

A Comparison of Some Recent Bayesian and Classical Procedures for Simultaneous Equation Models with Weak Instruments

Chuanming Gao, Kajal Lahiri*

Department of Economics, State University of New York at Albany,
Albany, NY 12222

Abstract: We compare the finite sample performance of a number of Bayesian and classical procedures for limited information simultaneous equations models with weak instruments by a Monte Carlo study. We consider recent Bayesian approaches developed by Chao and Phillips (1998, CP), Geweke (1996), Kleibergen and van Dijk (1998, KVD), and Zellner (1998). Amongst the sampling theory methods, OLS, 2SLS, LIML, Fuller's modified LIML, and the jackknife instrumental variable estimator (JIVE) due to Angrist, Imbens and Krueger (1999) and Blomquist and Dahlberg (1999) are also considered: Since the posterior densities and their conditionals in CP and KVD are non-standard, we propose a "Gibbs within Metropolis-Hastings" algorithm, which only requires the availability of the conditional densities from the candidate generating density. Our results show that in cases with very weak instruments, there is no single estimator that is superior to others in all cases. When endogeneity is weak, Zellner's MELO does the best. When the endogeneity is not weak and $\frac{1}{2}w_{12} > 0$; where $\frac{1}{2}$ is the correlation coefficient between the structural and reduced form errors, and w_{12} is the covariance between the unrestricted reduced form errors, BMOM outperforms all other estimators by a wide margin. When the endogeneity is not weak and $-\frac{1}{2} < 0$ (γ being the structural parameter), KVD approach seems to work very well. Surprisingly, the performance of JIVE was disappointing in all our experiments.

JEL classification: C30, C11, C13, C15

Keywords: Limited information estimation, Metropolis-Hastings algorithm, Gibbs sampler, Monte Carlo method

* Corresponding author. E-mail: KL758@cnsvox.albany.edu. Tel: (518) 442 4758; fax: (518) 442 4736. An earlier version of this paper was presented at the 1999 Joint Statistical Meetings and the 2000 Winter Meetings of the Econometric Society. We are grateful to Ingolf Dittmann, Terrence Kinal, Yuichi Kitamura, Roberto Mariano, Arnold Zellner, and Eric Zivot for many helpful comments and suggestions. However, the responsibility for any remaining errors and shortcomings is solely ours.

1 Introduction

Recent research on Bayesian analysis of the simultaneous equations models addresses a problem, raised initially by Maddala (1976), and now recognized as related to the problem of local non-identification when diffuse priors are used in traditional Bayesian analysis of such models, e.g., Dr̕ze (1976), Dr̕ze and Morales (1976), and Dr̕ze and Richard (1983).¹ In this paper, we will examine the approaches developed by Chao and Phillips (1998, hereafter CP), Geweke (1996), Kleibergen and van Dijk (1998, hereafter KVD), and Zellner (1998). The idea in KVD is to treat an overidentified simultaneous equations model (SEM) as a linear model with nonlinear parameter restrictions. While KVD focuses mainly on resolving the problem of local nonidentification, CP explores further the consequences of using a Jeffreys prior. By deriving the exact and (asymptotically) approximate representations for the posterior density of the structural parameter, they show that the use of a Jeffreys prior brings Bayesian inference closer to classical inference in the sense that this prior choice leads to posterior distributions which exhibit Cauchy-like tail behavior like the LIML estimator. Geweke (1996), being aware of the potential problem of local nonidentification, suggests a shrinkage prior such that the posterior density is properly defined for each parameter. In another novel approach, Zellner (1998) has developed a finite sample Bayesian method of moments (BMOM) procedure based on given data without specifying a likelihood function or introducing any sampling

¹Zellner (1998) provides the latest comprehensive review of the finite sample properties of SEM estimators, and emphasizes the need for finite sample optimal estimation procedure for such models.

assumptions.

For the Bayesian approaches considered, while Geweke (1996) proposes Gibbs sampling (GS) to evaluate the posterior density with a shrinkage prior, the posterior densities as well as their conditional densities resulting from CP and KVD are non-standard and cannot be readily simulated. In the category of "block-at-a-time" approach, we suggest a new MCMC procedure, which we call a "Gibbs within M-H" algorithm. The advantage of this algorithm is that it only requires the availability of the conditional densities from the candidate generating density. These conditional densities are used in a Gibbs sampler to simulate the candidate generating density, whose drawings, after convergence, are then weighted to generate drawings from the target density in a Metropolis-Hastings (M-H) algorithm. In this study, we will focus on weak instruments, where the classical approach has been particularly unsatisfactory.²

The main objective of the present paper is to compare the small sample performance of some Bayesian and classical approaches using Monte Carlo simulations. For the purpose of comparison, a number of classical methods including OLS, 2SLS, LIML, Fuller's modified LIML, and a recent jackknife instrumental variables estimator (JIVE) due to Angrist, Imbens and Krueger (1999) and Blomquist and Dahlberg (1999) are also computed from the generated data. Our simulation results from repeated sampling experiments provide some unambiguous guidelines for empirical practitioners.

The plan of the paper is as follows. In Section 2, we set up the model.

²There has been a growing interest in the estimation of LISEM with weak instruments. See Buse (1992), Bound, Jaeger and Baker (1995), Staiger and Stock (1997), Angrist, Imbens and Krueger (1999), Blomquist and Dahlberg (1999), among others.

Section 3 reviews in limited details the recent Bayesian approaches and JIVE. Section 4 suggests a new MCMC procedure for evaluating the posterior distributions for CP and KVD, and discusses the convergence diagnostics implemented. Section 5 presents simulation results and some discussions. Section 6 contains the main conclusions.

2 The Model

Consider the following limited information formulation of the m -equation simultaneous equations model (LISEM):

$$y_1 = \gamma_2' Y_2 + Z_1' \alpha + u; \quad (1)$$

$$Y_2 = Z_1' \beta_1 + Z_2' \beta_2 + V_2; \quad (2)$$

where $y_1 : (T \times 1)$ and $Y_2 : (T \times (m_j - 1))$ are the m included endogenous variables; $Z_1 : (T \times k_1)$ is an observation matrix of exogenous variables included in the structural equation (1); $Z_2 : (T \times k_2)$ is an observation matrix of exogenous variables excluded from (1); and u and V_2 are, respectively, a $T \times 1$ vector and a $T \times (m_j - 1)$ matrix of random disturbances to the system. We assume that $(u; V_2) \gg N(0; S - I_T)$, where the $m \times m$ covariance matrix S is positive definite symmetric (pds) and is partitioned conformably with the rows of $(u; V_2)$ as follows

$$S = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & S_{22} \end{pmatrix} :$$

The likelihood function for the model described by (1) and (2) can be

written as

$$L(\beta; \sigma; \Gamma_1; \Gamma_2; S|Y; Z) = (2\pi)^{-i} |S|^{-1} \exp\left\{-\frac{1}{2} \text{tr}[S^{-1}(u; V_2)'(u; V_2)]\right\} \quad (3)$$

where $Y = (y_1; Y_2)$ and $Z = (Z_1; Z_2)$:

The structural model described by (1) and (2) can alternatively be written in its reduced form

$$\begin{bmatrix} y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \quad (4)$$

where $\mu_1 = \beta_1 + \Gamma_1^{-1} \beta_2$, $\mu_2 = u + V_2^{-1} \beta_2$, $(\mu_1; \mu_2) \sim N(0; -I_T)$, $S = C' - C$, $C = \begin{bmatrix} 1 & 0 \\ 0 & I_{m_1} \end{bmatrix}$. The likelihood function corresponding to this alternative representation is:

$$L(\beta; \sigma; \Gamma_1; \Gamma_2; -jY; Z) = (2\pi)^{-i} |S|^{-1} \exp\left\{-\frac{1}{2} \text{tr}[-i^{-1}(\mu_1; \mu_2)'(\mu_1; \mu_2)]\right\} \quad (5)$$

The likelihood functions (3) and (5) are equivalent since the Jacobian between β and S is unity.

Geweke (1996) considers the following reduced rank regression specification³

$$Y = Z_1 A + Z_2 E + E; \quad (6)$$

where $A = (\beta_1; \mu_1)$, $E = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix}$ and $\epsilon = (I_{m_1}; -)$; $E = (V_2; \mu_1) \sim N(0; \underline{S} - I_T)$ with $\underline{S}^{-1} = \begin{bmatrix} S^{11} & S^{12} \\ S^{21} & S^{22} \end{bmatrix}$ partitioned conformably with the rows of $(V_2; \mu_1)$: Obviously, (6) is equivalent to (4) and the corresponding likelihood function is similar to (5).

Note that in the absence of restrictions on the covariance structure, (1) is fully identified if and only if $\text{rank}(\Gamma_2) = (m_1 - 1) \cdot k_2$.

³Geweke (1996) considered a more general specification. To facilitate comparison, for Geweke approach only, we have denoted $Y = (Y_2; y_1)$:

3 Review of some recent formulations

Among the most recent Bayesian approaches, Geweke (1996) used a shrinkage prior such that all parameters are identified (in the sense that a proper posterior distribution exists) even when β_2 has reduced rank. KVD treated overidentified SEMs as linear models with nonlinear parameter restrictions using the singular value decomposition. A diffuse or natural conjugate prior for the parameters of the embedding linear model results in the posterior for the parameters of the SEM having zero weight in the region of parameter space where β_2 has reduced rank. This is a feature of the Jacobian of transformation from the multivariate linear model to the SEM. CP used a prior by applying Jeffreys principle on the model described by (1) and (2) and the assumptions regarding the disturbances. An important quality of the Jeffreys prior in the present context is that it places no weight in the region of the parameter space where $\text{rank}(\beta_2) < (m_2 - 1)$ and relatively low weight in close neighborhoods of this region where the model is nearly unidentified.

3.1 Zellner's Bayesian method of moments approach (BMOM)

Among the various Bayesian treatments of SEM proposed by Zellner (1971, 1978, 1986, 1994, 1998), the recent Bayesian method of moments approach applies the principle of maximum entropy and generates optimal estimates which can be evaluated by double K-class estimators. Given the unrestricted reduced form equation $y_1 = Z\beta_1 + \epsilon_1$; Zellner (1998) considered a balanced

loss function,

$$L_b = \lambda L_g + (1 - \lambda)L_p$$

$$= \lambda (y_1' X \beta)^0 (y_1' X \beta) + (1 - \lambda) (Z' \beta_1 X \beta)^0 (Z' \beta_1 X \beta), \text{ for } 0 \leq \lambda \leq 1$$

where $X = (Y_2; Z_1)$; $\beta = (\beta_2; \beta_1)$; and $\hat{\beta}$ is an estimate of β . The BMOM estimate that minimizes $E L_b$; where the expectation is taken with respect to a probability density function of the β matrices appeared in unrestricted reduced form equations, is given by

$$\hat{\beta} = \begin{pmatrix} \beta_2 \\ \beta_1 \end{pmatrix} = \begin{pmatrix} Y_2' Y_2 & -K_1 \beta_2' \beta_2 \\ Z_1' Y_2 \end{pmatrix}^{-1} \begin{pmatrix} Y_2' Z_1 \\ Z_1' Y_1 \end{pmatrix} + \begin{pmatrix} (Y_2' Y_2 - K_2 \beta_2' \beta_2)^{-1} Y_2' \\ Z_1' Y_1 \end{pmatrix} \beta_1; \quad (7)$$

where

$$K_1 = 1 - \lambda; K_2 = \lambda; \text{ with } 0 \leq \lambda \leq 1;$$

$$\text{and } \beta_2 = (I - Z(Z'Z)^{-1}Z')Y_2;$$

BMOM estimate will vary depending on the value of λ . When $\lambda = 1$; it is the optimal estimate resulting from a "goodness of fit" loss function L_g . When $\lambda = 0$; it is the optimal estimate given by a precision of estimation loss function L_p . Meanwhile, the well-known minimum expected loss (MELO) estimator is derived using a precision of estimation loss function and may be evaluated as a K-class estimator with

$$K_1 = K_2 = 1 - \lambda; \text{ with } 0 \leq \lambda \leq 1;$$

3.2 The Geweke (1996) approach

Geweke (1996) assumes the following reference prior

$$j_S j_i^{(m+\nu+1)=2} \exp\left[-\frac{1}{2} \text{tr} S S^{-1}\right] \exp\left[-\frac{\lambda^2}{2} (\beta_2' + \text{tr} \beta_2' \beta_2 + \text{tr} A' A)\right]; \quad (8)$$

which is the product of an independent inverted Wishart distribution for \mathbb{S} with ν degrees of freedom and scale matrix S , and an independent $N(0; \lambda^2)$ shrinkage priors for each element of β and γ_2 : Geweke derived the respective conditional posterior distributions, which may be used to generate drawings through Gibbs sampling from the joint posterior distribution. Regarding the vector of parameters $(\mathbb{S}^{i-1}; \mathbf{A}; \gamma_2; \beta)$; we obtain the full conditional densities as follows:

(1) Conditional density of \mathbb{S}^{i-1}

$$\mathbb{S}^{i-1} | (\gamma_2; \beta; \mathbf{A}; \mathbf{Z}; \mathbf{Y}) \gg \text{Wishart}(T + \nu; \mathbf{G}^{i-1}); \quad (9)$$

where $\mathbf{G} = \mathbf{S} + (\mathbf{Y}_i - \mathbf{Z}_1 \mathbf{A}_i - \mathbf{Z}_2 \boldsymbol{\epsilon})^0 (\mathbf{Y}_i - \mathbf{Z}_1 \mathbf{A}_i - \mathbf{Z}_2 \boldsymbol{\epsilon})$:

(2) Conditional density of \mathbf{A}

$$\begin{aligned} & \text{vec}(\mathbf{A}) | (\gamma_2; \beta; \mathbb{S}^{i-1}; \mathbf{Z}; \mathbf{Y}) \\ & \gg N([\mathbb{S}^{i-1} - \mathbf{Z}_1^0 \mathbf{Z}_1 + \lambda^2 \mathbf{I}_{mk_1}]^{i-1} [\mathbb{S}^{i-1} - \mathbf{Z}_1^0 \mathbf{Z}_1] \text{vec}(\mathbf{A}); \\ & [\mathbb{S}^{i-1} - \mathbf{Z}_1^0 \mathbf{Z}_1 + \lambda^2 \mathbf{I}_{mk_1}]^{i-1}); \end{aligned} \quad (10)$$

where $\mathbf{A} = (\mathbf{Z}_1^0 \mathbf{Z}_1)^{i-1} \mathbf{Z}_1^0 (\mathbf{Y}_i - \mathbf{Z}_2 \boldsymbol{\epsilon})$:

(3) Conditional density of γ_2^4

$$\begin{aligned} & \text{vec}(\gamma_2) | (\beta; \mathbb{S}^{i-1}; \mathbf{A}; \mathbf{Z}; \mathbf{Y}) \\ & \gg N([\mathbb{S}^{11} - \mathbf{Z}_2^0 \mathbf{Z}_2 + \lambda^2 \mathbf{I}_{k_2(m_i-1)}]^{i-1} [\mathbb{S}^{11} - \mathbf{Z}_2^0 \mathbf{Z}_2] \text{vec}(\gamma_2); \\ & [\mathbb{S}^{11} - \mathbf{Z}_2^0 \mathbf{Z}_2 + \lambda^2 \mathbf{I}_{k_2(m_i-1)}]^{i-1}); \end{aligned} \quad (11)$$

⁴The expressions for the conditional densities of γ_2 and β given in Geweke (1996, expressions (11) and (13)) contain some typographical errors and are corrected here in (11) and (12).

where $\ln L(\mu; Y; Z)$ is the log-likelihood function as specified in (3), and $Q_X = I_T \otimes P_X$, $P_X = X(X'X)^{-1}X'$. As first noted by Poirier (1996), the prior in (13) places no weight where $\text{rank}(\Sigma_2) < (m - 1)$ through the factor $j^{\frac{1}{2}} Z_2' Q_{Z_1} Z_2 j^{1-2}$.

The joint posterior of the parameters of LISEM (1) and (2) is constructed as proportional to the product of the prior (13) and the likelihood function (3),

$$\begin{aligned}
 p(\beta; \sigma; \Sigma_1; \Sigma_2; S; Y; Z) &= p(\beta; \sigma; \Sigma_1; \Sigma_2; S) L(\beta; \sigma; \Sigma_1; \Sigma_2; S; Y; Z) \\
 &= j^{\frac{1}{4} 11} j^{(k_2 - m + 1) = 2} j^{\frac{1}{2}} j^{(T + k + m + 1) = 2} j^{\frac{1}{2}} Z_2' Q_{Z_1} Z_2 j^{1-2} \\
 &\quad \propto \exp\left\{-\frac{1}{2} \text{tr}[S^{-1}(u; V_2)'(u; V_2)]\right\} g; \quad (14)
 \end{aligned}$$

where $(u; V_2)$ is defined in (1) and (2). Note that (14) or its conditionals do not belong to any standard class of probability density functions.

3.4 The Kleibergen and van Dijk (1998) approach

To solve the problem of local nonidentification and also to avoid the so called Borel-Kolmogorov paradox [see Billingsley (1986) and Poirier (1995)], KVD considered (4) as a multivariate linear model with nonlinear parameter restrictions:

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \begin{bmatrix} \mu \\ \sigma \end{bmatrix} \begin{bmatrix} \frac{1}{A_1} \\ \otimes_2 \end{bmatrix} + \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}; \quad (15)$$

where A_1 is a $k_2 \times 1$ vector, \otimes_2 is a $k_2 \times (m - 1)$ matrix. Denote $\otimes = (A_1; \otimes_2)$. The reduced form model (4) is obtained if a reduced rank restriction is imposed on the linear model (15) such that $\text{rank}(\otimes) = (m - 1)$ instead of

m. Using a singular value decomposition (SVD) of \odot , they show that (15) is identical to the so-called unrestricted reduced form (URF) model,⁵

$$y_1 - Y_2 = Z_1 \beta_1 + Z_2 B + Z_2 B_\perp + \epsilon_1 - V_2 \epsilon_2; \quad (16)$$

where $B = \beta_1 - I_{m_i-1}$ is a $(k_2 - m + 1) \times 1$ vector. Z_2 and B_\perp are the orthogonal complements of Z_1 and B respectively, such that $Z_2' Z_1 = 0$, $B B_\perp' = 0$, and $Z_2' Z_2 = I_{k_2 - m + 1}$, $B_\perp B_\perp' = 1$ (i.e. $Z_2 = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$, $\beta_1 = (k_2 - m + 1) \times 1$, $\beta_2 = (k_2 - m + 1) \times (m_i - 1)$, where $\beta_2 = \begin{bmatrix} \beta_{21} \\ \beta_{22} \end{bmatrix}$, $\beta_{21} : (m_i - 1) \times (m_i - 1)$; $\beta_{22} : (k_2 - m + 1) \times (m_i - 1)$, and $B_\perp = (1 - \beta_{21}^{-1}) \beta_{21}^{-1} \beta_{22}$):

There is one-to-one correspondence between the parameters in (15) and (16). The SVD of \odot is,

$$\odot = U S V'; \quad (17)$$

where $U : k_2 \times k_2$, $U'U = I_{k_2}$; $V : m \times m$, $V'V = I_m$; and $S : k_2 \times m$ is a rectangular matrix containing the (nonnegative) singular values (in decreasing order) on its main diagonal ($= (s_{11}; s_{22}; \dots; s_{mm})$). Rewrite

$$U = \begin{bmatrix} U_{11} & U_{12} \\ U_{21} & U_{22} \end{bmatrix}, S = \begin{bmatrix} s_1 & 0 \\ 0 & s_2 \end{bmatrix} \text{ and } V = \begin{bmatrix} v_{11} & v_{12} \\ v_{21} & v_{22} \end{bmatrix}; \quad (18)$$

where $U_{11}, s_1, v_{21} : (m_i - 1) \times (m_i - 1)$; $v_{12} : 1 \times 1$; $v_{11}, v_{22} : (m_i - 1) \times 1$; $U_{12} : (m_i - 1) \times (k_2 - m + 1)$; $U_{21} : (k_2 - m + 1) \times (m_i - 1)$; $U_{22} : (k_2 - m + 1) \times (k_2 - m + 1)$; $s_2 : (k_2 - m + 1) \times 1$, then the following relationship between $(Z_2'; \epsilon_2)$ and $(U; S; V)$ results,

$$Z_2' = \begin{bmatrix} U_{11} \\ U_{21} \end{bmatrix} S_1 V_{21}', \epsilon_2 = V_{21}^{0'} v_{11}', \text{ and}$$

⁵Note that this formulation or the singular value decomposition does not change the identification status of the LISEM specified by (1) and (2). If $\text{rank}(Z_2) < (m_i - 1)$, ϵ_2 is locally nonidentified.

$$s_2 = (U_{22}U_{22}^0)^{i-1}U_{22}S_2V_{12}^0(V_{12}V_{12}^0)^{i-1} \quad (19)$$

Note that s_2 is obtained through pre- and postmultiplication of s_2 by orthogonal matrices while s_2 contains the smallest singular values of Θ and is invariant with respect to the ordering of variables contained in Y and Z_2 :

According to KVD, the above shows that the model described by (1) and (2) can be considered as equivalent to the linear model (16) with a nonlinear (reduced rank) restriction $s_2 = 0$ on the parameters. Therefore the priors and posteriors of the parameters of the LISEM (1) and (2) may be constructed as proportional to the priors and posteriors of the parameters of the linear model (16) evaluated at $s_2 = 0$:

A diffuse (Jeffreys) prior for the parameters $(\beta_1; \beta_2; \Theta; -)$ of the linear model⁶

$$p(\beta_1; \beta_2; \Theta; -) \propto j^{-j} i^{(k+m+1)=2} \\ \propto j^{-j} i^{(m+1)=2} j^{-i-1} - Z^0 Z j^{1=2}; \quad (20)$$

where $k = k_1 + k_2$, implies the prior for the parameters $(\beta; \beta_1; \beta_2; -)$ of the LISEM (4)

$$p(\beta; \beta_1; \beta_2; -) \\ \propto p(\beta_1; \beta_2; \Theta(\beta_2; \beta; -); -) j_{s=0} J(\Theta; (\beta_2; \beta; -)) j_{s=0} \\ \propto j^{-j} i^{(m+1)=2} j^{-i-1} - Z^0 Z j^{1=2} j(\Theta; (\beta_2; \beta; -)) j_{s=0} \\ \propto j^{-j} i^{(m+1)=2} j^{-i-1} - Z^0 Z j^{1=2} \\ \propto \epsilon^{-i} B^0 - I_{k_2} e_1 - \beta_2 B_2^0 - \beta_2; \quad (21)$$

⁶This is the prior suggested in Dräpe (1976). Zellner (1971) and Zellner, Bauwens and van Dijk (1988) used a similar prior with $i(m+1)=2$ in the exponent.

where \mathbf{X}_{jive} is the $T \times (m_j - 1 + k_1)$ matrix with t -th row defined by

$$\mathbf{Z}_{t \setminus i} \mathbf{b}_{j \setminus i} = \mathbf{Z}_t (\mathbf{Z}_{i \setminus t}^0 \mathbf{Z}_{i \setminus t})^{-1} (\mathbf{Z}_{i \setminus t}^0 \mathbf{X}_{i \setminus t}) = \frac{\mathbf{Z}_t \mathbf{b}_{j \setminus i} - h_t \mathbf{X}_{i \setminus t}}{1 - h_t};$$

$\mathbf{Z}_{i \setminus t}$ and $\mathbf{X}_{i \setminus t}$ are $(T - 1) \times k$ and $(T - 1) \times (m_j - 1 + k_1)$ matrices obtained after eliminating the t -th rows of \mathbf{Z} and \mathbf{X} matrices respectively, $\mathbf{b}_{j \setminus i} = (\mathbf{Z}^0 \mathbf{Z})^{-1} \mathbf{Z}^0 \mathbf{X}$, and $h_t = \mathbf{Z}_t (\mathbf{Z}^0 \mathbf{Z})^{-1} \mathbf{Z}_t^0$. In JIVE, the instrument is independent of the disturbances even in finite samples, which is achieved by using a 'leave-one-out' jackknife-type fitted value in place of the usual unrestricted reduced form predictions.

Angrist, Imbens and Krueger (1999) also proposed a second jackknife estimator that is a slight modification of (23). Similar to their study, we found that its performance is very similar to JIVE, and is not reported here.

4 Posterior simulator: "Gibbs within M-H" algorithm

Given the full conditional densities in (9) through (12) for the four blocks of parameters, evaluating the joint posterior densities in Geweke (1996) by Gibbs sampling is straightforward, see Geweke (1996) for a detailed description. Although Geweke's (1996) shrinkage prior does not meet the argument in KVD that the implied prior/posterior on the parameters of an embedding linear model should be well-behaved, we found that the use of Geweke's shrinkage prior does not lead to a reducible Markov Chain. With the specification of a shrinkage prior, when β_2 has reduced rank, the joint posterior density still depends on β and will not exhibit any asymptotic cusp. In the following we only discuss the posterior simulation for CP and KVD.

KVD suggested two simulation algorithms for the posterior (22): an Importance sampler and an Metropolis-Hastings algorithm. We found that their M-H algorithm performs unsatisfactorily with low acceptance rate even for reasonable parameter specifications. As mentioned earlier, since the posteriors (14) and (22) as well as their conditional posteriors do not belong to any standard class of probability density functions, Gibbs sampling can not be used. In this section, we suggest an alternative simulation algorithm which combines Gibbs sampling [see Casella and George (1992), and Chib and Greenberg (1996)] and Metropolis-Hastings algorithm [see Metropolis et al. (1953), Hastings (1970), Smith and Roberts (1993), Tierney (1994), Chib and Greenberg (1995)]. Our algorithm is different from the "M-H within Gibbs" algorithm, and can find its usefulness in other applications as well.

To generate drawings from the target density $p(x)$, we use a candidate-generating density $r(x)$. An Independence sampler, which is a special case of the M-H sampler, in algorithmic form is as follows:

0. Choose starting values x^0
1. Draw x^i from $r(x)$
2. Accept x^i with probability

$$\alpha(x^{i-1}; x^i) = \min \left\{ \frac{p(x^i)r(x^{i-1})}{p(x^{i-1})r(x^i)}, 1 \right\}, \quad \begin{array}{l} \text{if } p(x^{i-1})r(x^i) > 0 \\ \text{if } p(x^{i-1})r(x^i) = 0 \end{array} \quad (24)$$

otherwise $x^i = x^{i-1}$;

3. $i = i + 1$. Go to 1.

It is generally not feasible to draw all elements of the vector x simultaneously. A block-at-a-time possibility was first discussed in Hastings (1970, sec. 2.4) and then in Chib and Greenberg (1995) along with an example. Chib

and Greenberg (1995) considered applying the M-H algorithm in turn to sub-blocks of the vector x , which presumes that the target density $p(x)$ may be manipulated to generate full conditional densities for each of the subblocks of x , conditioning on other elements of x . However the full conditionals are sometimes not readily available from the target density for empirical investigators. The posteriors (14) and (22) happen to fall in this category. In this latter case, problems come up at step 1 while trying to generate drawings from the joint marginal density $r(x)$. Note that these drawings, whether accepted or rejected at step 2, satisfy the necessary reversibility condition if step 1 is performed successfully.

To simplify the notation, we consider a vector x which contains two blocks, $x = (x_1; x_2)$. KVD used the fact that

$$r(x_1; x_2) = r(x_1)r(x_2|x_1) \quad (25)$$

and suggested to draw x_1^i from $r(x_1)$ and then draw x_2^i from $r(x_2|x_1^i)$. The pair $(x_1^i; x_2^i)$ is then taken as a drawing from $r(x)$. It turns out that this strategy gives very low acceptance rate at step 2 in simulation studies for various reasonable parameter values. Sometimes the move never take place and the posterior has all its mass at the parameter values of the first drawing. The reason for the failure is that information is not updated at subsequent drawings and the transition kernel of (25) is static.

If the full conditionals $r(x_1|x_2)$ and $r(x_2|x_1)$ are available, which is usually true for many standard densities, we propose to use them in a Gibbs sampler to make independent drawings from the invariant density $r(x)$ after the Markov chain has converged.

The combined algorithm is thus as follows, which we call "Gibbs within M-H":

0. Choose starting values $x^0 = (x_1^0; x_2^0)$:
1. Draw x_1^i from $r(x_1|x_2^{i-1})$, draw x_2^i from $r(x_2|x_1^i)$:
2. Accept $x^i = (x_1^i; x_2^i)$ with probability $\min(1, \frac{q(x^{i-1}|x^i)}{q(x^i|x^{i-1})})$ as defined in (24), otherwise $x^i = x^{i-1}$:
3. $i = i + 1$. Go to 1.

As explained, step 2 is the Gibbs step and step 3 is the M-H step in our combined algorithm. In the following subsections, we describe the steps for implementing the above procedure to generate drawings from the posteriors (14) and (22).⁷

4.1 Implementing the CP approach

Note that the posterior in the CP approach is proportional to the product of the prior, which is uniformly bounded, and the likelihood function, which can be sampled by a Gibbs sampler. Therefore we choose the candidate-generating density the way suggested by Chib and Greenberg (1995): we use the likelihood function, $L(\beta; \alpha; \lambda_1; \lambda_2; S; Y; Z)$, as the candidate generating density for the posterior (14). Using precision matrix S^{i-1} ; the simulation steps are as follows,

0. Choose starting values $(\beta^0; \alpha^0; \lambda_1^0; \lambda_2^0; S^{i-1:0})$
1. Draw $S^{i-1:i}$ from $p(S^{i-1}|j^{i-1}; \alpha^{i-1}; \lambda_1^{i-1}; \lambda_2^{i-1}; Y; Z)$
 Draw $(\beta^i; \alpha^i; \lambda_1^i; \lambda_2^i)$ from $p(\beta; \alpha; \lambda_1; \lambda_2 | S^{i-1:i}; Y; Z)$

⁷The algorithm has been illustrated with a simple labor supply model in Gao and Lahiri (2000).

Similar to the way we implemented the CP approach, it is more convenient to work with the precision matrix Σ^{-1} in the conditional densities. Applying the procedure outlined above, the steps involved in constructing the Markov chain for the posterior (22) are summarized as follows,

0. Choose starting values $(\Theta^0; \Sigma^{-1;0})$
1. Draw $\Sigma^{-1;i}$ from $p(\Sigma^{-1;i} | \Theta^{i-1}; Y; Z)$
Draw Θ^i from $p(\Theta^i | \Sigma^{-1;i}; Y; Z)$
2. Perform a singular value decomposition of $\Theta^i = U^i S^i V^{i0}$
3. Compute $\Sigma^{-1; \Sigma^i; \frac{1}{2}; \frac{1}{2}}$ according to (18)-(19)
4. Compute $w(\Sigma^{-1; \Sigma^i; \frac{1}{2}; \frac{1}{2}; \Theta^i; \Sigma^{-1;i})$ according to (29)-(30)
5. Draw $(\frac{1}{2}; \frac{1}{2})$ from $p(\frac{1}{2}; \frac{1}{2} | \Sigma^{-1;i}; \Theta^i; \Sigma^{-1; \Sigma^i; \frac{1}{2}; \frac{1}{2}}; Y; Z)_{j_s=0}$
6. Accept $(\Sigma^{-1; \frac{1}{2}; \frac{1}{2}; \frac{1}{2}; \Sigma^{-1;i})$ as a drawing from the posterior with

probability,

$$\min \frac{w(\Sigma^{-1; \Sigma^i; \frac{1}{2}; \frac{1}{2}; \Theta^i; \Sigma^{-1;i})}{w(\Sigma^{-1;i-1; \Sigma^i; \frac{1}{2}; \frac{1}{2}; \Theta^{i-1}; \Sigma^{-1:(i-1)})}; 1;$$

otherwise, $(\Sigma^{-1; \frac{1}{2}; \frac{1}{2}; \frac{1}{2}; \Sigma^{-1;i}) = (\Sigma^{-1;i-1; \Sigma^i; \frac{1}{2}; \frac{1}{2}; \Theta^{i-1}; \Sigma^{-1:(i-1)})$;

7. $i = i + 1$. Go to 1.

Note that the conditional densities used in the first step are as follows:

$$p(\Sigma^{-1;i} | \Theta^i; Y; Z) \propto j^{-i} j^{(T+k_2i-m_i-1)/2} \exp\left[-\frac{1}{2} \text{tr}(\Sigma^{-1;i} G)\right]; \quad (31)$$

which follows a Wishart distribution $W_m(T + k_2; G; \Sigma^{-1;i})$ with $(T + k_2)$ degrees of freedom, where $G = Y^0 Q_Z Y + (\Theta^i; \mathbf{b})^0 Z_2^0 M_{Z_1} Z_2 (\Theta^i; \mathbf{b})$, and $\mathbf{b} = (Z_2^0 M_{Z_1} Z_2)^{-1} Z_2^0 M_{Z_1} Y$. In addition,

$$p(\Theta^i | \Sigma^{-1;i}; Y; Z) \propto j^{-i} j^{k_2/2} \exp\left[-\frac{1}{2} \text{tr}[\Sigma^{-1;i} (\Theta^i; \mathbf{b})^0 Z_2^0 M_{Z_1} Z_2 (\Theta^i; \mathbf{b})]\right]; \quad (32)$$

which is a matrix-variate normal density.

The conditional density used in step 5 is

$$p(\mu_1; \sigma_1^2 | j^{-1}; \mu_2; \sigma_2^2; \mu_3; \sigma_3^2; Y; Z) \propto j^{-i} j^{k_1=2} \exp\left\{ -\frac{1}{2} \text{tr}[-i^{-1}(\alpha_i \quad \mathbf{b})^0 Z_1^0 Z_1(\alpha_i \quad \mathbf{b})]g\right\}; \quad (33)$$

evaluated at $\mu_s = 0$; where $\alpha = (\mu_1 \quad \sigma_1^2)$; $\mathbf{b} = (Z_1^0 Z_1^0)^{-1} Z_1^0 (Y \quad Z_2^0)$:

4.3 Convergence Diagnosis

One important implementation issue associated with MCMC methods is that of determining the number of iterations required. There are various informal or formal methods for the diagnosis of convergence, see Cowles and Carlin (1996) and Brooks and Roberts (1999) for recent comprehensive reviews and recommendations. Since the posterior densities in (14) and (22) resulting from CP and KVD do not have moments of any positive integer order, most of the methods proposed in the MCMC literature which require the existence of at least the first moment (posterior mean) are ruled out. We are left with a very few alternatives that can be used in our context.

First, the popular Raftery and Lewis (1992) method has been recognized as the best for estimating the convergence rate of the Markov chain if quantiles of the posterior density are of major interest, although the method does not provide any information as to the convergence rate of the chain as a whole. Because we are interested in the posterior modes and medians for τ associated with the Bayesian approaches, we may largely rely on Raftery and Lewis' method to determine the number of burn-ins, and the subsequent number of iterations required to attain specified accuracy (e.g., estimating the 0.50 quantile in any posterior within $\$0.05$ with probability 0.95). But

we do not adopt their suggested skip-interval. MacEachern and Berliner (1994) showed that estimation quality is always degraded by discarding samples. We once experimented with using the skip-intervals and found that the results are basically the same if a sufficient number of iterations are run. This seems to be inefficient and sometimes infeasible in terms of computation time.

For each specification in our Monte Carlo study with repeated experiments, we determined the number of burn-ins and subsequent number of iterations by running the publicly available Fortran code `gibbsit` on MCMC output of 10,000 iterations from three or more testing replications. For KVD and CP approaches, the numbers of burn-ins for both the GS step and the M-H algorithm were estimated. It was found that the number of burn-ins in the GS step is negligible for most cases. However, we discarded more iterations as the transient phase than the estimated number of burn-ins.⁹ The estimated number of subsequent iterations across testing replications was stable for the Gibbs sampler (in both Geweke approach and the GS step for KVD and CP approaches), but it varied a lot for the M-H procedures, which is also demonstrated by the variation in acceptance rates over repeated experiments. We used a generous value for the number of subsequent iterations when feasible.

Second, for MCMC output from each testing replication, we also applied other convergence diagnostic methods, including percentiles derived from every quarter of the long chain, Yu and Mykland (1994)'s CUSUM plot, and

⁹In practice, there is often a concern about possible underestimation of true length of the burn-in period using the Raftery and Lewis method if the quantile of interest is not properly pre-prescribed, see Brooks and Roberts (1999).

Brooks' (1996) D sequence statistic. While the CUSUM partial sums actually involve averaging over sampling drawings, the computation of Brooks' statistic is justified on the basis that it is designed to measure the frequency of back-forth movement in the MCMC algorithm. However, these diagnostics may sometimes provide contradictory outcomes so that one has to be extra careful in interpreting them before making a judgment on convergence.

5 Simulation results and discussions

In this section, we present results of Monte Carlo experiments and discuss some of the findings. As mentioned before, for the purpose of comparison, we also computed a number of single K-class estimators including OLS, 2SLS, LIML and Fuller's modified LIML. In summary, the set of K-class estimator for the structural coefficients in model (1) and (2) is given by:

$$\begin{bmatrix} \mu \\ \mathbf{b} \end{bmatrix} = \begin{bmatrix} Y_2^0 Y_2^0 - K_1 \mathbf{b}_2^0 \mathbf{b}_2^0 & Y_2^0 Z_1^0 \\ Z_1^0 Y_2^0 & Z_1^0 Z_1^0 \end{bmatrix}^{-1} \begin{bmatrix} (Y_2^0 - K_2 \mathbf{b}_2^0) y_1^0 \\ Z_1^0 y_1^0 \end{bmatrix}; \quad (7)$$

where $\mathbf{b}_2^0 = Q_Z Y_2$:

The following LISEM estimators have been considered:

(1) Ordinary least squares (OLS)

$$K_1 = K_2 = 0:$$

(2) Two stage least squares (2SLS)

$$K_1 = K_2 = 1:$$

(3) Zellner's (1978) Bayesian minimum expected loss estimator (MELO)

within M-H" algorithm for the CP approach since the likelihood function in (3) is used as the candidate-generating density to explore the CP posterior.

(10) Posterior mode and median from CP approach using \Gibbs within M-H" algorithm

(11) Posterior mode and median from KVD approach using \Gibbs within M-H" algorithm

For the recent Bayesian approaches and LIML-GS, we report both (posterior) mode and median to show possible asymmetry in the marginal densities of β . Any preference for one over the other will depend on the researcher's loss function. We obtain 16 estimates for each generated data set. The data are generated from the model,

$$\begin{aligned} y_1 &= Y_2 \beta + u; \\ Y_2 &= Z_2 \beta + V_2; \end{aligned} \tag{34}$$

where y_1, Y_2 are $T \times 1$ such that $m = 2$; and $Z_2 : T \times k_2$. We further specify $\beta = 1$ and

$$S = \begin{pmatrix} \mu & \frac{1}{2} \\ \frac{1}{2} & 1 \end{pmatrix} \tag{35}$$

for $j = 0.20, 0.60, \text{ and } 0.95$.¹² Z_2 is simulated from a $N(0; I_{k_2} - I_T)$ distribution and $(u; V_2)$ from a $N(0; S - I_T)$ distribution. A constant term is added in each equation, i.e., Z_1 is a $T \times 1$ vector of 1's.

The simulation results are reported in Table 1 through Table 13. Tables 1 to 12 are for cases with $\beta > 0$; each table reporting results for one speci-

¹²We do not report cases with $\beta = 0.99$ or 1. As pointed out by Maddala and Jeong (1992), when the instruments are weak and β is very close to one, the exact finite sample distribution of IV estimator is bimodal. Our experiments show that the marginal posterior density of β from the recent Bayesian approaches exhibits a similar pattern.

ification. Tables 13 summarizes the results for cases with $\lambda < 0$ for BMOM and KVD for whom negative λ made a surprising difference. As mentioned before, we focus on the estimates of the structural parameter β . Specifically, we analyze the sensitivity of the various estimates of β with respect to the strength of the instrumental variables Z , the degree of overidentification ($k_2 \geq m + 1$), the degree of endogeneity (λ), and the sample size (T). Also, we will examine whether the performance of an estimator is symmetric with respect to the sign of parameter λ , an issue generally overlooked in the literature.¹³

Note that the strength of the instrumental variables for the included endogenous variable Y_2 is measured in terms of the adjusted R^2 by regressing Y_2 on $Z = (Z_1; Z_2)$. In the data generating process, we controlled \bar{R}^2 to be within $\pm 2.5\%$ of the specified value to reduce unnecessary variation. We did not experiment with extremely small \bar{R}^2 (say, 0.01 or less). In these cases the mean values of all estimators approached the point of concentration $\beta_{12} = \beta_{22}$; which is equal to $(\beta + \lambda)$ for our data generating process (DGP).

For each specification, the number of replications is 400. The number of burn-ins (nburn_GS and nburn_MH), and subsequent number of iterations (n) determined at the convergence diagnosis step are reported in the footnotes

¹³Denote $\beta = \begin{pmatrix} \beta_{11} & \beta_{12} \\ \beta_{21} & \beta_{22} \end{pmatrix}$. Using $S = C^0 - C$, we have $\beta_{11} = w_{11} + \lambda w_{12} + \beta_{22} w_{22}$; $\beta_{12} = w_{12} + \lambda w_{22}$; and $\beta_{22} = w_{22}$. Letting $\lambda = \beta_{12} = \frac{\beta_{11} - \beta_{22}}{\beta_{11} - \beta_{22}}$; the second relationship may be rewritten as:

$$\beta_{11} - \beta_{22} = \lambda \frac{\beta_{11} - \beta_{22}}{\beta_{11} - \beta_{22}}$$

If S is normalized as in (35) with $\beta_{11} = \beta_{22} = 1$, then $w_{12} = \beta + \lambda$. Therefore, in our context, given $\beta = 1$, the sign and magnitude of λ (or w_{12}) has a special significance.

to each table. The average acceptance rate and its standard deviation (in parentheses) across replications for each M-H routine are reported as well. To evaluate alternative estimators, we computed mean, standard deviation (Std), root of mean squared errors (RMSE), and mean absolute deviation (MAD) over repeated experiments for all the estimators considered.¹⁴ Since LIML, posterior densities for CP and KVD, as well as 2SLS in the just-identified case do not have finite moments of positive order in finite samples, one should interpret the computed mean, standard deviation and RMSE across replications for these estimators with caution. In this sense, the MAD across replications is a preferred measure to consider.

We will first look at cases reported in Tables 1 to 12 with $\frac{1}{2} > 0$. In Table 1, we consider a case ($T = 50$, $\frac{1}{2} = 0.60$, $k_2 = 4$) with moderately strong instruments ($\bar{R}^2 = 0.40$). It is found that with reasonably strong instruments all estimators designed for simultaneous equations perform reasonably well. As expected, OLS is seriously biased. BMOM has a slight edge over others in terms of RMSE and MAD. For all recent Bayesian approaches and LIML-GS the medians perform a little better than modes, and CP over KVD, in terms of bias, RMSE and MAD. Notice that the classical LIML estimates are different from LIML-GS (mode or median). As noted by Dr̄z̄e (1976), from a Bayesian viewpoint, LIML produces an estimate of β conditionally on the over-identifying restrictions, the modal values of all the remaining parameters, and a uniform prior. In other words, the concentrated likelihood function of β after concentrating out (i.e., maximizing

¹⁴Medians were also calculated. Since they were very close to the corresponding means in all our experiments, we did not report them in this paper.

with respect to) other reduced-form and nuisance parameters is a conditional density. However, LIML-GS is a marginal density with all other parameters being integrated out. Due to possible asymmetry in the distribution of the nuisance parameters, the modal/median values of LIML-GS may not coincide with classical LIML estimates. In all our experiments, we find that the median-unbiasedness property of (conditional) LIML does not carry over to the marginal LIML (i.e., LIML-GS); however, the former generally has a much larger standard deviation than the latter. In a way, LIML-GS brings the classical LIML estimator close to its Bayesian counterpart for the purpose comparison.

It is interesting to note that across all our tables, the difference between LIML-GS and CP can only be attributed to the importance of Jeffrey's prior. Compared to LIML-GS, typically CP has a smaller bias, but slightly larger standard deviation, even though the differences are very small. In some cases, however, the use of Jeffrey's prior reduces the bias in CP quite substantially. For example, in Table 4 with $T = 50$ and a high degree of overidentification, the bias is reduced from 0.36 to 0.25.

A simple case when the structural model is just identified ($k_2 = 1$) is reported in Table 2. For this case it is well known that classical LIML coincides with 2SLS. KVD approach does not accommodate the case of just-

identification since (15) requires $k_2 > (m_j - 1)$.¹⁵ In this case, we find that CP-Mode produces results closer to LIML-GS-Mode than to LIML. CP (1998) showed that for a two-equation just-identified SEM in orthonormal canonical form, the posterior density of β with Jeffreys prior has precisely the same functional form as the density of the finite sample distribution of the corresponding LIML estimator as obtained by Mariano and McDonald (1979). Our simulation results show that the assumption of orthonormal canonical form is crucial for their exact correspondence, which cannot be extended to a general SEM.¹⁶ In general, the Bayesian marginal density is not the same as the classical conditional density. Interestingly, JIVE is considerably more biased and has larger standard deviation than 2SLS. Also, CP-Median and LIML-GS-Median perform significantly worse than their modes. This is because in an exactly identified model with weak instruments the probability of local non-identification is substantial, and the resulting non-standard marginal density exhibits a very high variance. The same result holds true for Geweke-Median, but to a lesser extent. Thus, for exactly identified SEMs

¹⁵When $k_2 = (m_j - 1)$, a diffuse prior in (20) for the linear model implies that the prior for the parameters of the LISEM (4) is

$$p(\beta; \Sigma_1; \Sigma_2; -) \propto |\Sigma_2|^{-j} |\Sigma_1|^{-(k+m+1)-2j};$$

and the prior for the parameters of the LISEM (1) and (2) is

$$p(\beta; \Sigma_1; \Sigma_2; S) \propto |S| |\Sigma_2|^{-(k+m+1)-2j};$$

which is identical to the Jeffreys prior; see also expressions (22) and (42) in CP.

¹⁶Note that the relationship between the standardized parameter vector and the original parameter vector involves the nuisance parameters, cf. Phillips (1983). However, when a SEM is in orthonormal canonical form (i.e., the exogenous regressors are orthonormal and the disturbance covariance matrix Σ is an identity matrix), both the density of random parameter β from the CP approach and the probability density of the classical LIML estimator for β are conditional on these information.

with very weak instruments, mode of the marginal density is a more dependable measure of β . We should point out that in all other cases in this study the medians generally turned out to be more preferable than the modes in terms of bias, RMSE, and MAD (see Tables 11 and 12, for instance).

Results reported in Tables 3 through 12 consider cases with general over-identification and weak instruments. As noted in the literature, OLS and 2SLS are median-biased in the direction of the correlation coefficient ρ , and the bias in 2SLS grows with the degree of over-identification, and decreases as sample size increases. Results in Tables 3 through 10 confirm these results. Since MELO is a single K-class estimator with $0 < K < 1$, it is always between OLS and 2SLS estimates. The bias in MELO shows the same pattern as that of 2SLS. With moderate simultaneity, the median-bias in 2SLS can be as large as about 40% of the true value (see Table 8). We note that MELO, LIML-GS-Mode, and KVD-Mode or KVD-Median are also median-biased in the direction of ρ . But the bias in JIVE is consistently in the opposite direction of ρ . Classical LIML is remarkably median-unbiased when the instrumental variables are not very weak, which is well documented in the literature. We find that LIML is median-biased in the direction of ρ when the instruments are very weak (Table 8), which is consistent with the finding in Staiger and Stock (1998) using local-to-zero asymptotic theory. Even in this situation, the bias of LIML is much smaller than that of any other estimator, except BMOM.

The MAD of OLS is very close to its bias (i.e., relatively small Std) across all cases and it implies that OLS method is robust in the sense that it does not suffer from heavy tails or outlying estimates, see Zellner (1998). In this

sense, MELO and BMOM are all robust with relatively small standard deviations across replications. However, OLS exhibits large bias in the presence of simultaneity and is not so appealing. It is known that for a degree of overidentification strictly less than 7, 2SLS would have a smaller asymptotic mean squared error (AMSE) than LIML, cf. Mariano (1982). In cases with weak instruments the situation gets more complicated in finite samples. In our experiments, LIML has larger RMSE and MAD than 2SLS except in Tables 11 and 12 where λ was 0.95: Note that the degree of over-identification is 8 in Tables 4, 6, 8 and 10.

Among classical estimators, JIVE turns out to be least appealing. Monte Carlo simulations in Angrist, Imbens and Krueger (1999) showed that JIVE has slight median bias in the opposite direction of λ (but less than 2SLS) and have heavier tails than LIML. Our Table 6 is comparable to panel 2 of their Table I, and the results are also similar. Our other experiments show that JIVE may also have large absolute bias (larger than LIML) in the case with weak instruments, sometimes even greater than 2SLS (see Table 2). Generally, JIVE has slightly less bias than 2SLS, but this gain is overshadowed by enlarged standard deviation such that in finite samples it has no advantage over 2SLS in terms of MAD and RMSE. We also find that JIVE has greater RMSE and MAD than LIML. Blomquist and Dahlberg (1999) experimented with much larger sample sizes than ours. Comparing our Table 4 with Table 6 and with an unreported simulation with a sample size of 500, we found that the relative gain in JIVE is more than other estimators as sample size increases, even though its relative low standing remains valid.

Fuller's modified LIML estimators are included because Fuller1 is de-

signed to minimize the median-bias, and Fuller4 to minimize the mean-squared error. It seems that this conclusion is also problematic in the presence of weak instruments. Between the two, Fuller1 has smaller median-bias, and Fuller4 has smaller standard deviations across replications. However, in terms of RMSE or MAD, Fuller4 shows no advantage over Fuller1 in most of the cases.

Because all the estimators except OLS are consistent and their asymptotic distributions are also the same, results in Tables 3 through 6 confirm that their bias and dispersion decrease as sample size increases. But if the instruments are very weak (see Tables 7 and 8), their bias and dispersion may remain significant, a point emphasized forcefully by Zellner (1998). However, when the endogeneity is not strong (see Tables 9 and 10), their bias and dispersion may not be a big concern for some of the estimators.

Across all cases, we find that the bias in BMOM is small if λ is not too small and the structural equation (1) is overidentified. As sample size increases or degree of over-identification rises, the observed bias in BMOM decreases. The most striking feature of BMOM is that it exhibits the smallest MAD and Std when λ is not too small. MELO shows slightly smaller MAD and Std than BMOM if λ is small (see Tables 9-10). In cases with very weak instruments and high degree of over-identification, the MAD of BMOM is only one-fourth of that of other estimators (see Table 8). These are in accordance with Tsurumi (1990)'s finding that in many cases, ZEM has the least relative mean absolute deviation. Meanwhile, if λ is very small and the structural equation is overidentified, the bias in BMOM can be large; 2SLS, LIML-GS, Geweke, and CP perform remarkably well in these situations.

Next, we examine in more detail the recent Bayesian approaches. Overall, the median bias resulting from these approaches exhibits the same pattern as the bias of 2SLS, it increases with the degree of over-identification, and decreases as sample size rises. The Geweke (1996) approach used a shrinkage prior but its performance is comparable with LIML-GS and CP. The median-bias from PMOD-Geweke is the same or slightly less than that of LIML-GS-Mode, and the bias from Geweke-Median is always slightly less than that of LIML-GS-Median. Similar relationships are observed for MADs. These reflect the impact of the (informative) shrinkage prior on the posterior density.

For each specification, the acceptance rate in the M-H algorithm using CP approach is stable while that using KVD approach shows huge variation across replications. The acceptance rate for CP is generally above 40% except when sample size is small and the degree of overidentification is high. This shows that the posterior of CP is largely dominated by the likelihood function (3) and the Jeffreys prior generally carries little information. Second, in terms of the computed standard deviations (Stds) of the estimates across replications, CP-Mode has larger dispersion than LIML-GS-Mode, and CP-Median has larger dispersion than LIML-GS-Median. These also shed light on the notion that Jeffreys prior is less informative than a uniform prior. However, between the Jeffreys prior (13) used by CP and the implied prior (21) resulting from diffuse/Jeffreys prior on a linear model used by KVD, it is not clear which one is less informative.

As for the KVD (1998) approach, we observe that it performs as well as any other estimator if the instruments are not weak (see Table 1). But

when the instruments are weak, and $\frac{1}{2}$ is positive, KVD shows more bias and higher MAD than those from CP. In Tables 4 with $T = 50$ and high degree of overidentification, KVD performs as bad as OLS.

Next we consider cases with negative $\frac{1}{2}$; and the results are summarized in Table 13. We replicate each case in Tables 1 - 12 with the same specification except $\frac{1}{2}$ being negative. Since the performance of all estimators except BMOM and KVD were basically the same with respect to the sign of $\frac{1}{2}$; we only report results on these two in Table 13. We find that when $\frac{1}{2}$ changes sign, the bias of BMOM does not change sign and even increases in magnitude. Also note that the computed Stds for BMOM when $\frac{1}{2} < 0$ are close to the respective ones when $\frac{1}{2} > 0$. Therefore, for cases with $\frac{1}{2} < 0$, BMOM has large RMSEs/MADs and loses its attraction. Note that BMOM is the same as the double K-class estimator (DKC) with K values fixed. This asymmetry in the performance in DKC is not well recognized in the literature. However, the observed asymmetry in its bias with respect to $\frac{1}{2}$ in our experiments is readily explained by examining an expression for the mean of double K-class estimator (DKC) in Dwivedi and Srivastava (1984, Theorem 1). We can express \mathbf{b}_{DKC} as:

$$\mathbf{b}_{DKC} = \mathbf{b}_{K_1} + \frac{Y_2^0 Y_2^1}{Z_1^0 Y_2^1} (K_1 - K_2) \frac{Y_2^0 Z_1^1}{Z_1^0 Z_1^1} \mathbf{b}_2; \quad (36)$$

where \mathbf{b}_{K_1} is a single K-class estimator with characterizing scalar K_1 : When $Z_1^0 Z_2^1 = 0$, which is satisfied in our experimental specifications, a double K-class estimator of β may be written as

$$\mathbf{b}_{DKC} = \mathbf{b}_{K_1} + (K_1 - K_2) \frac{Y_2^0 Q_Z y_1}{Y_2^0 \Phi Y_2}$$

where $\Phi = (1 - K_1)Q_{z_1} + K_1P_{z_2}$: Observe that for $0 < K_1 < 1$; b_{K_1} is biased in the direction of $\frac{1}{2}$, as noted in Mariano (1982). Note also that $Y_2^0 \Phi Y_2 > 0$; and $Y_2^0 Q_{z_1} y_1$ provides an estimate of w_{12} . Although Dwivedi and Srivastava (1984) explored the dominance of double K-class over K-class using the exact MSE criterion, their guidelines for the selection of K_2 for a given K_1 are not entirely valid, because the conditions were derived from a small Monte Carlo simulation with cases with positive w_{12} and negative $\frac{1}{2}$ only. Since $K_1 < K_2$ for BMOM, when $\frac{1}{2}$ and w_{12} have the opposite sign, the second term in b_{DKC} will be of the same sign as the bias of b_{K_1} , therefore b_{DKC} (hence BMOM) will exhibit large bias. Otherwise, when $\frac{1}{2}w_{12} > 0$; the bias is mitigated. Based on our simulation results, we found that the sign of $\frac{1}{2}$ has no effect on the standard deviation of BMOM. This finding shows that the greater RMSE of BMOM when $\frac{1}{2}w_{12} < 0$ is due to the aggravated bias. For the specification corresponding to table 4 in Table 13 (i.e., $T = 50$, $\frac{1}{2} = -0.60$, $K_2 = 4$, $\bar{R}^2 = 0.10$), we find that for given $K_1 = 0.947$; RMSE is minimized if K_2 is chosen to be 0.829, which is much less than K_1 , and also less than $K_2 = 0.987$ used in BMOM.

In tables 3 - 12 we found that KVD with $\frac{1}{2} > 0$ performs very poorly, often with substantial bias and relatively high RMSE and MAD. CP uniformly dominates KVD in these cases. However, with $\frac{1}{2} < 0$ the picture turns around remarkably well in favor of KVD. As we see in Table 13, across all cases the bias tends to be negative and relatively small. With other parameter values being the same, KVD with $\frac{1}{2} < 0$ has significantly less RMSE and MAD than cases when $\frac{1}{2} > 0$, and performs unequivocally the best among all estimators when endogeneity is strong. However, since this

observed asymmetry is essentially a finite sample problem with KVD, the improved performance when $\frac{1}{2} < 0$ becomes less significant when the sample size increases from 50 to 100. With $\frac{1}{2} < 0$ the overall performance of KVD is very comparable to that of CP, if not slightly better in some cases.

After experimenting with widely different negative and positive values of $\frac{1}{2}$ and $\frac{1}{2}$, we found out that the performance of KVD is dependent on the sign of $\frac{1}{2}$, rather than on the sign of $\frac{1}{2}$: When $\frac{1}{2} > 0$; it performs very unsatisfactorily as documented in Tables 3-12. Kleibergen and Zivot (1998) have recently derived exact analytical expressions for the conditional densities of $\frac{1}{2}$ given $\frac{1}{2}$ for both the KVD and CP posteriors. They show that the difference between the two is in the Jacobian relating the unrestricted linear multivariate model to the restricted reduced form model. We expect that this additional term may account for the asymmetry in KVD with respect to $\frac{1}{2}$. In our experiments, we found that in finite samples, when $\frac{1}{2} > 0$; the reduced rank restriction using singular value decomposition shifts the marginal posterior for KVD away from the marginal posterior of the linear multivariate model. However, when the sample size gets large, the problem seems to go away.

6 Conclusions

This paper examines the relative merits of some recent developments in the Bayesian and classical analysis of limited information simultaneous equations models in situations where the instruments are very weak. Since the posterior densities and their conditionals in the Bayesian approaches devel-

oped by Chao and Phillips (1998) and Kleibergen and van Dijk (1998) are non-standard, we proposed and implemented a "Gibbs within Metropolis-Hastings" algorithm, which only requires the availability of the conditional densities from the candidate generating density. These conditional densities are used in a Gibbs sampler (GS) to simulate the candidate generating density, whose drawings, after convergence, are then weighted to generate drawings from the target density in a Metropolis-Hastings (M-H) algorithm. We rely on Raftery and Lewis (1992) method to determine the number of burn-ins, and the subsequent number of required iterations in order to ensure convergence. Through a MCMC simulation study, our results provide useful guidelines for empirical practitioners.

The first comforting result is that with reasonably strong instruments (marginal \bar{R}^2 in excess of 0.40) all estimators perform equally well in finite samples. In cases with very weak instruments (marginal \bar{R}^2 less than 0.10), there is no single estimator that is superior to others in all cases. When endogeneity is weak ($\frac{1}{2}$ less than 0.20), Zellner's MELO does the best. When the endogeneity is relatively strong ($\frac{1}{2}$ in excess of 0.60) and $\frac{1}{2}w_{12} > 0$, BMOM outperforms all other estimators by wide margins. When the endogeneity is strong but $\frac{1}{2} < 0$, KVD approach seems to get very appealing; but, otherwise, its performance is surprisingly very poor. With $\frac{1}{2} > 0$, as the sample size gets larger, the performance of KVD improves rapidly. Fortunately, the Geweke and CP approaches exhibit no such asymmetry and their performances based on bias, RMSE, and MAD are very similar. Based on the medians of marginal posteriors, their performance ranking is consistently a distant second. The record of JIVE is quite disappointing across all our

experiments, and is not recommended in practice. Even though JIVE is slightly less biased than 2SLS in most cases, its standard deviation is considerably higher, particularly in small samples. The most remarkable result in this study is that poor instruments can affect the performance of different estimators differently, depending on the signs and magnitudes of certain key parameters of the model. Given the finding that even in finite samples with very weak instruments BMOM and KVD perform so remarkably well on certain parts of the parameter space, more research is needed to understand the reasons for the asymmetry and find ways to fix the problem.

References

- Angrist, J., G.W. Imbens and A. Krueger (1999). Jackknife instrumental variables estimation. *Journal of Applied Econometrics* 14, 57-67.
- Billingsley, P. (1986). *Probability and Measure* (Wiley, New York).
- Blomquist, S. and M. Dahlberg (1999). Small sample properties of LIML and jackknife IV estimators: experiments with weak instruments. *Journal of Applied Econometrics* 14, 69-88.
- Bound, J., D.A. Jaeger, and R.M. Baker (1995). Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak. *Journal of the American Statistical Association* 90, 443-450.
- Brooks, S. and G.O. Roberts (1999). Assessing convergence of Markov chain Monte Carlo algorithms. *Statistics and Computing* (forthcoming).
- Brooks, S. (1996). Quantitative convergence diagnosis for MCMC via CUSUMS. Technical report, University of Bristol.
- Buse, A. (1992). The bias of instrumental variable estimators. *Econometrica* 60, 173-180.
- Casella, G. and E. George (1992). Explaining the Gibbs sampler. *The American Statistician* 46, 167-174.
- Chao, J.C. and P.C.B. Phillips (1998). Posterior distributions in limited information analysis of the simultaneous equations model using the Jeffreys prior. *Journal of Econometrics* 87, 49-86.
- Chib, S. and E. Greenberg (1995). Understanding the Metropolis-Hastings algorithm. *The American Statistician* 49, 327-335.
- Chib, S. and E. Greenberg (1996). Markov chain Monte Carlo simulation methods in econometrics. *Econometric Theory* 12, 409-431.
- Cowles, M.K. and B.P. Carlin (1996). Markov chain Monte Carlo convergence diagnosis: a comparative review. *Journal of the American Statistical Association* 91, 883-904.

- Dręze, J.H. (1976). Bayesian limited information analysis of the simultaneous equation model. *Econometrica* 44, 1045-1075.
- Dręze, J.H. and J.A. Morales (1976). Bayesian full information analysis of simultaneous equations. *Journal of American Statistical Association* 71, 329-354.
- Dręze, J.H. and J.-F. Richard (1983). Bayesian analysis of simultaneous equation systems. In: Z. Griliches, and M. Intriligator, eds., *Handbook of Econometrics* (North Holland, Amsterdam).
- Dwivedi, T.D, and V.K. Srivastava (1984). Exact finite sample properties of double k-class estimators in simultaneous equations. *Journal of Econometrics* 25, 263-283.
- Fuller W.A. (1977). Some properties of a modification of the limited information estimator. *Econometrica* 45, 939-953.
- Gao, C. and K. Lahiri (1999). Further consequences of viewing LIML as an iterated Aitken estimator. *Journal of Econometrics* (forthcoming).
- Gao, C. and K. Lahiri (2000). MCMC algorithms for two recent Bayesian limited information estimators. *Economics Letters* 66, 121-126.
- Geweke, J. (1996). Bayesian reduced rank regression in econometrics. *Journal of Econometrics* 75, 121-146.
- Hastings, W.K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* 57, 97-109.
- Kleibergen, F. (1997). Equality restricted random variables: densities and sampling algorithms. *Econometric Institute Report 9662/A*. Erasmus University Rotterdam.
- Kleibergen, F. (1998). Conditional densities in econometrics (1998). Discussion paper.
- Kleibergen, F. and H.K. van Dijk (1998). Bayesian simultaneous equation analysis using reduced rank structures. *Econometric Theory* 14, 701-743.

- Kleibergen, F. and E. Zivot (1998). Bayesian and classical approaches to instrumental variable regression. Econometric Institute Report 9835/A. Erasmus University Rotterdam.
- MacEachern, S.N. and L.M. Berliner (1994). Subsampling the Gibbs sampler. *The American Statistician* 48, 188-190.
- Maddala, G.S. (1976). Weak priors and sharp posteriors in simultaneous equation models. *Econometrica* 44, 345-351.
- Maddala, G.S. and J. Jeong (1992). On the exact small sample distribution of the instrumental variable estimator. *Econometrica* 60, 181-183.
- Mariano, R.S. and J.B. McDonald (1979). A note on the distribution functions of LIML and 2SLS structural coefficient in exactly identified case. *Journal of the American Statistical Association* 74, 847-848.
- Metropolis, N., A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, and E. Teller (1953). Equations of state calculations by fast computing machines. *Journal of Chemical Physics* 21, 1087-1092.
- Pagan, A. (1979). Some consequences of viewing LIML as an iterated Aitken estimator, *Economics Letters* 3, 269-372.
- Percy, D.F. (1992). Prediction for seemingly unrelated regressions. *Journal of the Royal Statistical Society B* 54, 243-252.
- Poirier, D. (1995). *Intermediate Statistics and Econometrics* (M.I.T. Press, Cambridge, Mass.).
- Poirier, D. (1996). Prior beliefs about τ . In: J.M. Bernardo, J.O. Berger, A.P. Dawid, and A.F.M. Smith, eds., *Bayesian Statistics 5* (Clarendon Press, Oxford).
- Raftery, A.E. and S.M. Lewis (1992). How many iterations in the Gibbs sampler? In: J.M. Bernardo, A.F.M. Smith, A.P. Dawid and J.O. Berger, eds., *Bayesian Statistics 4* (Oxford University Press).
- Staiger, D. and J.H. Stock (1997). Instrumental variables regression with weak instruments. *Econometrica* 65, 557-586.

- Smith, A.F.M. and G.O. Roberts (1993). Bayesian computation via the Gibbs sampler and related Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society B* 55, 3-23.
- Tsurumi, H. (1990). Comparing Bayesian and non-Bayesian limited information estimators. In: Geisser, S., Hodges, J.S., Press, S.J., Zellner, A., eds., *Bayesian and Likelihood Methods in Statistics and Econometrics* (North-Holland, Amsterdam).
- Tierney, L. (1994). Markov chains for exploring posterior distributions. *Annals of Statistics* 22, 1701-1767.
- Yu, B. and P. Mykland (1994). Looking at Markov samplers through Cusum path plots: a simple diagnostic idea. Technical Report 413, University of California at Berkeley, Dept. of Statistics.
- Zellner, A. (1971). *An Introduction to Bayesian Inference in Econometrics* (Wiley, New York).
- Zellner, A. (1978). Estimation of functions of population means and regression coefficients: a minimum expected loss (MELO) approach. *Journal of Econometrics* 8, 127-158.
- Zellner, A. (1986). Further results on Bayesian minimum expected loss (MELO) estimates and posterior distributions for structural coefficients. In: Slottje, D., eds., *Advances in Econometrics*, Vol. 5, pp. 171-182.
- Zellner, A. (1994). Bayesian and Non-Bayesian estimation using balanced loss functions. In: Gupta, S.S., Berger, J.O., eds., *Statistical Decision Theory and Related Topics*, Vol. V. Springer, New York, Chapter 28, pp. 377-390.
- Zellner, A. (1998). The finite sample properties of simultaneous equations' estimates and estimators: Bayesian and non-Bayesian approaches. *Journal of Econometrics* 83, 185-212.
- Zellner, A., L. Bauwens, and H.K. van Dijk (1988). Bayesian specification analysis and estimation of simultaneous equation models using Monte Carlo methods. *Journal of Econometrics* 38, 39-72.