



# Bayesian Instrumental Variables: Priors and Likelihoods

Hedibert F. Lopes & Nicholas G. Polson

To cite this article: Hedibert F. Lopes & Nicholas G. Polson (2014) Bayesian Instrumental Variables: Priors and Likelihoods, *Econometric Reviews*, 33:1-4, 100-121, DOI: [10.1080/07474938.2013.807146](https://doi.org/10.1080/07474938.2013.807146)

To link to this article: <http://dx.doi.org/10.1080/07474938.2013.807146>



Accepted author version posted online: 24 May 2013.



[Submit your article to this journal](#)



Article views: 213



[View related articles](#)



[View Crossmark data](#)



Citing articles: 2 [View citing articles](#)

## BAYESIAN INSTRUMENTAL VARIABLES: PRIORS AND LIKELIHOODS

Hedibert F. Lopes and Nicholas G. Polson

*Booth School of Business, University of Chicago, Chicago, Illinois, USA*

□ *Instrumental variable (IV) regression provides a number of statistical challenges due to the shape of the likelihood. We review the main Bayesian literature on instrumental variables and highlight these pathologies. We discuss Jeffreys priors, the connection to the errors-in-the-variables problems and more general error distributions. We propose, as an alternative to the inverted Wishart prior, a new Cholesky-based prior for the covariance matrix of the errors in IV regressions. We argue that this prior is more flexible and more robust than the inverted Wishart prior since it is not based on only one tightness parameter and therefore can be more informative about certain components of the covariance matrix and less informative about others. We show how prior-posterior inference can be formulated in a Gibbs sampler and compare its performance in the weak instruments case for synthetic as well as two illustrations based on well-known real data.*

**Keywords** Angrist–Krueger data; Bayesian learning; Cholesky decomposition; Demand for cigarettes; Errors-in-variables; Fat-tails; Inverted Wishart; IV regression.

**JEL Classification** C01; C11; C15; C26.

### 1. INTRODUCTION

Simultaneous equation models (SEMs) (Zellner, 1971; Chetty, 1966) and instrumental variables (IV) are fundamental tools in statistics and econometrics. Additional Bayesian work on SEM are Drèze (1976), Drèze and Morales (1976), Drèze and Richard (1983), and Kleibergen and van Dijk (1998), amongst others. IV regression has been tackled with a plethora of methods including: Bayesian approaches (Kleibergen and van Dijk, 2007; Lindley and El Sayyad, 1968; Zellner, 1971), Bayes-Stein shrinkage (Zellner and Vandaele, 1975), decision-theoretic methods (Chamberlain, 2007), method of moments (Zellner et al., 2007), semiparametric Dirichlet mixtures (Conley et al., 2008; Florens and

Address correspondence to Hedibert F. Lopes, Booth School of Business, University of Chicago, 5807 S. Woodlawn Ave, Chicago, IL 60637, USA; E-mail: hlopes@chicagobooth.edu

Simoni, 2010) and Monte Carlo simulation (Zellner et al., 1988), to name but a few. Comparisons between Bayesian and classical approaches have been discussed in Lindley and El Sayyad (1968) and Kleibergen and Zivot (2003). This problem is intertwined with that of “errors-in-the-variables” models (Minka, 1999; Zellner, 1971), co-integration (Kleibergen and Paap, 2002; Koop et al., 2010; Strachan, 2003; Villani, 2005), and reduced rank regression (Geweke, 1996). For a discussion of Bayesian approaches to IV and many examples in Economics and Marketing, see Lancaster (2004) and Rossi et al. (2005).

In this article, we will revisit and discuss key formulation, identification and estimation aspects when performing Bayesian inference in the instrumental variable regression model based on several of the above alternative representations. As opposed to simple linear regression models, Sims (2007), for instance, highlights the many issues with a priori assumptions in modeling the IV system as inferences can be very sensitivity to these specifications in the reduced form model. In addition, a common feature of many of these IV problem is that you can obtain unexpected “sharp” (and possibly ill-behaved) posterior distributions from “weak” prior distributions, mainly due to the heavy tails of the likelihood or the nonlinearity of the parameters of interest or both (Maddala, 1976; Hoogerheide et al., 2007; Hoogerheide and van Dijk, 2008a,b; Hoogerheide et al., 2008; Zellner, 1971).

We propose a new Cholesky-based prior for the covariance matrix of the errors in IV regressions, as an alternative to the inverted Wishart prior. Rossi et al. (2005) highlight the importance of priors on the error covariance-matrix  $\Sigma$  as one goes from the structural model to the reduced form (see also Lancaster, 2004). Rossi et al. (2005, p. 29) point out to several of the drawbacks of the Wishart distribution: “The most important is that the Wishart has only one tightness parameter. This means that we cannot be very informative on some elements of the covariance matrix and less informative on others.” We argue that our Cholesky-based prior is more flexible and avoids such drawbacks.

The remainder of the paper is organized as follows. The basic IV regression setup and its Bayesian solution are outlined in Section 2. Important departures, such as more general, noninformative prior specifications and more general error distributions, are presented in Section 3, where we also describe the prior sensitivity issues raised in Sims (2007). Section 4 makes the connection between IV regressions and the likelihood-based approach of Zellner and Minka in the “errors-in-the-variables” models (EVM). Our Cholesky-based prior for the simultaneous error covariance is introduced and discussed in Section 5. The section includes three illustration: one based on simulated data and two based on well-known IV studies. Section 6 concludes.

## 2. IV REGRESSION

The basic setup is as follows and is drawn from Rossi et al. (2005). Let  $y_i$  be the response variable and  $x_i$  the (endogenous) regressor obeying, for  $i = 1, \dots, n$ , the system of equations

$$x_i = z_i' \delta + \varepsilon_{1i} \quad (1)$$

$$y_i = \gamma + \beta x_i + \varepsilon_{2i}, \quad (2)$$

with  $z_i$  a  $p$ -dimensional vector of instruments, related to  $x_i$  but independent of  $\varepsilon_{2i}$ . For simplicity an intercept is included in  $z$ , such that there are, in fact, only  $p - 1$  IV in the above structure. We will assume that  $\varepsilon_i = (\varepsilon_{1i}, \varepsilon_{2i})'$  are independent and identically distributed (i.i.d.)  $N(0, \Sigma)$ , i.e., a bivariate normal distribution with zero mean vector and  $\Sigma$  variance-covariance matrix, where  $\Sigma$  has diagonal components  $\sigma_{11}$  and  $\sigma_{22}$  and off-diagonal component  $\sigma_{12} = \rho(\sigma_{11}\sigma_{22})^{1/2}$ . More general error structures are discussed in Section 3.1.

The key distinction between the above system of equations to a standard bivariate regression is the possible correlation between the error terms  $\varepsilon_{1i}$  and  $\varepsilon_{2i}$  (and, therefore, between  $x_i$  and  $\varepsilon_{2i}$ ). This leads to the well-known “endogeneity” bias when learning  $\beta$  from Eq. (2); that is the information of  $x_i$  that is correlated with  $\varepsilon_{2i}$  should not be used when learning about the regression parameter  $\beta$ .

The reduced form representation of Eqs. (1) and (2) is

$$x_i = z_i' \pi_x + v_{1i} \quad (3)$$

$$y_i = \gamma + z_i' \pi_y + v_{2i}, \quad (4)$$

where  $\pi_x = \delta$ ,  $v_{1i} = \varepsilon_{1i}$ ,  $\pi_y = \beta\delta$  and  $v_{2i} = \beta\varepsilon_{1i} + \varepsilon_{2i}$ . The relation between  $\varepsilon_i$  and  $v_i$  is

$$v_i = \begin{pmatrix} 1 & 0 \\ \beta & 1 \end{pmatrix} \varepsilon_i = B \varepsilon_i. \quad (5)$$

The model is not identified in the limiting case of  $\delta = 0$ . More generally, if the instruments  $z$  explain only a small portion of the variability of  $x$  (weak instrument case), then the likelihood function is concentrated around  $\beta + \sigma_{12}/\sigma_{11} = c$  for some estimable constant  $c$ .

Finally, it is worth noting that the interpretation of  $\rho$  as a measure of endogeneity needs to be extended when inference is performed from a Bayesian viewpoint. More specifically, the unconditional distribution of  $\varepsilon_i$ , once  $\Sigma$  is integrated out, might still exhibit dependence even when  $\rho = 0$ . Unconditionally, i.e., integrating out  $\Sigma$ , the dependence between  $\varepsilon_{1i}$  and  $\varepsilon_{2i}$  is actually higher than the conditional bivariate normal with a small

$\rho$  and potentially much higher when degrees of freedom for the prior is small. In this article, we will focus on what can be called “weak endogeneity in the frequentist sense,” or conditional endogeneity, to make it explicit we are not considering the case where  $\Sigma$  is integrated out.

## 2.1. Posterior Inference

Recall that  $\varepsilon_i$ s are i.i.d.  $N(0, \Sigma)$ , such that the distribution of the reduced form errors  $v_i$  is also bivariate normal  $N(0, \Omega)$  where

$$\Omega = B\Sigma B' = \begin{pmatrix} \sigma_{11} & \beta\sigma_{11} + \sigma_{12} \\ \beta\sigma_{11} + \sigma_{12} & \beta^2\sigma_{11} + 2\beta\sigma_{12} + \sigma_{22} \end{pmatrix} \quad (6)$$

so that  $\beta$  and  $\Sigma$  are intertwined in the reduced form and independent priors for both parameters would be counter-intuitive.

We will start, instead, with the prior specification for the structural parameters from Eqs. (1) and (2) discussed in Rossi et al. (2005):

$$\delta \sim N(d_0, D_0) \quad (7)$$

$$(\gamma, \beta)' \sim N(b_0, B_0) \quad (8)$$

$$\Sigma \sim IW(v_0, \Sigma_0), \quad (9)$$

for known hyperparameters  $d_0$ ,  $D_0$ ,  $b_0$ ,  $B_0$ ,  $v_0$ , and  $\Sigma_0$ .  $IW(v_0, \Sigma_0)$  stands for the inverted-Wishart distribution with parameters  $v_0$  (prior degrees of freedom) and  $\Sigma_0$  (prior scale matrix), whose density is  $p(\Sigma) \propto |\Sigma|^{-(v_0+p+1)/2} \exp\{-\frac{1}{2}\text{tr}\Sigma_0\Sigma^{-1}\}$  and  $v_0 > p$ . If  $v_0 \geq p+2$ , then  $E(\Sigma) = \Sigma_0/(v_0 - p - 1)$ . In this article,  $p = 2$ .

**Sampling  $\Sigma$ .** The full conditional distributions under the above specification are straightforward. First, the error variance has posterior

$$(\Sigma | \gamma, \beta, \delta, \text{data}) \sim IW(v_0 + n, \Sigma_0 + S),$$

where  $S = \sum_{i=1}^n \varepsilon_i \varepsilon_i'$  and  $\text{data} = \{(x_i, y_i, z_i); i = 1, \dots, n\}$ .

**Sampling  $(\gamma, \beta)$ .** Secondly, the regression parameters  $(\gamma, \beta)$  have a joint distribution of the form

$$(\gamma, \beta | \delta, \Sigma, \text{data}) \sim N(b_1, B_1),$$

where

$$B_1^{-1} = B_0^{-1} + \sum_{i=1}^n \tilde{x}_i \tilde{x}_i' \quad \text{and} \quad B_1^{-1} b_1 = B_0^{-1} b_0 + \sum_{i=1}^n \tilde{x}_i \tilde{y}_i,$$

for

$$\begin{aligned}\tilde{x}_i &= (1, x_i)' / \sigma_{2|1}^{1/2} \\ \tilde{y}_i &= (y_i - (x_i - z_i' \delta) \sigma_{12} / \sigma_{11}) / \sigma_{2|1}^{1/2} \\ \sigma_{2|1} &= \sigma_{22}(1 - \rho^2).\end{aligned}$$

**Sampling  $\delta$ .** Finally, the regression parameters for the instrument,  $\delta$ , have a distribution of the form

$$(\delta | \gamma, \beta, \Sigma, \text{data}) \sim N(d_1, D_1),$$

where

$$D_1^{-1} = D_0^{-1} + \sum_{i=1}^n \tilde{z}_i \tilde{z}_i' \quad \text{and} \quad D_1^{-1} d_1 = D_0^{-1} d_0 + \sum_{i=1}^n \tilde{z}_i \tilde{x}_i,$$

for

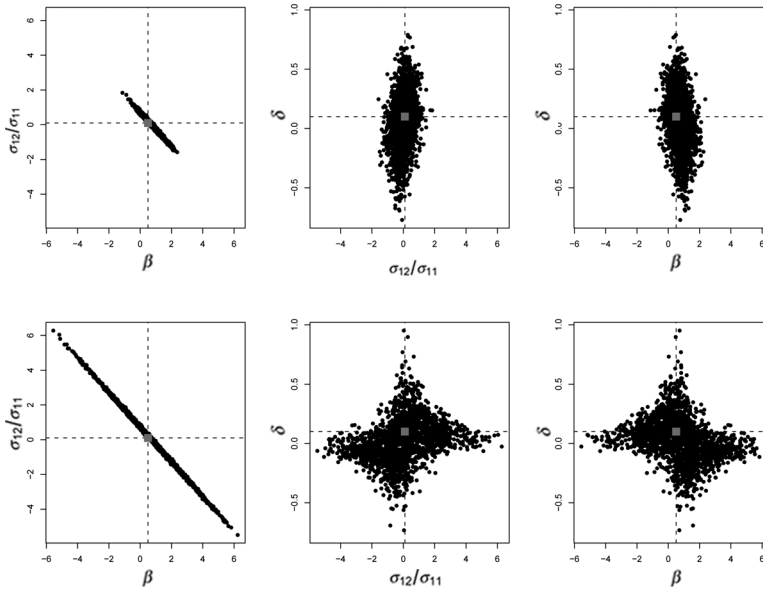
$$\begin{aligned}\tilde{x}_i &= (x_i - (y_i - \gamma - \beta x_i) \sigma_{12} / \sigma_{22}) / \sigma_{1|2}^{1/2} \\ \tilde{z}_i &= z_i / \sigma_{1|2}^{1/2} \\ \sigma_{1|2} &= \sigma_{11}(1 - \rho^2).\end{aligned}$$

## 2.2. Illustration

The performance of the above Gibbs sampler is illustrated with a data set with  $n = 200$  observations simulated from Eqs. (1) and (2) with  $\gamma = 1.0$ ,  $\beta = 0.5$ ,  $\sigma_{11} = \sigma_{22} = 1$ ,  $\delta = (1.0, 0.1)'$  (one weak instrument), and  $\rho = 0.1$  (low degree of endogeneity).

The prior hyperparameters are  $d_0 = b_0 = 0$  and  $D_0 = B_0 = 25I_2$ , suggesting relatively low prior information about  $\delta$ ,  $\beta$ , and  $\gamma$ . For  $\Sigma$  two priors are entertained:  $v_0 = 3$  and  $\Sigma_0 = 3I_2$  (relatively vague prior) and  $v_0 = 0.00001$  and  $\Sigma_0 = 0.00001I_2$  (non-informative prior). The distinction between “relatively vague” and “non-informative” is arbitrary and its only purpose is to distinguish two prior specification for  $\Sigma$ . As it can be seen, the “non-informative” prior tries to mimic the improper Jeffreys prior for covariance matrices, while the “relatively vague” prior is proper and has prior mean but no prior variance.

Figure 1 summarizes our findings and shows that the variability of  $\sigma_{12}/\sigma_{11}$  and of  $\beta$  are greatly affected by the choice of the prior on  $\Sigma$ . Posterior inference based on the relatively vague prior (top row) turns out to be rather informative when compared to posterior inference based on the non-informative prior (bottom row). Notice that  $\beta$  and  $\delta$



**FIGURE 1** Inverted Wishart prior:  $\Sigma \sim IW(v_0, \Sigma_0)$ . Top row:  $v_0 = 3$  and  $\Sigma_0 = 3I_2$  (relatively vague prior). Bottom row:  $v_0 = 0.00001$  and  $\Sigma_0 = 0.00001I_2$  (non-informative prior). The independent priors on  $\gamma$ ,  $\beta$ , and the components of  $\delta$  are all  $N(0, 25)$ . The Gibbs sampler is run for 1,000,000 draws (after discarding the initial 10,000 draws) with every 1,000th kept for posterior summaries.

(or  $\sigma_{12}/\sigma_{11}$ ) become relatively unrelated (top middle and right panels), suggesting again that the relatively vague prior turns out to be too informative regarding these linear dependences. Meanwhile, under the non-informative prior,  $\beta$  (similar for  $\sigma_{12}/\sigma_{11}$ ) becomes more diffuse a posteriori (varies between  $(-6, 6)$ ) when  $\delta$  approaches zero (or the instrument becomes innocuous) highlighting the sharp behavior of the likelihood function in the vicinity of non-identifiability (bottom middle and right panels). The pattern repeats itself, as expected, regardless of the sample size, since the (non-identified) likelihood dominates the non-informative prior even for fairly small sample sizes. These plots resemble the bivariate distributions with the conditional normal example of Arnold and Strauss (1991) and the example of uncorrelated models with normal marginals of Ebrahimi et al. (2010). Comparisons of the top and bottom panels show that the parameters of Wishart prior can have profound distributional effects. See also the discussion on Section 5 of Hoogerheide and van Dijk (2008a) where similar plots for  $(\beta, \delta)$  are exhibit.

To sum up, the trade-off between more precise (and less accurate) and less precise (and more accurate) posterior inference is an important, problem-specific part of the modeling process that needs to be dealt with on a case-by-case basis. In the next section, we discuss alternatives to the above model and prior specifications.

### 3. MORE ON IV REGRESSION

A number of authors have proposed the use of Jeffreys prior in the IV regression. From the reduced form (Eqs. (3) and (4)), let  $\pi = (\pi_x, \pi_y) = (\delta, \beta\delta)$  a  $(p \times 2)$  matrix of rank one. The Jeffreys prior is then given by

$$\left| \frac{\partial \pi}{\partial(\beta, \delta)} \left( \frac{\partial \pi}{\partial(\beta, \delta)} \right)' \right| = \|\delta\| (1 + \beta^2)^{\frac{1}{2}}, \quad (10)$$

which is suggestive of using priors with polynomial tails, such as a Cauchy (Chao and Phillips, 1998; Sims, 2007).

The reduced form IV model where inference for  $\pi_y = \beta\delta$  is required (see text after Eq. (4)) is related to the analysis of a product of two normal means (Lindley and El Sayyad, 1968; Berger and Bernardo, 1989). The issue is two-fold: we have a high dimensional parameter vector that requires prior regularisation and we have to learn how much shrinkage to employ. There is a fundamental trade-off in the fit of the two regression equation described above. This was illustrated by the example in the previous section and presented in Fig. 1.

Sims (2007) makes a number of important remarks about prior sensitivity in IV problems. Two of them are as follows:

- 1) *As the likelihood does not go to zero as  $\beta \rightarrow \infty$ , with  $\Sigma$  fixed, no matter what the sample size is there is an issue of prior sensitivity.*

Assume we combine a likelihood with a prior that has Gaussian tails. As the likelihood mode (maximum likelihood estimator) moves away from the prior mean (and mode), the posterior mean (and mode) will move away from the prior mean (and mode) at first (as expected), but then reverses the direction, coming back to settle at the prior mean (or mode) even when the likelihood mode and prior mean (or mode) are well separated. These effects of fat-tailed combinations of likelihoods and priors has been documented in the linear Bayes regression model by West (1984) and exploited to find sparse estimators by Carvalho et al. (2010).

- 2) *Put another way, in a sample where the posterior is highly non-Gaussian and has substantial, slowly declining tails, even apparently weak prior information  $(\delta, \beta) \sim N(0, 100I)$  can substantially affect the inference.*

One approach then is to make the prior flexible enough to accommodate fat-tails and to *learn* from the data how thick these tails should be, see Lopes and Polson (2011). There will be an interaction with the specification of the prior magnitude on  $\Sigma$ . See the discussion in Lindley and El Sayyad (1968) in their Bayesian treatment. Our approach will be to free up the prior on  $\Sigma$  and use a Cholesky-based prior rather than the standard inverse Wishart.



### 3.1. Dirichlet Process Prior

Conley et al. (2008) develop a Bayesian semiparametric approach to the IV regression problem. They use a normal-based Dirichlet process prior to jointly model structural and instrumental variable equations errors (the errors in our Eqs. (1) and (2)). More specifically,  $\varepsilon_i \sim N(0, \Sigma_i)$  replaces the standard  $\varepsilon_i \sim N(0, \Sigma)$ , with  $\Sigma_i$  now i.i.d. from the discrete random distribution  $G$ , which in turn is modeled by a Dirichlet process with concentration parameter  $\alpha$  and base distribution  $G_0$ , commonly denoted by  $G \sim DP(\alpha, G_0)$ . The marginal distribution of  $\Sigma_i$  (integrating out  $G$ ) is continuous and is called a mixture of Dirichlet Processes (MDP). See Escobar and West (1995, 1998) for Bayesian posterior inference via Markov chain Monte Carlo (MCMC) for this class of models.

Conley et al. performed extensive sampling experiments and showed that, when the errors are actually normally distributed, their prior specification leads to more efficient results when compared to standard Bayesian (Section 2) and classical methods. They say, in their conclusion, that “Our Bayesian semi-parametric procedure produces credibility regions which are dramatically shorter than confidence intervals based on the weak instrument asymptotics. The shorter intervals from our method are produced by more efficient use of sample information.” They go on and finish the paper saying that “... our nonparametric Bayesian method dominates Bayesian methods based on normal errors and may be preferable to methods from the recent weak instruments literature if the investigator is willing to trade-off lower coverage for dramatically smaller intervals.” To sum up, their Bayesian methodology should be the preferred one in many IV problems from here on.

### 3.2. Bayesian Model Averaging

Recent research has focussed on applying Bayesian model averaging across sets of instruments, exogeneity restrictions, the validity of identifying restrictions, and the set of exogenous regressors are explored by Eicher et al. (2009) and Koop et al. (2011). Another avenue, which we do not explore here, is to assume a flexible fat-tailed alternative such as a mixture of  $t_v$  distributions (mixed over  $v$ ). Our use of Cholesky priors for  $\Sigma$  follows through in these settings as well. As we will see, there can be interesting sensitivity issues to the prior specification.

## 4. ERRORS-IN-THE-VARIABLES MODELS

Minka (1999) proposed a proper Bayes approach to linear regression with errors in both equations. This builds on the well-known work of Zellner (1971, Chapter 5). Prior specification is an important feature of this

problem, as Minka observes *Statisticians have worked on regression but they often shoot themselves in the foot by using priors that are too weak or too strong*. Our flexible cholesky-based prior on  $\Sigma$  will allieviate the problem.

The issues can be seen in the simple one dimensional case. The EVM is given by

$$x_i = z_i + \varepsilon_{1i} \quad (11)$$

$$y_i = \beta z_i + \varepsilon_{2i}, \quad (12)$$

where we take  $\delta = 1$  and  $\gamma = 0$  from the previous specification. The  $z_1, \dots, z_n$ s are now latent unobserved variables, or Zellner's "incidental parameters," which we endow with a normal prior  $z_i \sim N(0, \tau^2)$ . The covariance of the error term is again  $\Sigma$ , which, for the purpose of illustration, will be considered diagonal, that is  $\Sigma = \text{diag}(\sigma_{11}^2, \sigma_{22}^2)$ .

Let  $(x, y) = (x_1, \dots, x_n, y_1, \dots, y_n)$  be the observed data and  $S = \sum_{i=1}^n (x_i, y_i)'(x_i, y_i)/n$  be the sample covariance matrix. Zellner (1971) was the first to calculate the Bayes marginal likelihood given by

$$p(x, y | \theta, \Sigma, \tau^2) \propto |\Sigma|^{-\frac{n}{2}} (1 + \tau^2 \theta' \Sigma^{-1} \theta)^{-\frac{1}{2}} \exp\{-0.5 \text{tr}(GS)\}, \quad (13)$$

where  $\theta = (1, \beta)'$  and  $G = \Sigma^{-1} - \Sigma^{-1} \theta (\theta' \Sigma^{-1} \theta + \tau^{-2})^{-1} \theta' \Sigma^{-1}$ .

Zellner's Bayes marginal likelihood approach can differ dramatically when compared to the conditional two stage least squares (2SLS) approach. Here, if  $q$  is the principal eigenvector of  $S\Sigma^{-1}$ , then the Maximum likelihood estimator (MLE) is  $\hat{\beta} = q_2/q_1$  and the eigenvalues play a minor role conditionally. In the marginal Bayes likelihood, however, they have a very influential effect and it is typical for the marginal likelihood to have two finite local maxima. The case where they agree is when  $\tau \rightarrow \infty$  although this is rarely the case in practice. Another special case of interest is when the largest eigenvalue is small. Then the marginal likelihood has one maximum and it is near zero and all the variation in the data can be explained by noise and  $\hat{\beta}$  is driven towards zero.

## 5. CHOLESKY-BASED PRIOR

This section will focus on one choice equation and one outcome equation, although it can directly be extended to multivariate scenarios. The use of Cholesky-based priors is not new and they have had success in many situations; see, for example, Lopes et al. (2011b) for an application to high dimensional stochastic volatility modeling (see also Pourahmadi, 1999, for longitudinal models).

The key idea is the following. Instead of modeling  $\Sigma$  via an inverted Wishart distribution with parameters  $v_0$  and  $\Sigma_0$ , i.e.,  $\Sigma \sim IW(v_0, \Sigma_0)$ ,

we will model the components of the recursive conditional regressions that arises from the Cholesky decomposition of  $\Sigma$ . More precisely, recall from Section 2.1 that  $\varepsilon_i = (\varepsilon_{1i}, \varepsilon_{2i})'$  are i.i.d.  $N(0, \Sigma)$ , and let

$$\Sigma = AHA' \quad (14)$$

be the Cholesky decomposition of  $\Sigma$  such that  $A$  is lower triangular with ones in the main diagonal and lower triangular component given by  $a_{21} = \sigma_{12}/\sigma_{11}$  and  $H = \text{diag}(\sigma_{11}, \sigma_{2|1})$ . The reverse transformation is given by  $\sigma_{12} = a_{21}\sigma_{11}$  and  $\sigma_{22} = \sigma_{2|1} + \sigma_{12}^2/\sigma_{11}$ . Therefore,

$$A^{-1}\varepsilon_i \sim N(0, H), \quad (15)$$

and  $\varepsilon_i \sim N(0, \Sigma)$  can be rewritten by the following recursive conditional regressions (or simply *triangular regressions*)

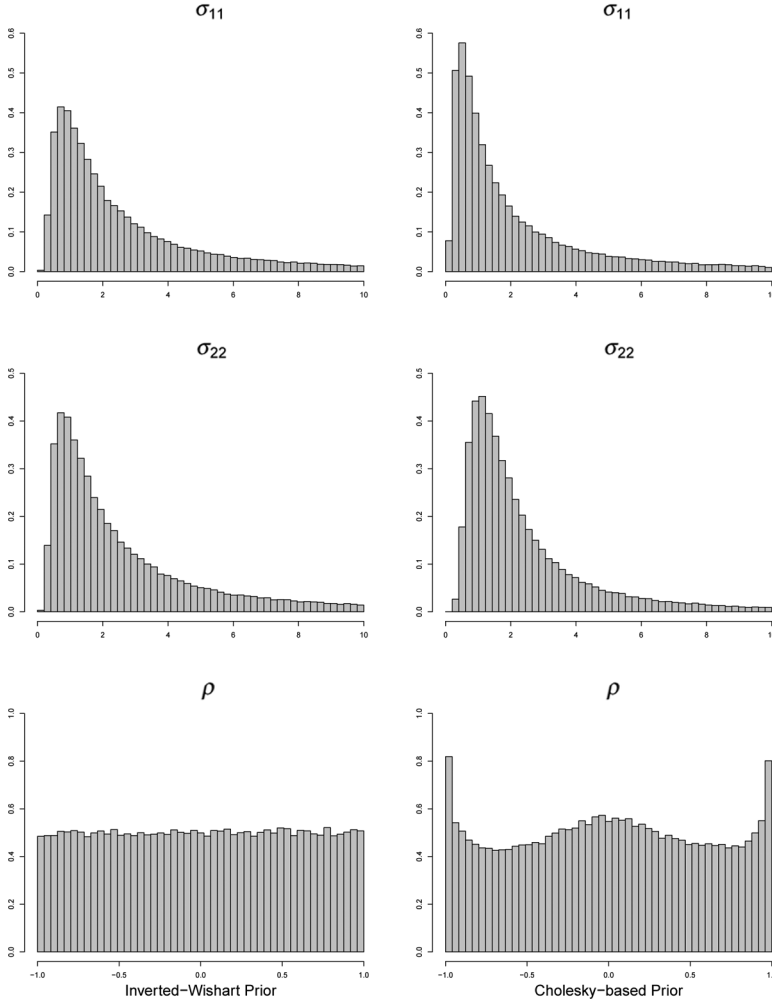
$$\varepsilon_{1i} \sim N(0, \sigma_{11}) \quad (16)$$

$$\varepsilon_{2i}|\varepsilon_{1i} \sim N(a_{21}\varepsilon_{1i}, \sigma_{2|1}). \quad (17)$$

The parameter  $a_{21}$  measures the strength of the correlation between  $\varepsilon_{1i}$  and  $\varepsilon_{2i}$ , while  $\sigma_{2|1}$  is the conditional residual variance. We then specify independent prior distributions for  $\sigma_{11}$ ,  $a_{21}$ , and  $\sigma_{2|1}$ . The implied prior for  $\Sigma$  can be directly obtained, either analytically or via Monte Carlo simulation. More specifically,  $\sigma_{11}$  is learned from Eq. (16),  $\sigma_{12}$  is learned from  $\sigma_{11}$  and  $a_{21}$  ( $= \sigma_{12}/\sigma_{11}$ ), see Eq. (17), and  $\sigma_{22}$  is learned from  $\sigma_{11}$ ,  $\sigma_{12}$  and  $\sigma_{2|1}$  ( $= \sigma_{22} - \sigma_{12}^2/\sigma_{11}$ ), also from Eq. (17).

The main attractiveness of the Cholesky-based prior is its relative freedom to independently quantify the uncertainty for the individual components of  $\Sigma$ , which turns out to be one of the major constraints of the Wishart and inverted Wishart distributions. See Rossi et al. (2005) and Lopes et al. (2011b) for additional discussion on the limitations of the Wishart distribution.

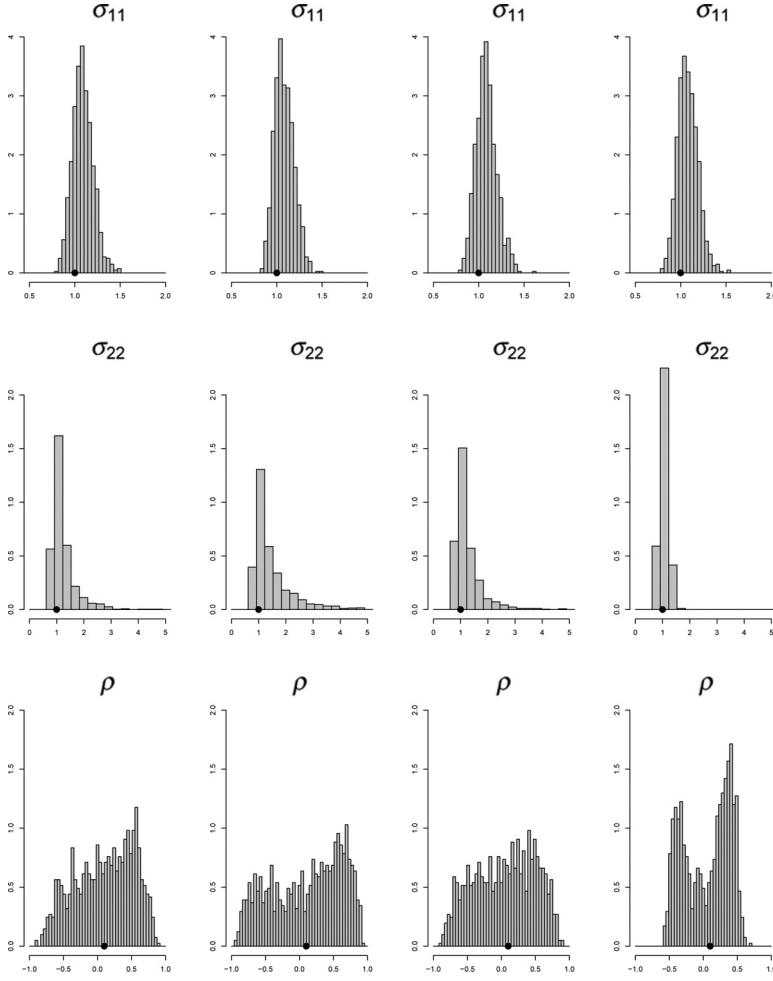
The MCMC scheme for the IV regression model with Cholesky-based prior is the same as with the inverted-Wishart prior for all parameters but the covariance ones, that is  $\sigma_{11}$  and  $\sigma_{2|1}$  and  $a_{21}$  (see Section 2.1). For  $\sigma_{11}$  and  $\sigma_{2|1}$ , we assign independent standard inverted gamma priors, while a normal prior is specified for  $a_{21}$ . These priors are combined with Eqs. (16) and (17) and, conditional on  $\gamma, \beta$ , and  $\delta$ , leads to standard Gibbs updates for  $\sigma_{11}$  and  $\sigma_{2|1}$  and  $a_{21}$ . The next section presents three illustrative applications.



**FIGURE 2** Inverted Wishart and Cholesky-based priors.  $\Sigma \sim IW(3, 3I_2)$  (left column) and  $\sigma_{11} \sim IG(0.75, 0.75)$ ,  $\sigma_{2|1} \sim IG(3, 3)$ , and  $a_{21} \sim N(0, 0.7)$  (right column).

### 5.1. Illustration 1: Synthetic Data

We continue in the context of weak instrument and low endogeneity introduced in Section 2.2, that is  $\delta = (1.0, 0.1)'$  and  $\rho = 0.1$ . Figure 2 compares the inverted Wishart prior,  $\Sigma \sim IW(3, 3I_2)$ , with the implied prior obtained from the Cholesky-based prior of  $(\sigma_{11}, a_{21}, \sigma_{2|1})$ , where  $\sigma_{11} \sim IG(0.75, 0.75)$ ,  $\sigma_{2|1} \sim IG(3, 3)$  and  $a_{21} \sim N(0, 0.7)$ . The hyperparameters were selected to make both prior distributions on  $\Sigma$  as similar as possible. In addition, two prior specifications for  $(\beta, \delta)$  are entertained. In the first case,  $(\beta, \delta) \sim N(0, 25I_3)$ , which represents diffuse but proper prior



**FIGURE 3** Marginal posteriors. Columns are histograms approximations to the marginal posterior distributions based on different prior specifications for  $(\beta, \delta)$  and  $\Sigma$ . From left to right:  $(\beta, \delta) \sim N(0, 25I_3)$  and  $\Sigma \sim IW(3, 3I_2)$ ;  $(\beta, \delta) \sim N(0, 25I_3)$  and Cholesky-based prior for  $\Sigma$  (see caption of Fig. 2);  $p(\beta, \delta) \propto \|\delta\|(1 + \beta^2)^{\frac{1}{2}}$  (see Eq. (10)) and  $\Sigma \sim IW(3, 3I_2)$ ;  $p(\beta, \delta) \propto \|\delta\|(1 + \beta^2)^{\frac{1}{2}}$  and Cholesky-based prior for  $\Sigma$ .

distributions. In the second case,  $p(\beta, \delta) \propto \|\delta\|(1 + \beta^2)^{\frac{1}{2}}$  (see Eq. 10), which can be thought of as a Jeffreys-type prior. This leads to four prior specifications for  $(\beta, \delta, \Sigma)$ . As before, the prior for  $\gamma$  is  $N(0, 25)$ .

The marginal posterior distributions for  $\sigma_{11}$ ,  $\sigma_{22}$ ,  $\rho$  and the joint marginal for  $(\sigma_{12}/\sigma_{11}, \sigma_{22})$  and  $(\beta, \delta)$  appear in Figs. 3 and 4, respectively. The learning of  $\sigma_{11}$  and  $\sigma_{22}$  is similar across prior specifications, with the forth prior specification (Jeffreys for  $(\beta, \delta)$  and Cholesky-based for  $\Sigma$ ) leading to more informative posterior for  $\sigma_{22}$  as well as  $\rho$  and  $\beta$ .

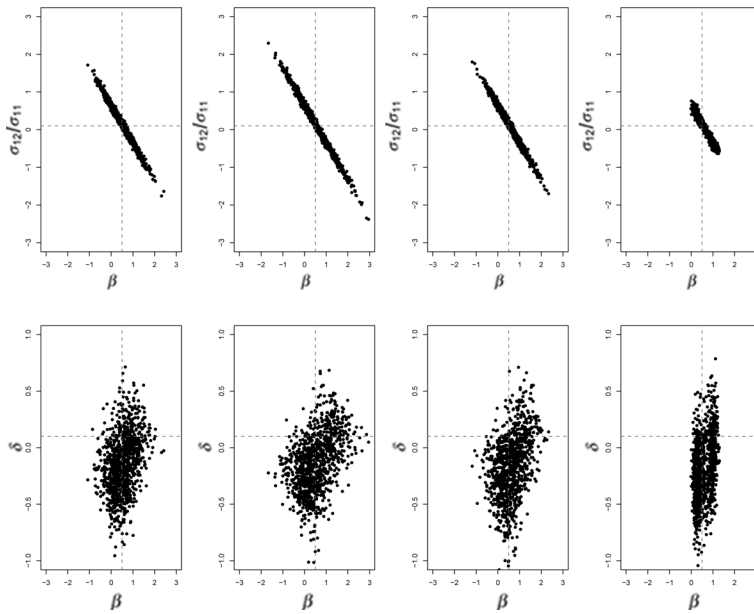
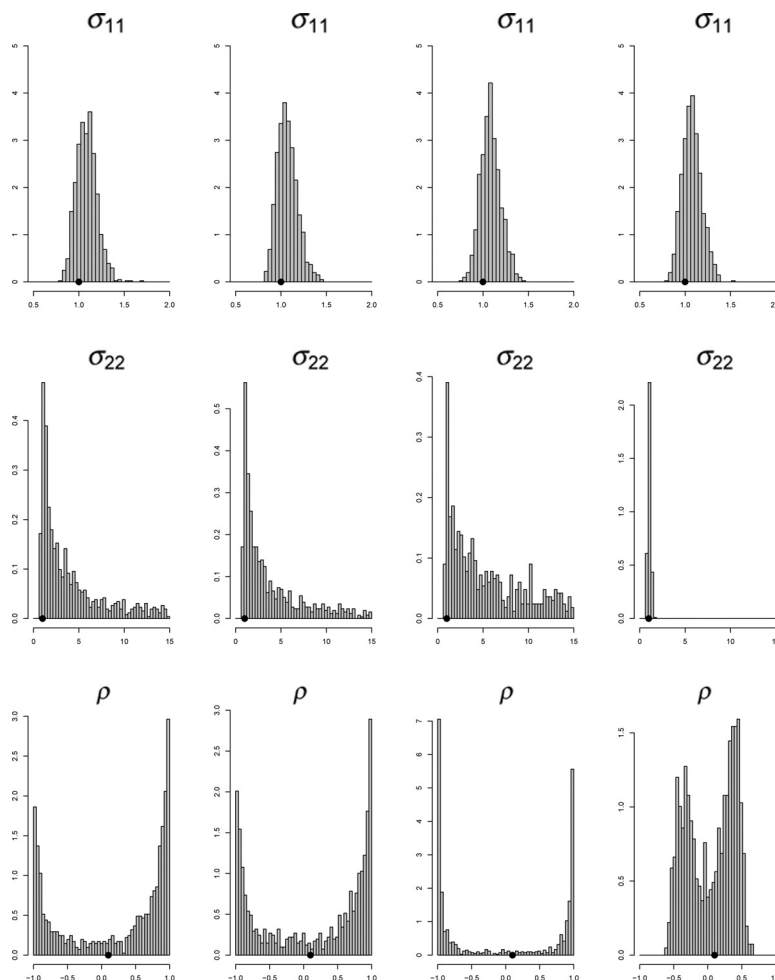


FIGURE 4 Joint posterior of  $(\beta, \sigma_{12}/\sigma_{11})$  and  $(\beta, \delta)$ . Columns are as in Fig. 3.

Similar results are found when the priors on  $\Sigma$  or on  $(\sigma_{11}, \sigma_{2|1}, a_{21})$  become extremely noninformative, that is, when  $\Sigma \sim IW(0.00001, 0.00001I_2)$  or  $\sigma_{11} \sim IG(0.000001, 0.000001)$ ,  $\sigma_{2|1} \sim IG(0.000001, 0.000001)$  and  $a_{21} \sim N(0, 1000)$ . See Figs. 5 and 6. It appears that the Jeffreys prior for  $(\beta, \delta)$  when combined with a diffuse Cholesky-based prior for  $\Sigma$  leads to unstable results, while the other three combination produce relatively similar results.

## 5.2. Illustration 2: Demand for Cigarettes

Here we revisited the demand for cigarettes illustration presented in Chapter 10 of Stock and Watson's (2003) textbook *Introduction to Econometrics*, pp. 339–341. The goal is to study the effect of price changes on the demand for cigarettes when using sales tax as an instrumental variable. The data set consists of annual data for the 48 continental United States for the year of 1995 and the proxy for price is the logarithm of the average real price per pack of cigarettes including all taxes ( $x$  in our notation). In addition, the proxy for consumption is the logarithm of the number of packs of cigarettes sold per capita in the state ( $y$  in our notation), while the proxy for sales tax is the portion of the tax on cigarettes arising from the general sales tax, measured in dollars per pack in real dollars, deflated by the consumer price index ( $z$  in our notation). The previous

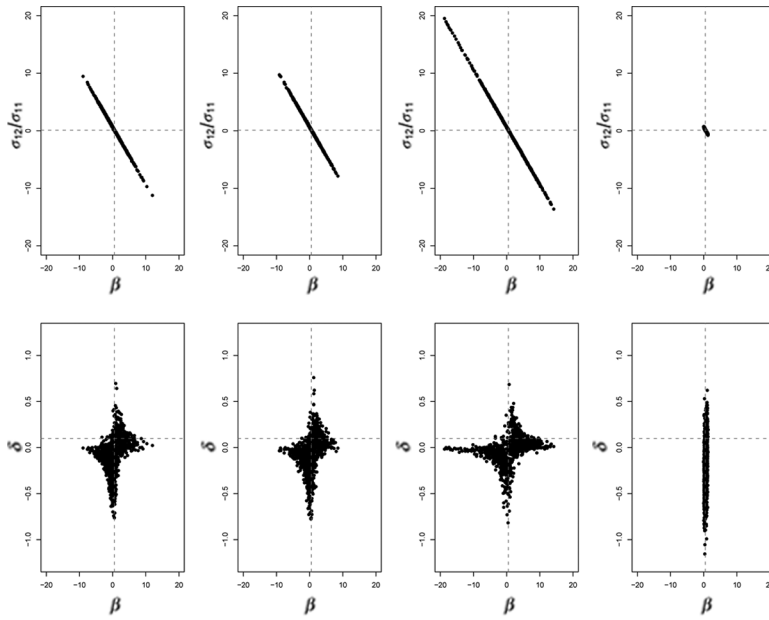


**FIGURE 5** Marginal posteriors. Same as Fig. 3, but based on more diffuse prior specification. For the inverted Wishart prior,  $\Sigma \sim IW(0.00001, 0.00001I_2)$ , while for the Cholesky-based prior,  $\sigma_{11} \sim IG(0.000001, 0.000001)$ ,  $\sigma_{2|1} \sim IG(0.000001, 0.000001)$ , and  $a_{21} \sim N(0, 1000)$ .

description was taken from the webpage for Stock and Watson's book at [http://wps.aw.com/aw\\_stock\\_ie\\_2/50/13016/3332253.cw/index.html](http://wps.aw.com/aw_stock_ie_2/50/13016/3332253.cw/index.html).

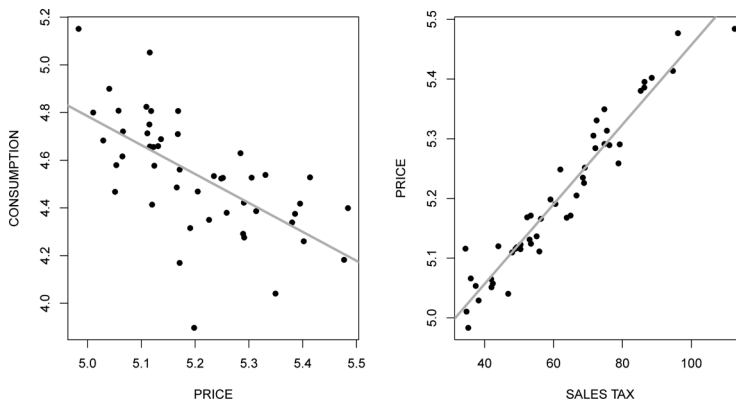
Figure 7 presents the data with ordinary least squares estimates given by  $(\gamma_{ols}, \beta_{ols}) = (10.850, -1.213)$  and  $\delta_{ols} = (4.79006, 0.00667)$ . The sample coefficient of correlation between the residuals of both OLS fits is around  $-0.1775306$ . These preliminary results suggest a scenario of a low degree of endogeneity and a relatively weak instrument, somewhat similar to the previous simulation exercise.

As with the simulation exercise, a fairly vague prior specification is used for all the model parameters, that is,  $\Sigma \sim IW(a_0, a_0I_2)$  or  $\sigma_{11} \sim$



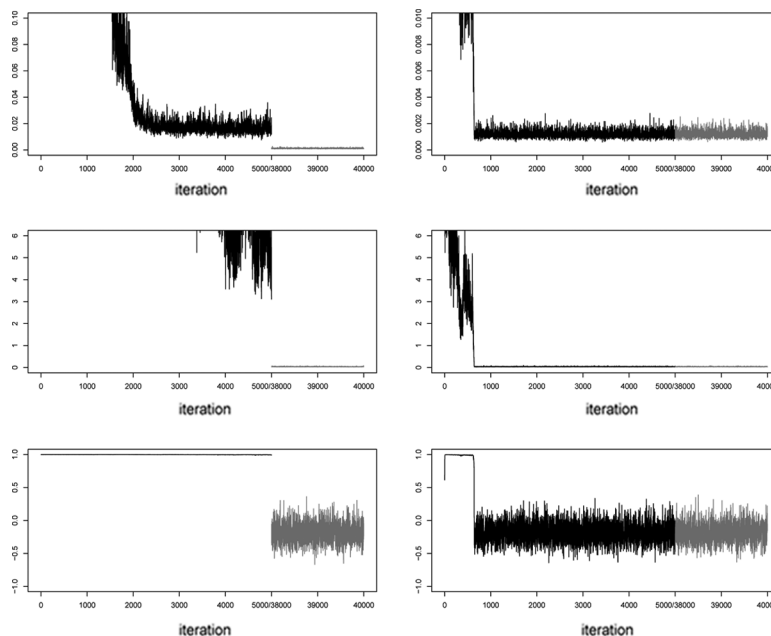
**FIGURE 6** Joint posterior of  $(\beta, \sigma_{12}/\sigma_{11})$  and  $(\beta, \delta)$ . Same as Fig. 4, but based on the prior specification described in Fig. 5.

$IG(a_0, a_0)$ ,  $\sigma_{2|1} \sim IG(a_0, a_0)$ , where  $a_0 = 0.0000001$ , and  $a_{21} \sim N(0, \sigma_0^2)$ ,  $(\beta, \delta) \sim N(0, \sigma_0^2 I_3)$ , and  $\gamma$  is  $N(0, \sigma_0^2)$ , where  $\sigma_0^2 = 1,000,000$ . The initial values for the parameters are  $\sigma_{11}^{(0)} = \sigma_{2|1}^{(0)} = 1$ ,  $a_{21}^{(0)} = 0$ ,  $\gamma^{(0)} = \beta^{(0)} = 0$ , and  $\delta^{(0)} = (0, 0)$ . The performance of both MCMC schemes are presented in



**FIGURE 7** Demand for cigarettes. The data set consists of annual data for the 48 continental U.S. states for the year of 1995. PRICE is the logarithm of the average real price per pack of cigarettes including all taxes. CONSUMPTION is the logarithm of the number of packs of cigarettes sold per capita in the state. SALES TAX is the portion of the tax on cigarettes arising from the general sales tax, measured in dollars per pack (in real dollars, deflated by the Consumer Price Index).





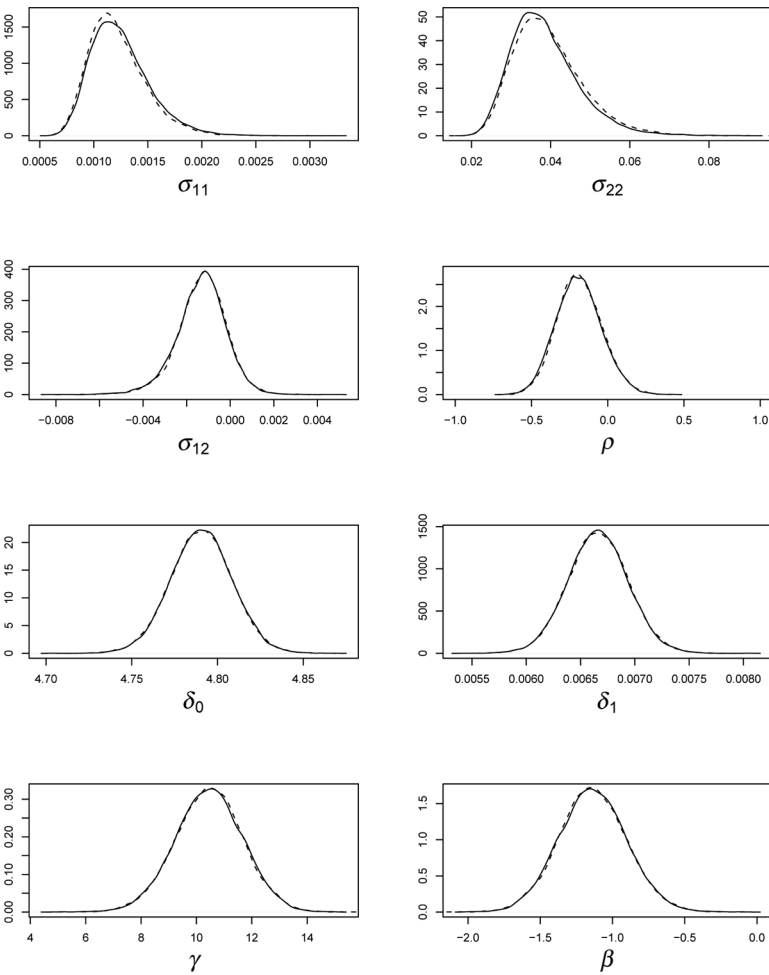
**FIGURE 8** Demand for cigarettes. Trace plots of MCMC based on diffuse priors. Inverted Wishart prior (left panels) and Cholesky-based priors (right panels). After 5,000 draws the MCMC based on the Inverted Wishart prior has not yet converged, while the MCMC based on the Cholesky-based prior converges after 1,000 draws.

Fig. 8. Both algorithms converge after several thousand draws, with the Cholesy-based one performing slightly better than the inverted Wishart one.

Posterior inference is summarized in Fig. 9 and Table 1. The posterior median for the parameter that measure the degree of endogeneity,  $\rho$ , is  $-0.1934$  while the posterior probability that  $\rho > 0$  is about 10% (under both prior specifications). On the other hand, the instrumental variable, sales tax, has a nonzero effect of price since the 95% credibility interval for  $\delta_1$ ,  $(0.0061, 0.0072)$ , is away from zero. Based on the posterior of  $\beta$ , one can argue that an increase in the price of 1% reduces consumption by 0.67% to 1.59%.

### 5.3. Illustration 3: Return to Education

Here we revisited the return to education illustration presented in Chapter 8 of Lancaster's textbook *An Introduction to Modern Bayesian Econometrics*, pp. 325–334. The goal is to study the effect of (years of) education on (log) wages when using quarter of birth as an instrumental variable. This is a fairly well-known study and was first proposed by Angrist and Krueger (1991). We follow Lancaster's simplification and focus only

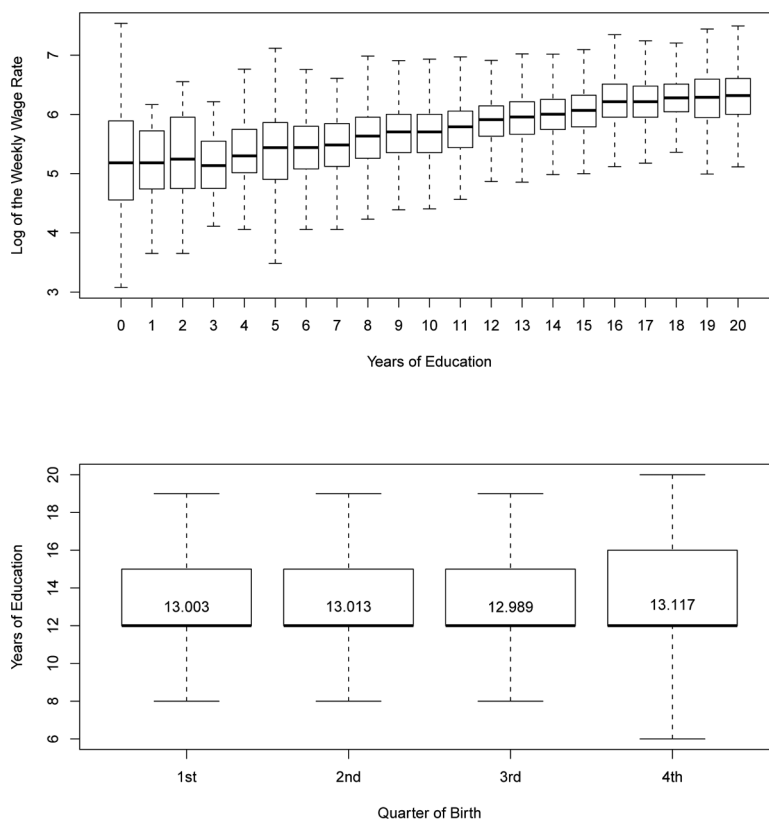


**FIGURE 9** Demand for cigarettes. Marginal posterior densities based on diffuse priors. Inverted Wishart prior (solid lines) and Cholesky-based priors (dashed lines).

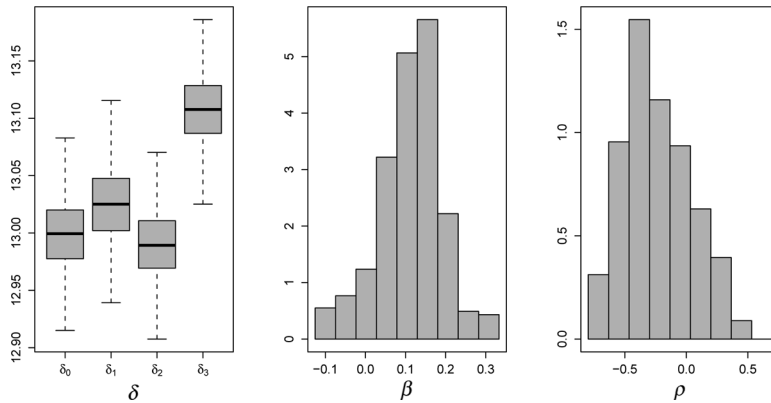
**TABLE 1** Demand for cigarettes. Summaries of the marginal posterior distributions

Parameter	Median	95% Credible interval
$\sigma_{11}$	0.0012	(0.0008, 0.0019)
$\sigma_{22}$	0.0374	(0.0254, 0.0584)
$\sigma_{12}$	−0.0013	(−0.0038, 0.0008)
$\rho$	−0.1934	(−0.4655, 0.1111)
$\delta_0$	4.7907	(4.7554, 4.8251)
$\delta_1$	0.0067	(0.0061, 0.0072)
$\gamma$	10.4560	(8.0198, 12.8142)
$\beta$	−1.1373	(−1.5904, −0.6694)

on men born in 1939, which is the last year of the 10-year study of Angrist and Krueger (1991). They argued that quarter of birth might be related to years of education due to age-related regulations to both enter and leave school. More precisely, children whose birthdays fall in the 4th quarter of the calendar year will enter (elementary) school the fall of that same year or within one or two months from their birthdays, while children whose birthdays fall, say, in the 1st quarter of the calendar year will enter school at least six months after their birthday. In addition, compulsory schooling laws require students to remain in school until a predetermined age (usually sixteen or seventeen). Figure 10 shows the summary of both relationships. As it can be seen by the differences between the average years of education per quarter of birth, the instrument (quarter of birth) is relatively weak.



**FIGURE 10** Return to education. Sample size of  $n = 35,805$  men born in 1939. A subset of the data set analyzed by Angrist and Krueger (1991). Across quarters of birth, the median years of education is 12, corresponding to completion of high school. Means are slightly increasing from 1st to 4th quarter of birth, with the difference between 13.117 (4th quarter) and the other quarters ranges between 5.5 and 6.5 weeks of education.



**FIGURE 11** Return to education. Marginal posterior densities based on diffuse priors. Posterior means of  $\delta_0$ ,  $\delta_1$ ,  $\delta_2$ , and  $\delta_3$  are 12.999, 13.025, 12.990, and 13.108, respectively, while  $Pr(\beta < 0 | \text{data}) = 7.7\%$  and  $Pr(\rho < 0 | \text{data}) = 73.4\%$ .

We perform Bayesian inference based on our Cholesky-based prior for the components of  $\Sigma$ . The prior hyperparameters specified as in the previous illustration, which leads to fairly vague prior information. Posterior summaries are presented in Fig. 11. As expected, apart from  $\delta_3$ , all  $\delta$ s are quite similar, corroborating with the initial suggestion that quarter of birth is a weak instrument for this illustration. Nonetheless, the posterior mean and median of the percentage marginal return,  $100\beta$ , are 11.2% and 11.9%, respectively, while the 95% credibility interval is (6.57%, 16.5%). Finally, the degree of endogeneity can be measure by  $\rho$ , whose posterior mean and median are  $-0.245$  and  $-0.191$ , respectively, while  $Pr(\rho < 0 | \text{data}) = 73.4\%$ . Lancaster (p. 333) says that “the posterior suggests that the structural form errors are negatively correlated and this is a bit surprising on the hypothesis that a major element of both  $\varepsilon_1$  and  $\varepsilon_2$  is *ability* and this variable tends to affect positively both education and wages. But the evidence is very far from conclusive.” We agree and claim that this example, as well as the previous ones, illustrates the inferential difficulties when combining weak instruments and low degree of endogeneity.

## 6. DISCUSSION

Instrumental variable likelihoods and their errors-in-the-variables cousin have challenging likelihood surfaces. With that comes the issue of sensitivity to prior specification. It has long been known that the Bayesian solution to the identifiability problem is attractive (Conley et al., 2008; Zellner, 1971) and that, given a prior, inference can be based on marginal likelihoods.

Here we introduce the use of Cholesky-based priors, which are more flexible than the traditional normal inverse-Wishart regression priors, and more realistic than an “uninformative” Jeffreys prior. Recent work in standard regression problems has been addressed with heavy-tailed Cauchy priors (Gelman et al., 2008). We show how prior-posterior inference can be formulated in a Gibbs sampler and compare its performance in the weak instruments case. Given modern-day computational methods for Bayesian inference (Gamerman and Lopes, 2006; Lopes et al., 2011a), these complicated likelihoods seem ripe for more discussion of prior sensitivity.

## ACKNOWLEDGMENTS

The authors would like to thank The University of Chicago Booth School of Business for providing financial support for our research. The authors are grateful to the Guest Editor, Ehsan Soofi, and the two anonymous referees, whose invaluable comments and suggestions significantly improved the presentation and the quality of the paper. Finally, we would like to thank our beloved friend Arnold Zellner for being invariably enthusiastic about Bayesian statistics and an important source of inspiration and support.

## REFERENCES

- Angrist, J. D., Krueger, A. B. (1991). Does compulsory school attendance affect schooling and earnings. *Quarterly Journal of Economics* 106:979–1014.
- Arnold, B. C., Strauss, D. (1991). Bivariate distributions with conditionals in prescribed exponential families. *Journal of the Royal Statistical Society, Series B* 53:365–375.
- Berger, J. O., Bernardo, J. M. (1989). Estimating the product of means: Bayesian analysis with reference priors. *Journal of American Statistical Association* 84:200–207.
- Carvalho, C. M., Polson, N. G., Scott, J. G. (2010). The Horseshoe estimator of sparse signals. *Biometrika* 97(2):465–480.
- Chamberlain, G. (2007). Decision theory applied to an instrumental variables model. *Econometrica* 75(3):609–652.
- Chao, J. C., Phillips, P. C. B. (1998). Posterior distribution in limited information analysis of the simultaneous equations model using Jeffreys prior. *Journal of Econometrics* 87:49–86.
- Chetty, V. K. (1966). Bayesian analysis of some simultaneous equation models and specification errors. *Unpublished PhD Thesis*, University of Wisconsin, Madison.
- Conley, T., Hansen, C., McCulloch, R. E., Rossi, P. E. (2008). A semi-parametric Bayesian approach to the instrumental variable problem. *Journal of Econometrics* 144:276–305.
- Drèze, J. H., Morales, J. A. (1976). Bayesian full information analysis of simultaneous equations. *Journal of American Statistical Association* 71:919–923.
- Drèze, J. H. (1976). Bayesian limited information analysis of the simultaneous equations model. *Econometrica* 44:1045–1075.
- Drèze, J. H., Richard, J. F. (1983). Bayesian analysis of simultaneous equations systems. In: Griliches, Z., Intrilligator, M. D., eds. *Handbook of Econometrics*. Vol. 1. Amsterdam: Elsevier Science.
- Ebrahimi, N., Hamedani, G. G., Soofi, E. S., Volkmer, H. (2010). A class of models for uncorrelated random variables. *Journal of Multivariate Analysis* 101:1859–1871.

- Eicher, T. S., Lenkoski, A., Raftery, A. E. (2009). Bayesian model averaging and endogeneity under model uncertainty: An application to development determinants. *Technical report*, University of Washington.
- Escobar, M. D., West, M. (1995). Bayesian density estimation and inference using mixtures. *Journal of the American Statistical Association* 90:577–588.
- Escobar, M. D., West, M. (1998). Computing non-parametric hierarchical models. In: Dey, D., Müller, P., and Sinha, D., eds. *Practical Nonparametric and Semiparametric Bayesian Statistics*. New York: Springer, pp. 1–22.
- Florens, J.-P., Simoni, A. (2010). Nonparametric estimation of an instrumental regression: A quasi-Bayesian approach based on regularized posterior. *Technical report*, Toulouse School of Economics.
- Gamerman, D., Lopes, H. F. (2006). *Markov Chain Monte Carlo: Stochastic Simulation for Bayesian Inference*. Baton-Rouge: Chapman & Hall/CRC.
- Gelman, A., Jakulin, A., Pittau, M. G., Su, Y.-S. (2008). A weakly informative default prior distribution for logistic and other regression models. *The Annals of Applied Statistics* 2:1360–1383.
- Geweke, J. (1996). Bayesian reduced rank regression. *Journal of Econometrics* 75:121–146.
- Hoogerheide, L. F., Kaashoek, J. F., van Dijk, H. K. (2007). On the shape of posterior densities and credible sets in instrumental variable regression models with reduced rank. *Journal of Econometrics* 139:154–180.
- Hoogerheide, L. F., Kleibergen, F., van Dijk, H. K. (2008). Natural conjugate priors for the instrumental variables regression model applied to the Angrist-Krueger data. *Journal of Econometrics* 138:63–103.
- Hoogerheide, L. F., van Dijk, H. K. (2008a). Possibly Ill-behaved posteriors in econometric models. *Technical report*, Tinbergen Institute.
- Hoogerheide, L. F., van Dijk, H. K. (2008b). Simulation-based Bayesian econometric inference. *Technical report*, Tinbergen Institute.
- Kleibergen, F., Paap, R. (2002). Priors, posteriors, and Bayes factors for a Bayesian analysis of co-integration. *Journal of Econometrics* 111:223–249.
- Kleibergen, F., van Dijk, H. K. (1998). Bayesian simultaneous equations analysis using reduced rank structures. *Econometric Theory* 14:701–743.
- Kleibergen, F., van Dijk, H. K. (2007). Natural conjugate priors for the instrumental variables regression model applied to the Angrist-Krueger data. *Journal of Econometrics* 138:63–103.
- Kleibergen, F., Zivot, E. (2003). Bayesian and Classical approaches to Instrumental Variable regression. *Journal of Econometrics* 114:29–72.
- Koop, G., Leon-Gonzalez, R., Strachan, R. (2010). Efficient posterior simulation for cointegrated models with priors on the cointegration space. *Econometric Reviews* 29:224–242.
- Koop, G., Leon-Gonzalez, R., Strachan, R. (2011). Bayesian Model averaging in the Instrumental variable regression model. *Technical report*, University of Strathclyde.
- Lancaster, T. (2004). *An Introduction to Modern Bayesian Econometrics*. Blackwell Publishing.
- Lindley, D. V., El Sayyad, G. M. (1968). The Bayesian estimation of a linear functional relationship. *Journal of Royal Statistical Society, Series B* 30:190–202.
- Lopes, H. F., Carvalho, C. M., Polson, N. G., Johannes, M. (2011a). Particle learning for sequential Bayesian computation (with discussion). In: Bernardo, J. M., Bayarri, M. J., Berger, J. O., Dawid, A. P., Heckerman, D., Smith, A. F. M., West, M., eds. *Bayesian Statistics*. Vol. 9, pp. 317–360.
- Lopes, H. F., McCulloch, R. E., Tsay, R. E. (2011b). Cholesky stochastic volatility. *Technical report*, The University of Chicago Booth School of Business.
- Lopes, H. F., Polson, N. G. (2011). Particle learning for fat-tailed distributions. *Technical report*, The University of Chicago Booth School of Business.
- Maddala, G. S. (1976). Weak priors and sharp posteriors in simultaneous equation models. *Econometrica* 44:345–351.
- Minka, T. (1999). Linear regression with errors in both variables: a proper bayesian approach. *MIT Media Lab Note* (10/8/99).
- Pourahmadi, M. (1999). Joint mean-covariance models with applications to longitudinal data: unconstrained parameterisation. *Biometrika* 86:677–690.

- Rossi, P. E., Allenby, G. M., McCulloch, R. E. (2005). *Bayesian Statistics and Marketing*. New York: Wiley.
- Sims, C. (2007). Thinking about instrumental variables. *Technical report*, Princeton University.
- Stock, J. H., Watson, M. W. (2003). *Introduction to Econometrics*. Boston: Addison Wesley.
- Strachan, R. W. (2003). Valid Bayesian estimation of the cointegrating error correction model. *Journal of Business & Economic Statistics* 21:185–195.
- Villani, M. (2005). Bayesian reference analysis of cointegration. *Econometric Theory* 21:326–357.
- West, M. (1984). Outlier models and prior distributions in Bayesian linear regression. *Journal of Royal Statistical Society, Series B* 46:431–439.
- Zellner, A. (1971). *An Introduction to Bayesian Inference in Econometrics*. New York: Wiley.
- Zellner, A., Vandaale, W. (1975). Bayes-Stein estimators for  $k$ -means, regression and simultaneous equation models. In: Fienberg, S. E., Zellner, A., eds. *Studies in Bayesian Econometrics and Statistics*. Amsterdam: North-Holland, pp. 627–653.
- Zellner, A., Bauwens, L., van Dijk, H. K. (1988). Bayesian specification analysis and estimation of simultaneous equation models using Monte Carlo methods. *Journal of Econometrics* 38:39–72.
- Zellner, A., Tobias, J., Ryu, H. K. (1997). Bayesian method of moments (BMOM) analysis of parametric and semi-parametric regression models. *Technical report*, University of Chicago.