^2 = 1/3$ ,  $\mu_{z_1} = \mu_{z_2} = 1$ ,  $\sigma_{z_1}^2 = \sigma_{z_2}^2 = 1$ , and  $\sigma_{z_1 z_2} = 0$ . OLS denotes the usual ordinary least squares estimator, PML-SMN<sub>K</sub> denotes Pseudo-ML based on a bivariate scale mixture of  $K$  normals, SS denotes the spherically symmetric SP estimator, ML denotes MLE which exploit the information of the true distribution of the shocks, including the degrees of freedom, and MD denotes the optimum minimum distance estimator. We compute the expected value of the Hessian and variance of the score of the finite mixture-based PMLEs using the expressions in Fiorentini and Sentana (2021) evaluated at the true values of the mean and variance parameters in  $\theta$  and the pseudo true values of the shape parameters, which we numerically obtain from samples of 40 million simulated observations.

## Figures

Figure 1: Relative efficiency OLS/IV for  $\alpha$

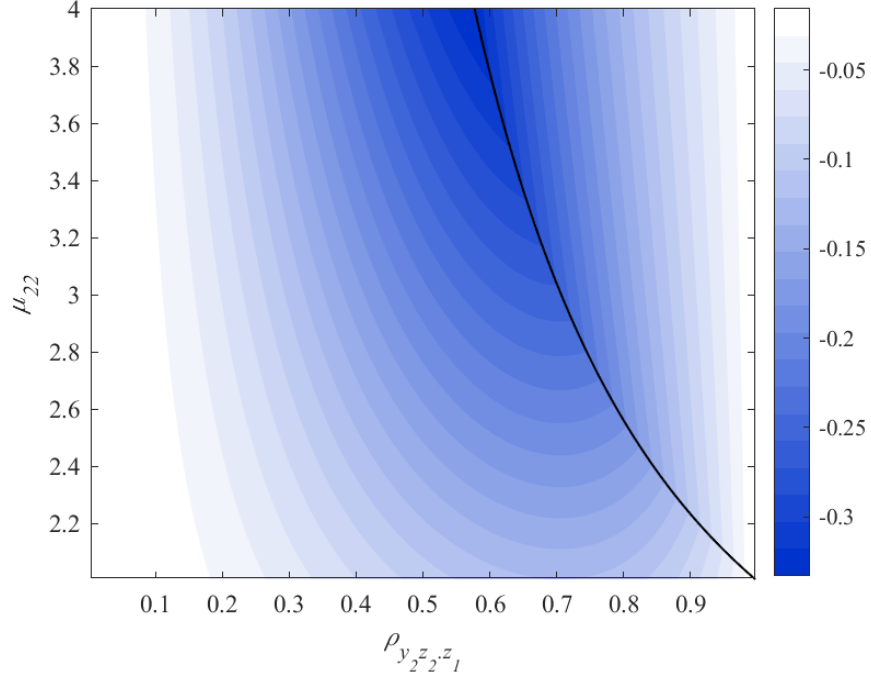


*Notes:* When the  $R^2$  of equation (2) coincides with  $\rho_{y_2 z_2, z_1}^2$ , the relative efficiency of the OLS/IV estimators of  $\alpha$  is given by

$$\frac{V(\hat{\alpha}_{LS})}{V(\hat{\alpha}_{IV})} = [(1 - \rho_{y_2 z_2, z_1}^2)\mu_{22} + \rho_{y_2 z_2, z_1}^2]\rho_{y_2 z_2, z_1}^2.$$

The solid line denotes the boundary line  $\mu_{22} = 1 + \rho_{y_2 z_2, z_1}^{-2}$  while the dotted line denotes the locus of  $(\rho_{y_2 z_2, z_1}, \mu_{22})$  combinations for which the IV estimator of  $\alpha$  reaches its maximum asymptotic efficiency relative to the corresponding OLS estimator, which is given by  $\rho_{y_2 z_2, z_1}^2 = \frac{1}{2}\mu_{22}/(\mu_{22} - 1)$ .

Figure 2: Relative efficiency MD/OLS-IV for  $\alpha$



*Notes:* When the  $R^2$  of equation (2) coincides with  $\rho_{y_2 z_2, z_1}^2$ , the relative efficiency of the MD/OLS and MD/IV estimators of  $\alpha$  are given by

$$\frac{V(\hat{\alpha}_{MD})}{V(\hat{\alpha}_{LS})} = \frac{\mu_{22}}{[1 + (\mu_{22} - 1)\rho_{y_2 z_2, z_1}^2][(1 - \rho_{y_2 z_2, z_1}^2)\mu_{22} + \rho_{y_2 z_2, z_1}^2]}$$

and

$$\frac{V(\hat{\alpha}_{MD})}{V(\hat{\alpha}_{IV})} = \frac{\mu_{22}\rho_{y_2 z_2, z_1}^2}{1 + (\mu_{22} - 1)\rho_{y_2 z_2, z_1}^2},$$

respectively. The solid line denotes the boundary line  $\mu_{22} = 1 + \rho_{y_2 z_2, z_1}^{-2}$ .

Figure 3: Relative efficiency Student  $t$  ML/MD for  $\alpha$  and  $\beta$

Figure 3a: Relative efficiency ML/MD for  $\alpha$

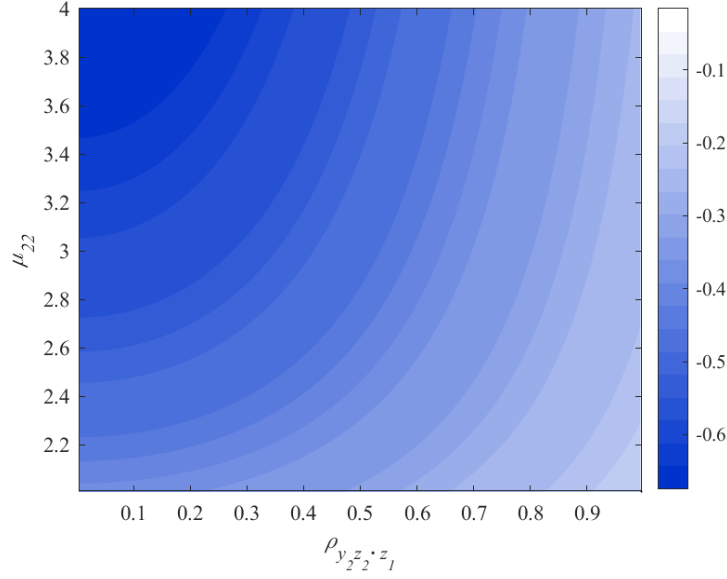
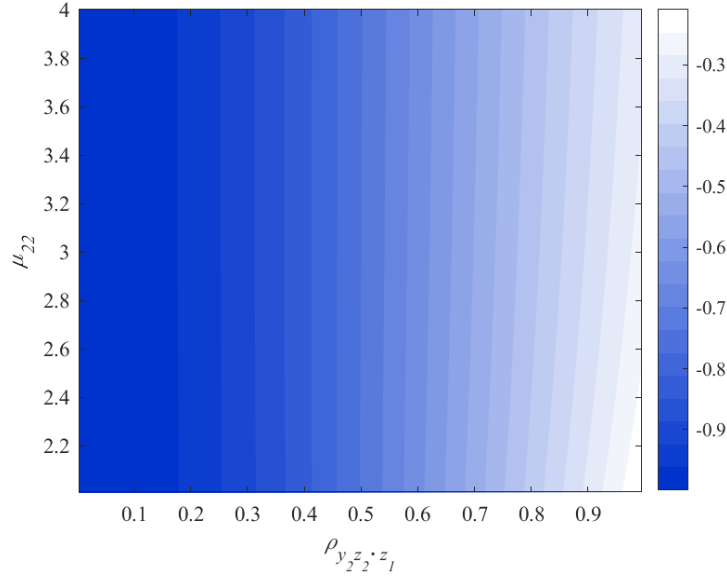


Figure 3b: Relative efficiency ML/MD for  $\beta$



*Notes:* When the  $R^2$  of equation (2) coincides with  $\rho_{y_2 z_2 \cdot z_1}^2$ , the relative efficiency of the MD/OLS and MD/IV estimators of  $\alpha$  and  $\beta$  are, respectively, given by

$$\frac{AVar(\sqrt{n}\hat{\alpha}_{MLt})}{AVar(\sqrt{n}\hat{\alpha}_{MD})} = \frac{1 + (\mu_{22} - 1)\rho_{y_2 z_2 \cdot z_1}^2}{[(1 - \rho_{y_2 z_2 \cdot z_1}^2)M_{ss} + \rho_{y_2 z_2 \cdot z_1}^2 M_{ll}]\mu_{22}}$$

and

$$\frac{AVar(\sqrt{n}\hat{\beta}_{MLt})}{AVar(\sqrt{n}\hat{\beta}_{MD})} = \frac{\rho_{y_2 z_2 \cdot z_1}^2}{[(1 - \rho_{y_2 z_2 \cdot z_1}^2)M_{ss} + \rho_{y_2 z_2 \cdot z_1}^2 M_{ll}]\mu_{22}},$$

and where  $M_{ll} = \nu(2 + \nu)/[(\nu - 2)(\nu + 4)]$  and  $M_{ss} = (\nu + 2)/(\nu + 4)$  with  $\nu = 2(2\mu_{22} - 1)/(\mu_{22} - 1)$ .

Figure 4: Monte Carlo spherical data generating processes versus Gaussian distribution

Figure 4a: Bivariate standard Gaussian distribution

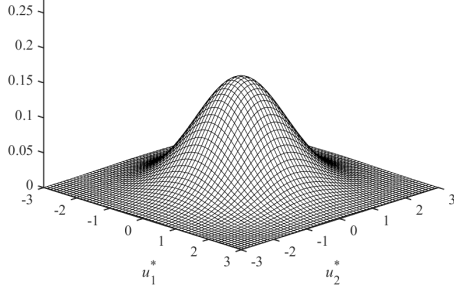


Figure 4b: Contours of a bivariate standard Gaussian distribution

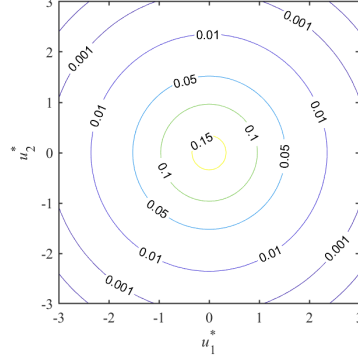


Figure 4c: Bivariate standard Student  $t$  distribution density

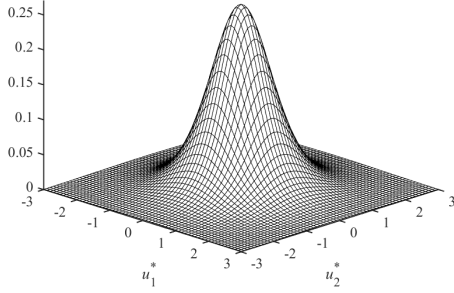


Figure 4d: Contours of a bivariate standard Student  $t$  density

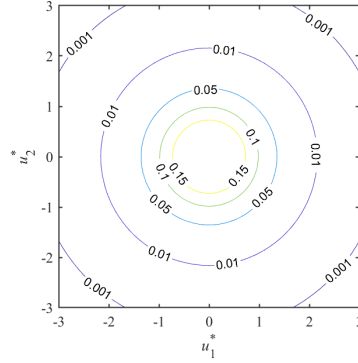


Figure 4e: Bivariate standard scale mixture of two normals

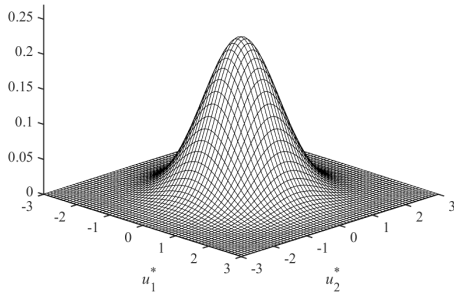
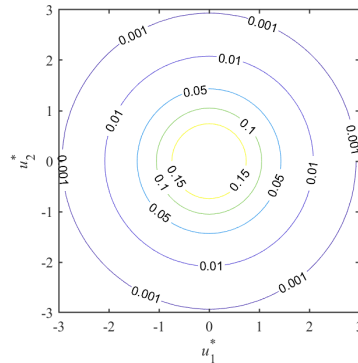


Figure 4f: Contours of a bivariate standard scale mixture of two normals



*Notes:* In all panels,  $E(u_{1i}^*) = E(u_{2i}^*) = 0$ ,  $V(u_{1i}^*) = V(u_{2i}^*) = 1$ , and  $cov(u_{1i}^*, u_{2i}^*) = 0$ . Panels c-d: Student  $t$  distribution with  $\nu = 5$  degrees of freedom. Panels e-f: Scale mixture of two normals with scale parameter  $\varkappa = 0.09$  and mixing probability  $\lambda = 0.05$ .

Figure 5: Monte Carlo non-spherical data generating processes versus Gaussian distribution

Figure 5a: Bivariate standard Gaussian distribution density

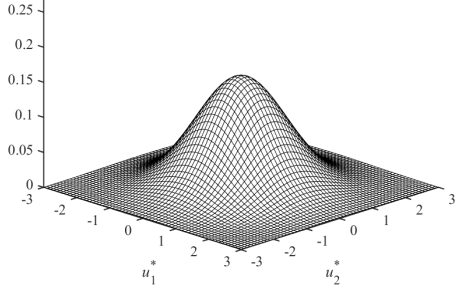


Figure 5b: Contours of a bivariate standard Gaussian distribution

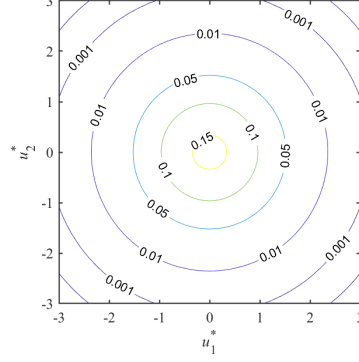


Figure 5c: Bivariate standard asymmetric Student  $t$  distribution density

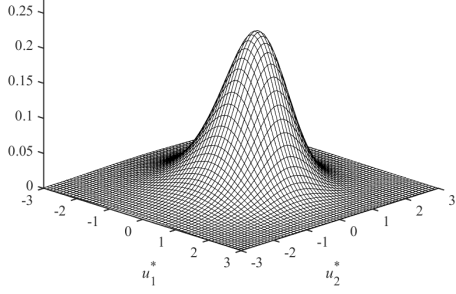


Figure 5d: Contours of a bivariate standard asymmetric Student  $t$  distribution

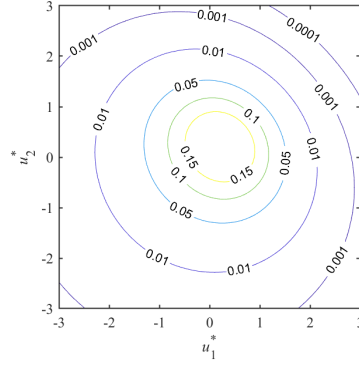


Figure 5e: Bivariate standard location-scale mixture of two normals

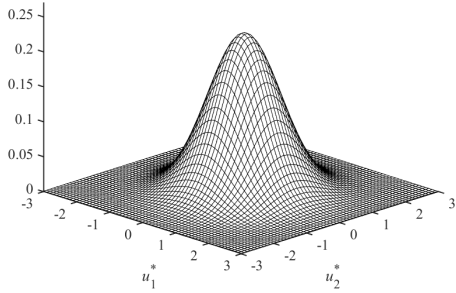
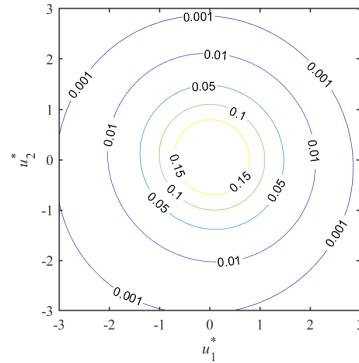


Figure 5f: Contours of a bivariate standard location-scale mixture of two normals



*Notes:* In all panels,  $E(u_{1i}^*) = E(u_{2i}^*) = 0$ ,  $V(u_{1i}^*) = V(u_{2i}^*) = 1$ , and  $cov(u_{1i}^*, u_{2i}^*) = 0$ . Panels c-d: Asymmetric Student  $t$  density with  $\nu = 9.65$  degrees of freedom, skewness parameters  $b_i = -1$ . Panels e-f: Location-scale mixture of two normals with mixing probability  $\lambda = 0.05$ , location vector  $\delta = -(1.01, 1.06)'$  and scale parameter  $\varkappa = 0.32$  (see Appendix D for details).

Figure 6: Monte Carlo results:  $T = 250$ ,  $\mu_{22} = 3$  and  $\rho_{y_2 z_2, z_1} = (\mu_{22} - 1)^{-1}$  (Tie)

Figure 6a: Student  $t$  innovations

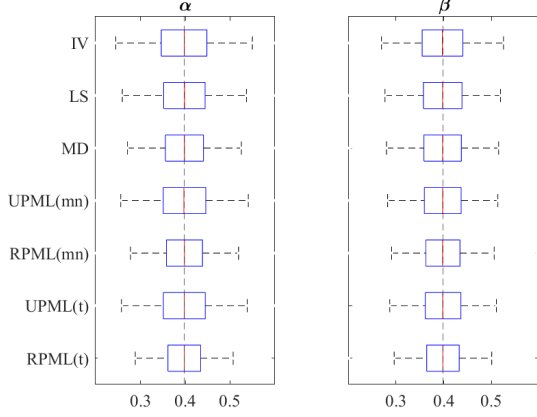


Figure 6b: Scale mixture of two normals

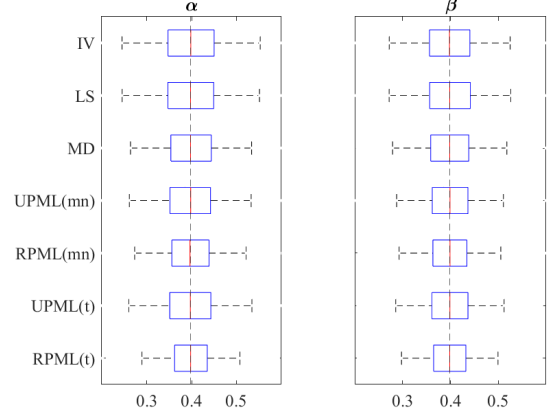


Figure 6c: Asymmetric Student  $t$  innovations

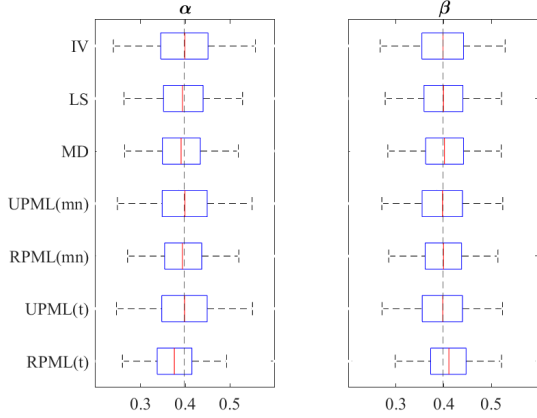
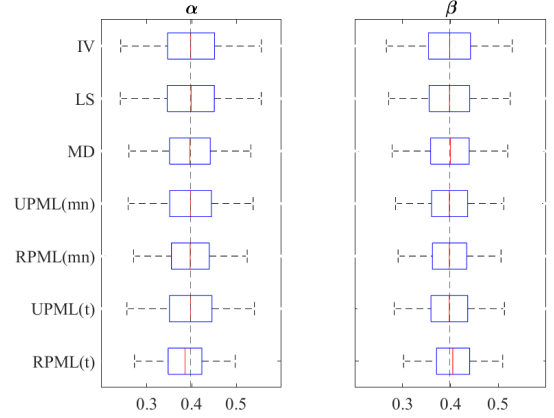


Figure 6d: Location-scale mixture of two normals



*Notes:* IV denotes the instrumental variables estimator, LS denotes the ordinary least squares estimator, MD denotes the optimum minimum distance estimator, UPML(mn) and RPML(mn) [UPML(t) and RPML(t)] denote the unrestricted and restricted ( $\sigma_{12} = 0$ ) Pseudo-ML estimators based on the mixture of two normals (mn) [Student- $t$  (t)] structural innovations. DGPs: Panel a: Student  $t$  distribution with  $\nu = 5$  degrees of freedom; Panel b: scale mixture of two normals with scale parameter  $\varkappa = 0.09$  and mixing probability  $\lambda = 0.05$ ; Panel c: asymmetric Student  $t$  density with  $\nu = 9.65$  degrees of freedom, skewness parameters  $b_i = -1$ ; and Panel d: location-scale mixture of two normals with mixing probability  $\lambda = 0.05$ , location vector  $\delta = -(1.01, 1.06)'$  and scale parameter  $\varkappa = 0.32$  (see Appendix D for details). In all DGPs, we set  $\sigma_1^2 = 1$  so that  $R_1^2 = 1/2$ . In order to have a tie between IV and LS, we set  $\rho_{y_2 z_2, z_1} = 1/\sqrt{2}$  so that  $\sigma_2^2 = 1/3$  and, therefore,  $\gamma = 0.20$ ,  $\alpha = \beta = 0.40$ ,  $\mu_0 = 0.15$  and  $\mu_1 = \mu_2 = 0.58$ .

Figure 7: Monte Carlo results:  $T = 250$ ,  $\mu_{22} = 3$  and  $\rho_{y_2 z_2, z_1} = \frac{1}{2}\mu_{22}(\mu_{22} - 1)^{-1}$  (Max)

Figure 7a: Student  $t$  innovations

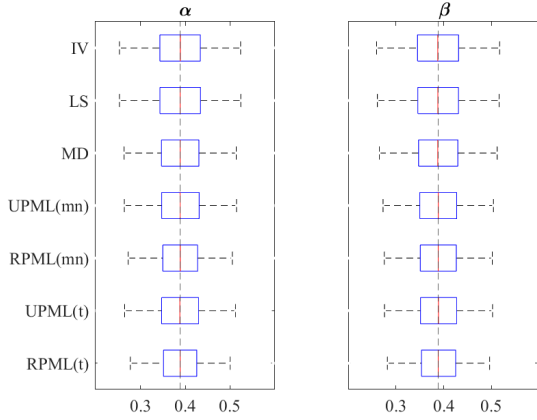


Figure 7b: Scale mixture of two normals

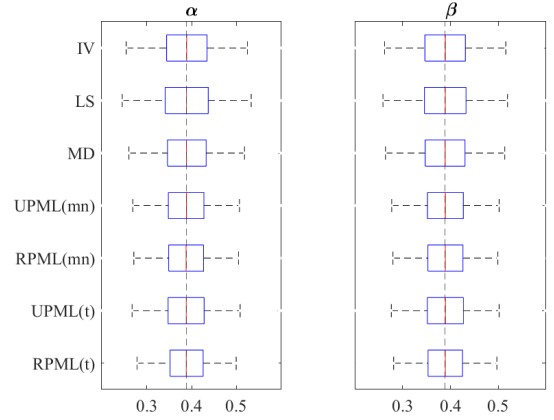


Figure 7c: Asymmetric Student  $t$  innovations

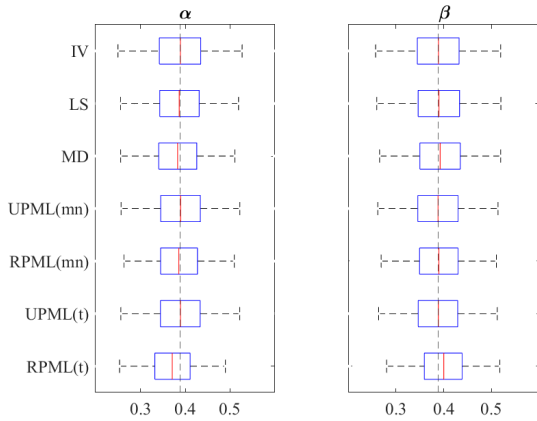
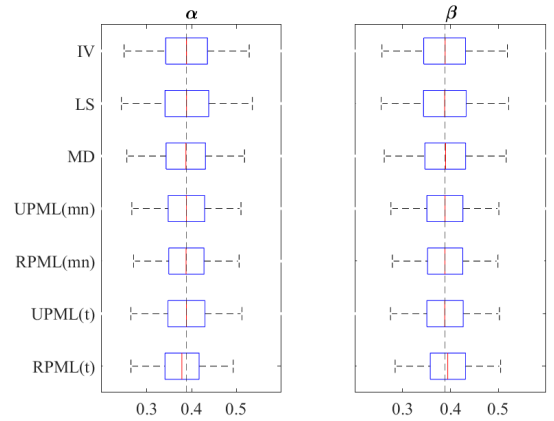


Figure 7d: Location-scale mixture of two normals



*Notes:* IV denotes the instrumental variables estimator, LS denotes the ordinary least squares estimator, MD denotes the optimum minimum distance estimator, UPML(mn) and RPML(mn) [UPML(t) and RPML(t)] denote the unrestricted and restricted ( $\sigma_{12} = 0$ ) Pseudo-ML estimators based on the mixture of two normals (mn) [Student-t (t)] structural innovations. DGPs: Panel a: Student  $t$  distribution with  $\nu = 5$  degrees of freedom; Panel b: scale mixture of two normals with scale parameter  $\varkappa = 0.09$  and mixing probability  $\lambda = 0.05$ ; Panel c: asymmetric Student  $t$  density with  $\nu = 9.65$  degrees of freedom, skewness parameters  $b_i = -1$ ; and Panel d: location-scale mixture of two normals with mixing probability  $\lambda = 0.05$ , location vector  $\delta = -(1.01, 1.06)'$  and scale parameter  $\varkappa = 0.32$  (see Appendix D for details). In all DGPs, we set  $\sigma_1^2 = 1$  so that  $R_1^2 = 1/2$ . In order to have a maximum relative efficiency of IV versus LS, we set  $\rho_{y_2 z_2, z_1} = \sqrt{3}/2$  so that  $\sigma_2^2 = 1/7$  and, therefore,  $\gamma = 0.22$ ,  $\alpha = \beta = 0.39$ ,  $\mu_0 = 0.31$  and  $\mu_1 = \mu_2 = 0.65$ .

Figure 8: Monte Carlo results:  $T = 250$ ,  $\mu_{22} = 7/3$  and  $\rho_{y_2 z_2, z_1} = (\mu_{22} - 1)^{-1}$  (Tie)

Figure 8a: Student  $t$  innovations

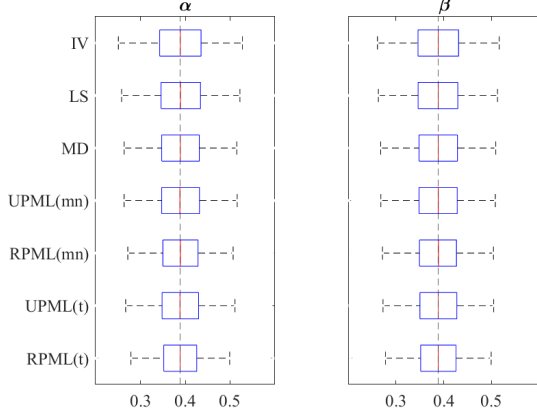


Figure 8b: Scale mixture of two normals

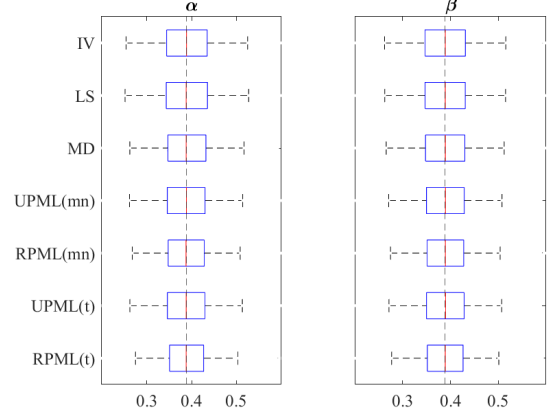


Figure 8c: Asymmetric Student  $t$  innovations

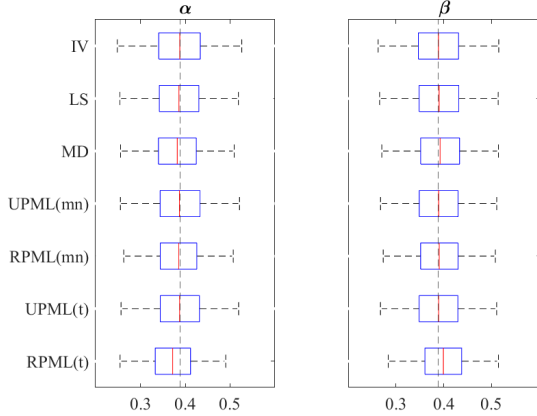
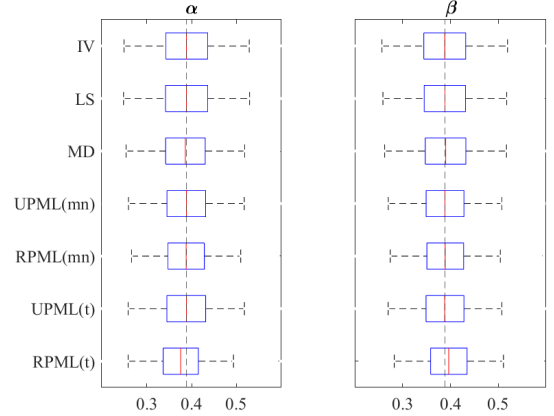


Figure 8d: Location-scale mixture of two normals



*Notes:* IV denotes the instrumental variables estimator, LS denotes the ordinary least squares estimator, MD denotes the optimum minimum distance estimator, UPML(mn) and RPML(mn) [UPML(t) and RPML(t)] denote the unrestricted and restricted ( $\sigma_{12} = 0$ ) Pseudo-ML estimators based on the mixture of two normals (mn) [Student-t (t)] structural innovations. DGPs: Panel a: Student  $t$  distribution with  $\nu = 11/2$  degrees of freedom; Panel b: scale mixture of two normals with scale parameter  $\varkappa = 0.12$  and mixing probability  $\lambda = 0.05$ ; Panel c: asymmetric Student  $t$  density with  $\nu = 10.38$  degrees of freedom, skewness parameters  $b_i = -1$ ; and Panel d: location-scale mixture of two normals with mixing probability  $\lambda = 0.05$ , location vector  $\delta = -(1.16, 1.24)'$  and scale parameter  $\varkappa = 0.38$  (see Appendix D for details). In all DGPs, we set  $\sigma_1^2 = 1$  so that  $R_1^2 = 1/2$ . In order to have a tie between IV and LS, we set  $\rho_{y_2 z_2, z_1} = \sqrt{3}/2$  so that  $\sigma_2^2 = 1/7$  and, therefore,  $\gamma = 0.22$ ,  $\alpha = \beta = 0.39$ ,  $\mu_0 = 0.31$  and  $\mu_1 = \mu_2 = 0.65$ .

Figure 9: Monte Carlo results:  $T = 1,000$ ,  $\mu_{22} = 3$  and  $\rho_{y_2 z_2, z_1} = (\mu_{22} - 1)^{-1}$  (Tie)

Figure 9a: Student  $t$  innovations

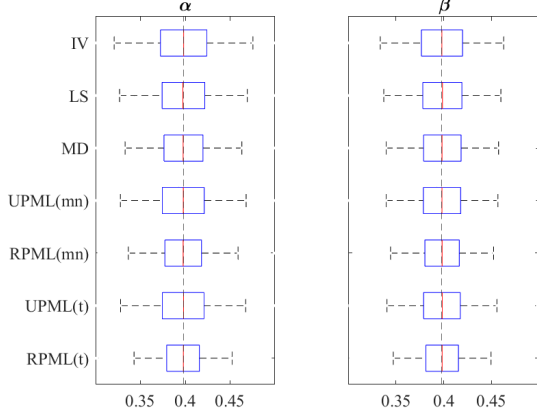


Figure 9b: Scale mixture of two normals

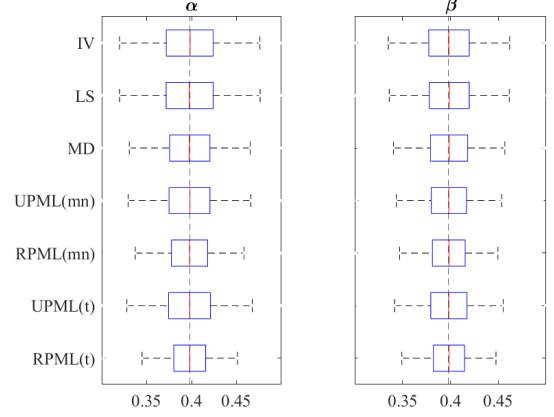


Figure 9c: Asymmetric Student  $t$  innovations

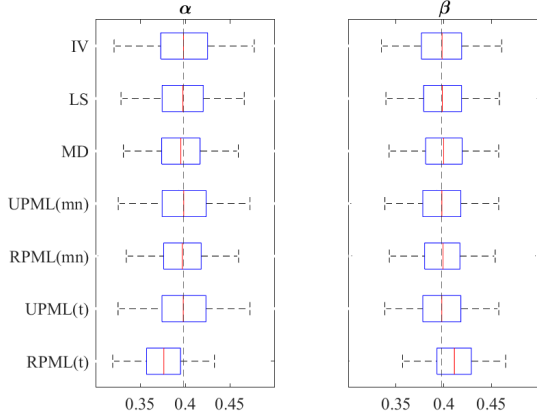
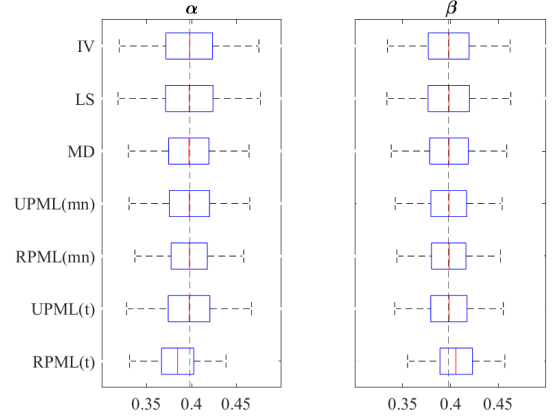


Figure 9d: Location-scale mixture of two normals



*Notes:* IV denotes the instrumental variables estimator, LS denotes the ordinary least squares estimator, MD denotes the optimum minimum distance estimator, UPML(mn) and RPML(mn) [UPML(t) and RPML(t)] denote the unrestricted and restricted ( $\sigma_{12} = 0$ ) Pseudo-ML estimators based on the mixture of two normals (mn) [Student-t (t)] structural innovations. DGPs: Panel a: Student  $t$  distribution with  $\nu = 5$  degrees of freedom; Panel b: scale mixture of two normals with scale parameter  $\varkappa = 0.09$  and mixing probability  $\lambda = 0.05$ ; Panel c: asymmetric Student  $t$  density with  $\nu = 9.65$  degrees of freedom, skewness parameters  $b_i = -1$ ; and Panel d: location-scale mixture of two normals with mixing probability  $\lambda = 0.05$ , location vector  $\delta = -(1.01, 1.06)'$  and scale parameter  $\varkappa = 0.32$  (see Appendix D for details). In all DGPs, we set  $\sigma_1^2 = 1$  so that  $R_1^2 = 1/2$ . In order to have a tie between IV and LS, we set  $\rho_{y_2 z_2, z_1} = 1/\sqrt{2}$  so that  $\sigma_2^2 = 1/3$  and, therefore,  $\gamma = 0.20$ ,  $\alpha = \beta = 0.40$ ,  $\mu_0 = 0.15$  and  $\mu_1 = \mu_2 = 0.58$ .

Figure 10: Monte Carlo results:  $T = 1,000$ ,  $\mu_{22} = 3$  and  $\rho_{y_2 z_2 \cdot z_1} = \frac{1}{2}\mu_{22}(\mu_{22} - 1)^{-1}$  (Max)

Figure 10a: Student  $t$  innovations

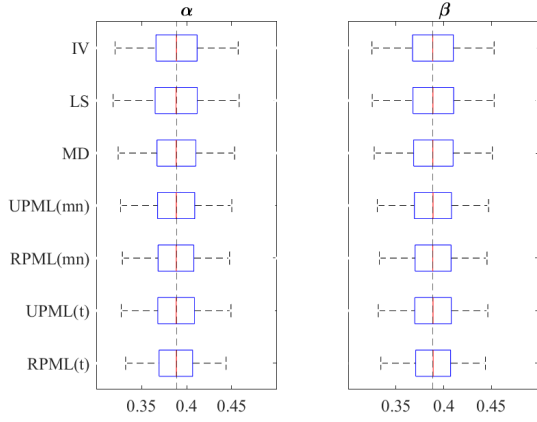


Figure 10b: Scale mixture of two normals

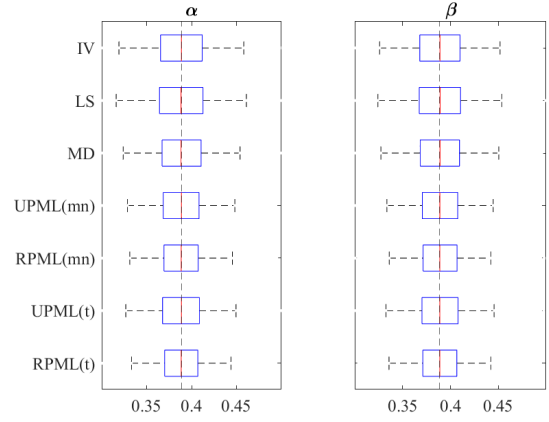


Figure 10c: Asymmetric Student  $t$  innovations

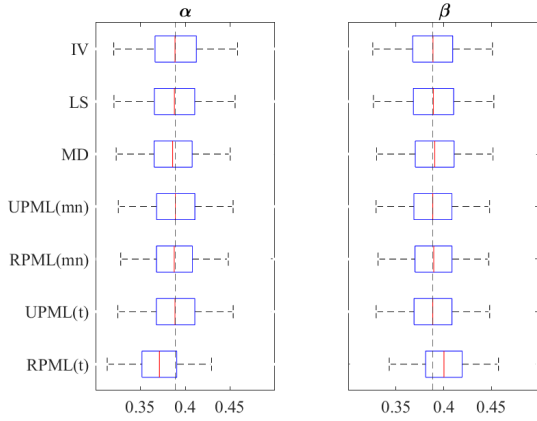
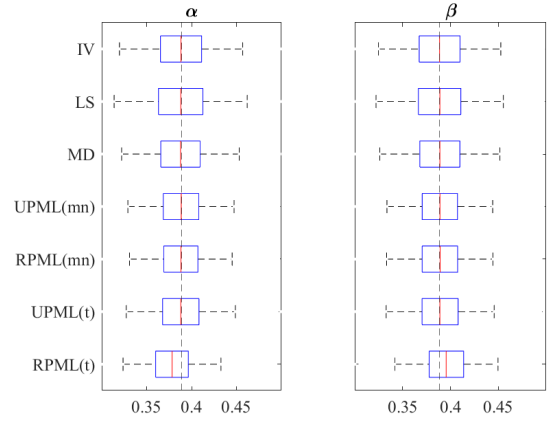


Figure 10d: Location-scale mixture of two normals



*Notes:* IV denotes the instrumental variables estimator, LS denotes the ordinary least squares estimator, MD denotes the optimum minimum distance estimator, UPML(mn) and RPML(mn) [UPML(t) and RPML(t)] denote the unrestricted and restricted ( $\sigma_{12} = 0$ ) Pseudo-ML estimators based on the mixture of two normals (mn) [Student-t (t)] structural innovations. DGPs: Panel a: Student  $t$  distribution with  $\nu = 5$  degrees of freedom; Panel b: scale mixture of two normals with scale parameter  $\varkappa = 0.09$  and mixing probability  $\lambda = 0.05$ ; Panel c: asymmetric Student  $t$  density with  $\nu = 9.65$  degrees of freedom, skewness parameters  $b_i = -1$ ; and Panel d: location-scale mixture of two normals with mixing probability  $\lambda = 0.05$ , location vector  $\delta = -(1.01, 1.06)'$  and scale parameter  $\varkappa = 0.32$  (see Appendix D for details). In all DGPs, we set  $\sigma_1^2 = 1$  so that  $R_1^2 = 1/2$ . In order to have a maximum relative efficiency of IV versus LS, we set  $\rho_{y_2 z_2 \cdot z_1} = \sqrt{3}/2$  so that  $\sigma_2^2 = 1/7$  and, therefore,  $\gamma = 0.22$ ,  $\alpha = \beta = 0.39$ ,  $\mu_0 = 0.31$  and  $\mu_1 = \mu_2 = 0.65$ .

Figure 11: Monte Carlo results:  $T = 1,000$ ,  $\mu_{22} = 7/3$  and  $\rho_{y_2 z_2, z_1} = (\mu_{22} - 1)^{-1}$  (Tie)

Figure 11a: Student  $t$  innovations

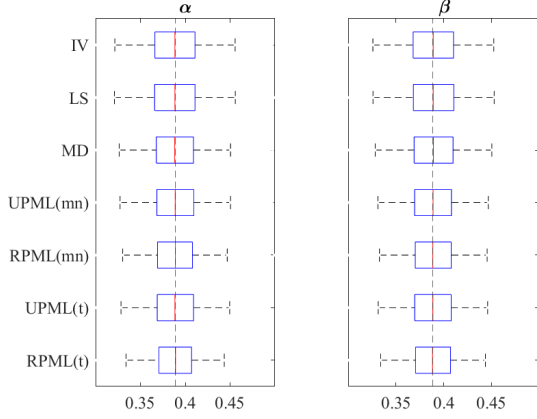


Figure 11b: Scale mixture of two normals

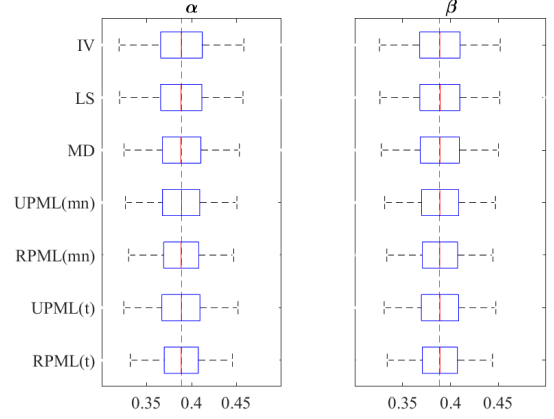


Figure 11c: Asymmetric Student  $t$  innovations

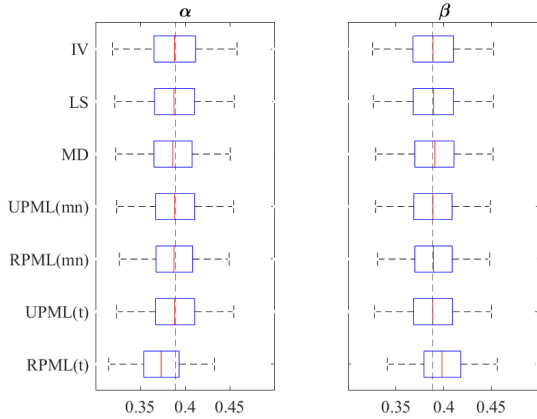
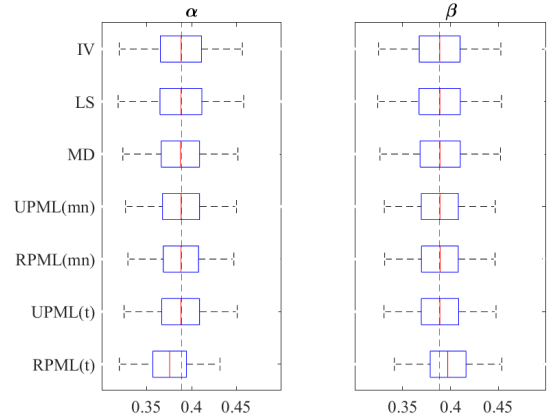


Figure 11d: Location-scale mixture of two normals



*Notes:* IV denotes the instrumental variables estimator, LS denotes the ordinary least squares estimator, MD denotes the optimum minimum distance estimator, UPML(mn) and RPML(mn) [UPML(t) and RPML(t)] denote the unrestricted and restricted ( $\sigma_{12} = 0$ ) Pseudo-ML estimators based on the mixture of two normals (mn) [Student-t (t)] structural innovations. DGPs: Panel a: Student  $t$  distribution with  $\nu = 11/2$  degrees of freedom; Panel b: scale mixture of two normals with scale parameter  $\varkappa = 0.12$  and mixing probability  $\lambda = 0.05$ ; Panel c: asymmetric Student  $t$  density with  $\nu = 10.38$  degrees of freedom, skewness parameters  $b_i = -1$ ; and Panel d: location-scale mixture of two normals with mixing probability  $\lambda = 0.05$ , location vector  $\delta = -(1.16, 1.24)'$  and scale parameter  $\varkappa = 0.38$  (see Appendix D for details). In all DGPs, we set  $\sigma_1^2 = 1$  so that  $R_1^2 = 1/2$ . In order to have a tie between IV and LS, we set  $\rho_{y_2 z_2, z_1} = \sqrt{3}/2$  so that  $\sigma_2^2 = 1/7$  and, therefore,  $\gamma = 0.22$ ,  $\alpha = \beta = 0.39$ ,  $\mu_0 = 0.31$  and  $\mu_1 = \mu_2 = 0.65$ .