

Bayesian order selection in heterogeneous hidden Markov models

Journal:	<i>Journal of Computational and Graphical Statistics</i>
Manuscript ID	JCGS-22-216
Manuscript Type:	Original Article
Keywords:	Bayesian Methods, Longitudinal Data, Markov Chain Monte Carlo (MCMC)

SCHOLARONE™
Manuscripts

Bayesian order selection in heterogeneous hidden Markov models

Yudan Zou, Yiqi Lin and Xinyuan Song*

Department of Statistics, The Chinese University of Hong Kong

Abstract

Hidden Markov models (HMMs) are valuable tools for analyzing longitudinal data due to their capability to describe dynamic heterogeneity. Conventional HMMs typically assume that the number of hidden states (i.e., the order of HMMs) is known or predetermined through criterion-based methods. However, prior knowledge about the order is often unavailable, and a pairwise comparison using criterion-based methods becomes increasingly tedious and computationally demanding when the model space enlarges. A few studies have considered simultaneously performing order selection and parameter estimation under the frequentist framework. Still, they focused only on homogeneous HMMs and thus cannot accommodate situations where potential covariates affect the between-state transition. This study proposes a Bayesian double-penalized (BDP) procedure to conduct a simultaneous order selection and parameter estimation for heterogeneous HMMs. We develop a novel Markov chain Monte Carlo algorithm coupled with an efficient adjust-bound reversible jump scheme to address the challenges in updating the order. Simulation studies show that the proposed BDP procedure considerably outperforms the commonly used criterion-based methods. An application to the Alzheimer’s Disease Neuroimaging Initiative study further confirms the utility of the proposed method.

Keywords: Bayesian method; double penalization; dynamic heterogeneity; longitudinal data; order selection

*Department of Statistics, The Chinese University of Hong Kong, Shatin, Hong Kong, E-mail: xysong@cuhk.edu.hk

1 Introduction

Hidden Markov Models (HMMs) have broad applications in medical, behavioral, social, and psychological sciences, wherein heterogeneous longitudinal data are frequently collected and analyzed. HMMs consist of two parts: a transition model to characterize the dynamic transition process between hidden states and a conditional regression (emission) model to examine state-specific covariate effects on the response of interest.

Conventional HMMs typically assume that the number of hidden states (i.e., order of HMM) is known or predetermined through criterion-based methods, such as the Akaike's information criterion (AIC, Akaike, 1974) and Bayesian information criterion (BIC, Schwarz, 1978). However, despite their successful applications in many substantive studies (e.g., Celeux and Durand, 2008; Ip et al., 2013; Song et al., 2017), these criterion-based methods conduct pairwise comparisons among candidate models, which could become increasingly tedious and computationally intensive when the model space is ample. Moreover, these procedures perform estimation in two stages: choosing the order in the first stage and estimating the parameter of the selected model in the second stage, and thus may not be as effective as single-stage approaches.

Penalization methods are valuable alternatives to their criterion-based counterparts in estimating HMMs with unknown order. For example, Mackay (2002) proposed to penalize small state proportions and obtained a consistent order estimate of HMMs. Chen and Khalili (2008) pointed out the necessity of preventing the second type of overfitting induced by similar-density components in finite mixture models and suggested a double penalization procedure to avoid the two types of overfitting simultaneously. Hung et al. (2013) introduced this double penalized method to non-regression Gaussian HMMs. Lin and Song (2022) further adapted the double penalization idea into regression-based HMMs with an extended expectation-maximization (EM) algorithm and established order selec-

tion consistency and algorithm convergence. In addition, Zhou et al. (2020) considered continuous-time HMMs with unknown order and proposed a modified penalized maximum likelihood estimation approach. Apart from the penalization methods developed under the frequentist framework, Liu and Song (2020) also considered a Bayesian approach by regarding the order of HMMs as a random variable and updating it together with other parameters using the reversible jump Markov chain Monte Carlo (MCMC) algorithm (Green, 1995). Nevertheless, the available methods in the frequentist and Bayesian frameworks are either criterion-based two-stage approaches or focused only on finite mixture models or homogeneous HMMs.

This study proposes a novel Bayesian double penalized (BDP) procedure to conduct a simultaneous order selection and parameter estimation for heterogeneous HMMs. The procedure includes two penalties. The first is a lower bound imposed on the summation of mixing proportions to prevent states with near-zero initial probabilities. The second is a least absolute shrinkage and selection operator (lasso)-type penalty introduced to the distance between regression coefficients to avoid states with nearly identical parameters. We develop a hybrid MCMC algorithm that integrates the data augmentation, Gibbs sampler, forward filtering backward sampling (FFBS, Baum et al., 1970), and the Metropolis-Hastings (MH) algorithm. In particular, we offer an efficient adjust-bound reversible jump (ABRJ) sampling scheme to address the challenges of updating the order in implementing the MCMC algorithm. Simulation studies in Section 5 demonstrate that the proposed BDP procedure considerably outperforms the commonly used AIC and BIC in order selection accuracy. In addition, by setting a sizable upper bound of the order (e.g., 200), the proposed method allows sufficient flexibility in estimating the order of HMMs and thus can accommodate the case where many states exist. Last but not least, the BDP procedure accomplishes order selection and parameter estimation in a single stage. By contrast, criterion-based approaches perform pairwise comparison and parameter estimation on a two-stage basis,

and the related computational burden dramatically increases when the candidate model space enlarges. Therefore, the proposed BDP procedure is also superior to the existing methods in terms of computational efficiency.

The rest of this article is organized as follows. Section 2 describes the model and related identifiability issues. Sections 3 and 4 present the BDP procedure and specific sampling schemes. Section 5 evaluates the empirical performance of the proposed method through simulation studies, and Section 6 reports an application to the Alzheimer's Disease Neuroimaging Initiative (ADNI) study. Section 7 concludes the paper. The technical details are provided in the Supplementary Material.

2 Model

Let $\mathbf{Y} = (\mathbf{y}_1', \dots, \mathbf{y}_n')'$, where $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})'$, and y_{it} is the response of subject i at time t ; $\mathbf{X} = (\mathbf{X}_1', \dots, \mathbf{X}_n')'$, where $\mathbf{X}_i = (\mathbf{x}_{i1}', \dots, \mathbf{x}_{iT}')'$, and \mathbf{x}_{it} is the covariate vector of subject i at time t ; $\mathbf{D} = (\mathbf{D}_1', \dots, \mathbf{D}_n')'$, where $\mathbf{D}_i = (\mathbf{d}_{i1}', \dots, \mathbf{d}_{iT}')'$, and \mathbf{d}_{it} is another covariate vector of subject i at time t , and the elements of \mathbf{d}_{it} can be distinct or overlapped with those of \mathbf{x}_{it} ; $\mathbf{Z} = (\mathbf{Z}_1', \dots, \mathbf{Z}_n')'$, where $\mathbf{Z}_i = (Z_{i1}, \dots, Z_{iT})'$, and Z_{it} is the hidden state of subject i at occasion t , which follows a first-order Markov chain and takes the values of $\{1, \dots, K\}$. Given the hidden state Z_{it} , y_{it} is assumed to be independent for all i and t , and is formulated through the conditional regression model as follows:

$$[y_{it}|Z_{it} = s] = \beta_s' \mathbf{x}_{it} + \delta_{it}, \quad (1)$$

where $\mathbf{x}_{it} = (1, x_{it1}, \dots, x_{it,p-1})'$ be the $p \times 1$ vector of covariates, β_s is the $p \times 1$ vector of state-specific regression coefficients, δ_{it} is the residual term independent of \mathbf{x}_{it} , and $[\delta_{it}|Z_{it} = s] \sim N(0, \psi_s)$.

Given that hidden states typically have a natural ranking and real meanings in most practical situations, we assume that hidden states $\{1, \dots, K\}$ are ordered. The hidden

transition process is then formulated by $Z_{i1} \sim \text{multinomial}(\pi_1, \dots, \pi_K)$ such that $0 \leq \pi_s \leq 1$ and $\sum_{s=1}^K \pi_s = 1$, and a continuation-ratio logit model (Agresti, 2003) as follows: for $t = 2, \dots, T, s = 1, \dots, K-1, u = 1, \dots, K$:

$$\log \left(\frac{P_{itus}}{P_{itu,s+1} + \dots + P_{ituK}} \right) = \zeta_{us} + \boldsymbol{\alpha}' \mathbf{d}_{it}, \quad (2)$$

where $P_{itus} = P(Z_{it} = s | Z_{i,t-1} = u)$, ζ_{us} is a transition-specific intercept, $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_q)'$ is a $q \times 1$ vector of regression coefficients. Let $\vartheta_{itus} = P(Z_{it} = s | Z_{it} \geq s, Z_{i,t-1} = u)$. Then, we can easily check that $\log \left(\frac{P_{itus}}{P_{itu,s+1} + \dots + P_{ituK}} \right) = \log \left(\frac{P(Z_{it}=s | Z_{i,t-1}=u)}{P(Z_{it} > s | Z_{i,t-1}=u)} \right) = \text{logit}(\vartheta_{itus})$, which is the log conditional odds of transitioning to the s th state instead of a higher state given $Z_{i,t-1} = u$. Therefore, the transition model (2) can be equivalently rewritten as $\text{logit}(\vartheta_{itus}) = \zeta_{us} + \boldsymbol{\alpha}' \mathbf{d}_{it}$.

Equations (1) and (2) define a heterogeneous HMM, under which some time-variant or baseline covariates affect the between-state transition. Let $\boldsymbol{\theta}$ be the vector containing all the regression and variance parameters in the proposed model. Then, the complete-data log-likelihood function of the proposed model is given by

$$\begin{aligned} \log p[\mathbf{Y}, \mathbf{X}, \mathbf{D}, \mathbf{Z} | \boldsymbol{\theta}] &= \sum_{i=1}^n [\log p(\mathbf{y}_i | \mathbf{X}_i, \mathbf{Z}_i, \boldsymbol{\theta}) + \log p(\mathbf{Z}_i | \mathbf{D}_i, \boldsymbol{\theta})] \\ &= \sum_{i=1}^n \sum_{t=1}^T \log p(y_{it} | \mathbf{x}_{it}, Z_{it}, \boldsymbol{\theta}) + \sum_{i=1}^n \sum_{t=2}^T \log p(Z_{it} | Z_{i,t-1}, \mathbf{d}_{it}, \boldsymbol{\theta}) + \sum_{i=1}^n \log p(Z_{i1} | \boldsymbol{\theta}) \quad (3) \\ &= \sum_{i=1}^n \sum_{t=1}^T \log p(y_{it} | \mathbf{x}_{it}, Z_{it}, \boldsymbol{\theta}) + \sum_{i=1}^n \sum_{t=2}^T \log(P_{itZ_{i,t-1}Z_{it}}) + \sum_{i=1}^n \log(P_{i10Z_{i1}}) \end{aligned}$$

where

$$\begin{aligned} P_{i10s} &= \pi_s, \quad s = 1, \dots, K, \\ P_{itu1} &= \frac{\exp(a_{itu1})}{1 + \exp(a_{itu1})}, \quad P_{ituK} = \prod_{j=1}^{K-1} \frac{1}{1 + \exp(a_{ituj})}, \\ P_{itus} &= \frac{\exp(a_{itus})}{1 + \exp(a_{itus})} \prod_{j=1}^{s-1} \frac{1}{1 + \exp(a_{ituj})}, \quad s = 2, \dots, K-1 \end{aligned} \quad (4)$$

with $a_{itus} = \zeta_{us} + \boldsymbol{\alpha}' \mathbf{d}_{it}$, for $t = 2, \dots, T, u = 1, \dots, K, s = 1, \dots, K-1$.

The proposed model is unidentifiable due to the label switching problem. Label switching arises because a random permutation of state labels does not change the likelihood function, which leads to a multi-modal posterior under a symmetric prior distribution. We address the problem by introducing the cluster ordering procedure proposed by Zhou et al. (2020) to sort the multidimensional parameters in the conditional regression model (1), which satisfies the atom property mentioned in Manole and Khalili (2021). We define the cluster ordering procedure in the context of the proposed model as follows.

Definition 2.1 A cluster ordering procedure is a mapping $\alpha_\beta: \{\beta_1, \dots, \beta_K\} \rightarrow \{\beta_{(1)}, \dots, \beta_{(K)}\}$, such that

$$\begin{aligned} \beta_{(1)} &= \arg \max_{\beta_j: j=1, \dots, K} \|\beta_j\|_2 \\ \beta_{(k)} &= \arg \min_{\beta_j \neq \beta_{(i)}, i=1, \dots, k-1} \|\beta_j - \beta_{(k-1)}\|_2, \quad k = 2, \dots, K, \end{aligned} \quad (5)$$

where $\|\cdot\|_2$ denotes the L_2 norm.

The cluster ordering procedure guarantees that the state labels are uniquely determined and induces a set of differences $\eta_1 = \beta_{(1)}$, and $\eta_k = \beta_{(k)} - \beta_{(k-1)}$ for $k = 2, \dots, K$. Recall that the hidden states are assumed ordered; we can then rewrite the conditional regression model (1) as follows:

$$[y_{it}|Z_{it} = s] = \sum_{k=1}^s (\eta'_k x_{it}) + \delta_{it}. \quad (6)$$

Hence, by constructing the bijective mapping between β_k and η_k , the complete-data log-likelihood function can be formulated as

$$\sum_{i=1}^n \sum_{t=1}^T (y_{it} - \sum_{k=1}^s \eta'_k x_{it})^2 - \sum_{i=1}^n \sum_{t=2}^T \log(P_{itus}) - \sum_{i=1}^n \log(P_{i10s}), \quad (7)$$

which facilitates the double penalization in the next section.

3 Bayesian double penalized (BDP) procedure

Two types of over-fitting exist in performing parameter estimation based on (7). The first type arises when some hidden states are almost empty and thus leading to near-zero mixing probabilities. The second type appears when two or more states have similar emission densities resulting in nearly identical parameter values. Chen and Khalili (2008) proposed a double penalized method to address the aforementioned over-fitting problems in the context of a finite mixture model. Ye et al. (2019) extended their method to a finite mixture of varying coefficient models. Manole and Khalili (2021) developed the Group-Sort-Fuse (GSF) procedure for order selection and parameter estimation in multidimensional finite mixture models. For HMMs, Hung et al. (2013) first considered the order selection problem and proposed a new double penalized procedure. Lin and Song (2022) extended the GSF procedure to regression-based HMMs. However, all the above methods focused either on finite mixture models or homogeneous HMMs and thus did not apply to the present heterogeneous HMMs. Moreover, they are developed in the frequentist framework, and their Bayesian versions have never been considered in the literature. This study aims to fill the gap and proposes a novel double penalized method under the Bayesian framework for the simultaneous order selection and parameter estimation of heterogeneous HMMs.

To prevent the first type of overfitting, we impose a lower bound on π_k to ensure the existence of a proper partition and avoid nearly empty states or near-zero mixing proportions. To address the second type of overfitting, we impose penalization on the norm of the discrepancy between different coefficient vectors. Manole and Khalili (2021) pointed out that the ordering procedure considerably outperforms the naive approach that penalizes all pairwise differences between β_k when many hidden states exist. Therefore, instead of naively penalizing all $\binom{K}{2}$ pairwise differences between β_k , $k = 1, \dots, K$, we merely penalize the L_2 -norm of $K - 1$ consecutive differences η_k . Specifically, the double

penalty is introduced to the complete-data log-likelihood function as follows:

$$\sum_{i=1}^n \sum_{t=1}^T (y_{it} - \sum_{k=1}^s \boldsymbol{\eta}'_k \mathbf{x}_{it})^2 - \sum_{i=1}^n \sum_{t=2}^T \log(P_{itus}) - \sum_{i=1}^n \log(P_{i10s}) - P(\boldsymbol{\theta}), \quad (8)$$

where

$$P(\boldsymbol{\theta}) = (1 - c_K) \sum_{k=1}^K \log(\pi_k) + \sum_{k=2}^K \gamma_k \|\boldsymbol{\eta}_k\|_2, \quad (9)$$

in which c_K and $\{\gamma_2, \dots, \gamma_K\}$ are tuning parameters.

In Equation (9), the first term introduces a penalty to $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)'$, whereas the second is a lasso-type penalty. $\boldsymbol{\eta}_k$ is a vector of the regression parameters in the k th state, and the tuning parameter γ_k is state-specific and penalizes the entire vector $\boldsymbol{\eta}_k$ rather than its elements. Thus, the second penalty is adaptive and group-wise, denoted as a modified adaptive group lasso (MAGlasso). Lin and Song (2022) considered a similar form of double penalization consisting of a lower bound penalty for mixing probabilities and the smoothly clipped absolute deviations (SCAD) penalty for the regression parameters in the context of regression-based HMMs. However, their method focused on homogeneous HMMs, and their EM-type algorithm poses theoretical and computational challenges in analyzing heterogeneous HMMs. The Bayesian method, however, enables us to implement the double penalization by introducing appropriate priors to relevant parameters. First, we assign a symmetric Dirichlet prior distribution to $\boldsymbol{\pi}$, denoted as $(\pi_1, \dots, \pi_K) \sim \text{Dir}(c_1, \dots, c_K)$, where the concentration parameter $c_K = c \frac{n}{K}$ and $c > 0$ is a preassigned constant, to avoid near-zero mixing proportions. With such a prior specification, we have the following proposition:

Proposition 3.1 *Suppose $Z_i \sim \text{categorical}(\pi_1, \dots, \pi_K)$, $i = 1, \dots, n$, with $0 \leq \pi_s \leq 1$, $\sum_{s=1}^K \pi_s = 1$, and $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K) \sim \text{Dir}(c_K, \dots, c_K)$, where $c_K = c \frac{n}{K}$, $c > 0$ is a constant. Then, we have*

$$E(\pi_s | \mathbf{Z}) \geq \frac{c}{c+1} \frac{1}{K}, \quad s = 1, \dots, K. \quad (10)$$

The derivation of Proposition 3.1 is provided in Appendix 1 of Supplementary Material. This proposition ensures that the conditional mean of each element of $\boldsymbol{\pi}$ is lower bounded by $\frac{c}{c+1} \frac{1}{K}$, thereby preventing near-zero probabilities or nearly empty states. The constant c can be determined according to the degree of penalty required for specific problems. The lower bound of $E(\pi_s|\mathbf{Z})$ is close to $\frac{1}{K}$ when c increases while it tapers off when c approaches zero. Typically, c around 0.5 can effectively prevent near-zero π_s . Alternatively, one can regard c as another tuning parameter and update it in the MCMC algorithm. However, our numerical results show that this data-driving method increases computational complexity but performs similarly to the approach that fixes c in the interval of $(0, 1)$.

Second, we use MAGlasso to tackle the second type of overfitting. Park and Casella (2008) introduced the Bayesian lasso to achieve the shrinkage on regression coefficients in a fully Bayesian framework. The basic idea of the Bayesian lasso is to penalize $\boldsymbol{\eta}_k$ by imposing a conditional Laplace prior on $\|\boldsymbol{\eta}_k\|_2$ as follows:

$$P(\boldsymbol{\eta}_k|\psi_k) = \frac{\gamma_k}{2\sqrt{\psi_k}} \exp\left(-\frac{\gamma_k}{\sqrt{\psi_k}}\|\boldsymbol{\eta}_k\|_2\right), \quad k = 2, \dots, K. \quad (11)$$

Then, the proposed model can be formulated through the following hierarchical representation: for $s = 1, \dots, K$,

$$\begin{aligned} [y_{it}|Z_{it} = s, \boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_s, \psi_s] &\sim N\left(\sum_{k=1}^s (\boldsymbol{\eta}'_k \mathbf{x}_{it}), \psi_s\right), \\ [\boldsymbol{\eta}_s|\psi_s, \tau_s^2] &\sim N(\mathbf{0}, \psi_s \tau_s^2 \mathbf{I}_p), \quad \psi_s^{-1} \stackrel{\text{ind}}{\sim} \text{Gamma}(a_{\psi s0}, b_{\psi s0}), \\ \tau_s^2 &\stackrel{\text{ind}}{\sim} \text{Gamma}\left(\frac{p+1}{2}, \frac{\gamma_s^2}{2}\right), \quad \gamma_k^2 \stackrel{\text{ind}}{\sim} \text{Gamma}(a_{\gamma k0}, b_{\gamma k0}), \quad k = 2, \dots, K, \end{aligned} \quad (12)$$

where $\mathbf{0}$ is the vector of zero elements, \mathbf{I}_p is the p -dimensional identity matrix, $a_{\psi s0}, b_{\psi s0}, a_{\gamma k0}$, and $b_{\gamma k0}$, are hyperparameters whose values are prespecified.

Proposition 3.2 *Under the hierarchical model (12), the conditional prior distribution of $\boldsymbol{\eta}_k$ has the form of (11).*

The derivation of Proposition 3.2 is provided in Appendix 1 of Supplementary Material. The MAGlasso procedure aims to update the tuning parameters by exploiting the data,

thereby automatically imposing large penalties on unimportant coefficients. This target can be naturally achieved by introducing dispersed priors with small hyperparameters $a_{\gamma k0}$ and $b_{\gamma k0}$. The degree of dispersion of the gamma priors determines the magnitudes of penalties imposed on unimportant components (Park and Casella, 2008). Typically, setting $a_{\gamma k0}$ to a positive integer (e.g., 1) and $b_{\gamma k0}$ to a small value (e.g., 0.1 or 0.01) can induce a dispersed gamma prior (Guo et al., 2012; Kang et al., 2019). With this prior specification, we can derive the posterior distribution of the tuning parameters, which have the following forms:

$$[\tau_s^{-2}|\cdot] \sim \text{Inverse-Gaussian} \left\{ \sqrt{\frac{\gamma_s^2 \psi_s}{\|\boldsymbol{\eta}_s\|_2^2}}, \gamma_s^2 \right\}, \quad (13)$$

$$[\gamma_s^2|\cdot] \sim \text{Gamma}(a_{\gamma s0} + \frac{p+1}{2}, b_{\gamma s0} + \frac{\tau_s^2}{2}). \quad (14)$$

If $\|\boldsymbol{\eta}_s\|_2$ is significant, τ_s^2 tends to be large based on (13). Then, the corresponding γ_s is dominated by τ_s^2 based on (14). On the contrary, if $\|\boldsymbol{\eta}_s\|_2$ is insignificant, τ_s^2 tends to be small, and the related γ_s is then dominated by the dispersed prior.

Considering that the Bayesian lasso does not shrink coefficients precisely to zero, we need a criterion to quantify the closeness of $\|\boldsymbol{\eta}_s\|_2$ to a zero vector. Based on the specification of (11), we can show that $P(\boldsymbol{\eta}_s|\mathbf{Y}, \mathbf{Z}, \mathbf{Z}, \boldsymbol{\theta}) \sim N(\boldsymbol{\eta}_s^*, \boldsymbol{\Sigma}_s^*)$, where $\boldsymbol{\eta}_s^*$ and $\boldsymbol{\Sigma}_s^*$ are provided in Appendix 2 of Supplementary Material. Therefore, the squared Mahalanobis distance $d_s^2 = (\boldsymbol{\eta} - \boldsymbol{\eta}_s^*)' \boldsymbol{\Sigma}_s^{*-1} (\boldsymbol{\eta} - \boldsymbol{\eta}_s^*) \sim \chi_p^2$ determines a hyper-ellipse density contour centered at $\boldsymbol{\eta}_s^*$. In this study, we adopt the 95% highest posterior credible region (HPCR) criterion (Harper and Hooker, 1976) or equivalently the smallest region covering 95% of posterior probability mass. $\boldsymbol{\eta}_s$ is regarded as redundant if its 95% HPCR covers $\mathbf{0}$. Alternatively, we can transform the decision rule to a direct comparison of the squared Mahalanobis distance between $\mathbf{0}$ and $\boldsymbol{\eta}_s^*$ with a critical value of χ_p^2 , i.e., if $\boldsymbol{\eta}_s^{*'} \boldsymbol{\Sigma}_s^{*-1} \boldsymbol{\eta}_s^* \leq \chi_{p,0.05}^2$, then $\boldsymbol{\eta}_s$ is redundant; otherwise, it is significant.

4 Posterior Sampling

4.1 Prior specification

Based on (11), the conditional prior distribution of $\boldsymbol{\eta}_s$ given ψ_s is

$$P(\boldsymbol{\eta}_s|\psi_s) \propto \exp\left(-\frac{\gamma_s}{\sqrt{\psi_s}}\|\boldsymbol{\eta}_s\|_2\right), \quad s = 2, \dots, K. \quad (15)$$

This conditional Laplace prior can be represented as the scale mixture of normals with an exponential mixing density, leading to a hierarchical representation of the full model as follows: for $s = 1, \dots, K$,

$$\begin{aligned} [y_{it}|Z_{it} = s, \boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_s, \psi_s] &\sim N\left(\sum_{k=1}^s (\boldsymbol{\eta}'_k \mathbf{x}_{it}), \psi_s\right), \\ [\boldsymbol{\eta}_s|\psi_s, \tau_s^2] &\sim N(\mathbf{0}, \psi_s \tau_s^2 \mathbf{I}_p), \quad \psi_s^{-1} \stackrel{ind}{\sim} \text{Gamma}(a_{\psi s 0}, b_{\psi s 0}), \\ \tau_s^2 &\stackrel{ind}{\sim} \text{Gamma}\left(\frac{p+1}{2}, \frac{\gamma_s^2}{2}\right). \end{aligned} \quad (16)$$

For the tuning parameter γ_k , we assign the following dispersed gamma prior:

$$\gamma_k^2 \stackrel{ind}{\sim} \text{Gamma}(a_{\gamma k 0}, b_{\gamma k 0}), \quad k = 2, \dots, K, \quad (17)$$

where $a_{\gamma k 0}$ and $b_{\gamma k 0}$ are hyperparameters whose values are prespecified to achieve a highly dispersed prior. The degree of dispersion determines the magnitude of the penalty on unimportant regression parameters.

For the mixing probabilities in $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)'$, we assign a Dirichlet prior as follows:

$$(\pi_1, \dots, \pi_K) \sim \text{Dir}(c_K, \dots, c_K), \quad c_K = c \frac{n}{K}, \quad (18)$$

where c is a hyperparameter. Typically, assigning a moderate value in $(0, 1)$ can avoid a too small or too large penalty on the mixing probabilities. Our simulation study shows that a value around 0.5 performs satisfactorily.

For other parameters involved in the transition model (2), we assign conjugate priors:

$$\zeta_{us} \stackrel{ind}{\sim} N(\zeta_{us 0}, \sigma_{us 0}^2), \quad \boldsymbol{\alpha}_k \stackrel{ind}{\sim} N(\boldsymbol{\alpha}_{k 0}, \boldsymbol{\Sigma}_{\alpha k 0}), \quad k = 1, \dots, q, \quad (19)$$

where ζ_{us0} , σ_{us0}^2 , α_{k0} , and $\Sigma_{\alpha k0}^2$ are hyperparameters with prespecified values. The common practice is to set ζ_{us0} and the elements of α_{k0} to zero and assign σ_{us0}^2 and the diagonal elements of $\Sigma_{\alpha k0}$ to large values to induce vague priors if the preliminary information about ζ_{us} and α_k is unavailable.

4.2 MCMC algorithm

Unlike conventional HMMs that prespecify K , this study regards K as another unknown parameter and updates it with other model parameters in θ . The Bayesian estimate of (θ, K) can be obtained through the mean of the posterior samples drawn from $P(\theta, K | \mathbf{Y}, \mathbf{X}, \mathbf{D})$. However, $P(\theta, K | \mathbf{Y}, \mathbf{X}, \mathbf{D})$ involves unknown hidden states, leading to intractable sampling from $P(\theta, K | \mathbf{Y}, \mathbf{X}, \mathbf{D})$. Using the data augmentation technique, we instead work on $P(\mathbf{Z}, \theta, K | \mathbf{Y}, \mathbf{X}, \mathbf{D})$. However, the joint posterior distribution $P(\mathbf{Z}, \theta, K | \mathbf{Y}, \mathbf{X}, \mathbf{D})$ is still complex. Thus, the Gibbs sampler is employed to iteratively update each component through sampling from its full conditional distribution as follows: (a) update hidden states by sampling \mathbf{Z} from $P(\mathbf{Z} | \mathbf{Y}, \mathbf{X}, \mathbf{D}, \theta, K)$, (b) update the model parameters by sampling θ from $P(\theta | \mathbf{Y}, \mathbf{X}, \mathbf{D}, \mathbf{Z}, K)$, and (c) update the order K by sampling from $P(K | \mathbf{Y}, \mathbf{X}, \mathbf{D}, \mathbf{Z}, \theta)$. Owing to the transitioning features of hidden states and nonlinearity of the transition model (2), steps (a) and (b) require MCMC techniques, such as the FFBS and MH algorithms. The full conditional distributions involved in steps (a) and (b) are derived in Appendix 2 of Supplementary Material. Step (c) is the so-called ABRJ step, which allows K to be updated at each MCMC iteration as follows.

Let (K_{min}, K_{max}) be the lower and upper bounds of K , and $(K_{min}^{(0)}, K_{max}^{(0)})$ and $(K_{min}^{(j)}, K_{max}^{(j)})$ be their values at the initial stage and j th iteration of the MCMC algorithm. Typically, we set $K_{min}^{(0)} = 2$ and $K_{max}^{(0)}$ to a relatively large positive integer (e.g., 100 or 200) to allow sufficient flexibility in updating K . At the $(j + 1)$ th iteration, $K^{(j+1)}$ can remain unchanged, increase, or decrease by 1. To update $K^{(j)}$, we first locate a state s_* , such that

$s_* = \operatorname{argmin}_{s=1,\dots,K^{(j)}} \|\boldsymbol{\eta}_s^{(j)}\|_2$. Then, we calculate $d_{s_*}^2 = \boldsymbol{\eta}_{s_*}^{(j)'} \boldsymbol{\Sigma}_{s_*}^{(j)-1} \boldsymbol{\eta}_{s_*}^{(j)}$ and compare $d_{s_*}^2$ with $\chi_{p,0.05}^2$. If $d_{s_*}^2 \leq \chi_{p,0.05}^2$, we regard this component as redundant and update $K^{(j)}$ downward to $K^{(j+1)} = \max(K^{(j)} - 1, K_{\min}^{(j)})$. Meanwhile, we adjust $K_{\max}^{(j)}$ as $K_{\max}^{(j+1)} = \min(K_{\max}^{(j)}, K^{(j)})$. If $d_{s_*}^2 > \chi_{p,0.05}^2$, we regard this component as necessary. Then, we jump $K^{(j)}$ upward to $K^{(j+1)} = \min(K^{(j)} + 1, K_{\max}^{(j)} - 1)$. If $d_{s_*}^2 > \chi_{p,0.05}^2$ but $K^{(j)} = K_{\max}^{(j)} - 1$, then $K^{(j)}$ remains unchanged, i.e., $K^{(j+1)} = K^{(j)}$. Figure 1 shows the strategy of updating K in the ABRJ step. A pseudocode for implementing the MCMC algorithm is given below:

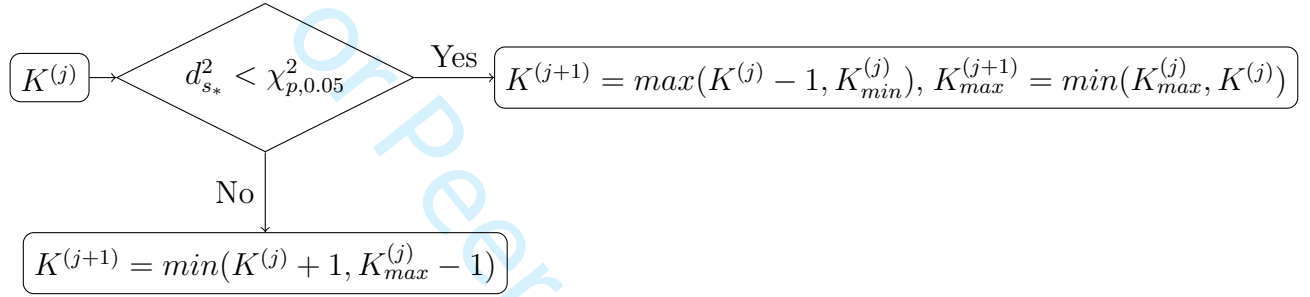


Figure 1: Strategy of updating K in the ABRJ step (c).

5 Simulation Study

This section includes three simulations to demonstrate the effectiveness of the proposed algorithm in order selection and parameter estimation under various scenarios. Simulation 1 evaluates the performance of the proposed method in the case of $K = 2$, Simulation 2 assesses estimation performance for HMMs with higher orders, and Simulation 3 focuses on order selection and compares the proposed method with AIC and BIC.

5.1 Simulation 1

This simulation considers a 2-state HMM with $p = 4$ and $q = 1$. Two sample sizes, $(n, T) = (50, 4)$, $(200, 4)$, are considered. In each setting, 100 datasets are generated from

Algorithm 1 MCMC algorithm for the estimation of heterogeneous HMMs

Data: $\mathbf{Y}, \mathbf{X}, \mathbf{D}, J, K_{min}^{(0)}, K_{max}^{(0)}$ $\triangleright J$ denotes the total number of iterations

- 1: $K^{(0)} = K_{min}^{(0)}$
- 2: **for** $j = 1$ to J **do**
- 3: Update $\mathbf{Z}^{(j)}$ by sampling from $P(\mathbf{Z}|\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}^{(j)}, K^{(j)})$ \triangleright FFBS algorithm
- 4: Update $\boldsymbol{\theta}^{(j)}$ by sampling from $P(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{X}, \mathbf{Z}, K^{(j)})$ \triangleright see details in Appendix B
- 5: $s_* = \operatorname{argmin}_{s=1, \dots, K^{(j)}} \|\boldsymbol{\eta}_s^{(j)}\|_2$
- 6: $\boldsymbol{\eta}_{s_*}^{(j)} = E(\boldsymbol{\eta}_{s_*}|\mathbf{Y}, \mathbf{X}, \mathbf{Z}, K^{(j)})$ \triangleright posterior mean vector
- 7: $\boldsymbol{\Sigma}_{s_*}^{(j)} = \operatorname{Var}(\boldsymbol{\eta}_{s_*}|\mathbf{Y}, \mathbf{X}, \mathbf{Z}, K^{(j)})$ \triangleright posterior covariance matrix
- 8: $d_{s_*}^2 = \boldsymbol{\eta}_{s_*}^{(j)'} \boldsymbol{\Sigma}_{s_*}^{(j)-1} \boldsymbol{\eta}_{s_*}^{(j)}$
- 9: **if** $d_{s_*}^2 < \chi_{p,0.05}^2$ **then**
- 10: $K^{(j+1)} = \max(K^{(j)} - 1, K_{min}^{(j)})$
- 11: $K_{max}^{(j+1)} = \min(K^{(j)}, K_{max}^{(j)})$
- 12: **else if** $d_{s_*}^2 \geq \chi_{p,0.05}^2$ **then**
- 13: $K^{(j+1)} = \min(K_{max}^{(j)} - 1, K^{(j)} + 1)$
- 14: **if** $K^{(j)} = K_{max}^{(j)} - 1$ **then**
- 15: $K^{(j+1)} = K^{(j)}$
- 16: **end if**
- 17: **end if**
- 18: $j = j + 1$
- 19: **end for**

the following model:

$$[y_{it}|Z_{it} = s] = \boldsymbol{\beta}_s' \mathbf{x}_{it} + \delta_{it}, \quad (20)$$

$$\operatorname{logit}(\vartheta_{itus}) = \zeta_{us} + \alpha d_{it},$$

where $\mathbf{x}_{it} = (1, x_{it1}, x_{it2}, x_{it3})'$, $x_{it1} \stackrel{\text{ind}}{\sim} N(0, 1)$, $x_{it2} \stackrel{\text{ind}}{\sim} U(-1, 1)$, $U(-1, 1)$ denotes the uniform distribution in $(-1, 1)$, $x_{it3} \stackrel{\text{ind}}{\sim} \operatorname{Bernoulli}(0.6)$, and $d_{it} \stackrel{\text{ind}}{\sim} N(0, 1)$. The true population values of the parameters are set as follows: $\boldsymbol{\beta}_1 = (2, 2, 1, 1)'$, $\boldsymbol{\beta}_2 = (0, 1, 2, -1)'$, $\psi_1 = \psi_2 = 0.25$, $\pi_1 = \pi_2 = 0.5$, $\boldsymbol{\zeta} = (\zeta_{11}, \zeta_{21}) = (-2, 2)'$, and $\alpha = -1$.

The hyperparameters of the prior distributions in (16)–(19) are specified as follows (Prior I): $a_{\psi s0} = 9$, $b_{\psi s0} = 4$, $a_{\gamma k0} = 1$, $b_{\gamma k0} = 0.1$, $c = 0.5$, $\alpha_{k0} = \zeta_{us0} = 0$, and $\sigma_{\alpha k0}^2 = \sigma_{us0}^2 = 1$. In implementing the MCMC algorithm, we impose a constraint described in Definition 2.1 (i.e., $\boldsymbol{\beta}_{(1)} > \boldsymbol{\beta}_{(2)}$) to each MCMC iteration to avoid label switching. Moreover, we set $K_{min}^{(0)} = 2$ and $K_{max}^{(0)} = 200$, which provides an extensive range for K . The

algorithm’s convergence is checked through the trace plots of the parameters. Figure 2(a) presents the trace plots of three MCMC chains of K starting from different initial values in an arbitrarily selected replication. The three MCMC chains of K mix rapidly and converge to the true value $K_0 = 2$ within a few iterations. Figure S1 of Supplementary Material presents the trace plots of three MCMC chains for other randomly selected parameters. Both figures indicate a fast convergence of the MCMC algorithm. To be conservative, we collect 10,000 posterior samples, discard the first 3000 iterations as burn-in, and calculate the bias and root mean square error (RMS) between the parameter estimates and their true values based on the remaining 7000 posterior samples corresponding to the selected order. Table 1 presents the estimation result. The bias and RMS are close to zero, and the performance improves when the sample size increases.

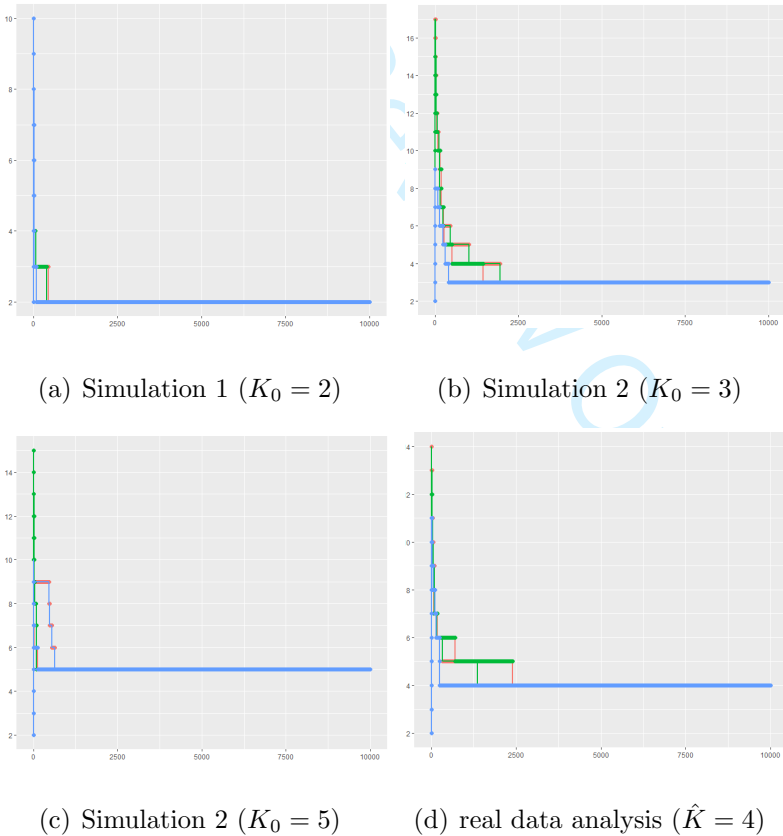


Figure 2: Trace plots of three MCMC chains of K in simulations and the ADNI study.

To reveal the sensitivity of Bayesian estimates to the prior input, we disturb the hy-

Table 1: Parameter estimates under Prior I in Simulation 1: $K_0 = 2$

$n = 50$						$n = 200$					
State 1			State 2			State 1			State 2		
Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS
Parameters in the conditional regression model											
β_{11}	0.034	0.049	β_{12}	-0.026	0.069	β_{11}	0.007	0.030	β_{12}	-0.006	0.040
β_{21}	-0.002	0.032	β_{22}	-0.001	0.035	β_{21}	-0.003	0.035	β_{22}	0.014	0.029
β_{31}	0.020	0.055	β_{32}	-0.018	0.051	β_{31}	0.001	0.019	β_{32}	-0.017	0.026
β_{41}	-0.020	0.083	β_{42}	0.045	0.070	β_{41}	-0.022	0.050	β_{42}	0.015	0.042
ψ_1	0.020	0.042	ψ_2	0.038	0.058	ψ_1	0.009	0.016	ψ_2	0.008	0.022
Parameters in the transition model											
ζ_{11}	0.030	0.168	ζ_{21}	-0.040	0.155	ζ_{11}	-0.008	0.096	ζ_{21}	-0.028	0.115
π_1	-0.007	0.052	π_2	0.007	0.052	π_1	-0.001	0.020	π_2	0.001	0.020
α_1	0.020	0.107				α_1	0.010	0.085			

perparameters as follows (Prior II): $a_{\psi s0} = 13$, $b_{\psi s0} = 6$, $a_{\gamma k0} = 1$, $b_{\gamma k0} = 0.01$, $c = 0.3$, $\alpha_{k0} = \zeta_{us0} = 2$, and $\sigma_{\alpha k0}^2 = \sigma_{us0}^2 = 100$. Table S1 of Supplementary Material reports the obtained results. The parameter estimates perform similarly to those in Table 1, indicating that the proposed Bayesian estimation is insensitive to the disturbed prior considered.

Furthermore, we check the sensitivity of Bayesian estimation to the misspecification of the distribution of δ_{it} by considering two nonnormal cases: (1) $\delta_{it} \sim U(-1, 1)$ and (2) $\delta_{it} \sim 0.4N(1, 1) + 0.6N(-1, 1)$. We simulate 100 datasets from the proposed model with $n = 200$ and δ_{it} drawn from case (1) or (2). The hyperparameters and other settings are the same as in Simulation 1. Table S2 of Supplementary Material presents the estimation results obtained under the two nonnormal cases. Except for the variance of δ_{it} and some parameters involved in the transition model, most parameter estimates are robust to the violation of the normality assumption of δ_{it} . Therefore, the impact of misspecifying the distribution of δ_{it} is mainly on estimating its variance.

5.2 Simulation 2

This simulation examines estimation performance for higher-order models. We first consider a 3-state HMM. Covariates \mathbf{x}_{it} is the same as in Simulation 1. For simplicity, we set $\mathbf{d}_{it} = \mathbf{x}_{it}^*$, where \mathbf{x}_{it}^* is the subvector of \mathbf{x}_{it} excluding 1. Two sample sizes, $(n, T) = (200, 6)$ and $(400, 6)$, are considered. The true population values of the parameters are set as $\beta_1 = (3, 3, 3, 3)'$, $\beta_2 = (0, 1, 2, 2)'$, $\beta_3 = (-4, 2, 1, 1)'$, $\psi = (0.25, 0.25, 0.25)'$, $\pi = (0.3, 0.4, 0.3)'$, $\zeta_1 = (-1, -1, -1)'$, $\zeta_2 = (1, 1, 1)'$, $\alpha = (1, -1, -1)'$. The prior specification and simulation settings are similar to Simulation 1, except that the hyperparameters for α are set as $\alpha_{k0} = \mathbf{0}$ and $\Sigma_{\alpha k0} = \mathbf{I}_3$. Figure 2(b) presents the trace plots of three MCMC chains of K starting from different initial values in an arbitrarily selected replication. Again, the MCMC chains mix and converge to the true value of $K_0 = 3$ rapidly. The trace plots of other parameters (Figure S2 of Supplementary Material) suggest that the algorithm converges within 3000 iterations. Therefore, we discard 3000 burn-in and use the remaining 7000 posterior samples to obtain the Bayesian estimates of the parameters. Table S3 of Supplementary Material shows the estimation results based on 100 replications under $(n, T) = (200, 6)$, indicating that the proposed method performs satisfactorily in bias and RMS. The results under $(n, T) = (400, 6)$ are further improved and not reported.

Next, we further increase the order to $K_0 = 5$. We consider a 5-state HMM with covariates $\mathbf{x}_{it} = (1, x_{it1}, x_{it2})'$, where $x_{it1} \stackrel{ind}{\sim} N(0, 1)$ and $x_{it2} \stackrel{ind}{\sim} U(-1, 1)$, and $\mathbf{d}_{it} = \mathbf{x}_{it}^*$. The other model setup is the same as above. The true population values of unknown parameters are set as $\beta_1 = (5, 5, 5)'$, $\beta_2 = (3, 4, 3)'$, $\beta_3 = (0, 3, 4)'$, $\beta_4 = (-2, 4, 3)'$, $\beta_5 = (-5, 5, 2)'$, $\psi = (0.25, 0.25, 0.25, 0.25, 0.25)'$, $\pi = (0.2, 0.2, 0.2, 0.2, 0.2)'$, $\zeta_1 = (-2, -2, -2, -2, -2)'$, $\zeta_2 = (-1, -1, -1, -1, -1)'$, $\zeta_3 = (1, 1, 1, 1, 1)'$, $\zeta_4 = (2, 2, 2, 2, 2)'$, $\alpha = (2, 1)'$. The prior and simulation settings are similar to those given above. Figure 2(c) presents the trace plots of three MCMC chains of K starting from different initial values in an arbitrarily se-

lected replication, showing that the iterative K quickly converges to its true value $K_0 = 5$. The trace plots of other parameters (Figure S3 of Supplementary Material) suggest that the MCMC chains mix well within 4000 iterations. Thus, we discard 4000 burn-in and use the remaining 6000 posterior samples to obtain the Bayesian estimates of the parameters. Table S4 of Supplementary Material presents the estimation results under $(n, T) = (200, 6)$, indicating that the proposed method performs satisfactorily in order selection and parameter estimation when K_0 increases to 5. The results under a larger size of $(n, T) = (400, 6)$ are better and not reported.

5.3 Simulation 3

This simulation assesses the performance of order selection. Considering that no existing methods can simultaneously estimate the order and model parameters of heterogeneous HMMs, we compare the proposed BDP procedure with criterion-based approaches, AIC and BIC, in order selection accuracy.

Table 2 presents the proportions of selected orders calculated based on 100 replications with $K_0 = \{3, 4, 5\}$ and $n = \{100, 200, 400\}$. The results show that the proposed BDP procedure consistently outperforms AIC and BIC in all the scenarios considered. In general, the performance of the three methods improves when the sample size increases but declines when K_0 increases. In particular, AIC and BIC perform poorly when $K_0 = 5$ regardless of the sample size; their correct selection proportions are below or around 0.5. By contrast, our proposed method performs much better, and its correct selection proportion attains 0.81 when $n = 400$. This comparison confirms that our procedure achieves higher accuracy than the conventional criterion-based approaches. In addition, unlike the two-stage methods that perform order selection and parameter estimation in two stages, our proposed method accomplishes the two tasks in a single stage. Moreover, our procedure guarantees that sampling only occurs when the states are necessary and retained. Hence, its advantages

Table 2: Proportions of selecting HMMs with different orders

K_0	\hat{K}	$n = 100$			$n = 200$			$n = 400$		
		AIC	BIC	BDP	AIC	BIC	BDP	AIC	BIC	BDP
3	2	0	0	0.13	0	0	0.10	0	0	0
	3	0.70	0.76	0.87	0.69	0.80	0.82	0.78	0.98	1
	4	0.18	0.22	0	0.25	0.20	0.08	0.20	0.02	0
	5	0.12	0.02	0	0.06	0	0	0.02	0	0
4	2	0	0	0.11	0	0	0.02	0	0	0
	3	0	0	0.16	0	0	0.21	0	0.12	0.09
	4	0.45	0.66	0.73	0.55	0.60	0.77	0.67	0.67	0.85
	5	0.30	0.31	0	0.30	0.25	0	0.21	0.09	0.06
5	6	0.25	0.03	0	0.15	0.15	0	0.12	0.12	0
	3	0	0	0.16	0	0	0.04	0	0	0.13
	4	0.24	0.28	0.23	0.18	0.21	0.15	0.07	0.14	0.06
	5	0.36	0.40	0.61	0.47	0.50	0.77	0.55	0.59	0.81
6	6	0.32	0.28	0	0.29	0.26	0.04	0.21	0.10	0
	7	0.08	0.04	0	0.06	0.03	0	0.17	0.17	0

in computational efficiency over existing ones become increasingly pronounced when the candidate model space enlarges.

6 Real Data Analysis

In this section, we applied the proposed method to the dataset extracted from the ADNI study to demonstrate the practical utility of the proposed method. ADNI is a longitudinal multicenter study that began in 2004, collecting various participants' imaging and clinical assessments. More information is referred to the official website: www.adni-info.org.

We focused on 616 subjects collected from the ADNI study with four follow-up visits, namely, the baseline, six months, 12 months, and 24 months. Alzheimer Disease Assessment Scale-Cognitive 13 (ADAS13), which reflects cognitive impairment in AD assessment, is treated as response y_{it} . Generally, a high ADAS13 score indicates low cognitive ability. In addition, some clinical and generic variables were considered as covariates. One is a

time-variant continuous variable, x_{it1} : the logarithm of the ratio of hippocampal volume over the whole brain volume (HIP). Other covariates include the genetic variable APOE- $\epsilon 4$, coded as 0, 1, 2, denoting the number of APOE- $\epsilon 4$ alleles and represented by x_{it2} ($x_{it2} = 1$ if carrying one allele and 0 otherwise) and x_{it3} ($x_{it3} = 1$ if carrying two alleles and 0 otherwise), patients' age at baseline, x_{it4} , and patients' gender, x_{it5} ($x_{it5} = 1$ if female). The main goal of this study is to simultaneously identify the number of hidden states and the state-specific relationship between ADAS 13 and its important risk factors.

Table 3: Parameter estimation results for ADNI study

State 1		State 2		State 3		State 4	
Par	Est(sd)	Par	Est(sd)	Par	Est(sd)	Par	Est(sd)
Parameters in the conditional regression model							
β_{11}	-0.803(0.037)	β_{12}	-0.191(0.059)	β_{13}	0.521(0.089)	β_{14}	1.559(0.114)
β_{21}	-0.164(0.036)	β_{22}	-0.281(0.048)	β_{23}	-0.301(0.042)	β_{24}	-0.331(0.090)
β_{31}	0.039(0.035)	β_{32}	0.135(0.056)	β_{33}	0.198(0.098)	β_{34}	0.253(0.116)
β_{41}	0.263(0.073)	β_{42}	0.542(0.136)	β_{43}	1.138(0.096)	β_{44}	1.906(0.202)
β_{51}	-0.070(0.094)	β_{52}	0.035(0.037)	β_{53}	0.061(0.048)	β_{54}	0.037(0.062)
β_{61}	-0.029(0.031)	β_{62}	0.046(0.051)	β_{63}	0.093(0.068)	β_{64}	0.399(0.101)
ψ_1	0.094(0.008)	ψ_2	0.101(0.007)	ψ_3	0.136(0.013)	ψ_4	0.421(0.049)
Parameters in the transition model							
ζ_{11}	2.863(0.255)	ζ_{21}	-2.217(0.234)	ζ_{31}	-3.867(0.450)	ζ_{41}	-3.537(0.479)
ζ_{12}	3.130(0.806)	ζ_{22}	3.652(0.462)	ζ_{32}	-1.701(0.320)	ζ_{42}	-3.343(0.431)
ζ_{13}	-0.762(0.852)	ζ_{23}	2.089(0.511)	ζ_{33}	3.941(0.542)	ζ_{43}	-2.271(0.441)
π_1	0.322(0.023)	π_2	0.314(0.019)	π_3	0.224(0.017)	π_4	0.140(0.013)
α_1	-0.139(0.099)	α_2	-0.538(0.229)	α_3	-0.742(0.350)	α_4	-0.019(0.104)
α_5	0.090(0.106)						

The prior specification and other settings are similar to the simulation study. We imposed constraint described in Definition 2.1 to each MCMC iteration to avoid label switching. The trace plots of K shown in Figure 2(d) indicate that the MCMC chains of K from different initial values quickly converge to $K = 4$, suggesting a 4-state HMM for the data. Figure S4 of Supplementary Material presents the trace plots of other parameters

involved in the selected model. The MCMC chains mixed well within 5000 iterations. Thus, we discarded 5000 burn-in iterations and used the remaining 5000 posterior samples to obtain the parameter estimates. Table 3 presents the parameter estimates of the selected 4-state HMM. Based on the results, we have the following observations.

First, the state-specific intercept β_{1s} exhibits a descending trend. Patients have the lowest ADAS mean score in state 1 and highest mean score in state 4. According to the existing literature (Kantarci et al., 2013), states 1 to 4 can be interpreted as CN, early mild cognitive impairment (MCI), late MCI, and AD accordingly.

Second, HIP (β_{2s}) exerts an adverse effect on ADAS13, implying that a sizable hippocampal volume is associated with a low ADAS13 score and thus high cognitive ability. Moreover, the magnitude of the HIP effect on ADAS13 increases from CN to AD, implying that hippocampal atrophy continuously impairs patients' cognitive ability during AD progression. The published medical reports (e.g., Dickerson and Wolk, 2013) also revealed that the loss of hippocampal volume significantly affects AD.

Third, the effects of APEP- $\epsilon 4$ (β_{3s} and β_{4s}) on ADAS13 are positive, suggesting that carrying APOE- $\epsilon 4$ increases AD risk, and such impact becomes increasingly pronounced with the disease progression. This finding is in line with the medical report (Risacher et al., 2015) that APOE- $\epsilon 4$ is a crucial biomarker of AD. Furthermore, the magnitude of β_{4s} is larger than β_{3s} for $s = 1, \dots, 4$, implying that carrying two alleles, in general, impairs cognitive function more significantly than carrying only one allele. Besides, patients' age and gender do not substantially affect ADAS13 when controlling hippocampal volume and APOE- $\epsilon 4$. An exception lies in $\beta_{64} = 0.399(0.101)$, which suggests that females suffer more severe cognitive decline than males in the late AD progression period. This result again agrees with the existing literature (e.g., Via et al., 2010; Kang et al., 2019).

Lastly, the transition pattern described by ζ exhibits a banding structure. That is, patients are likely to transit between adjacent states. Moreover, α_2 and α_3 are significant and

negative, implying that the transition pattern between hidden states exhibits heterogeneity. APOE- $\epsilon 4$ allele carriers are more likely to transit to a worse state rather than remain in the current one than noncarriers; carrying two alleles induces a higher risk of transitioning to a worse state than carrying one allele. This result is consistent with the existing finding (Eunjee et al., 2015) that APOE- $\epsilon 4$ alleles increase the risk of developing AD. However, other covariates, such as age and HIP, do not significantly affect the between-state transition given APOE- $\epsilon 4$. This result implies that conditional on APOE- $\epsilon 4$, the direct effects of age and hippocampal volume on the transition probability are weak.

7 Discussion

In this study, we have proposed a double penalized method to perform order selection and parameter estimation for heterogeneous HMMs under the Bayesian framework. In addition, we have developed a novel MCMC algorithm with an ABRJ sampling scheme to facilitate a joint updating of the order and model parameters. Multiple simulation studies and an application to the ANDI dataset demonstrate the superiority of the proposed method over existing ones and its utility in realistic settings. Furthermore, the proposed model can cope with general situations where specific covariates simultaneously influence the emission and transition processes.

The present work can be extended in several directions. First, we use a single indicator to represent a response or predictor. For example, we adopt ADAS13 to reflect cognitive ability in the ADNI data analysis. While in practice, multiple tests can be used to examine cognitive impairment, and their scores can be integrated into a univariate latent construct through factor analysis. Such an extension can accommodate latent responses or covariates, reduce model dimensionality, and improve interpretability. Second, our conditional regression model only accommodates a continuous variable. Given that complex data types are frequently encountered in medical, social, and psychological sciences, generalizing our

method to incorporate multivariate, functional, or image variables can considerably enhance model capability. Finally, this study mainly focuses on order selection. However, variable selection is potentially interesting in the presence of high-dimensional variables. Thus, we can consider additional penalties for simultaneous order and variable selection. These possible extensions may raise new theoretical and computational challenges.

Supplementary Material

Proof of the propositions in Section 3, full conditional distributions in Section 4, and additional numerical results in Sections 5 and 6 are provided in Supplementary Material.

References

Agresti, A. (2003). *Categorical Data Analysis*. John Wiley & Sons.

Akaike, H. (1974). New look at statistical-model identification. *IEEE Transactions on Automatic Control*, 19:716–723.

Baum, L. E., Petrie, T., Soules, G., and Weiss, N. (1970). A maximization technique occurring in the statistical analysis of probabilistic functions of markov chains. *The Annals of Mathematical Statistics*, 41(1):164–171.

Celeux, G. and Durand, J.-B. (2008). Selecting hidden markov model state number with cross-validated likelihood. *Computational Statistics*, 23:541–564.

Chen, J. and Khalili, A. (2008). Order selection in finite mixture models with a nonsmooth penalty. *Journal of the American Statistical Association*, 103:1674–1683.

Dickerson, B. and Wolk, D. (2013). Biomarker-based prediction of progression in mci: comparison of ad-signature and hippocampal volume with spinal fluid amyloid- β and tau. *Front Aging Neurosci*, 5:55.

- Eunjee, L., Hongtu, Z., Dehan, K., et al. (2015). BFLCRM: A bayesian functional linear cox regression model for predicting time to conversion to alzheimer's disease. *The Annals of Applied Statistics*, 9:2153–2178.
- Green, P. J. (1995). Reversible jump markov chain monte carlo computation and bayesian model determination. *Biometrika*, 82(4):711–732.
- Guo, R., Zhu, H., Chow, S.-M., and Ibrahim, J. G. (2012). Bayesian lasso for semiparametric structural equation models. *Biometrics*, 68(2):567–577.
- Harper, W. and Hooker, C. (1976). *Foundations of Probability Theory, Statistical Inference, and Statistical Theories of Science*.
- Hung, Y., Wang, Y., Zarnitsyna, V., Zhu, C., and Wu, C.-F. (2013). Hidden markov models with applications in cell adhesion experiments. *Journal of the American Statistical Association*, 108:1469–1479.
- Ip, E., Zhang, Q., Rejeski, J., Harris, T., and Kritchevsky, S. (2013). Partially ordered mixed hidden markov model for the disablement process of older adults. *Journal of the American Statistical Association*, 108(502):370–384.
- Kang, K., Song, X., Hu, X. J., and Zhu, H. (2019). Bayesian adaptive group lasso with semiparametric hidden markov models. *Statistics in Medicine*, 38:1634–1650.
- Kantarci, K., Gunter, J., et al. (2013). Focal hemosiderin deposits and β -amyloid load in the adni cohort. *Alzheimer's & Dementia : The Journal of the Alzheimer's Association*, 9:S116–S123.
- Lin, Y. and Song, X. (2022). Order selection for regression-based hidden markov model. *Journal of Multivariate Analysis*, to appear.

Liu, H. and Song, X. (2020). Bayesian analysis of hidden markov structural equation models with an unknown number of hidden states. *Econometrics and Statistics*, 18:29–43.

Mackay, R. (2002). Estimating the order of a hidden markov model. *Canadian Journal of Statistics*, 30:573–589.

Manole, T. and Khalili, A. (2021). Estimating the number of components in finite mixture models via the group-sort-fuse procedure. *The Annals of Statistics*, 49.

Park, T. and Casella, G. (2008). The Bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686.

Risacher, S. L. et al. (2015). Apoe effect on alzheimer’s disease biomarkers in older adults with significant memory concern. *Alzheimer’s & Dementia : The Journal of the Alzheimer’s Association*, 11,12:1417–1429.

Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6:461–464.

Song, X., Xia, Y., and Zhu, H. (2017). Hidden markov latent variable models with multivariate longitudinal data. *Biometrics*, 73:313–323.

Via, Jose, Lloret, and Ana (2010). Why women have more alzheimer’s disease than men: Gender and mitochondrial toxicity of amyloid- β peptide. *Journal of Alzheimer’s Disease*, 20:527–533.

Ye, M., Lu, Z., Li, Y., and Song, X. (2019). Finite mixture of varying coefficient model: Estimation and component selection. *Journal of Multivariate Analysis*, 171:452–474.

Zhou, J., Song, X., and Sun, L. (2020). Continuous time hidden markov model for longitudinal data. *Journal of Multivariate Analysis*, 179:104646.

Bayesian order selection in heterogeneous hidden Markov models

Supplementary Material

Appendix 1: Proof of propositions

Proof of Proposition 2.1

Suppose $Z_i \sim \text{categorical}(\pi_1, \dots, \pi_K)$ for $i = 1, \dots, n$, with $0 \leq \pi_s \leq 1$ and $\sum_{s=1}^K \pi_s = 1$, where $s = 1, \dots, K$. We assign a conjugate prior distribution $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K) \sim \text{Dir}(c_K, \dots, c_K)$, where $c_K = c \frac{n}{K}$, $c > 0$ is a non-negative constant. Given \mathbf{Z} , the posterior of $\boldsymbol{\pi}$ is

$$\begin{aligned} P(\boldsymbol{\pi}|\mathbf{Z}) &\propto P(\boldsymbol{\pi})P(\mathbf{Z}|\boldsymbol{\pi}) \\ &= \frac{\Gamma(Kc_K)}{\Gamma(c_K)^K} \prod_{s=1}^K \pi_s^{c_K-1} \prod_{s=1}^K \pi_s^{\sum_{i=1}^n I(Z_i=s)} \\ &\propto \prod_{s=1}^K \pi_s^{c_K+\sum_{i=1}^n I(Z_i=s)-1} \\ &\stackrel{D}{=} \text{Dir}(c_K + n_1, \dots, c_K + n_K). \end{aligned} \tag{S1}$$

Thus, $P(\boldsymbol{\pi}|\mathbf{Z}) \sim \text{Dir}(c_K + n_1, \dots, c_K + n_K)$, where $n_s = \sum_{i=1}^n I(Z_i = s)$. In addition,

$$\begin{aligned} E(\pi_s|\mathbf{Z}) &= \frac{c_K + n_s}{\sum_{j=1}^K (c_K + n_j)} = \frac{c_K + n_s}{(c+1)n} \\ &\geq \frac{c_K}{(c+1)n} = \frac{c}{c+1} \frac{1}{K}, \quad s = 1, \dots, K. \end{aligned} \tag{S2}$$

That is, the posterior mean of π_s is lower bounded by constant $\frac{c}{c+1} \frac{1}{K}$.

Proof of Proposition 2.2

The conditional Laplace distribution can be represented as the scale mixture of normals with an exponential mixing density

$$\frac{a}{2} \exp(-a|z|) = \int_0^\infty \frac{1}{\sqrt{2\pi s}} \exp\left(-\frac{z^2}{2s}\right) \frac{a^2}{2} \exp\left(-\frac{a^2}{2}s\right) ds, \quad a > 0. \quad (\text{S3})$$

Suppose we express our model by a set of hierarchical representation as follows:

$$\begin{aligned} [y_{it}|Z_{it} = s, \boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_s, \psi_s] &\sim N\left(\sum_{k=1}^s (\boldsymbol{\eta}'_k \mathbf{x}_{it}), \psi_s\right), \\ [\boldsymbol{\eta}_s|\psi_s, \tau_s^2] &\sim N(\mathbf{0}, \psi_s \tau_s^2 \mathbf{I}_p), \quad \psi_s^{-1} \stackrel{\text{ind}}{\sim} \text{Gamma}(a_{\psi s 0}, b_{\psi s 0}), \\ \tau_s^2 &\stackrel{\text{ind}}{\sim} \text{Gamma}\left(\frac{p+1}{2}, \frac{\gamma_s^2}{2}\right), \quad \gamma_s^2 \stackrel{\text{ind}}{\sim} \text{Gamma}(a_{\gamma s 0}, b_{\gamma s 0}), \quad s = 2, \dots, K. \end{aligned} \quad (\text{S4})$$

Then, the conditional distribution $\boldsymbol{\eta}_s|\psi_s$ can be obtained by integrating out τ_s^2 as follows:

$$\begin{aligned} P(\boldsymbol{\eta}_s|\psi_s) &= \int_0^\infty P(\boldsymbol{\eta}_s|\psi_s, \tau_s^2) P(\tau_s^2) d\tau_s^2 \\ &= \int_0^\infty (2\pi\psi_s\tau_s^2)^{-\frac{p}{2}} \exp\left(-\frac{\boldsymbol{\eta}'_s \boldsymbol{\eta}_s}{2\psi_s\tau_s^2}\right) \times \left(\frac{\gamma_s^2}{2}\right)^{\frac{p+1}{2}} (\tau_s^2)^{\frac{p-1}{2}} \exp\left(-\frac{\gamma_s^2}{2}\tau_s^2\right) d\tau_s^2 \\ &\propto \int_0^\infty (\tau_s^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}\left(\frac{\boldsymbol{\eta}'_s \boldsymbol{\eta}_s}{\psi_s\tau_s^2} + \gamma_s^2\tau_s^2\right)\right\} d\tau_s^2 \\ &= \int_0^\infty \frac{1}{\sqrt{2\pi\tau_s^2}} \exp\left(-\frac{\boldsymbol{\eta}'_s \boldsymbol{\eta}_s}{2\tau_s^2}\right) \frac{\gamma_s^2}{2} \exp\left(-\frac{\gamma_s^2}{2}\tau_s^2\right) d\tau_s^2 \\ &= \frac{\gamma_s}{2\sqrt{\psi_s}} \exp\left(-\frac{\gamma_s}{\sqrt{\psi_s}} \|\boldsymbol{\eta}_s\|_2\right) \end{aligned} \quad (\text{S5})$$

Thus, we obtain the conditional Laplace prior of $\boldsymbol{\eta}_s|\psi_s$ as follows:

$$P(\boldsymbol{\eta}_k|\psi_k) = \frac{\gamma_k}{2\sqrt{\psi_k}} \exp\left(-\frac{\gamma_k}{\sqrt{\psi_k}} \|\boldsymbol{\eta}_k\|_2\right), \quad k = 1, \dots, K. \quad (\text{S6})$$

Appendix 2: Full conditional distributions

1. Full conditional distribution of Z_{it}

Let $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})'$ denote the response of subject i , $\mathbf{x}_{it} = (x_{it1}, \dots, x_{itp})'$ and $\mathbf{d}_{it} = (d_{it1}, \dots, d_{itq})'$ denote covariate vectors in the conditional regression and transition models

respectively, and $\boldsymbol{\theta}$ denote the vector of all unknown parameters excluding K . Then,

$$\begin{aligned}
 P(Z_{it}|\cdot) &\propto P(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i, Z_{it}|\boldsymbol{\theta}) \\
 &= P(y_{i1}, \dots, y_{it}, \mathbf{x}_{i1}, \dots, \mathbf{x}_{it}, \mathbf{d}_{i1}, \dots, \mathbf{d}_{it}, Z_{it}|\boldsymbol{\theta}) \\
 &\times P(y_{i,t+1}, \dots, y_{iT}, \mathbf{x}_{i,t+1}, \dots, \mathbf{x}_{iT}, \mathbf{d}_{i,t+1}, \dots, \mathbf{d}_{iT}|\boldsymbol{\theta}, Z_{it}) \\
 &\doteq q_{it}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i, Z_{it}|\boldsymbol{\theta}) \times \tilde{q}_{it}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i|\boldsymbol{\theta}, Z_{it}).
 \end{aligned} \tag{S7}$$

Using the FFBS scheme, we first initialize $q_{i1}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i, Z_{i1}|\boldsymbol{\theta}) = P(y_{i1}, \mathbf{x}_{i1}, \mathbf{d}_{i1}, Z_{i1}|\boldsymbol{\theta})$
 $= P(y_{i1}, \mathbf{x}_{i1}, \mathbf{d}_{i1}|\boldsymbol{\theta}, Z_{i1})P(Z_{i1}|\boldsymbol{\theta})$ and calculate $q_{it}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i, Z_{it}|\boldsymbol{\theta})$ for $t = 2, \dots, T$ in a
 recursion manner as follows:

$$\begin{aligned}
 &q_{it}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i, Z_{it}|\boldsymbol{\theta}) \\
 &= P(y_{i1}, \dots, y_{it}, \mathbf{x}_{i1}, \dots, \mathbf{x}_{it}, \mathbf{d}_{i1}, \dots, \mathbf{d}_{it}, Z_{it}|\boldsymbol{\theta}) \\
 &= \sum_{u=1}^K P(y_{i1}, \dots, y_{it}, \mathbf{x}_{i1}, \dots, \mathbf{x}_{it}, \mathbf{d}_{i1}, \dots, \mathbf{d}_{it}, Z_{it}, Z_{i,t-1} = u|\boldsymbol{\theta}) \\
 &= \sum_{u=1}^K P(y_{i1}, \dots, y_{i,t-1}, \mathbf{x}_{i1}, \dots, \mathbf{x}_{i,t-1}, \mathbf{d}_{i1}, \dots, \mathbf{d}_{i,t-1}, Z_{i,t-1} = u|\boldsymbol{\theta}) \\
 &\times P(Z_{it}|Z_{i,t-1} = u, \mathbf{d}_{it}, \boldsymbol{\theta})P(y_{it}|\mathbf{x}_{it}, Z_{it}, \boldsymbol{\theta}) \\
 &= \sum_{u=1}^K [q_{i,t-1}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i, Z_{i,t-1} = u|\boldsymbol{\theta})P(Z_{it}|Z_{i,t-1} = u, \mathbf{d}_{it}, \boldsymbol{\theta})P(y_{it}|\mathbf{x}_{it}, Z_{it}, \boldsymbol{\theta})].
 \end{aligned} \tag{S8}$$

Similarly, we initialize $\tilde{q}_{iT}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i|\boldsymbol{\theta}, Z_{iT}) = 1$ and calculate $\tilde{q}_{it}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i|\boldsymbol{\theta}, Z_{it})$ for
 $t = T - 1, \dots, 1$ in a recursion manner as follows:

$$\begin{aligned}
 &\tilde{q}_{it}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i|\boldsymbol{\theta}, Z_{it}) \\
 &= P(y_{i,t+1}, \dots, y_{iT}, \mathbf{x}_{i,t+1}, \dots, \mathbf{x}_{iT}, \mathbf{d}_{i,t+1}, \dots, \mathbf{d}_{iT}|Z_{it}, \boldsymbol{\theta}) \\
 &= \sum_{u=1}^K P(y_{i,t+1}, \dots, y_{iT}, \mathbf{x}_{i,t+1}, \dots, \mathbf{x}_{iT}, Z_{i,t+1} = u|Z_{it}, \boldsymbol{\theta}) \\
 &= \sum_{u=1}^K P(y_{i,t+2}, \dots, y_{iT}, \mathbf{x}_{i,t+2}, \dots, \mathbf{x}_{iT}|Z_{i,t+1} = u, \boldsymbol{\theta})P(Z_{i,t+1} = u|Z_{it}, \mathbf{d}_{i,t+1}, \boldsymbol{\theta})
 \end{aligned}$$

$$\begin{aligned}
& \times P(y_{i,t+1} | \mathbf{x}_{i,t+1}, Z_{i,t+1} = u, \boldsymbol{\theta}) \\
& = \sum_{u=1}^K [\tilde{q}_{i,t+1}(\mathbf{y}_i, \mathbf{X}_i, \mathbf{D}_i | Z_{i,t+1} = u, \boldsymbol{\theta}) P(Z_{i,t+1} = u | Z_{it}, \mathbf{d}_{i,t+1}, \boldsymbol{\theta}) \\
& \quad \times P(y_{i,t+1} | \mathbf{x}_{i,t+1}, Z_{i,t+1} = u, \boldsymbol{\theta})].
\end{aligned}$$

2. Full conditional distributions of $\boldsymbol{\eta}_s$ and ψ_s^{-1}

$$[\boldsymbol{\eta}_s | \cdot] \sim N(\boldsymbol{\eta}_s^*, \boldsymbol{\Sigma}_s^*), \quad [\psi_s^{-1} | \cdot] \sim \text{Gamma}(a_{\psi_s}^*, b_{\psi_s}^*), \quad (\text{S9})$$

where $s_{\psi_s}^* = a_{\psi_s 0} + (n_s + p)/2$, and

$$\begin{aligned}
\boldsymbol{\Sigma}_s^* &= \left(\sum_{i=1}^n \sum_{t=1}^T \psi_s^{-1} \mathbf{x}_{it} \mathbf{x}_{it}' I(Z_{it} = s) + \psi_s^{-1} \tau_s^{-2} I_p \right)^{-1}, \\
\boldsymbol{\eta}_s^* &= \boldsymbol{\Sigma}_s^* \left[\sum_{i=1}^n \sum_{t=1}^T \psi_s^{-1} \mathbf{x}_{it} (y_{it} - \sum_{m=1}^{s-1} (\boldsymbol{\eta}_m' \mathbf{x}_{it})) I(Z_{it} = s) \right], \\
b_{\psi_s}^* &= b_{\psi_s 0} + \frac{1}{2} \left[\sum_{i=1}^n \sum_{t=1}^T (y_{it} - \sum_{k=1}^s (\boldsymbol{\eta}_k' \mathbf{x}_{it}))^2 I(Z_{it} = s) + \tau_s^{-2} \boldsymbol{\eta}_s' \boldsymbol{\eta}_s \right].
\end{aligned}$$

3. Full conditional distributions of π_s , ζ_{us} , and $\boldsymbol{\alpha}$

$$\begin{aligned}
[\boldsymbol{\pi} | \cdot] &\sim \text{Dir} \left(c_K + \sum_{i=1}^n I(Z_{i1} = 1), \dots, c_K + \sum_{i=1}^n I(Z_{i1} = K) \right), \\
[\zeta_{us} | \cdot] &\propto \exp \left\{ \sum_{v=s}^K \sum_{i=1}^n \sum_{t=2}^T \log (P_{ituv} \times I(Z_{it} = v, Z_{i,t-1} = u)) - \frac{(\zeta_{us} - \zeta_{us0})^2}{2\sigma_{\zeta_{us0}}^2} \right\}, \\
p[\boldsymbol{\alpha} | \cdot] &\propto \exp \left\{ \sum_{i=1}^n \sum_{t=2}^T \log (P_{itus} \times I(Z_{it} = s, Z_{i,t-1} = u)) - \frac{1}{2} (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)' \boldsymbol{\Sigma}_{\alpha}^{-1} (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) \right\},
\end{aligned}$$

where $\boldsymbol{\alpha}_0 = (\alpha_{10}, \dots, \alpha_{q0})'$ and $\boldsymbol{\Sigma}_{\alpha} = \text{diag}(\sigma_{\alpha 10}^2, \dots, \sigma_{\alpha q0}^2)$.

Appendix 3

This section provides additional numerical results. Tables S1–S4 present additional estimation results in Simulations 1 and 2. Figures S1, S2, and S3 show the trace plots of several parameters randomly selected from an arbitrary replication in Simulation 1 ($K_0 = 2$) and Simulation 2 ($K_0 = 3$ and 5). Figure S4 shows the trace plots of six parameters randomly selected in the real data example.

Table S1: Parameter estimates under Prior II in Simulation 1: $K_0 = 2$

$n = 50$						$n = 200$					
State 1			State 2			State 1			State 2		
Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS
Parameters in the conditional regression model											
β_{11}	-0.040	0.068	β_{12}	0.030	0.061	β_{11}	0.003	0.038	β_{12}	-0.014	0.053
β_{21}	-0.007	0.044	β_{22}	-0.001	0.053	β_{21}	0.018	0.026	β_{22}	0.013	0.028
β_{31}	0.013	0.038	β_{32}	-0.033	0.086	β_{31}	-0.006	0.027	β_{32}	-0.009	0.029
β_{41}	0.047	0.066	β_{42}	-0.046	0.089	β_{41}	-0.022	0.042	β_{42}	0.008	0.041
ψ_1	0.022	0.070	ψ_2	0.036	0.081	ψ_1	0.013	0.018	ψ_2	0.022	0.035
Parameters in the transition model											
ζ_{11}	-0.048	0.144	ζ_{21}	0.017	0.135	ζ_{11}	0.028	0.137	ζ_{21}	0.019	0.089
π_1	-0.010	0.052	π_2	0.010	0.052	π_1	-0.008	0.021	π_2	0.008	0.021
α_1	0.036	0.112				α_1	0.002	0.087			

Table S2: Parameter estimates under disturb residual in Simulation 1: $K_0 = 2$, $n = 200$

Case 1: $\delta_{it} \sim U(-1, 1)$						Case 2: $\delta_{it} \sim 0.4N(1, 1) + 0.6N(-1, 1)$					
State 1			State 2			State 1			State 2		
Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS
Parameters in the conditional regression model											
β_{11}	0.009	0.044	β_{12}	0.018	0.051	β_{11}	-0.091	0.096	β_{12}	0.061	0.073
β_{21}	0.007	0.027	β_{22}	-0.007	0.032	β_{21}	-0.012	0.021	β_{22}	-0.005	0.032
β_{31}	-0.006	0.024	β_{32}	-0.005	0.027	β_{31}	0.010	0.028	β_{32}	-0.011	0.044
β_{41}	-0.009	0.034	β_{42}	-0.027	0.079	β_{41}	-0.005	0.046	β_{42}	-0.036	0.079
ψ_1	0.091	0.103	ψ_2	0.108	0.163	ψ_1	0.256	0.267	ψ_2	0.272	0.276
Parameters in the transition model											
ζ_{11}	-0.021	0.092	ζ_{21}	0.047	0.230	ζ_{11}	-0.066	0.185	ζ_{21}	-0.013	0.124
π_1	-0.005	0.027	π_2	0.005	0.027	π_1	-0.002	0.029	π_2	0.002	0.029
α_1	-0.024	0.169				α_1	0.036	0.105			

Table S3: Parameter estimates in Simulation 2: $K_0 = 3$

State 1			State 2			State 3		
Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS
Parameters in the conditional regression model								
β_{11}	0.015	0.031	β_{12}	-0.006	0.030	β_{13}	0.001	0.037
β_{21}	0.013	0.039	β_{22}	-0.029	0.046	β_{23}	-0.006	0.046
β_{31}	0.012	0.029	β_{32}	0.021	0.036	β_{33}	0.011	0.026
β_{41}	-0.025	0.076	β_{42}	-0.001	0.041	β_{43}	-0.009	0.043
ψ_1	0.006	0.018	ψ_2	0.016	0.023	ψ_3	0.029	0.063
Parameters in the transition model								
ζ_{11}	0.016	0.086	ζ_{21}	0.014	0.110	ζ_{31}	-0.033	0.105
ζ_{12}	-0.023	0.136	ζ_{22}	-0.016	0.139	ζ_{32}	0.022	0.098
π_1	0.007	0.021	π_2	-0.009	0.029	π_3	0.002	0.023
α_1	0.036	0.073	α_2	0.023	0.102	α_3	0.006	0.145

Table S4: Parameter estimates in Simulation 2: $K_0 = 5$

State 1			State 2			State 3			State 4			State 5		
Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS	Par	Bias	RMS
Parameters in the conditional regression model														
β_{11}	0.028	0.044	β_{12}	-0.007	0.037	β_{13}	0.006	0.021	β_{14}	-0.023	0.040	β_{15}	0.016	0.049
β_{21}	-0.024	0.038	β_{22}	-0.017	0.046	β_{23}	-0.007	0.033	β_{24}	-0.024	0.045	β_{25}	-0.003	0.039
β_{31}	-0.031	0.073	β_{32}	0.037	0.059	β_{33}	0.029	0.056	β_{34}	0.040	0.086	β_{35}	0.021	0.046
ψ_1	0.006	0.021	ψ_2	0.004	0.020	ψ_3	0.016	0.025	ψ_4	0.029	0.044	ψ_5	0.038	0.091
Parameters in the transition model														
ζ_{11}	0.045	0.108	ζ_{21}	0.056	0.129	ζ_{31}	0.049	0.104	ζ_{41}	0.035	0.060	ζ_{51}	0.052	0.147
ζ_{12}	0.014	0.103	ζ_{22}	0.015	0.092	ζ_{32}	0.061	0.109	ζ_{42}	-0.025	0.134	ζ_{52}	-0.053	0.123
ζ_{13}	0.055	0.154	ζ_{23}	0.063	0.149	ζ_{33}	0.048	0.089	ζ_{43}	-0.069	0.110	ζ_{53}	0.010	0.058
ζ_{14}	-0.078	0.147	ζ_{24}	0.032	0.121	ζ_{34}	-0.058	0.160	ζ_{44}	-0.044	0.137	ζ_{54}	-0.074	0.129
π_1	-0.001	0.016	π_2	-0.001	0.023	π_3	0.003	0.012	π_4	-0.001	0.023	π_5	0.000	0.016
α_1	-0.009	0.050	α_2	-0.028	0.085									

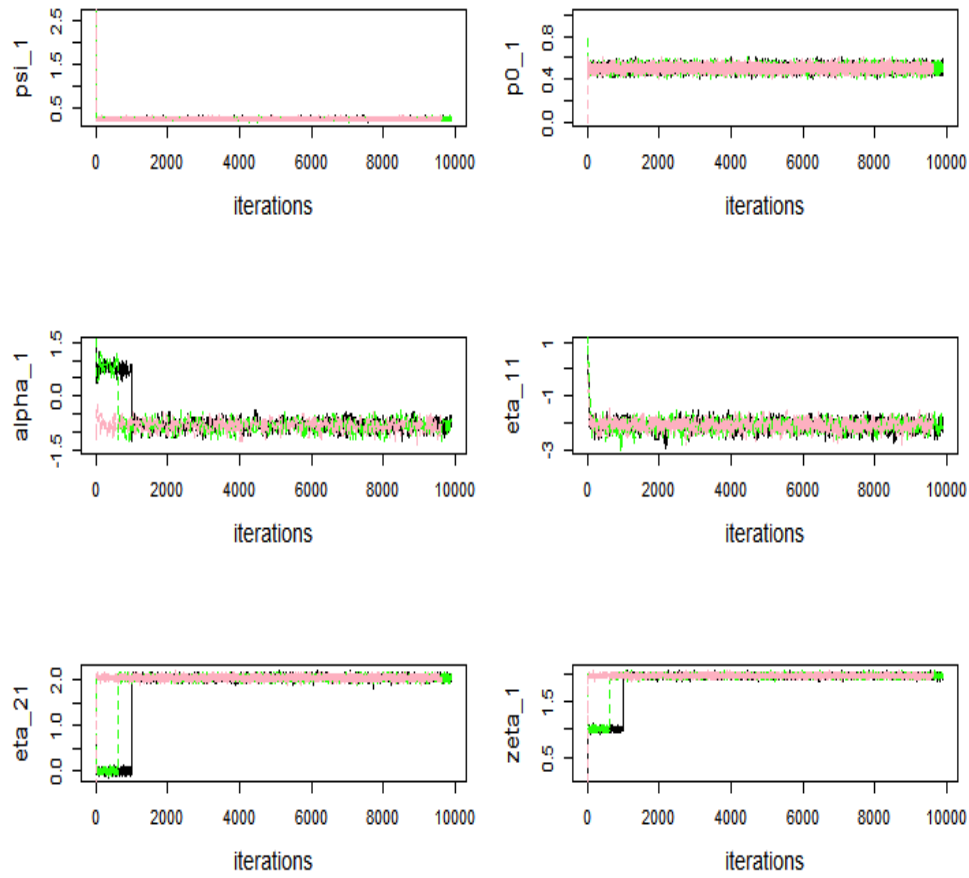


Figure S1: Trace plots of three MCMC chains of the randomly selected parameters in Simulation 1 ($K_0 = 2$).

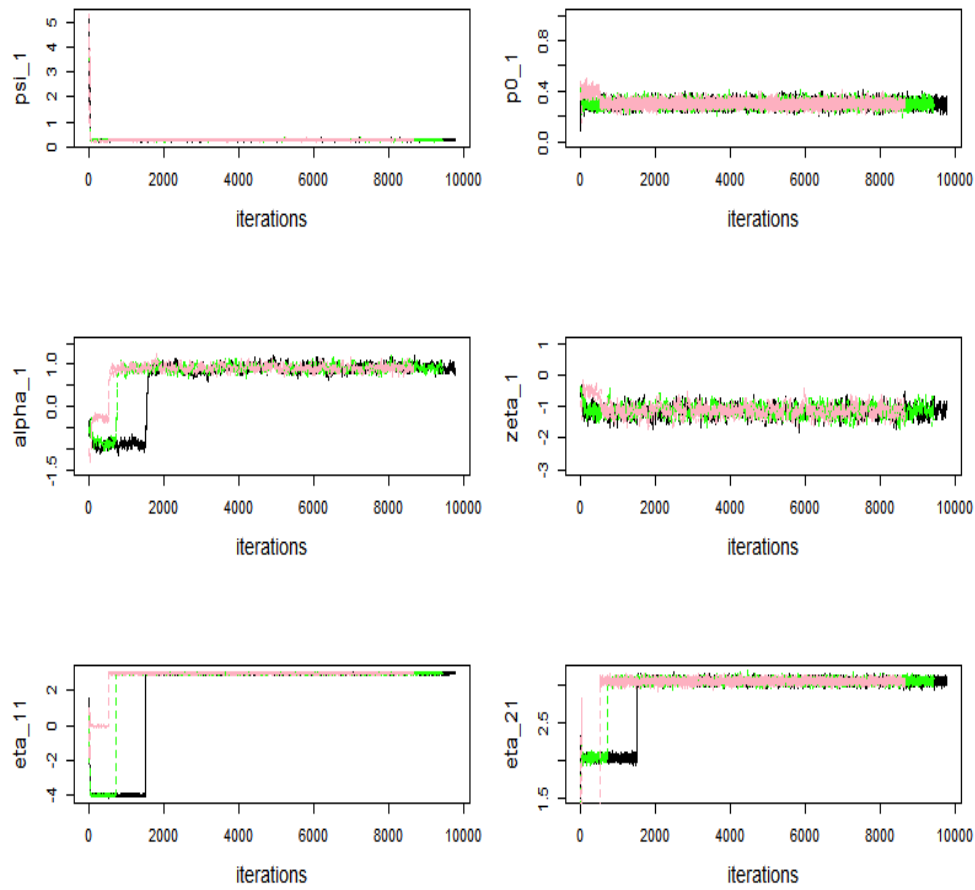


Figure S2: Trace plots of three MCMC chains of the randomly selected parameters in Simulation 2 ($K_0 = 3$).

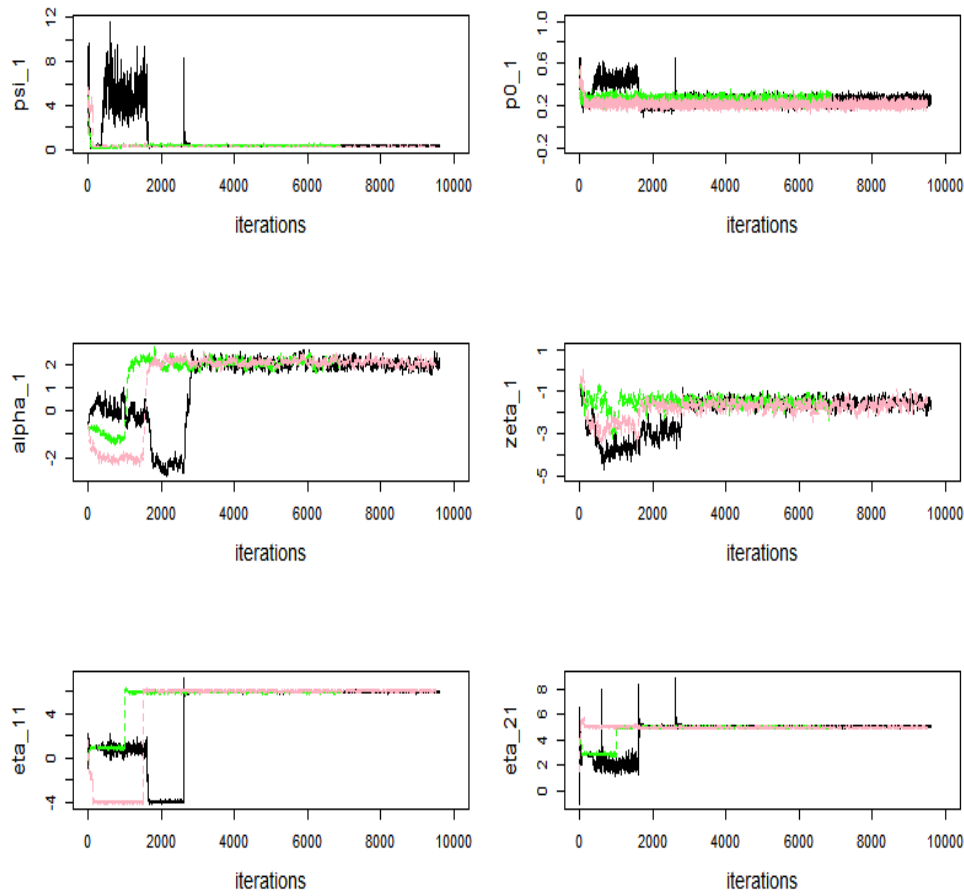


Figure S3: Trace plots of three MCMC chains of the randomly selected parameters in Simulation 2 ($K_0 = 5$).

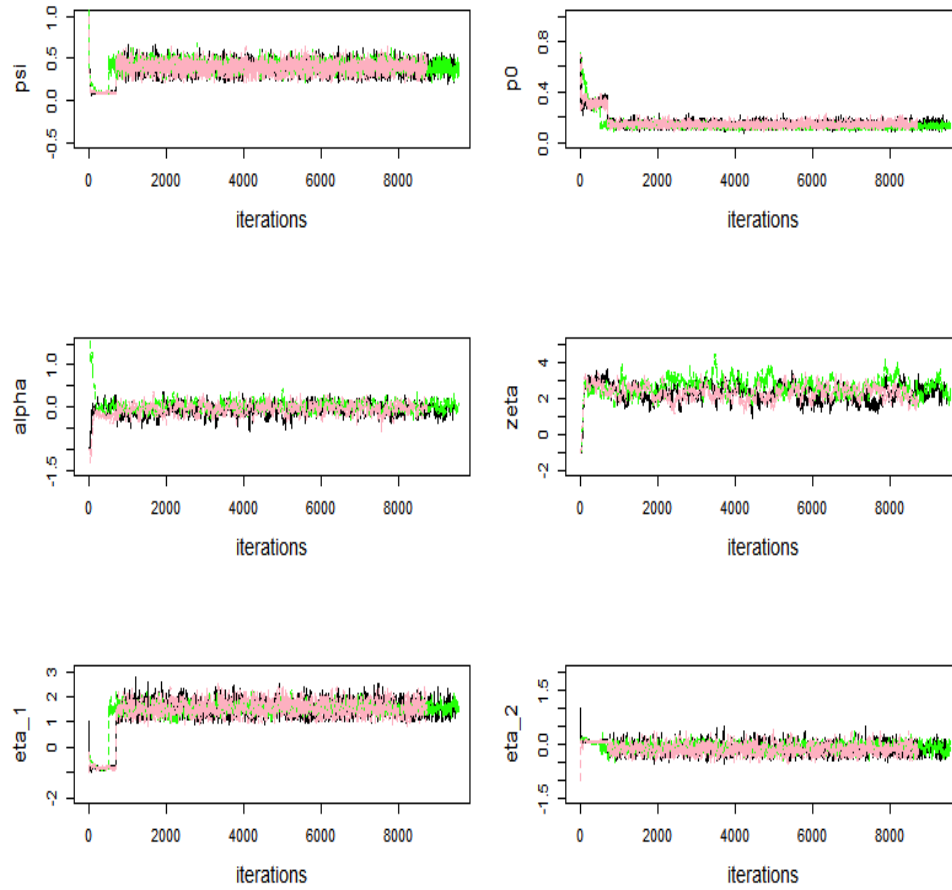


Figure S4: Trace plots of the three MCMC chains of the randomly selected parameters in the ADNI study.