

Markov-switching model selection using Kullback-Leibler divergence

Aaron Smith*

Department of Agricultural and Resource Economics
University of California, Davis

Prasad A. Naik

Graduate School of Management
University of California, Davis

Chih-Ling Tsai

Graduate School of Management
University of California, Davis
and
Guanghua School of Management
Peking University, P. R. China

Abstract

In Markov-switching regression models, we use Kullback-Leibler (KL) divergence between the true and candidate models to select the number of states and variables simultaneously. In applying Akaike information criterion (AIC), which is an estimate of KL divergence, we find that AIC retains too many states and variables in the model. Hence, we derive a new information criterion, Markov switching criterion (MSC), which yields a marked improvement in state determination and variable selection because it imposes an appropriate penalty to mitigate the over-retention of states in the Markov chain. MSC performs well in Monte Carlo studies with single and multiple states, small and large samples, and low and high noise. Furthermore, it not only applies to Markov-switching regression models, but also performs well in Markov-switching autoregression models. Finally, the usefulness of MSC is illustrated via applications to the U.S. business cycle and the effectiveness of media advertising.

JEL Classification: C22, C52.

Keywords: Advertising effectiveness, business cycles, EM algorithm, hidden Markov models, information criterion, Markov-switching regression.

*Corresponding author: Aaron Smith, Department of Agricultural and Resource Economics, University of California, One Shields Avenue, Davis, CA 95616. Phone 530-752-2138; Fax 530-752-5614; E-mail adsmith@ucdavis.edu.

1. Introduction

Economic systems often experience shocks that shift them from their present state into another state; for example, nations lurch into recession, government regimes change over time, and financial markets exhibit bubbles and crashes. These states tend to be stochastic and dynamic: if they occur once, they probably recur. To capture such probabilistic state transitions over time, Markov-switching models provide an analytical framework. In economics, Markov-switching models have been used for investigating the U.S. business cycle (Hamilton 1989), foreign exchange rates (Engel and Hamilton 1990), stock market volatility (Hamilton and Susmel 1994), real interest rates (Garcia and Perron 1996), corporate dividends (Timmermann 2001), the term structure of interest rates (Ang and Bekaert 2002a), and portfolio allocation (Ang and Bekaert 2002b), among others. Outside of economics, Markov-switching models find application in diverse fields such as computational biology (e.g., Durbin et al. 1998 for gene sequencing), computer vision (Bunke and Caelli 2001), and speech recognition (Rabiner and Juang 1993).

To estimate Markov-switching models, Baum and his colleagues (Baum and Petrie 1966, Baum et al. 1970) developed the forward-backward algorithm, which was extended to encompass general latent variable models under the expectation-maximization (EM) principle (see Dempster, Laird and Rubin 1977). If the number of states in Markov-switching models is known, the EM algorithm yields consistent parameter estimates, and statistical inference proceeds via standard maximum-likelihood theory (e.g., Bickel, Ritov and Rydén 1998). If the number of states is not known, however, the likelihood ratio test to infer the true number of states breaks down because regularity conditions do not hold (see Hartigan 1977, Hansen 1992, Garcia 1998).

The number of states is often not known a priori, so we propose applying Kullback-Leibler (KL) divergence to determine it. We note that KL divergence has been used in various model selection contexts (see, e.g., Sawa 1978, Leroux 1992, Sin and White 1996, Burnham and Anderson 2002). Specifically, Akaike's information criterion (AIC, see Akaike 1973) provides an estimate of KL distance but, in Markov-switching models, it misleads the users into selecting too many states (see Section 4.2). Consequently, one fits spurious regressions in nonexistent states; this misspecification results in incorrect inclusion of variables, which reduces the accuracy of estimated parameters and lowers the precision of model forecasts. Hence, the problem of simultaneous determination of the number of states to retain in the Markov chain and the variables to include in the regression model for each retained state remains open.

The objective of this paper is to develop a new information criterion for simultaneous selection of states and variables in Markov switching models. To accomplish this goal, we obtain an explicit approximation to the KL distance for the class of Markov switching regression models. The resulting Markov switching criterion (MSC) imposes an appropriate penalty, and so it mitigates the over-retention of states in the Markov chain and alleviates the tendency to over-fit the number of variables in each state. Moreover, in Monte Carlo studies, MSC performs well in single and multiple states, small and large samples, and low and high noise. Finally, it not only applies to Markov-switching regression models, but also performs well in Markov-switching autoregression models.

We present two empirical applications of MSC to understand (a) the business cycles in the US economy and (b) the effectiveness of media advertising. In the business cycle application, based on the minimum MSC value, we retain a three-state model for US GNP growth with one recessionary state and two expansionary states. The second expansionary state

occurred mostly after 1984, and it exhibits slower growth, lower volatility, and longer duration than the first one. This finding supports the notion of “great moderation” (see Kim and Nelson 1999a, McConnell and Perez-Quiros 2000, Stock and Watson 2003). In the advertising application, MSC suggests the retention of a two-state Markov-switching model for sales and advertising of the Lydia Pinkham brand; the results reveal new insights not discernible from the standard regression model.

We organize this paper as follows. In Section 2, we describe the model structure and estimation algorithm for multiple state Markov-switching models. We derive the information criterion in Section 3 and investigate its properties and performance under various conditions in Section 4. Section 5 presents empirical applications to business cycles and media advertising. Section 6 concludes the paper by identifying avenues for future research.

2. Estimating N -state Markov-switching models

We present the model structure, establish notation, and briefly describe the estimation of Markov-switching regressions, conditional on knowing the number of states N .

2.1 *Model structure*

Consider an N -state Markov chain. Let s_t denote an $N \times 1$ selection vector with elements $s_{ti} = 1$ or 0, according to whether the Markov chain resides in the state i ($i = 1, \dots, N$). The unobserved state vector s_t evolves according to an ergodic Markov chain with the transition probability matrix

$$P = \begin{bmatrix} p_{11} & \cdots & p_{1N} \\ \vdots & p_{ij} & \vdots \\ p_{N1} & \cdots & p_{NN} \end{bmatrix}, \quad (1)$$

where $p_{ij} = pr(s_{t+1,j} = 1 | s_{ti} = 1)$ and $\sum_{j=1}^N p_{ij} = 1$ for every $i = 1, \dots, N$. We define the ergodic probabilities of the Markov chain by the vector $\pi = (\pi_1, \dots, \pi_N)'$, where $\sum_{i=1}^N \pi_i = 1$.

At time t , when the chain is in state i (i.e., $s_{ti} = 1$), we observe the dependent variable y_t according to the regression model

$$y_t = x_t' \beta_i + \sigma_i \varepsilon_{ti}, \quad (2)$$

where $\varepsilon_{ti} \sim N(0, 1)$ is independently distributed over time $t = 1, \dots, T$, x_t contains K explanatory variables, and the $K \times 1$ vector β_i denotes their marginal impact when the chain is in the state i . If the chain moves to the state j , the marginal impact of exogenous variables is β_j with the corresponding level of noise σ_j^2 . To capture this “switching” in regression models, we rewrite (2) as follows:

$$y_t = x_t' \beta s_t + \sigma s_t \varepsilon_t \quad (3)$$

where $\beta = (\beta_1, \dots, \beta_N)$, $\sigma = (\sigma_1, \dots, \sigma_N)$, and the selection vector s_t indicates the state at time t . The matrix β and vector σ have dimensions $K \times N$ and $1 \times N$, respectively. Equations (1) and (3), together, constitute the N -state Markov-switching regression model. When x_t includes lagged values of y_t , we obtain the N -state Markov-switching autoregression model (e.g., Hamilton 1989). Next, we describe an EM algorithm to estimate this model.

2.2 EM algorithm

Suppose we observe the complete data, including the sequences of both the observed variables $Y = \{(y_t, x_t') : t = 1, \dots, T\}$ and the state variables $S = \{s_t : t = 1, \dots, T\}$. Then the complete data log-likelihood function L_c is

$$\begin{aligned}
L_c(\theta; Y, S) &= L(\beta, \sigma; Y | S) + L(P; S) \\
&= \sum_{t=1}^T \sum_{i=1}^N s_{ti} \log f_i(y_t; \beta_i, \sigma_i) + \sum_{t=1}^{T-1} \sum_{i=1}^N \sum_{j=1}^N s_{ti} s_{t+1,j} \log p_{ij} + \sum_{i=1}^N s_{1i} \log \pi_i,
\end{aligned} \tag{4}$$

where $f_i(y_t; \beta_i, \sigma_i) = (2\pi\sigma_i^2)^{-1/2} \exp(-\frac{1}{2}(y_t - x_t'\beta_i)^2 / \sigma_i^2)$ is the density of y_t conditional on $s_{ti} = 1$ (see McLachlan and Peel 2000, p. 329).

In the E-step, we evaluate the expectation of L_c with respect to the unobserved latent states S , given the observed data Y and provisional estimates of θ . Let θ^l denote the provisional estimates at the l -th iteration, and $Q(\theta; \theta^l) = E[L_c | Y, \theta^l]$. Because L_c is linear in s_{ti} , $s_{ti}s_{t+1,j}$, and s_{1i} , we obtain

$$Q(\theta, \theta^l) = \sum_{t=1}^T \sum_{i=1}^N \xi_{ti}^{(l)} \log f_i(y_t; \beta_i, \sigma_i) + \sum_{t=1}^{T-1} \sum_{i=1}^N \sum_{j=1}^N \tau_{ij}^{(l)} \log p_{ij} + \sum_{i=1}^N \xi_{1i}^{(l)} \log \pi_i, \tag{5}$$

where $\tau_{ij}^{(l)} = E(s_{ti}s_{t+1,j} | Y, \theta^l)$ and $\xi_{ti}^{(l)} = E(s_{ti} | Y, \theta^l)$. To compute $(\tau_{ij}^{(l)}, \xi_{ti}^{(l)})$, we apply the forward-backward algorithm (e.g., McLachlan and Peel 2000, p. 330), which yields

$$\tau_{ij}^{(l)} = \frac{a_{ti}^{(l)} p_{ij}^{(l)} f_i(y_{t+1}; \theta^l) b_{t+1,j}^{(l)}}{\sum_{i=1}^N \sum_{j=1}^N a_{ti}^{(l)} p_{ij}^{(l)} f_i(y_{t+1}; \theta^l) b_{t+1,j}^{(l)}} \tag{6}$$

and

$$\xi_{ti}^{(l)} = \sum_{j=1}^N \tau_{ij}^{(l)}. \tag{7}$$

The ‘‘forward’’ probabilities a_{ti} are given by the forward recursion

$$a_{t+1,i}^{(l)} = \left(\sum_{j=1}^N a_{tj}^{(l)} p_{ij}^{(l)} \right) f_i(y_{t+1}; \theta^l), \tag{8}$$

and the ‘‘backward’’ probabilities b_{ij} are given by the backward recursion

$$b_{ij}^{(l)} = \sum_{i=1}^N p_{ij}^{(l)} f_i(y_{t+1}; \theta^l) b_{t+1,i}^{(l)}. \tag{9}$$

We initialize these recursions by setting $a_{1i} = \pi_i^{(l)} f_i(y_1; \theta^l)$ and $b_{Tj} = 1$, where

$\pi^{(l)} = (\pi_1^{(l)}, \dots, \pi_N^{(l)})'$ is the principal eigenvector of $P^{(l)}\pi = \pi$.

In the M-step, we maximize $Q(\theta, \theta^l)$ with respect to $\theta = \text{vec}(\beta, \sigma, P)$ to obtain the closed form estimates for the $(l+1)$ -th iteration:

$$\beta_i^{(l+1)} = (X'W_i^{(l)}X)^{-1}X'W_i^{(l)}y, \quad (10)$$

$$(\sigma_i^{(l+1)})^2 = (y - X\beta_i^{(l+1)})'W_i^{(l)}(y - X\beta_i^{(l+1)})/T_i^{(l)}, \quad (11)$$

and

$$p_{ij}^{(l+1)} = \frac{\sum_{t=1}^{T-1} \tau_{ij}^{(l)}}{\sum_{t=1}^{T-1} \xi_{ii}^{(l)}}, \quad (12)$$

where $X = (x_1, \dots, x_T)'$, $y = (y_1, \dots, y_T)'$, $W_i^{(l)} = \text{diag}(\xi_i^{(l)})$, $\xi_i^{(l)} = (\xi_{1i}^{(l)}, \dots, \xi_{ii}^{(l)}, \dots, \xi_{Ti}^{(l)})'$, and $T_i^{(l)} = \text{tr}(W_i^{(l)})$. Using the provisional estimates θ^l , we obtain the new estimates $\theta^{(l+1)} = \text{vec}(\beta^{(l+1)}, \sigma^{(l+1)}, P^{(l+1)})$ via the equations (10) through (12). We iterate the E- and M-steps until the absolute difference $|\theta^{(l+1)} - \theta^l|$ decreases below a preset tolerance. The resulting vector $\hat{\theta} = \text{vec}(\hat{\beta}, \hat{\sigma}, \hat{P})$ converges to the maximum likelihood estimates, which are consistent and asymptotically normal (Bickel, Ritov and Rydén, 1998). For finite sample properties, see Psaradakis and Sola (1998). We close this section with two remarks.

Remark 1. We enhance the stability of this algorithm as follows. First, to avoid singularities in the likelihood function and reduce the chance of spurious local maxima, we follow Hathaway's (1985) suggestion to set a lower bound on the relative variances across states. Second, to prevent underflow of forward probabilities in (8), for each t and $i = 1, \dots, N$, we follow Leroux and Puterman's (1992) recommendation to scale a_{ti} with a constant r such that $10^{-r} \sum_{i=1}^N a_{ti}$ lies between 0.1 and 1.0 and then multiply it by 10^{-r} . Because a_{ti} appears in both the numerator and denominator of (6), the value of $\tau_{ij}^{(l)}$ does not change. Similarly, we prevent

underflow of backward probabilities in (9).

Remark 2. This EM algorithm enables the estimation of Markov-switching models with many observations because the forward-backward method is linear in T . Furthermore, because both the E- and M-steps are available in closed form, the EM algorithm is robust to numerical uncertainties encountered by quasi-Newton methods. For example, Hamilton (1990, pp. 40-41) notes that "...methods that seek to approximate the sample Hessian can easily go astray ...By contrast, the EM algorithm by construction finds an analytical interior solution to a particular subproblem." Nonetheless, like quasi-Newton methods, the EM does not guarantee convergence to global maxima (see McLachlan and Krishnan 1997, p. 34). Finally, the EM algorithm can also be used to obtain Bayesian modal values by augmenting the expected complete data likelihood with the logarithm of prior density; see Dempster, Laird and Rubin (1977, p. 6) for this connection between EM and Bayesian analysis and Kim and Nelson (1999b, Ch. 9) for implementation in Markov-switching models.

3. Deriving Markov-switching criterion

In the above estimation, the number of states N is assumed known, which need not be the case in practice. To determine the number of states, we approximate the true data generating process (DGP) using several candidate models, quantify the information loss between the DGP and each candidate model, and then choose the model that entails the minimum expected information loss (e.g., Burnham and Anderson 2002). Specifically, let $g(Y^*)$ denote the probability density function of the DGP and $f(Y^*; \theta)$ be the density function for a candidate model, where Y^* represents the data used for evaluating the model. As in Sawa (1978) and Sin and White (1996), we quantify information loss using the Kullback-Leibler (KL) divergence,

which is defined as

$$d_{KL}(g, f; \theta) = E_{Y^*} \left[\log \frac{g(Y^*)}{f(Y^*; \theta)} \right], \quad (13)$$

where $d_{KL} \geq 0$, and $E_{Y^*}(\cdot)$ denotes the expectation with respect to the data generating density g .

Equation (13) measures the divergence between the two densities g and f , indicating the information loss entailed when we approximate the DGP using a candidate model. Recently, Zellner (2002, p. 43) interprets d_{KL} as the difference in expected log heights of the two densities; for other divergence measures, see Rényi (1970) or Linhart and Zucchini (1986, p. 18).

The information loss in (13) depends on the model parameters θ . In practice, we evaluate (13) at $\hat{\theta}$ obtained by fitting the candidate model f with the observed sample Y . To remove the dependence of (13) on the particular sample Y , we adopt Akaike's (1985) approach to average d_{KL} across different independent samples Y drawn from the same DGP and choose a model that minimizes the expected information loss:

$$\begin{aligned} \bar{d}_{KL}(g, f; \hat{\theta}) &= E_Y \left(E_{Y^*} \left[\log \frac{g(Y^*)}{f(Y^*; \hat{\theta})} \right] \right) \\ &= E_Y(E_{Y^*}[\log(g(Y^*))]) - E_Y(E_{Y^*}[\log(f(Y^*; \hat{\theta}))]), \end{aligned}$$

where $E_Y(\cdot)$ indicates expectation with respect to the density g which generates the estimation sample, Y .

Because $E_Y(E_{Y^*}[\log(g(Y^*))])$ remains invariant across all candidate models (i.e., constant across different choices of f), it is sufficient to select the model that minimizes

$$\tilde{d}_{KL} = \tilde{d}_{KL}(g, f; \hat{\theta}) = -2E_Y(E_{Y^*}[\log(f(Y^*; \hat{\theta}))]), \quad (14)$$

where the dependence on g arises from the double expectation, and the multiplication by two is

for convenience. To derive an estimator for \tilde{d}_{KL} , we consider the Markov-switching regression model in (1) and (3) in which x_t does not contain lagged dependent variables. In the Appendix, we simplify (14) and obtain the *Markov-switching criterion*,

$$MSC = -2\log(f(Y, \hat{\theta})) + \sum_{i=1}^N \frac{\hat{T}_i(\hat{T}_i + \lambda_i K)}{\delta_i \hat{T}_i - \lambda_i K - 2}, \quad (15)$$

where $\log(f(Y, \hat{\theta}))$ is the maximized log-likelihood value, $\hat{T}_i = \text{tr}(\hat{W}_i)$, $\hat{W}_i = \text{diag}(\hat{\xi}_{1i}, \dots, \hat{\xi}_{Ti})$, $\delta_i = E[\pi_i^* / \hat{\pi}_i]$, $\lambda_i = E[(\pi_i^* / \hat{\pi}_i)^2]$, and π_i^* is the i -th element of the principal eigenvector of $P^* \pi^* = \pi^*$ for the best estimates $\theta^* = \text{vec}(\beta^*, \sigma^*, P^*) = \arg \min_{\theta} E_{Y^*}[-\log f(Y^*; \theta)]$. The subsequent remarks elaborate the properties of MSC and its implementation in practice.

Remark 3. The first term of MSC measures the lack of fit; its second term imposes a penalty for including redundant states and variables. Thus, MSC balances the trade-off between improving a model's fit to the data and achieving parsimony of the fitted model. To select the candidate model, we compute (15) for varying choices of states and variables (N, K) and retain the one that attains the smallest value.

Remark 4. In regression models without Markov switching, MSC is equivalent to both Hurvich and Tsai's (1989) criterion in finite samples and Akaike's (1973) criterion in large samples. Specifically, in regression models, $N = \delta = \lambda = 1$, and so $MSC_{N=1} = -2\log(f(Y, \hat{\theta})) + T(T + K)/(T - K - 2)$, which equals Hurvich and Tsai's (1989, p. 300) AIC_C criterion. Furthermore, by subtracting T from $MSC_{N=1}$, we obtain $-2\log(f(Y, \hat{\theta})) + 2(K + 1)\{T/(T - K - 2)\}$, which approaches Akaike's (1973) $AIC = -2\log(f(Y, \hat{\theta})) + 2(K + 1)$ in large samples. Thus, the proposed MSC generalizes the applicability of these criteria to N -state Markov-switching regression models.

Remark 5. When $N > 1$, MSC imposes penalty through $\delta_i = E[\pi_i^* / \hat{\pi}_i]$ and $\lambda_i = E[(\pi_i^* / \hat{\pi}_i)^2]$. Because the distribution of $\hat{\pi}_i$ is not known, to implement MSC, we investigate the behavior of $\bar{\delta}_i = E[\pi_i^* / \bar{\pi}_i]$ and $\bar{\lambda}_i = E[(\pi_i^* / \bar{\pi}_i)^2]$, where $\bar{\pi}_i = T^{-1} \sum_{t=1}^T s_{ti}$ and $E[\bar{\pi}_i] = \pi_i^*$. For $\bar{\delta}_i$, we invoke Jensen's inequality to obtain $\bar{\delta}_i = \pi_i^* E[1 / \bar{\pi}_i] \geq \pi_i^* / E[\bar{\pi}_i] = 1$. In other words, a lower bound for $\bar{\delta}_i$ is unity, which yields a larger value of MSC than would result from any other $\delta_i > 1$. For $\bar{\lambda}_i$, we applied Gabriel's (1959) formula for the distribution of $\bar{\pi}_i$ to compute $\bar{\lambda}_i$ for various $N \times N$ transition matrices P . These computations indicated that $\bar{\lambda}_i$ is an increasing function of the number of states N . Using these results, we set $\delta_i = 1$ and $\lambda_i = 1, N$, and N^2 to implement MSC. In Section 4, Monte Carlo simulations show that MSC with $\delta_i = 1$ and $\lambda_i = N$ performs satisfactorily.

Remark 6. The application of MSC in (15) is not specific to the EM algorithm; it can be used in conjunction with other estimation approaches. For example, one could obtain $\log(f(Y, \hat{\theta}))$ via quasi-Newton methods and find \hat{T}_i using the smoother in Hamilton (1990) or Kim (1994). Thus, the value of MSC in (15) can be computed to determine states and variables jointly.

Remark 7. We obtained average \bar{d}_{KL} to remove dependence of (13) on the estimation sample Y . Alternatively, we can consider the possibility of averaging by using a posterior density for θ and a predictive density for Y^* . This approach may provide better results in small samples, an issue that needs further investigation.

Remark 8. Bates and Granger (1969) and Leamer (1978) suggest combining multiple models rather than selecting the single best one. To this end, Burnham and Anderson (2004,

pp. 269-274) recommend computing $\Delta_k = MSC_k - MSC_{\min}$ for each model f_k relative to the model that yields the minimum MSC value, and then use the weights $w_k = \frac{\exp(-0.5\Delta_k)}{\sum_k \exp(-0.5\Delta_k)}$ to conduct multi-model inference. Furthermore, to assess degrees of confidence in alternative models, Burnham and Anderson (2002, p. 170) offer the following guidelines: Δ_k between 0-2 indicates a substantial empirical support for the model f_k ; Δ_k between 4-7 suggests considerably less support; $\Delta_k > 10$ implies essentially no empirical evidence in favor of that model (also see Raftery (1996, p. 252) for guidelines when using Bayes factors). Finally, alternative approaches for incorporating model uncertainty include forecast combinations (Timmermann 2005), Bayesian model averaging (e.g., Hoeting et al. 1999), frequentist model averaging (Hjort and Claeskens 2003), and adaptive mixing of methods (Yang 2001).

Remark 9. We note that model comparisons based on AIC are asymptotically equivalent to those based on Bayes factors when prior information is as precise as the likelihood (Kass and Raftery 1995, p. 790). When prior information is small relative to the information contained in data, the Bayesian information criterion (BIC) tends to select models with highest posterior probability. In investigating the number of states to retain in Markov-switching autoregressive models, Psaradakis and Spagnolo (2003, p. 246) conclude that BIC tends to underestimate the number of states. We encourage further research to investigate such comparisons using the proposed MSC.

Remark 10. Here we elucidate the theoretical justification for using Kullback-Liebler divergence in model selection. In information theory, Shannon's (1948) entropy is defined as $-\sum_x p(x)\log(p(x))$ for a discrete random variable with probability mass function $p(x)$.

Generalizing Shannon's entropy to two continuous density functions g and f , Kullback and

Leibler (1951) quantify “information” by defining $d_{KL} = \int g(x) \log(g(x)/f(x)) dx$ and by connecting it to R. A. Fisher’s notion of sufficient statistics. Akaike (1973, 1985) not only extends Kullback-Leibler information to quantify expected information loss (i.e., $-E[E[\log(f(x))]]$), but also deepens the connection with likelihood theory (see deLeeuw 1992) by showing that (a) the maximized log-likelihood value is a biased estimate of expected information loss, and (b) the magnitude of asymptotic bias equals the number of estimable parameters in the approximating model f . These theoretical findings furnish the justification for using KL divergence as a bridge between estimation theory and model selection, thereby unifying them under a common optimization framework (for further details, see Burnham and Anderson 2004, p. 268).

4. Monte Carlo studies

Here we describe the simulation settings as well as the model selection procedure, and then we present Monte Carlo results to illustrate the properties and performance of MSC. We also explore the applicability of MSC to Markov-switching autoregression models.

4.1 Simulation settings and model selection procedure

We investigate the following five settings:

- (i) *Markov-switching regression*: The true model consists of two states ($N^0 = 2$) and three variables in each state including an intercept. The true regression coefficients are $\beta^0 = (\beta_1^0, \beta_2^0)$, where $\beta_1^0 = (1, 2, 3)'$ and $\beta_2^0 = (4, 3, 2)'$. The explanatory variables are stored in the $T \times 3$ matrix X^0 , whose first column equals one and second and third columns are randomly drawn from a standard normal distribution. The $N^0 \times 1$ state

variable s_t^0 is a Markov chain with transition probabilities: $p_{ii}^0 = 0.95$ and $p_{ij}^0 = 0.05$ for each $i, j = 1, 2$. We obtain the dependent variable using the model in (3), $y_t = (x_t^0)' \beta^0 s_t^0 + \sigma^0 s_t^0 \varepsilon_t^0$, where x_t^0 denotes the t -th row of X^0 , $t = 1, 2, \dots, T = 250$, $\varepsilon_t^0 \sim N(0, 1)$, and $\sigma^0 = (\sigma_1^0, \sigma_2^0) = (0.5, 0.5)$.

In each state, we consider five candidate variables, which are stored in the matrix X of dimension $T \times 5$. The first three columns of X are the same as X^0 , and we randomly draw the last two columns from the standard normal distribution. We consider four candidate states (i.e., $N = 1, \dots, 4$), and the candidate regression models include up to five variables from X in a sequentially nested fashion. Thus, we have 20 possibilities (4 states by 5 variables) from which to choose the true model.

- (ii) *Markov-switching regression with small sample and high noise:* We consider two variations from the settings in (i). First, to study small sample performance, we conduct the above simulations using $T = 100$. Second, we set $\sigma_i^0 = 1$ for both $T = 100$ and 250 to understand the effect of a higher noise level.
- (iii) *Markov-switching autoregression:* We conduct the simulation in (i) for autoregressive models, where the t -th rows of X^0 and X contain $(1, y_{t-1}, y_{t-2})$, and $(1, y_{t-1}, y_{t-2}, y_{t-3}, y_{t-4})$, respectively, for $t = 5, 6, \dots, T$. The true coefficients, $\beta_1^0 = (1, 0.2, 0.3)'$ and $\beta_2^0 = (3, 0.5, 0.2)'$, satisfy the stationarity condition.
- (iv) *Markov-switching autoregression with small sample and high noise:* Analogous to (ii), we investigate two variations from the settings in (iii).
- (v) *Single state model:* We investigate the case with $N^0 = 1$ to examine whether MSC leads to spurious Markov-switching structure when the true model is a standard regression. For

fixed regressors, we use $\beta^0 = (1, 2, 3)'$; for autoregression, $\beta^0 = (1, 0.2, 0.3)'$.

We conduct 1000 repetitions in each of the above settings to assess how often MSC selects the true model.

We employ the following model selection procedure for each of the 20 state-variable combinations $\{(N, K): N = 1, \dots, 4, K = 1, \dots, 5\}$. First, we choose initial parameter values using the K -means method (MacQueen 1967) to classify observations in the matrix (y, X) into N states. Then, we apply the EM algorithm to estimate the Markov-switching regression model. Next, we compute MSC in (15). We also constrain the term $(\delta_i \hat{T}_i - \lambda_i K - 2)$ in (15) to exceed unity in each realization to ensure positive penalty. Finally, we select the model that yields the smallest MSC value across all the 20 state-variable combinations.

4.2 Monte Carlo results

Here we present one figure and five tables to illustrate the accuracy and performance of MSC. In addition, we substantiate the claim that AIC overestimates the number of states in Markov-switching models.

Accuracy of MSC. We assess the accuracy of MSC by computing its proximity to the true KL distance. To this end, we estimate the true KL distance \tilde{d}_{KL} in (14) using the three steps: (a) randomly draw an estimation sample Y to obtain the EM estimates $\hat{\theta}$; (b) draw a holdout sample Y^* to evaluate $\log(f(Y^*; \hat{\theta}))$; (c) perform 100 repetitions with different holdout samples Y^* 's to estimate $E_{Y^*}[\log(f(Y^*; \hat{\theta}))]$. We repeat the steps (a)-(c) 100 times for different estimation samples Y to evaluate the double-expectation in (14).

Figure 1 presents the proximity plots for the MSC values from (15) using $\lambda_i = N$ for the setting (i) in Section 4.1. Panel A presents the results for state selection. It shows that both MSC

and \tilde{d}_{KL} achieve their minimum at the true number of states, i.e., $N^0 = 2$. Furthermore, MSC and \tilde{d}_{KL} are close when $N \leq N^0$, while MSC exceeds \tilde{d}_{KL} when $N > N^0$. In other words, MSC approximates \tilde{d}_{KL} reasonably well and imposes a larger penalty when the number of states exceeds those in the data generating process. This larger penalty mitigates overestimation of the number of states. Panel B, which presents the results for variable selection, depicts that MSC and \tilde{d}_{KL} are uniformly close. Thus, for the purposes of model selection, the proposed MSC reasonably approximates the Kullback-Leibler distance.

Performance of MSC. We investigate the simultaneous selection of states and variables in Markov-switching regression models (see the setting (i) in Section 4.1). We assess the performance of a criterion by the relative frequency of selecting various states and variables, while the measure of accuracy is how often the criterion selects the correct number of states that were used in the DGP. Table 1 reports the frequency of correct state and variable selection using MSC with $\lambda = 1$, N and N^2 . (Note that the subscript i on λ is suppressed in the rest of the paper.) For $\lambda = 1$, Panel A shows that incorrect model selection is asymmetric. Specifically, the zeros in Panel A reveal that $MSC_{\lambda=1}$ never underestimates the number of states or variables. But, $MSC_{\lambda=1}$ correctly selects two states 360 times and three variables 666 times out of 1000 occasions. Consequently, the joint frequency of selecting the correct states and variables is only 30.9%. Despite this unsatisfactory performance, we note that the *conditional* frequency of variable selection $(0, 0, \frac{309}{360}, \frac{29}{360}, \frac{22}{360})$ is satisfactory. This finding can be explained using Panel B of Figure 1, which shows that MSC estimates the true KL distance accurately when the number of states is known. More importantly, this finding underscores the insight that the model selection performance can be improved if we determine the true states accurately. To this end,

we investigate the performance of MSC with $\lambda = N$ and N^2 as stated in Remark 5.

For $\lambda = N$, Panel B indicates a marked improvement in model selection performance. Specifically, $MSC_{\lambda=N}$ correctly selects the two-state model in each of the 1000 realizations. We explain this improvement using Panel A of Figure 1, which exhibits that $MSC_{\lambda=N}$ imposes larger penalty than the KL distance, thus mitigating the tendency to fit too many states. Moreover, we find diminishing returns to further increases in penalty via $\lambda = N^2$ because performance improves marginally beyond that due to $MSC_{\lambda=N}$ (see Panel C in Table 1).

Table 2 demonstrates the robustness of these findings via the simulation setting (ii). When we increase the noise level from $\sigma_i^0 = 0.5$ to 1, the performance of $MSC_{\lambda=1}$ further deteriorates. The joint frequency of correctly selecting both the states and variables decreases from 309 to 124. In contrast, $MSC_{\lambda=N}$ and $MSC_{\lambda=N^2}$ perform well, as evidenced by the small decrease in the joint frequency from 992 to 981 and from 1000 to 998, respectively. In other words, these small decreases indicate that the performance of both the criteria do not deteriorate substantially as the noise level increases. We observe qualitatively similar findings when the sample size decreases from $T = 250$ to 100. It is worth noting that $MSC_{\lambda=N}$ is less sensitive to noise level in small samples than $MSC_{\lambda=N^2}$. Specifically, as the noise level increases for $T = 100$, the joint selection frequency of $MSC_{\lambda=N}$ decreases by 4.6% (from 951 to 907) compared to 13.9% for $MSC_{\lambda=N^2}$ (from 861 to 741). In other words, $MSC_{\lambda=N}$ outperforms $MSC_{\lambda=1}$ and $MSC_{\lambda=N^2}$ when both the sample size is small and the signal is weak.

We repeat the above analyses for the Markov-switching *autoregression* models described in the setting (iii). Table 3 reports the joint selection frequency by MSC with $\lambda = 1, N$, and N^2 in 1000 realizations. As before, incorrect model selection is asymmetric; $MSC_{\lambda=1}$ never understates

the number of states and seldom underestimates the number of variables. $MSC_{\lambda=N}$ outperforms $MSC_{\lambda=1}$ with 979 correct selections out of 1000 occasions (see Panel B in Table 3). This superior performance is due to the penalty imposed by $MSC_{\lambda=N}$, which mitigates the tendency to fit excessive states. We can marginally improve this performance from 979 to 984 by using a stronger penalty via $\lambda = N^2$ (compare Panels B and C in Table 3).

Table 4 shows that these findings are robust to various scenarios in the setting (iv). As the noise level increases in large samples, $MSC_{\lambda=1}$ performs poorly, whereas $MSC_{\lambda=N}$ and $MSC_{\lambda=N^2}$ perform satisfactorily as evidenced by smaller decreases in the joint frequency. We obtain qualitatively similar results for the small sample case. Moreover, $MSC_{\lambda=N}$ is less sensitive to the noise level in small samples than $MSC_{\lambda=N^2}$; for example, the correct selection frequency of $MSC_{\lambda=N}$ decreases by 46% (from 744 to 402) compared to 99.4% for $MSC_{\lambda=N^2}$ (from 171 to 1). Thus, $MSC_{\lambda=N}$ outperforms $MSC_{\lambda=1}$ and $MSC_{\lambda=N^2}$ when both the sample size is small and the signal is weak.

Single-state model. While MSC detects Markov switching when it does exist, can MSC reject Markov switching when it does not exist? To answer this question, we examine the setting (v) and use MSC to select the number of states (but not variables). In Table 5, Panels A and B show the correct selection frequency for the fixed regressor and autoregression settings, respectively. We find that $MSC_{\lambda=1}$ performs poorly regardless of the noise level or the sample size. However, the last two columns indicate that $MSC_{\lambda=N}$ and $MSC_{\lambda=N^2}$ correctly select a single-state model more than 90% of the occasions. Thus, $MSC_{\lambda=N}$ and $MSC_{\lambda=N^2}$ do not yield spurious Markov-switching structure when the true model is a standard regression.

We close this section by substantiating the claim in the *Introduction* that the AIC-based

estimate of KL divergence retains too many states and variables. We compute $AIC = -2\log f(y, \hat{\theta}) + 2d$, where $d = (NK + N^2)$ denotes the number of free parameters in θ . For the sake of illustration, we use the low noise and large sample setting (ii), which is favorable for AIC. Table 6 reveals that AIC selects more states and variables than in the DGP and that the correct joint selection frequency is only 48.1%. Thus, by using AIC in practical applications, users stand about equal chance to retain a correct or an incorrect model; when it is the latter, they would fit spurious regressions in non-existent states. We next present two empirical examples to illustrate the usefulness of $MSC_{\lambda=N}$ in practice.

5. Empirical examples

We first study the business cycle in the US economy and then the effectiveness of media advertising.

5.1 *U.S. real GNP growth*

Hamilton (1989) was first to formulate the Markov-switching autoregression model to capture business cycles in real GNP. In his formulation, the mean GNP growth rate switches between two states: recessions and expansions. Hansen (1992) extends this model to allow both the mean growth rate and the autoregressive coefficients to switch between states. We study this extended model, which is given by equations (1) and (3), where $x_t = (1, y_{t-1}, y_{t-2}, y_{t-3}, y_{t-4})'$ and y_t is quarterly real GNP growth in chained 1996 dollars. We use seasonally adjusted data that span the period 1947:1 through 2002:4 (see <http://www.bea.doc.gov>). We exclude 16 quarterly observations (1999:1 to 2002:4) from the estimation sample and use these excluded observations to evaluate one-quarter-ahead forecasts. The estimation sample comprises $T = 203$ observations because we also exclude 5 observations for computing the growth rate and the initial lagged

values.

We apply the EM algorithm described in Section 2 to these data, and consider various state-variable combinations (N, K) , where $N = 1, \dots, 4$ and $K = 1, \dots, 5$. We estimate 20 different N -state Markov-switching autoregression models and compute the two estimates of KL divergence: AIC and $MSC_{\lambda=N}$. Based on the minimum AIC value, we would select a model with $N^* = 4$ and $K^* = 5$, which is the largest model in this set of 20 candidate models. This finding is consistent with the simulation evidence (see Table 6), which reveals AIC's tendency to select more states and variables than necessary.

On the other hand, the minimum value of $MSC_{\lambda=N}$ yields $N^* = 3$ and $K^* = 1$, indicating the retention of the three-state model with no autoregressive lags (i.e., intercepts only). Table 7 reports the parameter estimates for this retained model, which identifies one recessionary state ($i = 1$) and two expansionary states ($i = 2, 3$). The estimated decline in real GNP during recessions is -0.10% per quarter; the mean growth rates during the two expansion states are 1.50% and 0.85% per quarter.

In Figure 2, we present the estimated smoothed probability sequence $\hat{\xi}_i = (\hat{\xi}_{1i}, \dots, \hat{\xi}_{Ti})'$ based on (7) and overlay it with the recessionary periods (in gray bars) noted by the National Bureau of Economic Research. Panel A shows that the estimated probability of recession reasonably matches the actual recessions. Panels B and C display the two types of expansions. The first type occurred exclusively before 1984, while the second occurred mostly during the 80s and the 90s. Because $\hat{\sigma}_3 = 0.42 < \hat{\sigma}_2 = 0.91$, the recent expansionary state ($i = 3$) exhibits lower volatility than the previous one ($i = 2$). This finding supports the phenomenon of great moderation — first discovered by Kim and Nelson (1999a) and McConnell and Perez-Quiros (2000) — which is characterized by a reduction in the variance of economic growth since 1984.

We compare the forecasting performance of this retained model to that of a benchmark model that specifies $\ln(\text{GNP})$ as a random walk with constant drift. Over the period 1999-2002, the mean squared forecast errors are 0.351 and 0.433 for the retained model and the random walk model, respectively. In addition, the mean absolute forecast error was 0.539 for the retained model and 0.546 for the random walk. The retained three-state Markov-switching model performs well because it adapts to the recession in 2001, whereas the random walk model does not (see Figure 2).

5.2 Advertising effectiveness

In marketing, brand managers commonly use the advertising model, $y_t = \beta^{(0)} + \beta^{(1)}z_t + \beta^{(2)}y_{t-1} + \varepsilon_t$, to determine the effectiveness of advertising (Bucklin and Gupta 1999, p. 262), where y_t denotes brand sales at time t , z_t represents advertising spending, and ε_t is the normal error term. The coefficient $\beta^{(1)}$ measures the effectiveness of current advertising; the coefficient $\beta^{(2)}$, known as the carryover effect, captures the cumulative impact of past advertising reflected in the attained sales y_{t-1} (see, e.g., Palda 1964, p. 13). We extend this advertising model by incorporating regime shifts so that the parameter vector $\beta_i = (\beta_i^{(0)}, \beta_i^{(1)}, \beta_i^{(2)})'$ is specific to each regime $i = 1, \dots, N$. This extension marks the first application of Markov-switching models in the advertising literature (see Feichtinger, Hartl and Sethi 1994, Mantrala 2002, Naik and Raman 2003).

We apply this extended model to Lydia Pinkham company's annual sales and advertising data from 1914 through 1960 (Palda 1964). This classic data set exhibits a few unique features: relatively stable product design during this period; advertising primarily affects sales, given the absence of channel members or sales force; and the lack of close competitors. These market conditions comport with the above advertising model. Furthermore, after the second World War

ended, Lydia Pinkham management demonstrated the product's efficacy to the Federal Trade Commission (FTC), which permitted them to make stronger claims in their advertising copy. Moreover, they switched from pure newspaper advertising to a mix of multiple media, which comprised newspaper, magazine, radio, and even television. (See Palda 1964, pp. 25-26 for details.)

Given these changes in market conditions, we consider the possibility of a distinct post-war regime(s) by estimating various Markov-switching models with state-variable combinations (N, K) , for $N = 1, \dots, 4$ and $K = 1, \dots, 3$. Then we compute AIC and $MSC_{\lambda=N}$ for each combination. AIC selects a model with $N^* = 3$ states, which, given the simulation results in Table 6, is likely to be more than necessary. In contrast, $MSC_{\lambda=N}$ retains two states (i.e., $N^* = 2$). The smoothed probabilities $\hat{\xi}_1 = (\hat{\xi}_{11}, \dots, \hat{\xi}_{T1})'$ indicate that the first state persisted from 1914 through 1945, whereas the second state lasted from 1946-1960. This regime switch coincided with the FTC's approval of stronger copy and the beginning of multiple media spending.

Table 8 shows the different estimates of advertising effectiveness and carryover effects for the pre- and post-war regimes. Specifically, advertising is more effective in the post-war era ($\hat{\beta}_2^{(1)} = 1.17 > \hat{\beta}_1^{(1)} = 0.43$) due to stronger copy and multiple media. In addition, the carryover effect is smaller in the post-war era ($\hat{\beta}_2^{(2)} = 0.27 < \hat{\beta}_1^{(2)} = 0.53$), given the shorter duration for the impact of past advertising to accumulate. Thus, these new findings are not discernible from the standard regression model of advertising.

6. Concluding remarks

Markov-switching regression models provide an analytical framework to study both shifts in regimes and the differential impact of explanatory variables across regimes (or states).

