

GENERALIZED METHODS FOR RESTRICTED MARKOV-SWITCHING MODELS WITH INDEPENDENT STATE VARIABLES

CHRISTOPHER A. SIMS, DANIEL F. WAGGONER, AND TAO ZHA

ABSTRACT. We generalize the existing Bayesian method for estimating Markov-switching models by allowing for independent Markov processes and various restrictions on transition matrices. This generalized method is used to develop tools for estimating both identified and reduced-form multivariate dynamic models. Moreover, we develop (1) an efficient block-wise optimization algorithm for obtaining the maximum likelihood or posterior estimates, (2) a new modified harmonic means method to deal with the common situation in which the likelihood is zero and a Gaussian approximation is inadequate, and (3) an estimation method for Markov-switching VARs with restrictions on both time variation and the lag structure.

I. INTRODUCTION

This paper extends the methods of Hamilton (1989), Chib (1996), and Kim and Nelson (1999) to Markov-switching models with independent state variables and linear restrictions on transition matrices and develops generalized methods and efficient algorithms for both estimation and inference of such models.

We apply this generalized methodology to both identified and unrestricted vector autoregression (VAR) models. The application allows for four key elements: (1) simultaneity, (2) over-identifying restrictions on both contemporaneous coefficients as well as the lag structure, (3) separating changes in structural residual variances from changes in coefficients of structural equations, and (4) separating changes in coefficients in one structural equation (e.g., monetary policy) from those in another equation (e.g., fiscal policy). Our framework is particularly useful in addressing questions related to the current debate on whether monetary policy and the private sector's behavior have significantly changed in recent history.¹

Date: August 14, 2006.

Key words and phrases. Volatility, coefficient changes, discontinuous shifts, Lucas critique, independent Markov processes.

We thank Tim Cogley for helpful comments. Eric Wang provided excellent research assistance in computation on the Linux operating system. We acknowledge the technical support on parallel and grid computation from Computing College of Georgia Institute of Technology. The views expressed herein do not necessarily reflect those of the Federal Reserve Bank of Atlanta or the Federal Reserve System.

¹For this debate, consult Cogley and Sargent (2002), Canova and Gambetti (2004), Beyer and Farmer (2004), Cogley and Sargent (2005), Primiceri (2005), and Sims and Zha (2006).

There are two common problems associated with applied Bayesian analysis in the recent macroeconomic literature. First, the Markov-Chain Monte-Carlo (MCMC) algorithm often begins with an arbitrary starting point without searching for the maximum likelihood estimate (MLE) or the estimate at the peak of the posterior density function. In high-dimensional multivariate dynamic models, however, an arbitrary starting point is likely to be in the extremely low probability region and the posterior draws simulated from the MCMC algorithm may get stuck in this region. Our block-wise optimization method is designed to help find efficiently the MLE or the posterior estimate for a complicated dynamic model.

The second problem is that when posterior odds ratios are reported in the macroeconomic literature, there often has been no diagnostic analysis of how accurate the computed odds ratios are. For example, the likelihood for most dynamic stochastic general equilibrium (DSGE) models has zero values in the interior points of the parameter space. Therefore, the modified harmonic means (MHM) method widely used in the macroeconomic literature is likely to produce a very inaccurate estimate of the marginal data density. Our new way of implementing the MHM method is designed to deal with this problem explicitly.

The rest of the paper is organized as follows. Section II develops a generalized Bayesian method for estimating Markov-switching models with independent state variables and linear restrictions on transition matrices. Based on this generalized method, Section III develops tools for estimation and inference of both identified and unrestricted vector autoregression (VAR) models. In Section IV, we develop a block-wise optimization method for estimating the Markov-switching models. This method proves much more efficient than the existing expectation-maximization (EM) algorithm. In Section V, we develop a new implementation of the MHM method. A three-variable VAR application to the post-war US data is presented in Section VI. The conclusion is given in Section VII.

II. MARKOV-SWITCHING MODEL

II.1. Distributional assumptions. Let $(Y_t, Z_t, \theta, Q, S_t)$ be a collection of random variables where

$$\begin{aligned} Y_t &= (y_1, \dots, y_t) \in (\mathbb{R}^n)^t, \\ Z_t &= (z_1, \dots, z_t) \in (\mathbb{R}^m)^t, \\ \theta &= (\theta_i)_{i \in H} \in (\mathbb{R}^r)^h, \\ Q &= (q_{i,j})_{(i,j) \in H \times H} \in \mathbb{R}^{h^2}, \\ S_t &= (s_0, \dots, s_t) \in H^{t+1}, \\ S_{t+1}^T &= (s_{t+1}, \dots, s_T) \in H^{T-t}, \end{aligned}$$

and H is a finite set with h elements and is usually taken to be the set $\{1, \dots, h\}$. The vector y_t contains the endogenous variables and the vector z_t contains the exogenous variables. Our analysis, however, encompass the case in which there are no exogenous variables. The matrix Q is a Markov transition matrix and $q_{i,j}$ is the probability that s_t is equal to i

given that s_{t-1} is equal to j . The matrix Q is restricted to satisfy

$$q_{i,j} \geq 0 \text{ and } \sum_{i \in H} q_{i,j} = 1.$$

We shall follow the convention that if u and v are random vectors for which a density function exists, $p(u, v)$ denotes the density function. The marginal and conditional density functions are expressed as

$$p(v) = \int p(u, v) du,$$

and

$$p(u | v) = \frac{p(u, v)}{\int p(u, v) du}.$$

We assume that $p(u, v)$ is integrable. Hence, $p(u | v)$ and $p(v)$ will exist for almost all v . The objects θ and Q are parameters, Y_t and Z_t are observed data, and S_t can be considered either a sequence of unobserved variables or a vector of nuisance parameters. We assume that $(Y_t, Z_t, \theta, Q, S_t)$ has a joint density function $p(Y_t, Z_t, \theta, Q, S_t)$, where we use the Lebesgue measure² on $(\mathbb{R}^n)^t \times (\mathbb{R}^m)^t \times (\mathbb{R}^r)^h \times \mathbb{R}^{h^2}$ and the counting measure on H^{t+1} . This density satisfies the following conditions.

Condition 1.

$$p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q, S_{t-1}) = q_{s_t, s_{t-1}}$$

for $t > 0$.

Condition 2.

$$p(y_t | Y_{t-1}, Z_t, \theta, Q, S_t) = p(y_t | Y_{t-1}, Z_t, \theta, s_t)$$

for $t > 0$.

Condition 3.

$$p(z_t | Y_{t-1}, Z_{t-1}, \theta, Q, S_t) = p(z_t | Y_{t-1}, Z_{t-1}).$$

Condition 1 states formally that the sequence S_t evolves according to an exogenous Markov process with the transition matrix Q . Condition 2 is needed for obtaining a standard posterior density function of Q conditional on S_T .³ Condition 3 ensures that z_t is an exogenous variable.

II.2. Propositions. From Conditions 1 - 3, one can prove the following propositions (the proofs can be found in Hamilton (1989), Chib (1996), and Kim and Nelson (1999)). These propositions are used throughout the rest of this paper.

Proposition 1.

$$p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q) = \sum_{s_{t-1} \in H} q_{s_t, s_{t-1}} p(s_{t-1} | Y_{t-1}, Z_{t-1}, \theta, Q)$$

for $t > 0$.

²Instead of the Lebesgue measure, any sigma finite measure on \mathbb{R}^n and \mathbb{R}^m can be used as long as the product measure is used on $(\mathbb{R}^n)^t$ and $(\mathbb{R}^m)^t$.

³This tractable result no longer holds for most regime-switching rational expectations models (Farmer, Waggoner, and Zha, 2006). In that case, the Metropolis algorithm may be used instead.

Proposition 2.

$$p(s_t | Y_t, Z_t, \theta, Q) = \frac{p(y_t | Y_{t-1}, Z_t, \theta, s_t) p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q)}{\sum_{s_{t-1} \in H} p(y_t | Y_{t-1}, Z_t, \theta, s_t) p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q)}$$

for $t > 0$.

Proposition 3.

$$p(s_t | Y_t, Z_t, \theta, Q, s_{t+1}) = p(s_t | Y_T, Z_T, \theta, Q, s_{t+1}^T)$$

for $0 \leq t < T$.

Proposition 4.

$$p(y_t, z_t | Y_{t-1}, Z_{t-1}, \theta, Q, S_T) = (y_t, z_t | Y_{t-1}, Z_{t-1}, \theta, Q, S_t)$$

for $0 < t \leq T$.

II.3. Independent state variables and restrictions on Q . An important part of this paper is to consider ℓ independent state variables s_t^1, \dots, s_t^ℓ . Let $h = \prod_{k=1}^\ell h_k$, $H = \prod_{k=1}^\ell H_k$ where $H_k = \{1, \dots, h_k\}$, and $s_t = (s_t^1, \dots, s_t^\ell)$ where $s_t^k \in H_k$. The transition matrix Q is of the form

$$Q = Q_1 \otimes \dots \otimes Q_\ell$$

where $Q_k = (q_{i,j}^k)$ is an $h_k \times h_k$ matrix such that

$$q_{i,j}^k \geq 0 \text{ and } \sum_{i \in H_k} q_{i,j}^k = 1, \text{ for every } j \in H_k.$$

The tensor product representation of Q implies that if $i = (i_1, \dots, i_\ell) \in H$ and $j = (j_1, \dots, j_\ell) \in H$, then $q_{i,j} = \prod_{k=1}^\ell q_{i_k, j_k}^k$. Conditional on Q , the Markov process s_t consists of ℓ independent Markov processes s_t^k .

We wish to impose a wide range of linear restrictions on Q while maintaining the property that the posterior distribution for each column of parameters is of Dirichlet. For notational clarity, in the following analysis we suppress both the subscript k and the superscript k that indicate a particular independent Markov process under consideration.

For $1 \leq j \leq h$, let M_j be an $h \times o_j$ matrix ($o_j \leq h$) such that all the elements of M_j are non-negative, there is at most one positive element in each row of M_j , and the sum of the elements in each column of M_j equals one. Denote $w_j = [w_{1,j}, \dots, w_{o_j,j}]' \in \mathbb{R}^{o_j}$ with

$$w_{i,j} \geq 0 \text{ and } \sum_{i=1}^{o_j} w_{i,j} = 1.$$

Linear restrictions on the j^{th} column of Q can be expressed in the form

$$q_j = M_j w_j.$$

It is straightforward to verify that the conditions on M_j and w_j guarantee that

$$q_{i,j} \geq 0 \text{ and } \sum_{i=1}^h q_{i,j} = 1.$$

II.4. Examples. Our class of restrictions on Q includes most examples discussed in the literature. For example, Sims (1999) discusses a structural break with an irreversible regime change. In a two-state case where the second state is absorbing or irreversible, we have

$$M_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

In general, exclusion restrictions of the form $q_{i,j} = 0$ require that the i^{th} row of M_j be zero.

As a second example, a symmetric jumping among states considered by Sims (2001) introduces a parsimonious parameterization of Q that may be used to avoid over-parameterization.

The transition matrix studied by Sims (2001) has the following form

$$Q = \begin{bmatrix} \pi_1 & (1 - \pi_2)/2 & 0 \\ 1 - \pi_1 & \pi_2 & 1 - \pi_3 \\ 0 & (1 - \pi_2)/2 & \pi_3 \end{bmatrix}, \quad (1) \text{?eqn:Q3states?}$$

where π_1 , π_2 , and π_3 are free parameters to be estimated. These linear restrictions can be expressed as

$$M_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, M_2 = \begin{bmatrix} 0 & 1/2 \\ 1 & 0 \\ 0 & 1/2 \end{bmatrix}, M_3 = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

A third example pertains to incremental changes in the model parameters (Cogley and Sargent, 2005).⁴ This kind of parameter drifts can be approximated arbitrarily well by expanding the number of states while containing the elements of Q in a much smaller dimension. Our approach has advantage over that of Cogley and Sargent (2005) because it allows for occasional discontinuous shifts in regime as well as frequent, incremental changes in parameters. One way to achieve this task is to concentrate weight on the diagonal of Q (Zha, In press). Specifically, one can express incremental increases and discontinuous jumps among $n + 1$ states as

$$Q = \begin{bmatrix} \pi_1 & \alpha_2(1 - \pi_2) & \dots & \alpha_{n+1}^n(1 - \pi_{n+1}) \\ \alpha_1(1 - \pi_1) & \pi_2 & \dots & \alpha_{n+1}^{n-1}(1 - \pi_{n+1}) \\ \alpha_1^2(1 - \pi_1) & \alpha_2(1 - \pi_2) & \dots & \alpha_{n+1}^{n-2}(1 - \pi_{n+1}) \\ \dots & \dots & \dots & \dots \\ \alpha_1^n(1 - \pi_1) & \alpha_2^{n-1}(1 - \pi_2) & \dots & \pi_{n+1} \end{bmatrix},$$

where π_i is a free parameter and $0 < \alpha_i < 1$ is taken as a given. The restrictions can be written as

$$M_1 = \begin{bmatrix} 1 & 0 \\ 0 & \alpha_1 \\ 0 & \alpha_1^2 \\ \dots & \dots \\ 0 & \alpha_1^n \end{bmatrix}, M_2 = \begin{bmatrix} 0 & \alpha_2 \\ 1 & 0 \\ 0 & \alpha_2 \\ \dots & \dots \\ 0 & \alpha_2^{n-1} \end{bmatrix}, \dots, M_{n+1} = \begin{bmatrix} 0 & \alpha_{n+1}^n \\ 0 & \alpha_{n+1}^{n-1} \\ 0 & \alpha_{n+1}^{n-2} \\ \dots & \dots \\ 1 & 0 \end{bmatrix},$$

where the elements in each column of M_i must sum up to 1.

⁴See also Sims (1993); Cogley and Sargent (2002); Stock and Watson (2003); Canova and Gambetti (2004); Primiceri (2005).

II.5. Prior. We begin with the Dirichlet prior on Q_k where $k \in \{1, \dots, \ell\}$ indicates the k^{th} state variable. For $\alpha_{i,j}^k > 0$ for every $i, j \in H_k \times H_k$ where H_k has h_k elements, the prior on Q_k is assumed to be

$$p(Q_k) = \prod_{j \in H_k} \left(\frac{\Gamma\left(\sum_{i \in H_k} \alpha_{i,j}^k\right)}{\prod_{i \in H_k} \Gamma\left(\alpha_{i,j}^k\right)} \right) \times \prod_{(i,j) \in H_k \times H_k} \left(q_{i,j}^k\right)^{\alpha_{i,j}^k - 1}.$$

The prior on Q , denoted by $p(Q)$, follows directly from the tensor relationship between Q and Q_k and from the fact that $h = \prod_{k=1}^{\ell} h_k$.

Suppressing both the superscript k and the subscript k , one can obtain the prior on w_j from the mapping from w_j to q_j :

$$p(w_j) \propto \prod_{i_w=1}^{o_j} (w_{i_w,j})^{\beta_{i_w,j} - 1} \quad (2) \text{?eqn:priorofbj?}$$

where w_j 's are independent of one another and $\beta_{i_w,j}$ is defined by

$$\beta_{i_w,j} = 1 + \sum_{i \in \{i: M_j(i, i_w) > 0\}} (\alpha_{i,j} - 1),$$

where $i \in \{1, \dots, h\}$ and $M_j(i, i_w)$ is the element at the i^{th} row and the i_w^{th} column of M_j .

The joint prior density for θ, Q, S_T is

$$p(\theta, Q, S_T) = p(\theta, Q) p(s_0 | \theta, Q) \prod_{t=1}^T p(s_t | \theta, Q, S_{t-1})$$

By Condition 1, $p(s_t | \theta, Q, S_{t-1}) = q_{s_t, s_{t-1}}$. We assume that the prior on θ is independent of the prior on Q and that $p(s_0 | \theta, Q) = \frac{1}{h}$ for every $s_0 \in H$.⁵ The resulting prior has the following form

$$p(\theta, Q, S_T) = \frac{p(\theta) p(Q)}{h} \prod_{t=1}^T q_{s_t, s_{t-1}}. \quad (3) \text{?eqn:prior?}$$

II.6. Likelihood. Using Proposition 4 and Conditions 2 and 3, one can show that the joint density of Y_T and Z_T conditional on θ and Q is

$$p(Y_T, Z_T | \theta, Q) = \prod_{t=1}^T p(y_t, z_t | Y_{t-1}, Z_{t-1}, \theta, Q)$$

Note

$$\begin{aligned} p(y_t, z_t | Y_{t-1}, Z_{t-1}, \theta, Q) &= \sum_{s_t \in H} p(y_t, z_t, s_t | Y_{t-1}, Z_{t-1}, \theta, Q) \\ &= \sum_{s_t \in H} p(y_t, z_t | Y_{t-1}, Z_{t-1}, \theta, Q, s_t) p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q) \end{aligned}$$

⁵The conventional assumption for $p(s_0 | \theta, Q)$ is the ergodic distribution of Q , if it exists. This convention, however, makes the resulting conditional posterior distribution of Q an unknown and complicated distribution.

and

$$\begin{aligned} p(y_t, z_t | Y_{t-1}, Z_{t-1}, \theta, Q, s_t) \\ &= p(y_t | Y_{t-1}, Z_t, \theta, Q, s_t) p(z_t | Y_{t-1}, Z_{t-1}, \theta, Q, s_t) \\ &= p(y_t | Y_{t-1}, Z_t, \theta, s_t) p(z_t | Z_{t-1}), \end{aligned}$$

it follows that

$$\begin{aligned} p(Y_T, Z_T | \theta, Q) &= \prod_{t=1}^T p(z_t | Z_{t-1}) \prod_{t=1}^T \left[\sum_{s_t \in H} p(y_t | Y_{t-1}, Z_t, \theta, s_t) p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q) \right] \\ &= p(Z_T) \prod_{t=1}^T \left[\sum_{s_t \in H} p(y_t | Y_{t-1}, Z_t, \theta, s_t) p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q) \right] \end{aligned}$$

Conditional on the vector of exogenous variables Z_t , the likelihood of Y_T is

$$p(Y_T | Z_T, \theta, Q) = \prod_{t=1}^T \left[\sum_{s_t \in H} p(y_t | Y_{t-1}, Z_t, \theta, s_t) p(s_t | Y_{t-1}, Z_{t-1}, \theta, Q) \right] \quad (4) \text{?eqn:likelihood}$$

This likelihood can be evaluated recursively, using Propositions 1 and 2.

II.7. Posterior distribution. By the Bayes rule, it follows from (3) and (4) that the posterior distribution of (θ, Q) is

$$p(\theta, Q | Y_T, Z_T) \propto p(\theta, Q) p(Y_T | Z_T, \theta, Q). \quad (5) \text{?eqn:posterior?}$$

The posterior density $p(\theta, Q | Y_T, Z_T)$ is unknown and complicated; the MCMC simulation directly from this distribution can be inefficient and problematic. One can, however, use the idea of Gibbs sampling to obtain the empirical joint posterior density $p(\theta, Q, S_T | Y_T, Z_T)$ by sampling alternately from the following conditional posterior distributions:

$$\begin{aligned} p(S_T | Y_T, Z_T, \theta, Q), \\ p(Q | Y_T, Z_T, S_T, \theta), \\ p(\theta | Y_T, Z_T, Q, S_T). \end{aligned}$$

Simulation from the conditional posterior density $p(\theta | Y_T, Z_T, Q, S_T)$ is model-dependent, which we will discuss in Section III. In this section we study the first two conditional posterior distributions.

II.7.1. Conditional posterior distribution of S_T . The distribution of S_T conditional on Y_T, Z_T, θ , and Q is

$$\begin{aligned} p(S_T | Y_T, Z_T, \theta, Q) &= p(s_T | Y_T, Z_T, \theta, Q) p(S_{T-1} | Y_T, Z_T, \theta, Q, S_T^T) \\ &= p(s_T | Y_T, Z_T, \theta, Q) \prod_{t=0}^{T-1} p(s_t | Y_T, Z_T, \theta, Q, S_{t+1}^T) \end{aligned}$$

where $S_{t+1}^T = \{s_{t+1}, \dots, s_T\}$. From Proposition 3,

$$\begin{aligned} p(s_t | Y_T, Z_T, \theta, Q, S_{t+1}^T) &= p(s_t | Y_t, Z_t, \theta, Q, s_{t+1}) \\ &= \frac{p(s_t, s_{t+1} | Y_t, Z_t, \theta, Q)}{p(s_{t+1} | Y_t, Z_t, \theta, Q)} \\ &= \frac{p(s_{t+1} | Y_t, Z_t, \theta, Q, s_t) p(s_t | Y_t, Z_t, \theta, Q)}{p(s_{t+1} | Y_t, Z_t, \theta, Q)} \\ &= \frac{q_{s_{t+1}, s_t} p(s_t | Y_t, Z_t, \theta, Q)}{p(s_{t+1} | Y_t, Z_t, \theta, Q)} \end{aligned}$$

The conditional density $p(s_t | Y_T, Z_T, \theta, Q, S_{t+1}^T)$ is straightforward to evaluate according to Propositions 1 and 2. Starting with s_T and working backward, we can easily sample S_T from the posterior conditional on Y_T, Z_T, θ, Q by using the following fact

$$\begin{aligned} p(s_t | Y_T, Z_T, \theta, Q) &= \sum_{s_{t+1} \in H} p(s_t, s_{t+1} | Y_T, Z_T, \theta, Q) \\ &= \sum_{s_{t+1} \in H} p(s_t | Y_T, Z_T, \theta, Q, s_{t+1}) p(s_{t+1} | Y_T, Z_T, \theta, Q) \\ &= \sum_{s_{t+1} \in H} p(s_t | Y_t, Z_t, \theta, Q, s_{t+1}) p(s_{t+1} | Y_T, Z_T, \theta, Q). \end{aligned}$$

Note that this density can also be evaluated recursively.

II.7.2. Conditional posterior distribution of Q_k . The conditional posterior density of Q derives directly from the conditional posterior density of the free parameters w_j .⁶ It follows from Condition 1 and the prior (2) that

$$p(w_j | Y_T, Z_T, \theta, S_T) \propto \prod_{i_w=1}^{o_j} (w_{i_w, j})^{n_{i_w, j} + \beta_{i_w, j} - 1}$$

where $n_{i_w, j}$ is the number of transitions from $s_{t-1} = j$ to $s_t \in \{s_t : M_j(s_t, i_w) > 0\}$.

III. STRUCTURAL VAR MODELS

The methodology developed so far is used by Rubio-Ramírez, Waggoner, and Zha (2005) and Sims and Zha (2006) to study a class of simultaneous-equation multivariate dynamic models that are commonly used for policy analysis. In this section, we develop and detail the econometric methods specific to these types of models.

III.1. Likelihood. We consider a class of models of the following form:

$$y_t' A(s_t) = \sum_{i=1}^p y_{t-i}' A_i(s_t) + z_t' C(s_t) + \varepsilon_t' \Xi^{-1}(s_t), \text{ for } 1 \leq t \leq T, \quad (6) \text{?eqn:structural}$$

where

- p is a lag length;
- y_t is an n -dimensional column vector of endogenous variables at time t ;

⁶To be consistent with Section II.5, we suppress both the superscript k and the subscript k that indicate a particular Markov process under study.

- z_t is an m -dimensional column vector of exogenous and deterministic variables at time t ;
- ε_t is an n -dimensional column vector of unobserved random shocks at time t ;
- For $1 \leq k \leq h$, $A(k)$ is an invertible $n \times n$ matrix and $A_t(k)$ is an $n \times n$ matrix;
- For $1 \leq k \leq h$, $C(k)$ is an $m \times n$ matrix;
- For $1 \leq k \leq h$, $\Xi(k)$ is an $n \times n$ diagonal matrix.

For the rest of the paper we take the initial conditions y_0, \dots, y_{1-p} as given. Let

$$x_t = \begin{bmatrix} y_{t-1} \\ \vdots \\ y_{t-p} \\ z_t \end{bmatrix} \quad \text{and} \quad F(s_t) = \begin{bmatrix} A_1(s_t) \\ \vdots \\ A_p(s_t) \\ C(s_t) \end{bmatrix}.$$

Then (6) can be written in the compact form:

$$y'_t A(s_t) = x'_t F(s_t) + \varepsilon'_t \Xi^{-1}(s_t), \text{ for } 1 \leq t \leq T \quad (7) \text{?eqn:compactstr}$$

We introduce the following notation that will be used repeatedly in this paper:

$$A = (A(1), \dots, A(h)), \quad F = (F(1), \dots, F(h)), \quad \Xi = (\Xi(1), \dots, \Xi(h)), \\ \theta = (A, F, \Xi),$$

$$Y_t = \begin{bmatrix} y'_1 \\ \vdots \\ y'_t \end{bmatrix}_{t \times n}, \quad Z_t = \begin{bmatrix} z'_1 \\ \vdots \\ z'_t \end{bmatrix}_{t \times k}, \quad S_t = \begin{bmatrix} s_0 \\ \vdots \\ s_t \end{bmatrix}_{(t+1) \times 1}.$$

We assume that

$$p(\varepsilon_t | Y_{t-1}, Z_t, S_T, \theta, Q) = \text{normal}(\varepsilon_t | \mathbf{0}, I_n),$$

where $\mathbf{0}$ denotes a vector or matrix of zeros, I_n denotes the $n \times n$ identity matrix, and $\text{normal}(x | \mu, \Sigma)$ denotes the multivariate normal distribution of x with mean μ and variance Σ .⁷ This assumption is equivalent to

$$p(y_t | Y_{t-1}, Z_t, S_t, \theta, Q) = \text{normal}(y_t | \mu_t(s_t), \Sigma(s_t)) \quad (8) \text{?eqn:likelihood}$$

where

$$\mu_t(k) = (F(k)A^{-1}(k))' x_t$$

and

$$\Sigma(k) = (A(k)\Xi^2(k)A'(k))^{-1}$$

Let $a_j(k)$ be the j^{th} column of $A(k)$, $f_j(k)$ be the j^{th} column of $F(k)$, and $\xi_j(k)$ be the j^{th} diagonal element of $\Xi(k)$. Define

$$a(k) = \begin{bmatrix} a_1(k) \\ \vdots \\ a_n(k) \end{bmatrix}_{n^2 \times 1}, \quad f(k) = \begin{bmatrix} f_1(k) \\ \vdots \\ f_n(k) \end{bmatrix}_{(pn+m)n \times 1}, \quad \text{and} \quad \xi(k) = \begin{bmatrix} \xi_1(k) \\ \vdots \\ \xi_n(k) \end{bmatrix}_{n \times 1}$$

⁷The matrix Σ must be symmetric and non-negative semi-definite.

It follows from (8) that

$$\begin{aligned} p(y_t | Y_{t-1}, Z_t, S_t, \theta, Q) &= |\Sigma(s_t)|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (y_t - \mu(s_t))' \Sigma^{-1}(s_t) (y_t - \mu(s_t))\right) \\ &= |A(s_t) \Xi(s_t)| \exp\left(-\frac{1}{2} (y_t' A(s_t) - x_t' F(s_t)) \Xi^2(s_t) (A'(s_t) y_t - F'(s_t) x_t)\right) \\ &= |A(s_t)| \prod_{j=1}^n |\xi_j(s_t)| \exp\left(-\frac{\xi_j^2(s_t)}{2} (y_t' a_j(s_t) - x_t' f_j(s_t))^2\right). \end{aligned}$$

We consider the case where the state variable $s_t = [s_{1t} \ s_{2t}]$ is a composite one such that either $s_{1t} = s_{2t}$ or s_{1t} and s_{2t} are independent random variables. The analytical results for more complicated cases will follow directly. We let a_j and f_j depend on s_{1t} and ξ_j depend on s_{2t} . Thus, the conditional likelihood function $p(y_t | Y_{t-1}, Z_t, S_t, \theta, Q)$ is equal to

$$|A(s_{1t})| \prod_{j=1}^n |\xi_j(s_{2t})| \exp\left(-\frac{\xi_j^2(s_{2t})}{2} (y_t' a_j(s_{1t}) - x_t' f_j(s_{1t}))^2\right). \quad (9) \text{?eqn:likelihood}$$

Given (9), the likelihood of Y_T can be formed by following (4).

III.2. A priori restrictions.

III.2.1. *Restrictions on time variation.* If we let all parameters vary across states, the number of free parameters in the model becomes impractically high when the system of equations is large or the lag length is long. For a typical quarterly model with 5 lags and 6 endogenous variables, for example, the number of parameters in $F(s_{1t})$ is of order 180 for each state. Given the post-war macroeconomic data, however, it is not uncommon to have some states lasting for only a few years and thus the number of associated observations is far less than 180 quarters. It is therefore essential to simplify the model by restricting the degree of time variation in the model's parameters. Such a restriction entails complexity and difficulties that have not been dealt with in the simultaneous-equation literature.

To begin with, we rewrite F as

$$F(s_{1t}) = G(s_{1t}) + \bar{S} A(s_{1t}). \quad (10) \text{?eqn:FG?}$$

$\begin{matrix} m \times n \\ m \times n \\ m \times n \\ n \times n \end{matrix}$

where

$$\bar{S} = \begin{bmatrix} \mathbf{I}_n \\ \mathbf{0} \\ (m-n) \times n \end{bmatrix}.$$

We let G be a collection of all $G(k)$ for $k = 1, \dots, h_1$. If we place a prior distribution on $G(s_{1t})$ that has mean zero, the specification of \bar{S} is consistent with the reduced-form random walk feature implied by the existing Bayesian VAR models (Sims and Zha 1998). This type of prior tends to imply that greater persistence (in the sense of a tighter concentration of the prior on the random walk) is associated with smaller disturbance variances. This feature is reasonable, as it is consistent with the idea that beliefs about the unconditional variance of the data are *not* highly correlated with beliefs about the degree of persistence in the data.

Let $g_j(k)$ be the j^{th} column of $G(k)$. The time-variation restrictions imposed on $g_j(k)$ can be generally expressed by two components, one being time varying and the other being constant across states. Denote the first component by the $r_{g,j} \times 1$ vector $g_{\delta_j(k)}$ and the second component by the $h_1 r_{g,j} \times 1$ vector g_{ψ_j} , where the subscripts $\delta_j(k)$ and ψ_j will be discussed further in Section III.2.2. We express $g_j(k)$ for $k = 1, \dots, h_1$ in the form

$$\text{diag} \left([g_j(1)' \quad \dots \quad g_j(h_1)']' \right) = \text{diag} \left([g'_{\delta_j(1)} \quad \dots \quad g'_{\delta_j(h_1)}] \right) \text{diag} (g_{\psi_j}), \quad (11) \text{?eqn:gdecompose?}$$

where $\text{diag}(x)$ is the diagonal matrix with the diagonal being the column vector x . The long vector g_{ψ_j} is formed by stacking h_1 sub-vectors and the k^{th} sub-vector corresponds to the parameters in the k^{th} state.

In this paper we focus on the following three cases of restricted time variations for $a_j(s_{1t})$ and $g_j(s_{1t})$ for the j^{th} equation where $j \in \{1, \dots, n\}$, although our general method is capable of dealing with other time variation cases.

$$a_j(s_{1t}) \xi_j(s_{2t}), g_{ij,\ell}(s_{1t}) \xi_j(s_{2t}), c_j(s_{1t}) \xi_j(s_{2t}) = \begin{cases} a_j, g_{ij,\ell}, c_j & \text{Case I} \\ a_j \xi_j(s_{2t}), g_{ij,\ell} \xi_j(s_{2t}), c_j \xi_j(s_{2t}) & \text{Case II} \\ a_j(s_{1t}) \xi_j(s_{2t}), g_{\psi_{ij,\ell}} g_{\delta_{ij}(s_{1t})} \xi_j(s_{2t}), c_j(s_{1t}) \xi_j(s_{2t}) & \text{Case III} \end{cases}, \quad (12) \text{?eqn:tvcases?}$$

where $g_{ij,\ell}(s_{1t})$ is the element of $g_j(s_{1t})$ for the i^{th} variable at the ℓ^{th} lag and $c_j(s_{1t})$ is a vector of parameters corresponding to the exogenous variable z_t in equation j . The parameter $g_{\psi_{ij,\ell}}$ is the element of g_{ψ_j} for the i^{th} variable at the ℓ^{th} lag in any state; it is constant across states. The parameter $g_{\delta_{ij}(s_{1t})}$ is the element of $g_{\delta_j(s_{1t})}$ for the i^{th} variable in state s_{1t} at any lag. Thus, when the state s_{1t} changes, $g_{\delta_{ij}(s_{1t})}$ changes with variables but does not vary across lags. The variability across variables when the state changes is necessary to allow long run (policy) responses to vary over time, while the restriction on the time variation across lags is essential to prevent over-parameterization. The parameters a_j , $g_{ij,\ell}$, and c_j without the symbol (s_{1t}) mean that these parameters are restricted to be independent of state (i.e., constant across time).

In this setup, we include $c_j(k)$ in the stacked column vector g_{ψ_j} . In principle, one could include the time-varying parameter $c_j(k)$ as part of the time-varying component vector $g_{\delta_j(k)}$. With the normalization $c_j(1) = 1$, however, the likelihood function for $c_j(k)$ where $k \geq 2$ is so ill-behaved that our Gibbs sampler fails to work. Moreover, our reparameterization of grouping $c_j(k)$ in g_{ψ_j} preserves the prior correlations between $c_j(k)$ and the other lagged coefficients as implied by the Sims and Zha (1998) dummy-observation prior, an important part of the prior specification. It is important to note that the other elements of g_{ψ_j} are restricted to be invariant to state.

Case I represents a traditional constant-parameter VAR equation, which has been dealt with extensively in the literature and thus will not be a focal discussion of this paper. Case II represents a structural equation with time-varying disturbance variances only. In this case,

$\xi_j(s_{2t})$ measures volatility for the j^{th} structural equation. Case III represents a structural equation with time-varying coefficients.⁸

There are many applications that derive directly from various combinations of Case II and Case III for different equations. Some combinations, for example, enable one to distinguish regime shifts in policy behavior from their effects on private sector behavior — the practical lesson of the Lucas critique. The model with Case II for all equations suggests no structural break for both policy and the private sector; the model with Case II for the policy equation and Case III for all other equations hypothesizes that the policy rule is stable and structural breaks originate from the private sector. Both of these models, while consistent with rational expectations, take the view that the Lucas critique is unimportant in practice. On the other hand, the model with Case III for all equations is most consistent with the Lucas critique and if found to have a superior fit to the data, suggests that extrapolating the effects of policy changes from linear approximations may be misleading.⁹ The model with Case III for the policy equation and Case II for all other equations is an unconventional but quite interesting hypothesis. It is unconventional because it contradicts many theoretical examples delivered by rational expectations. Yet it implies that the Lucas critique may be practically unimportant because, despite regime shifts in policy, the private sector responds linearly to the history of policy variables.

III.2.2. *Identifying restrictions.* It is well known that the model (7) would be unidentified without further identifying restrictions. We follow the identified VAR literature and apply linear restrictions on A and F in the form of

$$\mathfrak{R}_j \begin{bmatrix} a_j \\ f_j \end{bmatrix} = 0, \quad (13) \text{?eqn:linearrest}$$

where \mathfrak{R}_j is an $(n + np + m) \times (n + np + m)$ and is not of full rank. Appendix A shows that the above restrictions are equivalent to the existence of an $n \times r_{b,j}$ matrix U_j with orthonormal columns, a $(pn + m) \times r_{g,j}$ matrix V_j with orthonormal columns, and a $(pn + m) \times n$ matrix \hat{W}_j with $V_j' \hat{W}_j = 0$ such that

$$a_j(k) = U_j b_j(k), \quad (14) \{?\}$$

$$f_j(k) = V_j g_j(k) - \hat{W}_j U_j b_j(k). \quad (15) \{?\}$$

The $r_{b,j} \times 1$ vector $b_j(k)$ and the $r_{g,j} \times 1$ vector $g_j(k)$ are free parameters to be estimated. If we replace \hat{W}_j in (15) with $W_j = \hat{W}_j + V_j \tilde{W}_j$ for any $r_{g,j} \times n$ matrix \tilde{W}_j , the underlying linear restrictions (13) will still hold, although $V_j' W_j \neq 0$ in general. For \bar{S} defined in (10), one can show that there exists \tilde{W}_j such that $W_j = \bar{S}$ where

$$\tilde{W}_j = V_j' (\bar{S} - \hat{W}_j).$$

⁸The reduced-form equation for Case III, however, has both time-varying coefficients and heteroscedastic disturbances. This feature reinforces the point that one should work directly on the structural form, not the reduced-form, of the model.

⁹Theoretical arguments for this view can be found in Cooley, LeRoy, and Raymon (1984), Sims (1987), and more recently Leeper and Zha (2003).

It follows from (9), (14), and (15) that

$$p(y_t | Y_{t-1}, Z_t, S_t, \theta, Q) = |A(s_{1t})| \left[\prod_{j=1}^n |\xi_j(s_{2t})| \exp \left(-\frac{\xi_j^2(s_{2t})}{2} ((y'_t + x'_t W_j) U_j b_j(s_{1t}) - x'_t V_j g_j(s_{1t}))^2 \right) \right]. \quad (16) \text{?eqn:Restricted}$$

In addition to the time-variation restrictions (12), the lag coefficient vector $g_j(k)$ for $k \in \{1, \dots, h_1\}$ may be further restricted. Specifically, one may impose linear restrictions directly on $g_{\delta_j(k)}$ and g_{ψ_j} through the affine transformation from $\mathbb{R}^{r_{\delta,j}}$ to $\mathbb{R}^{r_{g,j}}$

$$g_{\delta_j(k)} = \Delta_j \delta_j(k) + \bar{\delta}_j \quad (17) \text{?eqn:restriction}$$

and the affine transformation from $\mathbb{R}^{r_{\psi,j}}$ to $\mathbb{R}^{h_1 r_{g,j}}$

$$g_{\psi_j} = \Psi_j \psi_j, \quad (18) \text{?eqn:restriction}$$

where Δ_j is an $r_{g,j} \times r_{\delta,j}$ matrix, Ψ_j is an $h_1 r_{g,j} \times r_{\psi,j}$ matrix, $\bar{\delta}_j$ is an $r_{g,j} \times 1$ vector, $\delta_j(k)$ is an $r_{\delta,j} \times 1$ vector, and ψ_j is an $r_{\psi,j} \times 1$ vector. The vectors $\delta_j(k)$ and ψ_j are free parameters to be estimated, while the other vectors and matrices on the right hand sides of (17) and (18) are given by linear restrictions. We assume without loss of generality that Δ_j and Ψ_j have orthonormal columns so that both $\Delta_j' \Delta_j$ and $\Psi_j' \Psi_j$ are identity matrices.

Consider the most common situation in which the constant term is the only exogenous variable. As implied by (12), $r_{\delta,j}$ is much smaller than $r_{g,j}$ so that the time varying component has a small dimension. Similarly, the dimension $r_{\psi,j}$ is much smaller than $h_1 r_{g,j}$. For Case II, we set $\Delta_j = \mathbf{0}$ and $\bar{\delta}_j = \mathbf{1}$ where $\mathbf{1}$ denotes a vector or matrix of ones. In practice, therefore, there is no free parameter vector $\delta_j(k)$ to deal with. All the sub-vectors in g_{ψ_j} that correspond to different states are the same. Thus, the dimension $r_{\psi,j}$ is no greater than $r_{g,j}$. For Case III, we set

$$\bar{\delta}_j = \begin{bmatrix} \mathbf{0} \\ np \times 1 \\ 1 \end{bmatrix},$$

where the last element corresponds to the constant term in the j^{th} equation. The first np elements in the k^{th} sub-vector of g_{ψ_j} are restricted to be the same as those elements in any other sub-vector.

III.2.3. The prior. We begin with a prior imposed directly on $a_j(k)$, g_{ψ_j} , $\delta_j(k)$, and $\xi_j^2(k)$. The prior on the free parameters $b_j(k)$ and ψ_j will then be derived from the linear restrictions (14) and (18).

In order to use the reference prior in the VAR literature, we let the prior distributions of $a_j(k)$ and g_{ψ_j} take the Gaussian form:

$$p(a_j(k)) = \text{normal}(a_j(k) | \mathbf{0}, \bar{\Sigma}_{a_j}), \quad (19) \text{?eqn:aprior?}$$

$$p(g_{\psi_j}) = \text{normal}(g_{\psi_j} | \mathbf{0}, \tilde{\Sigma}_{g_{\psi_j}}), \quad (20) \text{?eqn:gpsiprior?}$$

for $k = 1, \dots, h_1$ and $j = 1, \dots, n$, where $\tilde{\Sigma}_{g_{\psi_j}} = I_{h_1} \otimes \tilde{\Sigma}_g$. The prior covariance matrices $\bar{\Sigma}_{a_j}$ and $\tilde{\Sigma}_g$ are the same as the prior covariance matrices specified by Sims and Zha (1998) for the contemporaneous and lagged coefficients in the constant-parameter VAR model.

Because these prior covariance matrices are the same across k , $a_j(k)$ has exactly the same prior distribution for different values of k so that k is essentially irrelevant for this prior.¹⁰ In other words, our prior is symmetric across states, for a priori knowledge of how they should differ is difficult to obtain through the prior distribution of this kind.

Following Sims and Zha (1998), we also incorporate into the model the $n + 1$ “dummy observations” formed from the initial observations as an additional part of the prior. These dummy observations, used as an additional prior component, express widely-held beliefs in unit roots and cointegration in macroeconomic series and play an indispensable role in improving out-of-sample forecast performance. Let Y_d be an $(n + 1) \times n$ matrix of dummy observations on the left hand side of system (7) and X_d be an $(n + 1) \times m$ matrix of dummy observations on the right hand side such that

$$Y_d A(k) = X_d (G_\psi + \bar{S}A(k)) + \tilde{E}_d, \quad (21) \text{?eqn:dummyprior}$$

where G_ψ is a $(pn + m) \times n$ matrix formed from g_{ψ_j} and \tilde{E}_d is an $(n + 1) \times n$ matrix of standard normal random variables. If we add the diffuse prior

$$p(\text{vec}(A(k))) \propto |A(k)|^{-(n+1)}$$

to correct the degrees of freedom for the overall prior of $A(k)$, it can be shown that combining the dummy prior (21) and the normal prior (19)-(20) leads to the following overall prior:¹¹

$$p(a_j(k)) = \text{normal}(a_j(k) | \mathbf{0}, \bar{\Sigma}_{a_j}), \quad (22) \text{?eqn:apriorfina}$$

$$p(g_{\psi_j}) = \text{normal}(g_{\psi_j} | \mathbf{0}, \bar{\Sigma}_{g_{\psi_j}}), \quad (23) \text{?eqn:gpsiprior}$$

where $\bar{\Sigma}_{g_{\psi_j}} = I_{h_1} \otimes \bar{\Sigma}_g$ and

$$\bar{\Sigma}_g = (X_d' X_d + \bar{\Sigma}_g^{-1})^{-1}.$$

¹⁰In our setup, the state variable s_{1t} for $A(s_{1t})$ and the state variable s_{2t} for $\Xi(s_{2t})$ are independently treated. In Sims and Zha (2006), the two state variables are the same. For the Case II model, therefore, $a_j(k)$ are restricted to be the same for all k 's under the Sims and Zha setup and we denote this vector by a_j^* . This restriction implies that the prior covariance matrix for a_j^* differs from $\bar{\Sigma}_{a_j}$. To see this point, consider two standard normal random variables x_1 and x_2 . With the restriction $x_1 = x_2$, one can show that

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix}' = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}' x^*,$$

where x^* is normally distributed with mean 0 and variance 2. Thus, the distribution of x^* is different from that of x_1 or x_2 . By analogy, $a_j(1)$ and $a_j(2)$ can be thought as x_1 and x_2 ; and a_j^* as x^* . For the examples we have studied, it turns out that the prior under our current setup gives a higher marginal data density with the hyperparameter values suggested by Sims and Zha (1998) and Robertson and Tallman (1999, 2001).

¹¹The proof follows directly from the fact (Sims and Zha, 1998) that

$$(X_d' X_d + \bar{\Sigma}_{g_{\psi_j}}^{-1})^{-1} (X_d' Y_d + \bar{\Sigma}_{g_{\psi_j}}^{-1} \bar{S}) = \bar{S},$$

$$Y_d' Y_d + \bar{\Sigma}_{a_j}^{-1} + \bar{S}' \bar{\Sigma}_{g_{\psi_j}}^{-1} \bar{S} - \bar{\Sigma}_{0j}^{-1} = \bar{\Sigma}_{a_j}^{-1},$$

where

$$\bar{\Sigma}_{0j}^{-1} = (Y_d' X_d + \bar{S}' \bar{\Sigma}_{g_{\psi_j}}^{-1}) (X_d' X_d + \bar{\Sigma}_{g_{\psi_j}}^{-1})^{-1} (X_d' Y_d + \bar{\Sigma}_{g_{\psi_j}}^{-1} \bar{S}).$$

Given the linear restrictions (14) and (18), one can derive from (22) and (23) that the implied prior distribution for $b_j(k)$ and ψ_j is

$$p(b_j(k)) = \text{normal}(b_j(k) | \mathbf{0}, \bar{\Sigma}_{b_j}), \quad (24) \text{?eqn:bprior?}$$

$$p(\psi_j) = \text{normal}(\psi_j | \mathbf{0}, \bar{\Sigma}_{\psi_j}), \quad (25) \text{?eqn:psiprior?}$$

where

$$\bar{\Sigma}_{b_j} = \left(U_j' \bar{\Sigma}_{a_j}^{-1} U_j \right)^{-1},$$

$$\bar{\Sigma}_{\psi_j} = \left(\Psi_j' \bar{\Sigma}_{g_{\psi_j}}^{-1} \Psi_j \right)^{-1}.$$

The prior distribution of $\delta_j(k)$ is assumed to be normal:

$$p(\delta_j(k)) = \text{normal}(\delta_j(k) | \mathbf{0}, \bar{\Sigma}_{\delta_j(k)}), \quad (26) \text{?eqn:deltaprior}$$

where $\bar{\Sigma}_{\delta_j(k)} = \sigma_\delta^2 I_{r_{\delta,j}}$ and $I_{r_{\delta,j}}$ is the $r_{\delta,j} \times r_{\delta,j}$ identity matrix.

The prior distribution of $\xi_j^2(k)$ is assumed to have the gamma density function:

$$p(\xi_j^2) = \gamma(\xi_j^2 | \bar{\alpha}_j, \bar{\beta}_j), \quad (27) \text{?eqn:xi2prior?}$$

where

$$\gamma(x | \alpha, \beta) = \frac{1}{\Gamma(\alpha)} \beta^\alpha x^{\alpha-1} e^{-\beta x}.$$

III.3. The posterior distribution. Given the likelihood function (16) and the prior density function (24)-(27), our objective is to obtain the conditional posterior density function $p(\theta | Y_T, Z_T, S_T, Q)$ by sampling alternately from the following conditional posterior distributions:

$$p(b_j(k) | Y_T, Z_T, S_T, G, \Xi, Q, b_i(k)),$$

$$p(\delta_j(k) | Y_T, Z_T, S_T, A, \Xi, Q, \psi_j),$$

$$p(\psi_j | Y_T, Z_T, S_T, A, \Xi, Q, \delta_j(k)),$$

$$p(\xi_j^2(k) | Y_T, Z_T, S_T, A, G, Q),$$

where $i \neq j$ and $i = 1, \dots, n$. We now discuss each of these four conditional density functions.

III.3.1. Conditional posterior density of $b_j(k)$. Combining the likelihood (16) and the prior (24) implies that the posterior density of $b_j(k)$, conditional on S_T, G, Ξ, Q , and $b_i(k)$ for $i \neq j$, is proportional to

$$\exp\left(-\frac{1}{2} b_j'(k) \bar{\Sigma}_{b_j}^{-1} b_j(k)\right) \prod_{t \in \{t: s_{1t}=k\}} \left[|A(k)| \exp\left(-\frac{\xi_j^2(s_{2t})}{2} (y_t' a_j(k) - x_t' f_j(k))^2\right) \right],$$

for $k = 1, \dots, h_1$. It is important to note that both $a_j(k)$ and $f_j(k)$ are affine functions of $b_j(k)$. To evaluate the above density kernel more efficiently, we sometimes use the

following functional form:

$$\exp\left(-\frac{1}{2}b_j'(k)\bar{\Sigma}_{b_j}^{-1}b_j(k)\right)|A(k)|^{T_{1,k}} \times \prod_{t \in \{t: s_{1t}=k\}} \left[\exp\left(-\frac{1}{2}(a_j'(k)\Sigma_{yy,k}a_j(k) - 2f_j'(k)\Sigma_{xy,k}a_j(k) + f_j'(k)\Sigma_{xx,k}f_j(k))\right) \right].$$

where $T_{1,k}$ is the number of t 's such that $s_{1t} = k$,

$$\begin{aligned} \Sigma_{yy,k} &= \sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t})y_t y_t', \\ \Sigma_{xy,k} &= \sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t})x_t y_t', \\ \Sigma_{xx,k} &= \sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t})x_t x_t'. \end{aligned}$$

Unlike the constant-parameter simultaneous-equation VAR models studied by Waggoner and Zha (2003a), the above conditional posterior density of $b_j(k)$ is nonstandard. We thus use a Metropolis algorithm with the following proposal density for the transition from $b_j(k)$ to $b_j^*(k)$

$$p(b_j^*(k) | b_j(k), Y_T, Z_T, S_T, b_1, \dots, b_{j-1}, b_{j+1}, \dots, b_n, G, \Xi, Q) = \text{normal}\left(b_j^*(k) | \mathbf{0}_{r_{b,j} \times 1}, \kappa_{b_j(k)} \Sigma_{b_j(k)}\right) \quad (28) \text{?eqn: jumping?}$$

where $b_j^*(k)$ is a proposal draw, $\kappa_{b_j(k)}$ is a scale factor that can be adjusted to keep the acceptance ratio optimal (e.g., between 25% and 40%), and

$$\Sigma_{b_j(k)}^{-1} = \bar{\Sigma}_{b_j(k)}^{-1} + U_j'(YY_k + W_j'XY_k + XY_k'W_j + W_j'XX_kW_j)U_j$$

where

$$\begin{aligned} YY_k &= \sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t})y_t y_t' \\ XY_k &= \sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t})x_t y_t' \\ XX_k &= \sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t})x_t x_t' \end{aligned}$$

III.3.2. *Conditional posterior densities of $\delta_j(k)$ and ψ_j .* As discussed in Section III.2.2, the long vector g_{ψ_j} is stacked from h_1 sub-vectors. It can be seen from (??) that the restriction matrix Ψ_j can be formed from h_1 corresponding sub-matrices. If we denote

$$g_{\psi_j} = \begin{bmatrix} g_{\psi_{j,1}} \\ \dots \\ g_{\psi_{j,k}} \\ \dots \\ g_{\psi_{j,h_1}} \end{bmatrix}, \quad \Psi_j = \begin{bmatrix} \Psi_{j,1} \\ \dots \\ \Psi_{j,k} \\ \dots \\ \Psi_{j,h_1} \end{bmatrix},$$

we have

$$g_{\psi_{j,k}} = \Psi_{j,k} \psi_j. \quad (29) \text{?eqn:restriction}$$

From the conditional likelihood (16), the prior distribution (26), and the restriction (17), one can obtain the posterior density kernel of $\delta_j(k)$ conditional on S_T, A, Ξ, Q , and ψ_j as

$$\prod_{k=1}^{h_1} \exp\left(-\frac{1}{2} \delta_j(k)' \tilde{\Sigma}_{\delta_j(k)}^{-1} \delta_j(k)\right) \times \prod_{t \in \{t: s_{1t}=k\}} \exp\left(-\frac{\xi_j^2(s_{2t})}{2} \left((y'_t + x'_t W_j) U_j b_j(k) - x'_t V_j \text{diag}(g_{\psi_{j,k}}) (\Delta_j \delta_j(k) + \bar{\delta}_j)\right)^2\right).$$

Rearranging the terms in the above equation leads to

$$p(\delta_j(k) | Y_T, Z_T, S_T, A, \Xi, Q, \psi_j) = \text{normal}(\delta_j(k) | \tilde{\mu}_{\delta_j(k)}, \tilde{\Sigma}_{\delta_j(k)}), \quad (30) \text{?eqn:posterior}$$

where

$$\begin{aligned} \hat{\Sigma}_{\delta_j(k)}^{-1} &= \Delta'_j \text{diag}(g_{\psi_{j,k}}) V'_j \Sigma_{xx,k} V_j \text{diag}(g_{\psi_{j,k}}) \Delta_j, \\ \tilde{\Sigma}_{\delta_j(k)}^{-1} &= \bar{\Sigma}_{\delta_j(k)}^{-1} + \hat{\Sigma}_{\delta_j(k)}^{-1}, \\ \hat{\mu}_{\delta_j(k)} &= \Delta'_j \text{diag}(g_{\psi_{j,k}}) V'_j \left[\sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t}) x_t (y'_t + x'_t W_j) \right] U_j b_j(k), \\ \tilde{\mu}_{\delta_j(k)} &= \tilde{\Sigma}_{\delta_j(k)} \left(\hat{\mu}_{\delta_j(k)} - \hat{\Sigma}_{\delta_j(k)}^{-1} \bar{\delta}_j \right). \end{aligned}$$

Similarly, from the conditional likelihood (16), the prior distribution (25), and the restriction (29), we obtain the posterior density kernel of ψ_j conditional on S_T, A, Ξ, Q , and δ_j as

$$\prod_{k=1}^{h_1} \exp\left(-\frac{1}{2} \psi'_j \tilde{\Sigma}_{\psi_j}^{-1} \psi_j\right) \times \prod_{t \in \{t: s_{1t}=k\}} \exp\left(-\frac{\xi_j^2(s_{2t})}{2} \left((y'_t + x'_t W_j) U_j b_j(k) - x'_t V_j \text{diag}(g_{\delta_j(k)}) \Psi_{j,k} \psi_j\right)^2\right).$$

Rearranging the terms in the above equation gives

$$p(\psi_j | Y_T, Z_T, S_T, A, \Xi, Q, \delta_j) = \text{normal}(\psi_j | \tilde{\mu}_{\psi_j}, \tilde{\Sigma}_{\psi_j}), \quad (31) \text{?eqn:posterior}$$

where

$$\begin{aligned} \hat{\Sigma}_{\psi_j}^{-1} &= \sum_{k=1}^{h_1} \Psi'_{j,k} \text{diag}(g_{\delta_j(k)}) V'_j \Sigma_{xx,k} V_j \text{diag}(g_{\delta_j(k)}) \Psi_{j,k}, \\ \tilde{\Sigma}_{\psi_j}^{-1} &= \bar{\Sigma}_{\psi_j}^{-1} + \hat{\Sigma}_{\psi_j}^{-1}, \\ \hat{\mu}_{\psi_j} &= \sum_{k=1}^{h_1} \Psi'_{j,k} \text{diag}(g_{\delta_j(k)}) V'_j \left[\sum_{t \in \{t: s_{1t}=k\}} \xi_j^2(s_{2t}) x_t (y'_t + x'_t W_j) \right] U_j b_j(k), \\ \tilde{\mu}_{\psi_j} &= \tilde{\Sigma}_{\psi_j} \hat{\mu}_{\psi_j}. \end{aligned}$$

III.3.3. *Conditional posterior density of $\xi_j^2(k)$.* Let $T_{2,k}$ be the number of elements in $\{t : s_{2t} = k\}$ for $k = 1, \dots, h_2$. It follows that

$$p(\xi_j^2(k) | Y_T, Z_T, S_T, A, G, Q) = \gamma(\xi_j^2(k) | \tilde{\alpha}_j(k), \tilde{\beta}_j(k)), \quad (32) \text{?eqn:posteriorx}$$

where

$$\begin{aligned} \tilde{\alpha}_j(k) &= \bar{\alpha}_j + \frac{T_{2,k}}{2}, \\ \tilde{\beta}_j(k) &= \bar{\beta}_j + \frac{1}{2} \sum_{t \in \{t: s_{2t}=k\}} (y'_t a_j(s_{1t}) - x'_t f_j(s_{1t}))^2. \end{aligned}$$

III.4. **Other types of Markov processes.** The previous analysis can be easily extended to other types of Markov processes. If we wish to synchronize the two state variables s_{1t} and s_{2t} into one state variable s_t , we simply need to replace these two independent state variables by this one common state variable s_t in the likelihood function. If we wish to have an independent Markov process for the coefficients in each equation, s_{1t} becomes a composite state variable consisting of $s_{j,1t}$ for $j = 1, \dots, n$. In this case, we simply replace s_{1t} by $s_{j,1t}$ for the time-varying coefficients in equation j in the likelihood function.

III.5. **Normalization.** To obtain accurate posterior distributions of functions of θ (such as long run responses and historical decompositions), we must normalize both the signs of structural equations and the labels of states; otherwise, the posterior distributions will be symmetric with multiple modes, making statistical inferences of interest meaningless. Such normalization is also essential to achieving efficiency in evaluating the marginal data density for model comparison.¹²

For both purposes, we normalize the signs of structural equations the same way. Specifically, we use the Waggoner and Zha (2003b) normalization rule to determine the column signs of $A(k)$ and $F(k)$ for any given $k \in \{1, \dots, h\}$. Since our original prior is unnormalized and symmetric around the origin, this prior density must be multiplied by 2^n when the marginal data density is estimated with MCMC draws that are normalized by the rule proposed by Waggoner and Zha (2003b).

The scale normalization on $\xi_j(k)$ and $\delta_j(k)$ imposes the restrictions $\xi_j(k) = 1$ and $\delta_j(k) = \mathbf{1}_{r_{\delta,j} \times 1}$ for all $j \in \{1, \dots, n\}$ and $k \in \{1, \dots, h_1\}$, where the notation $\mathbf{1}_{r_{\delta,j} \times 1}$ denotes the $r_{\delta,j} \times 1$ vector of 1's. In principle, one could use other normalization rules as suggested by Hamilton, Waggoner, and Zha (2004). The marginal data density, however, is invariant to this kind of normalization.

IV. BLOCKWISE OPTIMIZATION ALGORITHM

In spite of the complexity inherent in the multiple-equation models considered in this paper, it is essential to obtain the estimate of θ at the peak of the posterior distribution (5). The posterior estimate or the maximum likelihood estimate, serving as a starting point for our MCMC algorithm, ensures that an unreasonably long sequence of posterior draws

¹²Note that the marginal data density is invariant to the way parameters are normalized, as long as the Jacobian transformations of the parameters are taken into account explicitly.

do not get stuck in the low probability region. Moreover, used as a reference point in normalization, it helps avoid distorting the statistical inferences likely to be produced by inappropriate normalization. And the likelihood value conditional on the posterior estimate helps detect obvious errors in computing marginal data densities for posterior odds ratios.

Hamilton (1994) proposes an expectation-maximizing (EM) algorithm for a simple Markov-switching model. For multivariate dynamic models, however, the expectation step in general has no analytical form. Chib (1996) proposes a Monte Carlo EM (MCEM) algorithm in which the evaluation of the E-step of the EM algorithm is approximated by Monte Carlo simulations from the posterior distribution.

As shown in Sims and Zha (2006), these MC simulations can be very expensive computationally. When the number of parameters is small, one may obtain the posterior estimate of θ by simply finding the value of θ that maximizes the posterior density $p(\theta, Q | Y_T, Z_T)$ given by (5). Sims (2001) uses this approach for his single-equation model. But for a system of multivariate dynamic equations, the number of model parameters can be too large for a straight maximization routine to be reliable.

In this paper, we propose a different algorithm. We use the Gibbs-sampling idea to break the parameters θ, Q into two blocks of parameters θ and Q . In the multivariate dynamic models considered in this paper, we break the block of parameters θ further into three sub-blocks, one containing $b_j(k)$ for $k = 1, \dots, h_1$, one containing $g_j(k)$ for $k = 1, \dots, h_1$, and third sub-block containing $\xi_j^2(k)$ for $k = 1, \dots, h_2$. Given an initial guess of the values of the parameters, one can use the standard hill-climbing optimization routine (e.g., the Quasi-Newton BFGS algorithm) to find the values of each block of parameters that maximizes the posterior density while holding other blocks of parameters fixed at the previous values. Iterate this algorithm across blocks until it converges. For each iteration, we also employ a *constrained* optimization method to check whether there are boundary solutions associated with Q or any other model parameters.

V. NEW IMPLEMENTATION OF THE MHM METHOD

For structural models such as DSGE models and over-identified VARs, the modified harmonic mean (MHM) method of Gelfand and Dey (1994) is a widely used method to compute the marginal data density. In this section we discuss the potential problem with this method when used for multivariate dynamic models and propose a new way of implementing the MHM method to remedy this problem. For notational clarity, we restrict ourselves to the constant-parameter case, treat θ as a collection of all the free parameters in the model, and omit the exogenous variables Z_T . At the end of this section, we discuss how to handle the Markov-switching models.

We begin by denoting the likelihood function by $p(Y_T | \theta)$ and the prior density be $p(\theta)$, both of which must have proper probability densities instead of their kernels. Given these two objects, the marginal data density is defined as

$$p(Y_T) = \int p(Y_T | \theta) p(\theta) d\theta. \tag{33} \text{?eqn:mdd?}$$

The MHM method used to approximate (33) numerically is based on a theorem that states

$$p(Y_T)^{-1} = \int_{\Theta} \frac{h(\theta)}{p(Y_T | \theta)p(\theta)} p(\theta | Y_T) d\theta, \quad (34) \text{?eqn:mhm?}$$

where Θ is the support of the posterior probability density and $h(\theta)$, often called a *weighting function*, is any probability density whose support is contained in Θ . Denote

$$m(\theta) = \frac{h(\theta)}{p(Y_T | \theta)p(\theta)}.$$

A numerical evaluation of the integral on the right hand side of (34) can be accomplished in principle through the Monte Carlo (MC) integration

$$\hat{p}(Y_T)^{-1} = \frac{1}{N} \sum_{i=1}^N m(\theta^{(i)}), \quad (35) \text{?eqn:mhmmc?}$$

where $\theta^{(i)}$ is the i^{th} draw of θ from the posterior distribution $p(\theta | Y_T)$. If $m(\theta)$ is bounded above, the rate of convergence from this MC approximation is likely to be practical.

Geweke (1999) proposes an implementation with $h(\cdot)$ constructed from the posterior simulator. The sample mean $\bar{\theta}$ and sample covariance matrix $\bar{\Omega}$ can be calculated from draws of θ from the posterior simulator. The weighting function is chosen to be a truncated multivariate Gaussian density with mean $\bar{\theta}$ and covariance $\bar{\Omega}$. The Gaussian density is truncated to ensure that the support of the weighting function is contained in the support of posterior. This method has been used for many DSGE and identified VAR models. But there are three potential problems associated with use of this standard method. First, the posterior density may be quite small at the sample mean, especially when the posterior density has multiple peaks. Second, a truncated Gaussian density function may be a poor local approximation to the posterior density. Third, as one can see from (6), the likelihood tends to be zero in the interior points of the domain Θ . All these problems can make this particular choice of weighting function an inefficient way to implement the MHM method for many multivariate dynamic models.

To deal with these problems, we propose a more general class of distributions than the Gaussian family, center and scale these distributions differently, and truncate these distributions in a more sophisticated manner. The easiest of these to deal with is the centering and scaling. Instead of centering at the sample mean, we center at the posterior mode $\hat{\theta}$ and instead of scaling by the sample covariance matrix, we use

$$\hat{\Omega} = \frac{1}{N} \sum_{i=1}^N \left(\theta^{(i)} - \hat{\theta} \right) \left(\theta^{(i)} - \hat{\theta} \right)'$$

where $\theta^{(i)}$ denotes the i^{th} draw from the posterior simulator and N is the sample size. Of course computing the posterior mode is more expensive than computing the sample mean (see Section IV), but this greatly improves efficiency. Instead of the Gaussian family of distributions, we use elliptical distributions. An elliptical distribution centered at $\hat{\theta}$ and

scaled by $\hat{S} = \sqrt{\hat{\Omega}}$ has a density of the form

$$g(\theta) = \frac{\Gamma(k/2)}{2\pi^{k/2} |\det(\hat{S})|} \frac{f(r)}{r^{k-1}}$$

where k is the dimension of θ , $r = \sqrt{(\theta - \hat{\theta})' \hat{\Omega}^{-1} (\theta - \hat{\theta})}$ and f is any one-dimensional density defined on the positive reals. We note that the Gaussian is an example of an elliptical distribution. Given that we know how to make draws from the one dimensional density f , making draws for an elliptical distribution is easy. Simply draw x from a k -dimensional standard Gaussian distribution and r from the density f , and define

$$\theta = \frac{r}{\|x\|} \hat{S}x + \hat{\theta}.$$

The one-dimensional density f is chosen in the following manner. For each draw $\theta^{(i)}$ from the posterior, let

$$r^{(i)} = \sqrt{(\theta^{(i)} - \hat{\theta})' \hat{\Omega}^{-1} (\theta^{(i)} - \hat{\theta})}$$

From these simulated $r^{(i)}$, we can easily form an estimate of their cumulative density function. The density f should be chosen so that its cumulative density closely matches the estimated one. There are many ways of doing this, for instance, f could be chosen to be a step function so that the cumulative density is piecewise-linear approximation of the estimated cumulative density. We chose a somewhat simpler technique. The density f has support on $[a, b]$ and is defined by

$$f(r) = \frac{vr^{v-1}}{b^v - a^v}$$

The parameters a , b , and v are chosen as follows. Let c_1 , c_{10} , and c_{90} be chosen so that one percent of the $r^{(i)}$ are less than c_1 , ten percent of the $r^{(i)}$ are less than c_{10} , and ninety percent of the $r^{(i)}$ are less than c_{90} . The parameter a is equal to c_1 and b and v are so chosen so that roughly ten percent of the draws from f will be less than c_{10} and ninety percent of the draws will be less than c_{90} . This translates into

$$v = \frac{\log(1/9)}{\log(c_{10}/c_{90})} \text{ and } b = \frac{c_{90}}{0.9^{1/v}}.$$

We now turn to the method we use to truncate the elliptical distribution g . Let L be a positive number and Θ_L be the region defined by

$$\Theta_L = \{\theta : p(Y_T | \theta)p(\theta) > L\}.$$

The weighting function h is chosen to be the elliptical distribution truncated so that its support is Θ_L . If q_L is the probability that a draw from the elliptical distribution lies in Θ_L , then h is given by

$$h(\theta) = \frac{\chi_{\Theta_L}(\theta)}{q_L} g(\theta).$$

The value of q_L can be estimated from draws from the elliptical density g . Since we can take independent draws from an elliptical density, the estimate of q_L has a binomial distribution and its accuracy can be readily obtained. The higher the truncation value of

L is, the larger the effective sample size of a sequence of MCMC draws is, but the less acceptable the value of \hat{q}_L becomes. Therefore, there is a balance between having a high value of L and having a reasonable estimate of q_L . Our experience tells us that a good choice of L is a value such that 90% of the draws from the posterior lie in Θ_L .

The new MHM method described by (??) is computationally more demanding than the standard method (??), but it avoids the potential problem associated with the unboundedness of $m(\theta)$. Denote

$$k(\theta|Y_T) = p(Y_T | \theta)p(\theta).$$

Clearly, $k(\theta|Y_T)$ is the kernel of the posterior probability density $p(\theta|Y_T)$. The procedure for implementing our new MHM method is as follows.

- (1) Simulate a sequence of posterior draws $\theta^{(i)}$ and record the minimum and maximum values of $k(\theta|Y_T)$, denoted by k_{\min} and k_{\max} respectively. Let $k_{\min} < L < k_{\max}$.
- (2) Simulate random draws of θ from $g(\theta)$ and compute the proportion of these draws that belong to Θ_L . This proportion, denoted by \hat{q}_L , is the estimate of q_L . Note that \hat{q}_L has a binomial distribution and depends on the number of MCMC draws and the sample simulated from $h(\cdot)$. If $\hat{q}_L < 1.0e - 06$, this estimate is unreliable because three or four standard deviations will include the value zero. As a rule of thumb, we keep $\hat{q}_L \geq 1.0e - 04$.
- (3) For each value of L , estimate the marginal data density according to (35).

We have thus far discussed the constant-parameter case. For the Markov-switching models, the only difference is the treatment of the transition matrix Q in which w_j for $j = 1, \dots, h$ is a vector of free parameters as discussed in Section II.5. The transition matrix parameters w_j 's are treated separately from θ and we use a Dirichlet density instead of a truncated power density as the weighting function for w_j .

VI. APPLICATION

In this section we apply our method developed in the previous sections to a regime-switching three-variable VAR model with five lags. The three variables are those commonly used by recent DSGE models: log GDP (x_t), GDP-deflator inflation (π_t), and the federal funds rate (R_t). The data are quarterly from 1959:I to 2005:IV. Recent debate on changes in monetary policy has focused on whether the coefficients in the policy equation have changed or the variance sizes for structural shocks have changed. Using the notation in Section III.1, we let

$$y_t = [x_t \quad \pi_t \quad R_t]'$$

Following the identification of Christiano, Eichenbaum, and Evans (2005), we consider the upper triangular matrix $A(s_t)$ where the last equation is the interest rate equation. We study a large number of models and compare their fits to the data. The types of models are described as follows

- #v:** Each equation is of Case II with # states under one common Markov process. We call this type of model “variance-only.”

- #vm:** Each equation is of Case III with # states under one common Markov process. We call this type of model “all-change” (i.e., both variances and means changing).
- #vRm:** The interest rate (R) equation (i.e., the third equation in our application) is of Case III and the other two equations are of Case II with # states under one common Markov process. We call this type of model “policy-change” (i.e., both variances and coefficients in the policy equation changing).
- #₁v#₂m:** Each equation is of Case III, with #₁ states under one Markov process for $a_j(s_{1t})$ and $f_j(s_{1t})$ and with #₂ states under the other independent Markov process for $\xi_j(s_{2t})$, where $j = 1, \dots, n$. We call this type of model “variance-with-all-change.”
- #₁v#₂Rm:** The interest rate equation is of Case III and the other equations are of Case II, with #₁ states under one Markov process for $a_j(s_{1t})$ and $f_j(s_{1t})$ and with #₂ states under the other independent Markov process for $\xi_j(s_{2t})$, where $j = 1, \dots, n$. We call this type of model “variance-with-policy-change.”

For all these quarterly models, the tightness values for the BVAR reference prior are, in the notation of Sims and Zha (1998), $\lambda_0 = 1.0$, $\lambda_1 = 1.0$, $\lambda_2 = 1.0$, $\lambda_3 = 1.2$, $\lambda_4 = 0.1$, $\mu_5 = 1.0$, and $\mu_6 = 1.0$. These hyperparameter values determines the prior covariance matrices $\bar{\Sigma}_{b_j}$ and $\bar{\Sigma}_{\psi_j}$. For other prior settings, we follow Sims and Zha (2006) and set $\sigma_\delta = 50$, $\bar{\alpha}_j = 1.0$, and $\bar{\beta}_j = 1.0$. For the prior distribution of the transition probability q_j as discussed in Section II.5, we set $\alpha_{i,j} = 1$ for $i \neq j$ and

$$\alpha_{j,j} = \frac{q_d(h-1)}{1-q_d}.$$

This prior setting, differing from that of Sims and Zha (2006), not only implies that the expected value of the probability of staying in the previous state is q_d but also allows the posterior estimate to $q_{k,k}$ to be one (i.e., allowing the k^{th} state to be an absorbing one). For our quarterly data, we set $q_d = 0.85$, implying a prior belief that the average duration of staying in the same state is between 6 and 7 quarters.

Using the blockwise optimization algorithm described in Section IV, we obtain the posterior estimates of the model parameters.¹³ With this estimate or any random value near the neighborhood of this estimate as a starting point, we simulate a sequence of 20 million MCMC draws to compute the marginal data density using the new method described in Section V.¹⁴ For the case of 3 states, the restricted transition matrix takes the form of (1). For the case of 4 states, the transition matrix is restricted as

$$Q = \begin{bmatrix} \pi_1 & (1-\pi_2)/2 & 0 & 0 \\ 1-\pi_1 & \pi_2 & (1-\pi_3)/2 & 0 \\ 0 & (1-\pi_2)/2 & \pi_3 & 1-\pi_4 \\ 0 & 0 & (1-\pi_3)/2 & \pi_4 \end{bmatrix}.$$

¹³For our C program, this algorithm takes less than 1 minute while the EM algorithm takes about 9 hours on a Pentium-4 personal desktop computer.

¹⁴It takes about 20 minutes to simulate one million MCMC draws.

Table 1 reports log values of marginal data densities for nine different models. The MDDs are not sensitive to the cutoff value L for our new MHM method. In the table, we report only one value and the corresponding \hat{q}_L . For each sequence of MCMC draws, we use the software R to compute an effective sample size (ESS) (i.e., the sample size adjusted for serial correlation of MCMC draws) according to Plummer, Best, Cowles, and Vines (2005). For all the models studied in Table 1, the computed ESSs are near one million.¹⁵ Based on the ESS, thus, the numerical standard error on the estimated MDD is trivially small. On the similar magnitude, we obtain very small numerical standard errors based on the procedure of Newey and West (1987).

It is known, however, that these measures tend to deliver much smaller numerical standard errors than the actual ones. We propose a different measure by breaking a sequence of 20 million MCMC draws into 10 successive blocks with each block having 2 million draws. For each block, we compute log value of the inverse of the estimated mean of $m(\theta)$ (by a proper scaling to avoid an overflow in computation). The standard deviation of log MDD is then computed according to the differences of log MDD across blocks and reported in Table 1. These standard deviations imply that our computed MDDs are tightly estimated. Figures 1 and 2 plot the values of log MDD across blocks for the 4-state variance-only model and the 3v2R_m variance-with-policy-change model. As can be seen, the estimated log MDD is quite stable across blocks.

The best-fit model is 3-state or 4-state variance-only model, which clearly dominates all other models. Among the models with changing coefficients, the 3v2Rm variance-with-policy-change model is the best, which does not improve upon the 3-state variance-only model. The conclusion that the variance-only model dominates is also given by the Schwarz criterion applied to the posterior kernel.

Our exercises point to the fact that accurate calculation of the MDD is an extremely challenging task and give reasons why the conventional way of implementing the MHM method in the existing literature can be problematic.

VII. CONCLUSION

We have developed a generalized Bayesian method for Markov-switching models with independent Markov processes and restricted transition matrices. Based this generalized method, we have developed tools for estimating both identified and unrestricted VARs. Our proposed blockwise optimization algorithm proves much more efficient than the existing EM algorithm for obtaining the MLEs or posterior estimates of model parameters. Our new way of implementing the MHM method deals explicitly with the problem of zero likelihood in the interior points of the parameter space. Because of this problem, uncertainty about the estimated posterior odds ratios reported in the recent macroeconomic literature is likely to be far more than what is commonly thought. Our new method of computing the marginal data density seems essential to achieving a reliable estimate of marginal likelihood.

¹⁵Because of some memory management problems associated with the program R , the ESSs are estimated on the smaller sample thinned by every twenty MCMC draws.

TABLE 1. Marginal data densities by new MHM method

	2v		2vRm	2v2m	2v2Rm				
log(MDD)	1819.87		1832.14	1858.22	1836.38				
s.d. of log(MDD)	0.0509		0.0234	0.039	0.8125				
log(L)	1687.6 0		1699.33	1638.17	1697.58				
\hat{q}_L	4.9e-1		2.1e-3	3.6e-4	1.1e-4				
	2v3Rm	3v	3v2Rm	4v					
log(MDD)		1863.58	1861.80	1866.39					
s.d. of log(MDD)		0.4792	0.2696	0.3325					
log(L)		1711.76	1681.71	1713.22					
\hat{q}_L		2.1e-4	1.9e-4	1.3e-4					

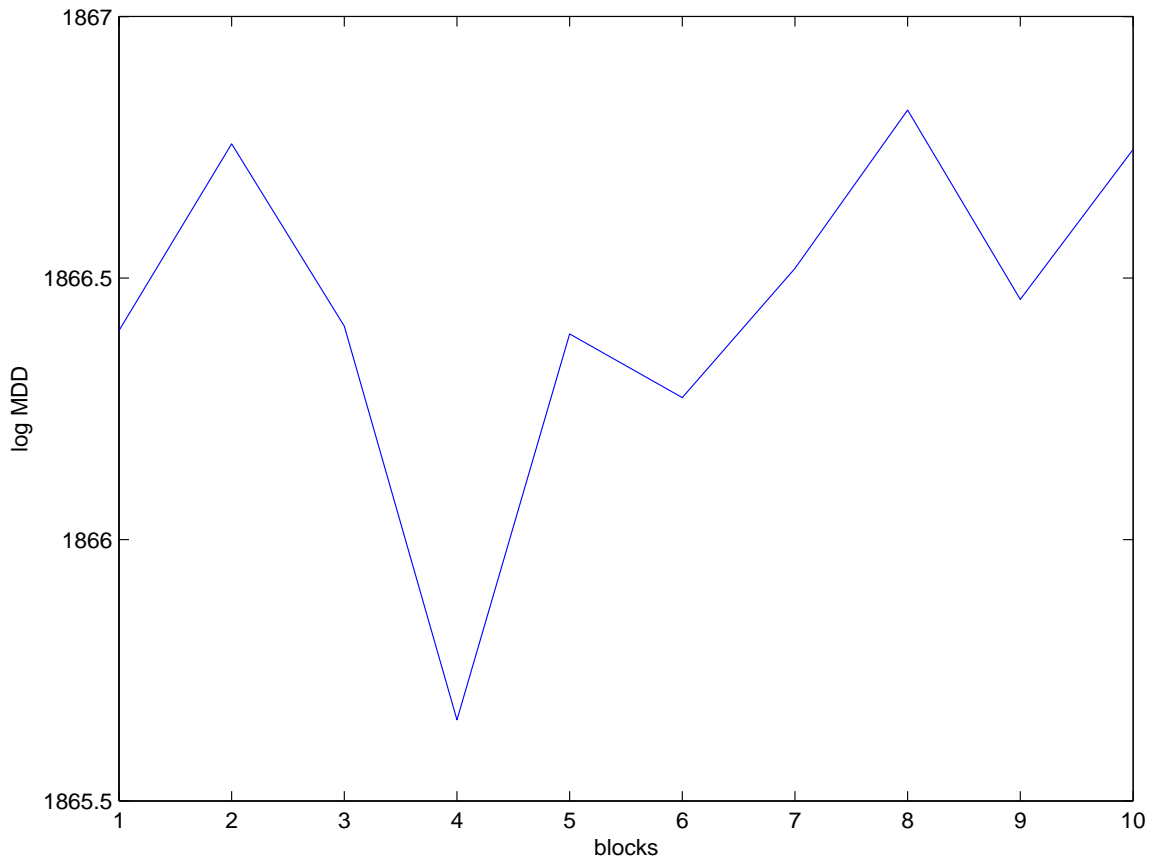


FIGURE 1. The 4-state variance-only (4v) model.

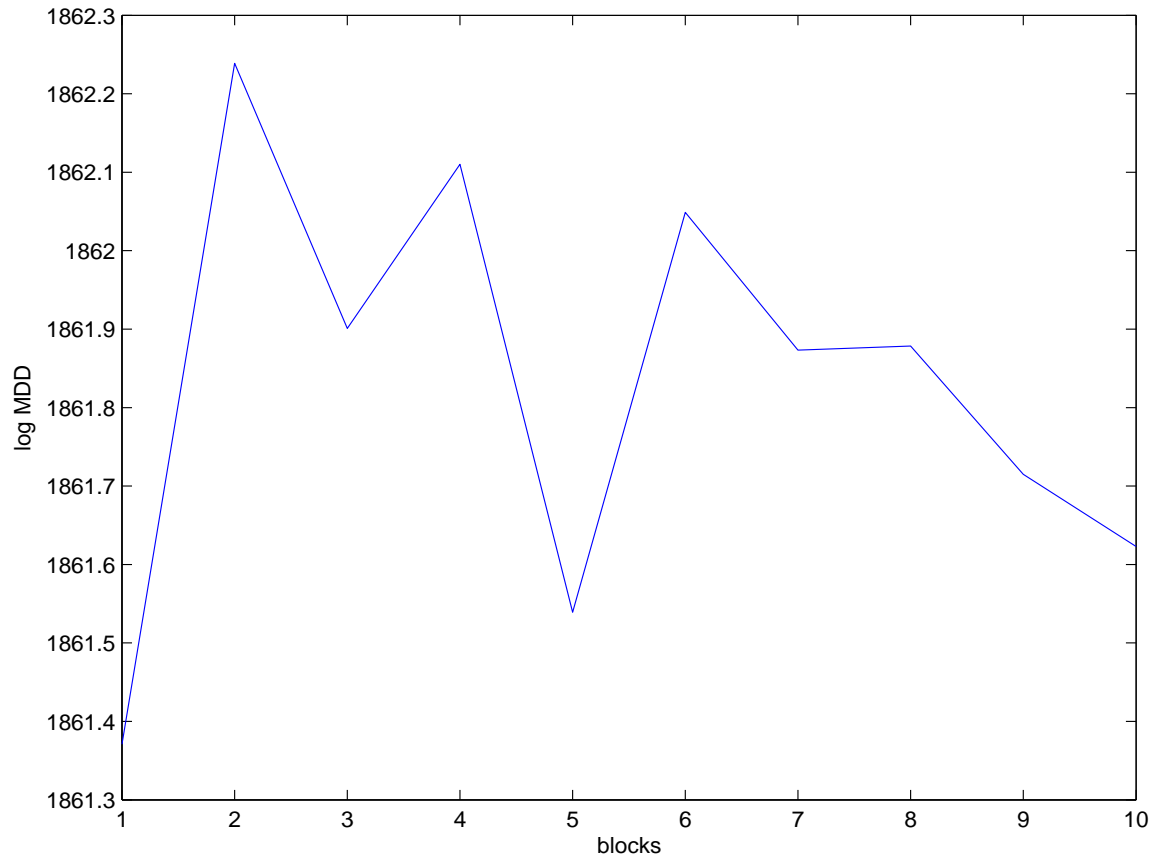


FIGURE 2. The 3v2Rm variance-with-policy-change model.

APPENDIX A. COMPUTING U_j , V_j , AND W_j

We assume that a_j and f_j satisfy linear restrictions of the form

$$Q_j \begin{bmatrix} a_j \\ f_j \end{bmatrix} = 0$$

where Q_j is a $(n+k) \times (n+k)$ with $k = np+m$. The matrix Q_j will not be of full rank. We show that there exist a $n \times q_j$ matrix U_j with orthogonal columns, a $(pn+m) \times r_j$ matrix V_j with orthogonal columns, and a such that $(pn+m) \times n$ matrix W_j with $W_j'V_j = 0$ such that

$$\begin{aligned} a_j(k) &= U_j b_j(k) \\ f_j(k) &= V_j g_j(k) - W_j U_j b_j(k) \end{aligned}$$

To prove this we rely on the following result:

Proposition 5. Given any $r \times s$ matrix X with $r \geq s$, there exist an invertible $r \times r$ matrix Y and a $s \times s$ orthogonal matrix $\begin{bmatrix} \hat{Z} & Z \end{bmatrix}$ where Z is a $s \times q$ matrix and \hat{Z} is a $s \times (s-q)$ such that

$$Y^{-1}X = \begin{bmatrix} \hat{Z}' \\ 0 \end{bmatrix}$$

Proof. This follows directly from the singular value decomposition of X . Let $X = UDV'$ where U is an $r \times r$ orthogonal matrix, V is a $s \times s$ orthogonal matrix, and D is a $r \times s$ diagonal matrix where the first $s-q$ diagonal elements are non-zero and the last q diagonal elements are zero. The first $s-q$ columns of V will be \hat{Z} , the last q columns of V will be Z , and $Y = UE$ where E is the $r \times r$ diagonal matrix whose first $s-q$ diagonal elements are the first $s-q$ diagonal elements of D , and the last $r-(s-q)$ diagonal elements are one. \square

Applying the above proposition to the last k columns of Q_j , there exists a $(n+k) \times (n+k)$ invertible matrix Y_1 and a $k \times k$ orthogonal matrix $\begin{bmatrix} \hat{V}_j & V_j \end{bmatrix}$ where \hat{V}_j is $k \times (k-r_j)$ and V_j is $k \times r_j$ such that

$$Y_1^{-1}Q_j = \begin{bmatrix} \hat{W}_j & \hat{V}_j' \\ \tilde{U}_j & 0 \end{bmatrix}$$

Now applying the above proposition to the $(n+r_j) \times n$ matrix \tilde{U}_j , there exists a $(n+r_j) \times (n+r_j)$ invertible matrix Y_2 and a $n \times n$ orthogonal matrix $\begin{bmatrix} \hat{U}_j & U_j \end{bmatrix}$ where \hat{U}_j is $n \times (n-q_j)$ and U_j is $n \times q_j$ such that

$$\begin{bmatrix} I_{k-r_j} & 0 \\ 0 & Y_2^{-1} \end{bmatrix} Y_1^{-1}Q_j = \begin{bmatrix} \hat{W}_j & \hat{V}_j' \\ \hat{U}_j' & 0 \\ 0 & 0 \end{bmatrix}$$

Thus a_j and f_j satisfy the restrictions if and only if

$$\begin{bmatrix} \hat{W}_j & \hat{V}_j' \\ \hat{U}_j' & 0 \end{bmatrix} \begin{bmatrix} a_j \\ f_j \end{bmatrix} = 0.$$

Since both $\hat{V}_j' \hat{V}_j$ and $\hat{U}_j' \hat{U}_j$ are equal to a identity matrix, writing $a_j = U_j b_j + \hat{U}_j c_j$ and $f_j = V_j g_j + \hat{V}_j h_j$ gives

$$\begin{bmatrix} \hat{W}_j & \hat{V}_j' \\ \hat{U}_j' & 0 \end{bmatrix} \begin{bmatrix} a_j \\ f_j \end{bmatrix} = \begin{bmatrix} \hat{W}_j U_j b_j + \hat{W}_j \hat{U}_j c_j + h_j \\ c_j \end{bmatrix}.$$

This is zero if and only if $c_j = 0$ and $h_j = -\hat{W}_j U_j b_j$. If we define $W_j = \hat{V}_j \hat{W}_j$, then the result follows.

REFERENCES

- [BF04] [1] BEYER, A., AND R. E. FARMER (2004): “What We Don’t Know About the Monetary Transmission Mechanism and Why We Don’t Know It,” Centre for Economic Policy Research Discussion Paper No. 4811.
- [CanovaG04] [2] CANOVA, F., AND L. GAMBETTI (2004): “Structural Changes in the US Economy: Bad Luck or Bad Policy,” Manuscript, Universitat Pompeu Fabra.
- [sC96] [3] CHIB, S. (1996): “Calculating Posterior Distributions and Model Estimates in Markov Mixture Models,” *Journal of Econometrics*, 75, 79–97.
- [CEE05] [4] CHRISTIANO, L., M. EICHENBAUM, AND C. EVANS (2005): “Nominal Rigidities and the Dynamics Effects of a Shock to Monetary Policy,” *Journal of Political Economy*, 113, 1–45.
- [CSCS02] [5] COGLEY, T., AND T. J. SARGENT (2002): “Evolving US Post-Wolrd War II Inflation Dynamics,” *NBER Macroeconomics Annual*, 16, 331–373.
- [CSCS05b] [6] ——— (2005): “Drifts and Volatilities: Monetary Policies and Outcomes in the Post WWII U.S.,” *Review of Economic Dynamics*, 8, 262–302.
- [CLR84] [7] COOLEY, T. F., S. F. LEROY, AND N. RAYMON (1984): “Econometric policy evaluation: Note,” *The American Economic Review*, 74, 467–470.
- [FWZ06] [8] FARMER, R. E., D. F. WAGGONER, AND T. ZHA (2006): “Minimal State Variable Solutions to Markov-Switching Rational Expectations Models,” Unpublished Manuscript.
- [GD94] [9] GELFAND, A. E., AND D. K. DEY (1994): “Bayesian Model Choice: Asymptotics and Exact Calculations,” *Journal of the Royal Statistical Society (Series B)*, 56, 501–514.
- [jG99] [10] GEWEKE, J. (1999): “Using Simulation Methods for Bayesian Econometric Models: Inference, Development, and Communication,” *Econometric Reviews*, 18(1), 1–73.
- [jH89] [11] HAMILTON, J. D. (1989): “A New Approach to the Economic Analysis of Nonstationary Time Series and the Business Cycle,” *Econometrica*, 57(2), 357–384.
- [jH94] [12] ——— (1994): *Times Series Analsis*. Princeton University Press, Princeton, NJ.
- [HWZ04] [13] HAMILTON, J. D., D. F. WAGGONER, AND T. ZHA (2004): “Normalization in Econometrics,” Federal Reserve Bank of Atlanta Working Paper 2004-13.
- [KN99] [14] KIM, C.-J., AND C. R. NELSON (1999): *State-Space Models with Regime Switching*. MIT Press, London, England and Cambridge, Massachusetts.
- [LZ03] [15] LEEPER, E. M., AND T. ZHA (2003): “Modest Policy Interventions,” *Journal of Monetary Economics*, 50(8), 1673–1700.
- [NW87] [16] NEWEY, W. K., AND K. K. WEST (1987): “A Simple Positive Semi-Definite Heteroskedasticity and Autocorrelation Consistent Covariance Matrix,” *Econometrica*, 55, 703–708.
- [coda05] [17] PLUMMER, M., N. BEST, K. COWLES, AND K. VINES (2005): “The coda Package,” Version 0.10-2, November, plummer@iarc.fr.
- [gP05] [18] PRIMICERI, G. (2005): “Time Varying Structural Vector Autoregressions and Monetary Policy,” *Review of Economic Studies*, 72, 821–852.
- [RT99] [19] ROBERTSON, J. C., AND E. W. TALLMAN (1999): “Vector Autoregressions: Forecasting and Reality,” *Federal Reserve Bank of Atlanta Economic Review*, First Quarter, 4–18.

- [RT01][20] ——— (2001): “Improving Federal-Funds Rate Forecasts in VAR Models Used for Policy Analysis,” *Journal of Business and Economic Statistics*, 19(3), 324–330.
- [RWZ05][21] RUBIO-RAMÍREZ, J. F., D. F. WAGGONER, AND T. ZHA (2005): “Markov-Switching Structural Vector Autoregressions: Theory and Application,” *Federal Reserve Bank of Atlanta, Working Paper 2005-27*.
- [cS87][22] SIMS, C. A. (1987): “A rational expectations framework for short-run policy analysis,” in *New approaches to monetary economics*, ed. by W. A. Barnett, and K. J. Singleton, pp. 293–308. Cambridge University Press, Cambridge, England.
- [cS93][23] ——— (1993): “A 9 Variable Probabilistic Macroeconomic Forecasting Model,” in *Business Cycles, Indicators, and Forecasting*, ed. by J. H. Stock, and M. W. Watson, vol. 28 of *NBER Studies in Business Cycles*, pp. 179–214. University of Chicago Press.
- [cS99][24] ——— (1999): “Drift and Breaks in Monetary Policy,” Manuscript, Princeton University.
- [cS01][25] ——— (2001): “Stability and Instability in US Monetary Policy Behavior,” Manuscript, Princeton University.
- [SZ98a][26] SIMS, C. A., AND T. ZHA (1998): “Bayesian Methods for Dynamic Multivariate Models,” *International Economic Review*, 39(4), 949–968.
- [SZ06][27] ——— (2006): “Were There Regime Switches in US Monetary Policy?,” *The American Economic Review*, 96, 54–81.
- [jSmW03][28] STOCK, J. H., AND M. W. WATSON (2003): “Has the Business Cycles Changed? Evidence and Explanations,” *Monetary Policy and Uncertainty: Adapting to a Changing Economy*, Federal Reserve Bank of Kansas City Symposium, Jackson Hole, Wyoming, August 28–30.
- [WZ03b][29] WAGGONER, D. F., AND T. ZHA (2003a): “A Gibbs Sampler for Structural Vector Autoregressions,” *Journal of Economic Dynamics and Control*, 28(2), 349–366.
- [WZ03a][30] ——— (2003b): “Likelihood Preserving Normalization in Multiple Equation Models,” *Journal of Econometrics*, 114(2), 329–347.
- [tZ06][31] ZHA, T. (In press): “Vector Autoregressions,” in *The New Palgrave Dictionary of Economics*, ed. by L. E. Blume, and S. Durlauf. Palgrave Macmillan.