

Bayesian Nonparametric Inference of Switching Dynamic Linear Models

Emily Fox, *Member, IEEE*, Erik Sudderth, *Member, IEEE*, Michael Jordan, *Fellow, IEEE*,
and Alan Willsky, *Fellow, IEEE*

Abstract—Many complex dynamical phenomena can be effectively modeled by a system that switches among a set of conditionally linear dynamical modes. We consider two such models: the switching linear dynamical system (SLDS) and the switching vector autoregressive (VAR) process. Our Bayesian nonparametric approach utilizes a hierarchical Dirichlet process prior to learn an unknown number of persistent, smooth dynamical modes. We additionally employ automatic relevance determination to infer a sparse set of dynamic dependencies allowing us to learn SLDS with varying state dimension or switching VAR processes with varying autoregressive order. We develop a sampling algorithm that combines a truncated approximation to the Dirichlet process with efficient joint sampling of the mode and state sequences. The utility and flexibility of our model are demonstrated on synthetic data, sequences of dancing honey bees, the IBOVESPA stock index, and a maneuvering target tracking application.

I. INTRODUCTION

LINEAR dynamical systems (LDSs) are useful in describing dynamical phenomena as diverse as human motion [3], [4], financial time-series [5]–[7], maneuvering targets [8], [9], and the dance of honey bees [10]. However, such phenomena often exhibit structural changes over time, and the LDS models which describe them must also change. For example, a ballistic missile makes an evasive maneuver; a country experiences a recession, a central bank intervention, or some national or global event; a honey bee changes from a *waggle*

to a *turn right* dance. Some of these changes will appear frequently, while others are only rarely observed. In addition, there is always the possibility of a new, previously unseen dynamical behavior. These considerations motivate us to develop a Bayesian nonparametric approach for learning *switching* LDS (SLDS) models. We also consider a special case of the SLDS—the switching vector autoregressive (VAR) model—in which direct observations of the underlying dynamical process are assumed available.

One can view the SLDS, and the simpler switching VAR process, as an extension of hidden Markov models (HMMs) in which each HMM state, or *mode*, is associated with a linear dynamical process. Within the signal processing community, such HMM-based models have received considerable attention and proven useful in modeling the complex time evolution of signals. Specifically, HMMs have a long history of signal processing applications, with major success stories in speech processing (see the early influential tutorial by Rabiner [11]). While the HMM makes a strong Markovian assumption that observations are conditionally independent given the mode, the SLDS and switching VAR processes are able to capture more complex temporal dependencies often present in real data. Applications of switching linear dynamical processes, with roots in the control and econometrics literature, have recently become more prevalent within signal processing [10], [12]–[14]. However, most existing methods for learning SLDS and switching VAR processes rely on either fixing the number of HMM modes, such as in the preceding papers, or considering a change-point detection formulation where each inferred change is to a new, previously unseen dynamical mode, such as in [15]. There is growing interest in expanding the modeling framework to remove the purely parametric assumption of these previous formulations. In this paper we show how one can, in a seamless manner, remain agnostic about the number of dynamical modes while still allowing for returns to previously exhibited dynamical behaviors.

The rapidly developing field of *Bayesian nonparametrics* provides a new direction for analyzing HMMs

E. Fox is with the Department of Statistical Science, Duke University, Durham, NC, 27708 USA e-mail: fox@stat.duke.edu. E. Sudderth is with the Department of Computer Science, Brown University, Providence, RI, 02912 USA e-mail: sudderth@cs.brown.edu. M. Jordan is with the Department of Electrical Engineering and Computer Science, and Department of Statistics, University of California, Berkeley, CA, 94720 USA e-mail: jordan@eecs.berkeley.edu. A. Willsky is with the Department of Electrical Engineering and Computer Science, MIT, Cambridge, MA, 02139 USA e-mail: willsky@mit.edu. This work was supported in part by MURIs funded through AFOSR Grant FA9550-06-1-0324 and ARO Grant W911NF-06-1-0076. E. Fox was partially supported by NSF Grant DMS-0903022. Preliminary versions (without detailed development or analysis) of this work have been presented at two conferences [1], [2].

with unknown state space cardinality. In particular, it has been shown that the hierarchical Dirichlet process (HDP) provides a useful prior on the HMM parameters [16], [17]. An alternative formulation of a Bayesian nonparametric HMM with application to music analysis has been presented in [18], though without the shared sparsity induced by the HDP. Another application of Bayesian nonparametrics to music analysis was presented in [19], where the authors propose Dirichlet process clustering of fixed-length segments of a time series, with each cluster modeling the dynamics of the given segments via a different finite HMM. See also [20] for a signal processing application of Dirichlet processes, specifically nonparametric modeling of excitations to a switching dynamical process. In this paper we make use of a variant of the HDP-HMM—the *sticky HDP-HMM* of [21]—to obtain improved control over the number of modes inferred; such control is crucial for the problems we examine. Our Bayesian nonparametric approach for learning switching dynamical processes extends the sticky HDP-HMM formulation to infer an unknown number of persistent dynamical modes and thereby capture a wider range of temporal dependencies. We then explore a method for learning which components of the underlying state vector contribute to the dynamics of each mode by employing *automatic relevance determination* (ARD) [22]–[24]. The resulting model allows for learning realizations of SLDS that switch between an unknown number of dynamical modes with possibly varying state dimensions, or switching VAR processes with varying autoregressive orders.

A. Previous System Identification Techniques

Paoletti et. al. [25] provide a survey of recent approaches to identification of switching dynamical models. The most general formulation of the problem involves learning: (i) the number of dynamical modes, (ii) the model order, and (iii) the associated dynamic parameters. For noiseless switching VAR processes, Vidal et. al. [26] present an exact algebraic approach, though relying on fixing a maximal mode space cardinality and autoregressive order. Psaradakis and Spagnolo [27] alternatively consider a penalized likelihood approach to identification of stochastic switching VAR processes.

For SLDS, identification is significantly more challenging, and methods typically rely on simplifying assumptions such as deterministic dynamics or knowledge of the mode space. Huang et. al. [28] present an approach that assumes deterministic dynamics and embeds the input/output data in a higher-dimensional space and finds the switching times by segmenting the data into distinct subspaces [29]. Kotsalis et. al. [30] develop a

balanced truncation algorithm for SLDS assuming the mode switches are i.i.d. within a fixed, finite set; the authors also present a method for model-order reduction of HMMs (see also [31]). In [32], a realization theory is presented for *generalized jump-Markov linear systems* in which the dynamic matrix depends both on the previous mode and current mode. Finally, when the number of dynamical modes is assumed known, Ghahramani and Hinton [33] present a variational approach to segmenting the data from a *mixture of experts* SLDS into the linear dynamical regimes and learning the associated dynamic parameters. For questions on observability and identifiability of SLDS in the absence of noise, see [34].

In the Bayesian approach that we adopt, we coherently incorporate noisy dynamics and uncertainty in the mode space cardinality. Our choice of prior penalizes more complicated models, both in terms of the number of modes and the state dimension describing each mode, allowing us to distinguish between the set of equivalent models described in [34]. Thus, instead of placing hard constraints on the model, we simply increase the posterior probability of simpler explanations of the data. As opposed to a penalized likelihood approach using *Akaike's information criterion* (AIC) [35] or the *Bayesian information criterion* (BIC) [36], our approach provides a model complexity penalty in a purely Bayesian manner.

In Sec. II, we provide background on the switching linear dynamical systems we consider herein, and previous Bayesian nonparametric methods of learning HMMs. Our Bayesian nonparametric switching linear dynamical systems are described in Sec. III. We proceed by analyzing a conjugate prior on the dynamic parameters, and a sparsity-inducing prior that allows for variable-order switching processes. The section concludes by outlining a Gibbs sampler for the proposed models. In Sec. IV we present results on synthetic and real datasets, and in Sec. V we analyze a set of alternative formulations that are commonly found in the maneuvering target tracking and econometrics literature.

II. BACKGROUND

A. Switching Linear Dynamic Systems

A state space (SS) model consists of an underlying state, $\mathbf{x}_t \in \mathbb{R}^n$, with dynamics observed via $\mathbf{y}_t \in \mathbb{R}^d$. A linear time-invariant (LTI) SS model is given by

$$\mathbf{x}_t = A\mathbf{x}_{t-1} + \mathbf{e}_t \quad \mathbf{y}_t = C\mathbf{x}_t + \mathbf{w}_t, \quad (1)$$

where \mathbf{e}_t and \mathbf{w}_t are independent Gaussian noise processes with covariances Σ and R , respectively.

An order r VAR process, denoted by $\text{VAR}(r)$, with observations $\mathbf{y}_t \in \mathbb{R}^d$, can be defined as

$$\mathbf{y}_t = \sum_{i=1}^r A_i \mathbf{y}_{t-i} + \mathbf{e}_t \quad \mathbf{e}_t \sim \mathcal{N}(0, \Sigma). \quad (2)$$

Every $\text{VAR}(r)$ process can be described in SS form, though not every SS model may be expressed as a $\text{VAR}(r)$ process for finite r [37].

The dynamical phenomena we examine in this paper exhibit behaviors better modeled as switches between a set of linear dynamical models. We define a *switching linear dynamical system* (SLDS) by

$$\begin{aligned} z_t \mid z_{t-1} &\sim \pi_{z_{t-1}} \\ \mathbf{x}_t = A^{(z_t)} \mathbf{x}_{t-1} + \mathbf{e}_t^{(z_t)} \quad \mathbf{y}_t &= C \mathbf{x}_t + \mathbf{w}_t. \end{aligned} \quad (3)$$

The first-order Markov process z_t with transition distributions $\{\pi_j\}$ indexes the mode-specific LDS at time t , which is driven by Gaussian noise $\mathbf{e}_t^{(z_t)} \sim \mathcal{N}(0, \Sigma^{(z_t)})$. One can view the SLDS as an extension of the classical hidden Markov model (HMM) [11], which has the same mode evolution, but conditionally *independent* observations:

$$\begin{aligned} z_t \mid z_{t-1} &\sim \pi_{z_{t-1}} \\ y_t \mid z_t &\sim F(\theta_{z_t}) \end{aligned} \quad (4)$$

for an indexed family of distributions $F(\cdot)$ where θ_i are the *emission parameters* for mode i .

We similarly define a *switching VAR*(r) process by

$$\begin{aligned} z_t \mid z_{t-1} &\sim \pi_{z_{t-1}} \\ \mathbf{y}_t &= \sum_{i=1}^r A_i^{(z_t)} \mathbf{y}_{t-i} + \mathbf{e}_t^{(z_t)}. \end{aligned} \quad (5)$$

B. Dirichlet Processes and the Sticky HDP-HMM

To examine a Bayesian nonparametric SLDS, and thus relax the assumption that the number of dynamical modes is known and fixed, it is useful to first analyze such methods for the simpler HMM. One can equivalently represent the finite HMM of Eq. (4) via a set of *transition probability measures* $G_j = \sum_{k=1}^K \pi_{jk} \delta_{\theta_k}$, where δ_{θ} is a mass concentrated at θ . We then operate directly in the parameter space Θ and transition between emission parameters with probabilities given by $\{G_j\}$. That is,

$$\begin{aligned} \theta'_t \mid \theta'_{t-1} &\sim G_{j: \theta'_{t-1} = \theta_j} \\ y_t \mid \theta'_t &\sim F(\theta'_t). \end{aligned} \quad (6)$$

Here, $\theta'_t \in \{\theta_1, \dots, \theta_K\}$ and is equivalent to θ_{z_t} of Eq. (4). A Bayesian nonparametric HMM takes G_j to

be *random*¹ with an infinite collection of atoms corresponding to the infinite HMM mode space.

The *Dirichlet process* (DP), denoted by $\text{DP}(\gamma, H)$, provides a distribution over discrete probability measures with an infinite collection of atoms

$$G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k} \quad \theta_k \sim H, \quad (7)$$

on a parameter space Θ that is endowed with a *base measure* H . The weights are sampled via a *stick-breaking construction* [38]:

$$\beta_k = \nu_k \prod_{\ell=1}^{k-1} (1 - \nu_{\ell}) \quad \nu_k \sim \text{Beta}(1, \gamma). \quad (8)$$

In effect, we have divided a unit-length stick into lengths given by the weights β_k : the k^{th} weight is a random proportion ν_k of the remaining stick after the previous $(k-1)$ weights have been defined. Letting $\beta = [\beta_1 \ \beta_2 \ \dots]$, we denote this distribution by $\beta \sim \text{GEM}(\gamma)$.

The DP has proven useful in many applications due to its clustering properties, which are clearly seen by examining the *predictive distribution* of draws $\theta'_i \sim G_0$. Because probability measures drawn from a DP are discrete, there is a strictly positive probability of multiple observations θ'_i taking identical values within the set $\{\theta_k\}$, with θ_k defined as in Eq. (7). For each value θ'_i , let z_i be an indicator random variable that picks out the unique value θ_k such that $\theta'_i = \theta_{z_i}$. Blackwell and MacQueen [39] introduced a Pólya urn representation of the θ'_i :

$$\begin{aligned} \theta'_i \mid \theta'_1, \dots, \theta'_{i-1} &\sim \frac{\gamma}{\gamma + i - 1} H + \sum_{j=1}^{i-1} \frac{1}{\gamma + i - 1} \delta_{\theta'_j} \\ &= \frac{\gamma}{\gamma + i - 1} H + \sum_{k=1}^K \frac{n_k}{\gamma + i - 1} \delta_{\theta_k}. \end{aligned} \quad (9)$$

Here, n_k is the number of observations θ'_i taking the value θ_k . From Eq. (9), and the discrete nature of G_0 , we see a reinforcement property of the DP that induces sparsity in the number of inferred mixture components.

A hierarchical extension of the DP, the hierarchical Dirichlet process (HDP) [16], has proven useful in defining a prior on the set of HMM transition probability measures G_j . The HDP defines a collection of probability measures $\{G_j\}$ on the same support points $\{\theta_1, \theta_2, \dots\}$ by assuming that each discrete measure

¹Formally, a random measure on a measurable space Θ with sigma algebra \mathcal{A} is defined as a stochastic process whose index set is \mathcal{A} . That is, $G(A)$ is a random variable for each $A \in \mathcal{A}$.

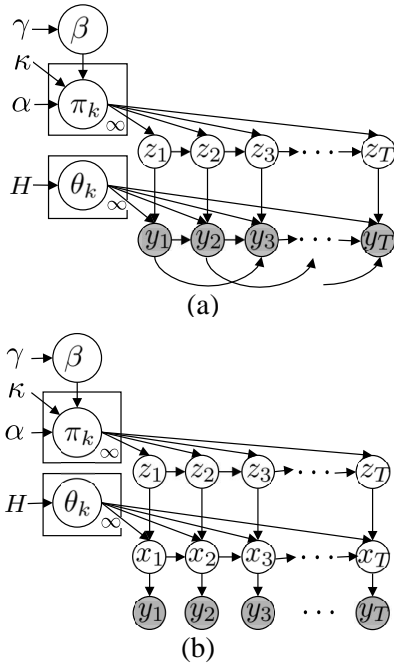


Fig. 1. Sticky HDP-HMM prior on (a) switching VAR(2) and (b) SLDS processes with the mode evolving as $z_{t+1} | \{\pi_k\}_{k=1}^\infty, z_t \sim \pi_{z_t}$ for $\pi_k | \alpha, \kappa, \beta \sim \text{DP}(\alpha + \kappa, (\alpha\beta + \kappa\delta_k)/(\alpha + \kappa))$. Here, $\beta | \gamma \sim \text{GEM}(\gamma)$ and $\theta_k | H \sim H$. The dynamical processes are as in Table I.

G_j is a variation on a global discrete measure G_0 . Specifically, the Bayesian hierarchical specification takes $G_j \sim \text{DP}(\alpha, G_0)$, with G_0 itself a draw from a $\text{DP}(\gamma, H)$. Through this construction, one can show that the probability measures are described as

$$\begin{aligned} G_0 &= \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k} & \beta | \gamma &\sim \text{GEM}(\gamma) \\ G_j &= \sum_{k=1}^{\infty} \pi_{jk} \delta_{\theta_k} & \pi_j | \alpha, \beta &\sim \text{DP}(\alpha, \beta) \end{aligned} \quad (10)$$

$$\theta_k | H \sim H.$$

Here, we use the notation $\pi_j = [\pi_{j1} \ \pi_{j2} \ \dots]$. Applying the HDP prior to the HMM, we obtain the *HDP-HMM* of Teh et. al. [16]. This corresponds to the model in Fig. 1(a), but without the edges between the observations.

By defining $\pi_j \sim \text{DP}(\alpha, \beta)$, the HDP prior encourages modes to have similar transition distributions. Namely, the mode-specific transition distributions are *identical* in expectation:

$$\mathbb{E}[\pi_{jk} | \beta] = \beta_k. \quad (11)$$

However, it does not differentiate self-transitions from moves between modes. When modeling dynamical processes with mode persistence, the flexible nature of the HDP-HMM prior allows for mode sequences with unrealistically fast dynamics to have large posterior probability. Recently, it has been shown [21] that one may mitigate this problem by instead considering a *sticky*

HDP-HMM where π_j is distributed as follows:

$$\pi_j | \beta, \alpha, \kappa \sim \text{DP}\left(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_j}{\alpha + \kappa}\right). \quad (12)$$

Here, $(\alpha\beta + \kappa\delta_j)$ indicates that an amount $\kappa > 0$ is added to the j^{th} component of $\alpha\beta$. This construction increases the expected probability of self-transition by an amount proportional to κ . Specifically, the expected set of weights for transition distribution π_j is a convex combination of those defined by β and mode-specific weight defined by κ :

$$\mathbb{E}[\pi_{jk} | \beta, \alpha, \kappa] = \frac{\alpha}{\alpha + \kappa} \beta_k + \frac{\kappa}{\alpha + \kappa} \delta(j, k). \quad (13)$$

Here, $\delta(j, k)$ denotes the discrete Kronecker delta. When $\kappa = 0$ the original HDP-HMM of Teh et. al. [16] is recovered. We place a prior on κ and learn the self-transition bias from the data. See [21] for details.

III. THE HDP-SLDS AND HDP-AR-HMM

We now consider a significant extension of the sticky HDP-HMM for both SLDS and VAR modeling, capturing dynamic structure underlying the observations by allowing switches among an unknown number dynamical modes. Our proposed Bayesian nonparametric approach aims to capture these uncertainties. Additionally, the methodology allows both learning the number of modes and estimating the dimensionality and associated parameterization of the system state process. Fig. 1(b) illustrates the *HDP-SLDS* model, while Fig. 1(a) illustrates the *HDP-AR-HMM* model (for the case of VAR(2)). The generative processes for these two models are summarized in Table I.

The prior on the underlying discrete-valued Markov process $\{z_t\}$ is just as in the sticky HDP-HMM. The question now is in determining an appropriate base measure H for the model parameters θ_k . For the HDP-SLDS, we place priors on the *dynamic parameters* $\{A^{(k)}, \Sigma^{(k)}\}$ and on the measurement noise covariance R and infer their posterior from the data. Note that we assume the dynamics of the latent state process are mode-specific, while the measurement mechanism is not. This assumption could be modified to allow for both a mode-specific measurement matrix $C^{(z_t)}$ and noise $w_t^{(z_t)} \sim \mathcal{N}(0, R^{(z_t)})$. However, such a choice is not always necessary nor appropriate for certain applications, and can have implications on the identifiability of the model. Based on a shared measurement matrix C , we fix $C = [I_d \ 0]$ without loss of generality, implying that it is the first d components of the state that are measured. Our choice of the state dimension n is, in essence, a choice of model order, and an issue we address in Sec. III-A2.

| | HDP-AR-HMM | HDP-SLDS |
|----------------------|---|--|
| Mode dynamics | $z_t \mid z_{t-1} \sim \pi_{z_{t-1}}$ | $z_t \mid z_{t-1} \sim \pi_{z_{t-1}}$ |
| Observation dynamics | $\mathbf{y}_t = \sum_{i=1}^r A_i^{(z_t)} \mathbf{y}_{t-i} + \mathbf{e}_t^{(z_t)}$ | $\mathbf{x}_t = A^{(z_t)} \mathbf{x}_{t-1} + \mathbf{e}_t^{(z_t)}$ $\mathbf{y}_t = C \mathbf{x}_t + \mathbf{w}_t$ |

TABLE I

DYNAMIC EQUATIONS FOR THE HDP-AR-HMM AND HDP-SLDS. HERE, π_j IS AS DEFINED IN EQ. (12) FOR THE STICKY HDP-HMM. THE ADDITIVE NOISE PROCESSES ARE DISTRIBUTED AS $\mathbf{e}_t^{(k)} \sim \mathcal{N}(0, \Sigma^{(k)})$ AND $\mathbf{w}_t \sim \mathcal{N}(0, R)$.

For the HDP-AR-HMM, we similarly place a prior on the dynamic parameters, which in this case consist of $\{A_1^{(k)}, \dots, A_r^{(k)}, \Sigma^{(k)}\}$. Our specific choice of priors is discussed in Sec. III-A.

A Gibbs sampling inference scheme for our models is derived in Sec. III-B. There is, of course, a difference between the steps required for the SLDS-based model (in which there is an unobserved continuous-valued state \mathbf{x}_t) and the AR-based model. In particular, for the HDP-SLDS the algorithm iterates among the following steps:

- 1) Sample the state sequence $\mathbf{x}_{1:T}$ given the mode sequence $z_{1:T}$ and SLDS parameters $\{A^{(k)}, \Sigma^{(k)}, R\}$.
- 2) Sample the mode sequence $z_{1:T}$ given the state sequence $\mathbf{x}_{1:T}$, HMM parameters $\{\pi_k\}$, and dynamic parameters $\{A^{(k)}, \Sigma^{(k)}\}$.
- 3) Sample the HMM parameters $\{\pi_k\}$ and SLDS parameters $\{A^{(k)}, \Sigma^{(k)}, R\}$ given the sequences $z_{1:T}$, $\mathbf{x}_{1:T}$, and $\mathbf{y}_{1:T}$.

For the HDP-AR-HMM, step (1) does not exist. Step (2) then involves sampling the mode sequence $z_{1:T}$ given the observations $\mathbf{y}_{1:T}$ (rather than $\mathbf{x}_{1:T}$), and step (3) involves conditioning solely on the sequences $z_{1:T}$ and $\mathbf{y}_{1:T}$ (not $\mathbf{x}_{1:T}$). Also, we note that step (2) involves a fairly straightforward extension of the sampling method developed in [21] for the simpler HDP-HMM model; the other steps, however, involve new constructs, as they require capturing and dealing with the temporal dynamics of the underlying continuous state models. Sec. III-A provides the structure of the posteriors needed to develop these steps.

A. Priors and Posteriors of Dynamic Parameters

We begin by developing a prior to regularize the learning of the dynamic parameters (and measurement noise) conditioned on a fixed mode assignment $z_{1:T}$. To make the connections between the samplers for the HDP-SLDS and HDP-AR-HMM explicit, we introduce the concept of *pseudo-observations* $\psi_{1:T}$ and rewrite the dynamic equation for both the HDP-SLDS and HDP-AR-HMM generically as

$$\psi_t = \mathbf{A}^{(k)} \bar{\psi}_{t-1} + \mathbf{e}_t^{(k)}, \quad (14)$$

where we utilize the definitions outlined in Table II.

For the HDP-AR-HMM, we have simply written the dynamic equation in Table I in matrix form by concatenating the lag matrices into a single matrix $\mathbf{A}^{(k)}$ and forming a *lag observation vector* $\bar{\psi}_t$ comprised of a series of previous observation vectors. For this section (for the HDP-SLDS), we assume a sample of the state sequence $\mathbf{x}_{1:T}$ (and hence $\{\psi_t, \bar{\psi}_t\}$) is available so that Eq. (14) applies equally well to both the HDP-SLDS and the HDP-AR-HMM. Methods for resampling this state sequence are discussed in Sec. III-B.

Conditioned on the mode sequence, one may partition this dynamic sequence into K different linear regression problems, where $K = |\{z_1, \dots, z_T\}|$. That is, for each mode k , we may form a matrix $\Psi^{(k)}$ with n_k columns consisting of the ψ_t with $z_t = k$. Then,

$$\Psi^{(k)} = \mathbf{A}^{(k)} \bar{\Psi}^{(k)} + \mathbf{E}^{(k)}, \quad (15)$$

where $\bar{\Psi}^{(k)}$ is a matrix of the associated $\bar{\psi}_{t-1}$, and $\mathbf{E}^{(k)}$ the associated noise vectors.

1) *Conjugate Prior on $\{\mathbf{A}^{(k)}, \Sigma^{(k)}\}$* : The *matrix-normal inverse-Wishart* (MNIW) prior [40] is conjugate to the likelihood model defined in Eq. (15) for the parameter set $\{\mathbf{A}^{(k)}, \Sigma^{(k)}\}$. Although this prior is typically used for inferring the parameters of a single linear regression problem, it is equally applicable to our scenario since the linear regression problems of Eq. (15) are independent conditioned on the mode sequence $z_{1:T}$. We note that while the MNIW prior does not enforce stability constraints on each mode, this prior is still a reasonable choice since each mode need not have stable dynamics for the SLDS to be stable [41], and conditioned on data from a stable mode, the posterior distribution will likely be sharply peaked around stable dynamic matrices.

Let $\mathbf{D}^{(k)} = \{\Psi^{(k)}, \bar{\Psi}^{(k)}\}$. The posterior distribution of the dynamic parameters for the k^{th} mode decomposes as

$$p(\mathbf{A}^{(k)}, \Sigma^{(k)} \mid \mathbf{D}^{(k)}) = p(\mathbf{A}^{(k)} \mid \Sigma^{(k)}, \mathbf{D}^{(k)}) p(\Sigma^{(k)} \mid \mathbf{D}^{(k)}). \quad (16)$$

The resulting posterior of $\mathbf{A}^{(k)}$ is straightforwardly derived to be (see [42])

$$p(\mathbf{A}^{(k)} \mid \Sigma^{(k)}, \mathbf{D}^{(k)}) = \mathcal{MN} \left(\mathbf{S}_{\psi\bar{\psi}}^{(k)} \mathbf{S}_{\bar{\psi}\bar{\psi}}^{-(k)}, \Sigma^{(k)}, \mathbf{S}_{\bar{\psi}\bar{\psi}}^{(k)} \right), \quad (17)$$

| | HDP-AR-HMM | HDP-SLDS |
|-------------------------|--|--|
| Dynamic matrix | $\mathbf{A}^{(k)} = [A_1^{(k)} \dots A_r^{(k)}] \in \mathbb{R}^{d \times (d \star r)}$ | $\mathbf{A}^{(k)} = A^{(k)} \in \mathbb{R}^{n \times n}$ |
| Pseudo-observations | $\boldsymbol{\psi}_t = \mathbf{y}_t$ | $\boldsymbol{\psi}_t = \mathbf{x}_t$ |
| Lag pseudo-observations | $\bar{\boldsymbol{\psi}}_t = [\mathbf{y}_{t-1}^T \dots \mathbf{y}_{t-r}^T]^T$ | $\bar{\boldsymbol{\psi}}_t = \mathbf{x}_{t-1}$ |

TABLE II

NOTATIONAL CONVENIENCES USED IN DESCRIBING THE GIBBS SAMPLER FOR THE HDP-AR-HMM AND HDP-SLDS.

with $\mathbf{B}^{-(k)}$ denoting $(\mathbf{B}^{(k)})^{-1}$ for a given matrix \mathbf{B} , $\mathcal{MN}(M, V, K)$ denoting a matrix-normal prior² for $\mathbf{A}^{(k)}$ with mean matrix M and left and right covariances K^{-1} and V , and

$$\begin{aligned} \mathbf{S}_{\psi\bar{\psi}}^{(k)} &= \bar{\boldsymbol{\Psi}}^{(k)} \bar{\boldsymbol{\Psi}}^{(k)T} + K & \mathbf{S}_{\psi\bar{\psi}}^{(k)} &= \boldsymbol{\Psi}^{(k)} \bar{\boldsymbol{\Psi}}^{(k)T} + MK \\ \mathbf{S}_{\psi\psi}^{(k)} &= \boldsymbol{\Psi}^{(k)} \boldsymbol{\Psi}^{(k)T} + MKM^T. \end{aligned} \quad (18)$$

The marginal posterior of $\Sigma^{(k)}$ is

$$p(\Sigma^{(k)} | \mathbf{D}^{(k)}) = \text{IW}\left(n_k + n_0, \mathbf{S}_{\psi|\bar{\psi}}^{(k)} + S_0\right), \quad (19)$$

where $\text{IW}(n_0, S_0)$ denotes an inverse-Wishart prior for $\Sigma^{(k)}$ with n_0 degrees of freedom and scale matrix S_0 , and is updated by data terms $\mathbf{S}_{\psi|\bar{\psi}}^{(k)} = \mathbf{S}_{\psi\psi}^{(k)} - \mathbf{S}_{\psi\bar{\psi}}^{(k)} \mathbf{S}_{\bar{\psi}\bar{\psi}}^{-(k)} \mathbf{S}_{\bar{\psi}\psi}^{(k)T}$ and $n_k = |\{t \mid z_t = k, t = 1, \dots, T\}|$.

2) *Alternative Prior — Automatic Relevance Determination*: The MNIW prior leads to full $\mathbf{A}^{(k)}$ matrices, which (i) becomes problematic as the model order grows in the presence of limited data; and (ii) does not provide a method for identifying irrelevant model components (i.e. state components in the case of the HDP-SLDS or lag components in the case of the HDP-AR-HMM.) To jointly address these issues, we alternatively consider *automatic relevance determination* (ARD) [22]–[24], which encourages driving components of the model parameters to zero if their presence is not supported by the data.

For the HDP-SLDS, we harness the concepts of ARD by placing independent, zero-mean, spherically symmetric Gaussian priors on the columns of the dynamic matrix $\mathbf{A}^{(k)}$:

$$p(\mathbf{A}^{(k)} | \boldsymbol{\alpha}^{(k)}) = \prod_{j=1}^n \mathcal{N}\left(\mathbf{a}_j^{(k)}; 0, \alpha_j^{-(k)} \mathbf{I}_n\right). \quad (20)$$

Each precision parameter $\alpha_j^{(k)}$ is given a $\text{Gamma}(a, b)$ prior. The zero-mean Gaussian prior penalizes non-zero columns of the dynamic matrix by an amount determined by the precision parameters. Iterative estimation of these hyperparameters $\alpha_j^{(k)}$ and the dynamic matrix $\mathbf{A}^{(k)}$ leads to $\alpha_j^{(k)}$ becoming large for columns whose evidence in the data is insufficient for overcoming the penalty induced by the prior. Having $\alpha_j^{(k)} \rightarrow \infty$ drives $\mathbf{a}_j^{(k)} \rightarrow 0$, implying that the j^{th} state component does not contribute

to the dynamics of the k^{th} mode. Thus, examining the set of large $\alpha_j^{(k)}$ provides insight into the order of that mode. Looking at the k^{th} dynamical mode alone, having $\mathbf{a}_j^{(k)} = 0$ implies that the realization of *that mode* is not minimal since the associated Hankel matrix

$$\mathcal{H} = \begin{bmatrix} C^T & CA^T & \dots & (CA^{d-1})^T \\ G & AG & \dots & A^{d-1}G \end{bmatrix} \equiv \mathcal{OR} \quad (21)$$

has reduced rank. However, the overall SLDS realization may still be minimal.

For our use of the ARD prior, we restrict attention to models satisfying the property that the state components that are observed are relevant to *all* modes of the dynamics:

Criterion 3.1: If for some realization \mathcal{R} a mode k has $\mathbf{a}_j^{(k)} = 0$, then that realization must have $\mathbf{c}_j = 0$, where \mathbf{c}_j is the j^{th} column of C . Here we assume, without loss of generality, that the observed states are the first components of the state vector.

This assumption implies that our choice of $C = [I_d \ 0]$ does not interfere with learning a sparse realization. We could avoid restricting our attention to models satisfying Criterion 3.1 by considering a more general model where the measurement equation is mode-specific and we place a prior on $C^{(k)}$ instead of fixing this matrix. However, this model leads to identifiability issues that are considerably less pronounced in the above case.

The ARD prior may also be used to learn variable-order switching VAR processes. Here, the goal is to “turn off” entire *lag blocks* $A_i^{(k)}$ (whereas in the HDP-SLDS we were interested in eliminating columns of the dynamic matrix.) Instead of placing independent Gaussian priors on each column of $\mathbf{A}^{(k)}$ as we did in Eq. (20), we decompose the prior over the lag blocks $A_i^{(k)}$:

$$p(\mathbf{A}^{(k)} | \boldsymbol{\alpha}^{(k)}) = \prod_{i=1}^r \mathcal{N}\left(\text{vec}(A_i^{(k)}); 0, \alpha_i^{-(k)} \mathbf{I}_{d^2}\right). \quad (22)$$

Since each element of a given lag block $A_i^{(k)}$ is distributed according to the same precision parameter $\alpha_i^{(k)}$, if that parameter becomes large the entire lag block will tend to zero.

In order to examine the posterior distribution on the dynamic matrix $\mathbf{A}^{(k)}$, it is useful to consider the Gaussian induced by Eq. (20) and Eq. (22) on a vectorization

²If $A \sim \mathcal{MN}(M, V, K)$, then $\text{vec}(A) \sim \mathcal{N}(\text{vec}(M), K^{-1} \otimes V)$, with \otimes denoting the Kronecker product.

of $\mathbf{A}^{(k)}$. Our ARD prior on $\mathbf{A}^{(k)}$ is equivalent to a $\mathcal{N}(0, \Sigma_0^{(k)})$ prior on $\text{vec}(\mathbf{A}^{(k)})$, where

$$\Sigma_0^{(k)} = \text{diag} \left(\alpha_1^{(k)}, \dots, \alpha_1^{(k)}, \dots, \alpha_m^{(k)}, \dots, \alpha_m^{(k)} \right)^{-1}. \quad (23)$$

Here, $m = n$ for the HDP-SLDS with n replicates of each $\alpha_i^{(k)}$, and $m = r$ for the HDP-AR-HMM with d^2 replicates of $\alpha_i^{(k)}$. (Recall that n is the dimension of the HDP-SLDS state vector \mathbf{x}_t , r the autoregressive order of the HDP-AR-HMM, and d the dimension of the observations \mathbf{y}_t .) To examine the posterior distribution of $\mathbf{A}^{(k)}$, we note that we may rewrite the state equation as,

$$\begin{aligned} \boldsymbol{\psi}_{t+1} &= \begin{bmatrix} \bar{\boldsymbol{\psi}}_{t,1} I_\ell & \bar{\boldsymbol{\psi}}_{t,2} I_\ell & \cdots & \bar{\boldsymbol{\psi}}_{t,\ell * r} I_\ell \end{bmatrix} \text{vec}(\mathbf{A}^{(k)}) \\ &\quad + \mathbf{e}_{t+1}^{(k)} \quad \forall t | z_t = k \\ &\triangleq \tilde{\Psi}_t \text{vec}(\mathbf{A}^{(k)}) + \mathbf{e}_{t+1}^{(k)}, \end{aligned} \quad (24)$$

where $\ell = n$ for the HDP-SLDS and $\ell = d$ for the HDP-AR-HMM. Using Eq. (24), we derive the posterior distribution as

$$\begin{aligned} p(\text{vec}(\mathbf{A}^{(k)}) | \mathbf{D}^{(k)}, \Sigma^{(k)}, \Sigma_0^{(k)}) \\ = \mathcal{N}^{-1} \left(\sum_{t|z_t=k} \tilde{\Psi}_{t-1}^T \Sigma^{-(k)} \boldsymbol{\psi}_t, \right. \\ \left. \Sigma_0^{-(k)} + \sum_{t|z_t=k} \tilde{\Psi}_{t-1}^T \Sigma^{-(k)} \tilde{\Psi}_{t-1} \right). \end{aligned} \quad (25)$$

See [42] for a detailed derivation. Here, $\mathcal{N}^{-1}(\vartheta, \Lambda)$ represents a Gaussian $\mathcal{N}(\mu, \Sigma)$ with information parameters $\vartheta = \Sigma^{-1}\mu$ and $\Lambda = \Sigma^{-1}$. Given $\mathbf{A}^{(k)}$, and recalling that each precision parameter is gamma distributed, the posterior of $\alpha_\ell^{(k)}$ is given by

$$p(\alpha_\ell^{(k)} | \mathbf{A}^{(k)}) = \text{Gamma} \left(a + \frac{|\mathcal{S}_\ell|}{2}, b + \frac{\sum_{(i,j) \in \mathcal{S}_\ell} a_{ij}^{(k)^2}}{2} \right) \quad (26)$$

The set \mathcal{S}_ℓ contains the indices for which $a_{ij}^{(k)}$ has prior precision $\alpha_\ell^{(k)}$. Note that in this model, regardless of the number of observations \mathbf{y}_t , the size of \mathcal{S}_ℓ (i.e., the number of $a_{ij}^{(k)}$ used to inform the posterior distribution) remains the same. Thus, the gamma prior is an informative prior and the choice of a and b should depend upon the cardinality of \mathcal{S}_ℓ (see Sec. IV-B for an example). For the HDP-SLDS, this cardinality is given by the maximal state dimension n , and for the HDP-AR-HMM, by the square of the observation dimensionality d^2 .

We then place an inverse-Wishart prior $\text{IW}(n_0, S_0)$ on $\Sigma^{(k)}$ and look at the posterior given $\mathbf{A}^{(k)}$:

$$p(\Sigma^{(k)} | \mathbf{D}^{(k)}, \mathbf{A}^{(k)}) = \text{IW} \left(n_k + n_0, \mathbf{S}_{\psi|\bar{\psi}}^{(k)} + S_0 \right), \quad (27)$$

where here, as opposed to in Eq. (19), we define

$$\mathbf{S}_{\psi|\bar{\psi}}^{(k)} = \sum_{t|z_t=k} (\boldsymbol{\psi}_t - \mathbf{A}^{(k)} \bar{\boldsymbol{\psi}}_{t-1})(\boldsymbol{\psi}_t - \mathbf{A}^{(k)} \bar{\boldsymbol{\psi}}_{t-1})^T. \quad (28)$$

3) *Measurement Noise Posterior*: For the HDP-SLDS, we additionally place an $\text{IW}(r_0, R_0)$ prior on the measurement noise covariance R . The posterior distribution is given by

$$p(R | \mathbf{y}_{1:T}, \mathbf{x}_{1:T}) = \text{IW}(T + r_0, S_R + R_0), \quad (29)$$

where $S_R = \sum_{t=1}^T (\mathbf{y}_t - C\mathbf{x}_t)(\mathbf{y}_t - C\mathbf{x}_t)^T$. Here, we assume that R is shared between modes. The extension to mode-specific measurement noise is straightforward.

B. Gibbs Sampler

For inference in the HDP-AR-HMM, we use a Gibbs sampler that iterates between sampling the mode sequence, $z_{1:T}$, and the set of dynamic and sticky HDP-HMM parameters. The sampler for the HDP-SLDS is identical with the additional step of sampling the state sequence, $\mathbf{x}_{1:T}$, and conditioning on this sequence when resampling dynamic parameters and the mode sequence. Periodically, we interleave a step that sequentially samples the mode sequence $z_{1:T}$ marginalizing over the state sequence $\mathbf{x}_{1:T}$ in a similar vein to that of Carter and Kohn [43]. We describe the sampler in terms of the pseudo-observations $\boldsymbol{\psi}_t$, as defined by Eq. (14), in order to clearly specify the sections of the sampler shared by both the HDP-AR-HMM and HDP-SLDS.

1) *Sampling Dynamic Parameters* $\{\mathbf{A}^{(k)}, \Sigma^{(k)}\}$: Conditioned on the mode sequence, $z_{1:T}$, and the pseudo-observations, $\boldsymbol{\psi}_{1:T}$, we can sample the dynamic parameters $\boldsymbol{\theta} = \{\mathbf{A}^{(k)}, \Sigma^{(k)}\}$ from the posterior densities of Sec. III-A. For the ARD prior, we then sample $\alpha^{(k)}$ given $\mathbf{A}^{(k)}$. In practice we iterate multiple times between sampling $\alpha^{(k)}$ given $\mathbf{A}^{(k)}$ and $\mathbf{A}^{(k)}$ given $\alpha^{(k)}$ before moving to the next sampling stage.

2) *Sampling Measurement Noise R (HDP-SLDS only)*: For the HDP-SLDS, we additionally sample the measurement noise covariance R conditioned on the sampled state sequence $\mathbf{x}_{1:T}$.

3) *Block Sampling $z_{1:T}$* : As shown in [21], the mixing rate of the Gibbs sampler for the HDP-HMM can be dramatically improved by using a *truncated* approximation to the HDP and jointly sampling the mode sequence using a variant of the forward-backward algorithm. In the case of our switching dynamical systems, we must account for the direct correlations in the observations in our likelihood computation. The variant of the forward-backward algorithm we use here then

involves computing backward messages $m_{t+1,t}(z_t) \propto p(\psi_{t+1:T}|z_t, \bar{\psi}_t, \pi, \theta)$ for each $z_t \in \{1, \dots, L\}$ with L the chosen truncation level, followed by recursively sampling each z_t conditioned on z_{t-1} from

$$p(z_t | z_{t-1}, \psi_{1:T}, \pi, \theta) \propto \pi_{z_{t-1}, z_t} p(\psi_t | \bar{\psi}_{t-1}, \mathbf{A}^{(z_t)}, \Sigma^{(z_t)}) m_{t+1,t}(z_t). \quad (30)$$

Joint sampling of the mode sequence is especially important when the observations are directly correlated via a dynamical process since this correlation further slows the mixing rate of the sequential sampler of Teh et. al. [16]. Note that using an order L weak limit approximation to the HDP still encourages the use of a sparse subset of the L possible dynamical modes.

4) *Block Sampling $\mathbf{x}_{1:T}$ (HDP-SLDS only)*: Conditioned on the mode sequence $z_{1:T}$ and the set of SLDS parameters $\theta = \{\mathbf{A}^{(k)}, \Sigma^{(k)}, R\}$, our dynamical process simplifies to a time-varying linear dynamical system. We can then block sample $\mathbf{x}_{1:T}$ by first running a backward Kalman filter to compute $m_{t+1,t}(\mathbf{x}_t) \propto p(\mathbf{y}_{t+1:T}|\mathbf{x}_t, z_{t+1:T}, \theta)$ and then recursively sampling each \mathbf{x}_t conditioned on \mathbf{x}_{t-1} from

$$p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}_{1:T}, z_{1:T}, \theta) \propto p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}^{(z_t)}, \Sigma^{(z_t)}) p(\mathbf{y}_t | \mathbf{x}_t, R) m_{t+1,t}(\mathbf{x}_t). \quad (31)$$

The messages are given in information form by $m_{t,t-1}(\mathbf{x}_{t-1}) \propto \mathcal{N}^{-1}(\vartheta_{t,t-1}, \Lambda_{t,t-1})$, where the information parameters are recursively defined as

$$\begin{aligned} \vartheta_{t,t-1} &= \mathbf{A}^{(z_t)^T} \Sigma^{-(z_t)} (\Sigma^{-(z_t)} + \Lambda_{t|t}^b)^{-1} \vartheta_{t|t}^b \\ \Lambda_{t,t-1} &= \mathbf{A}^{(z_t)^T} \Sigma^{-(z_t)} \mathbf{A}^{(z_t)} \\ &\quad - \mathbf{A}^{(z_t)^T} \Sigma^{-(z_t)} (\Sigma^{-(z_t)} + \Lambda_{t|t}^b)^{-1} \Sigma^{-(z_t)} \mathbf{A}^{(z_t)}. \end{aligned} \quad (32)$$

The standard $\vartheta_{t|t}^b$ and $\Lambda_{t|t}^b$ updated information parameters for a backward running Kalman filter are given by

$$\begin{aligned} \Lambda_{t|t}^b &= C^T R^{-1} C + \Lambda_{t+1,t} \\ \vartheta_{t|t}^b &= C^T R^{-1} y_t + \vartheta_{t+1,t}. \end{aligned} \quad (33)$$

See [42] for a derivation and for a more numerically stable version of this recursion.

5) *Sequentially Sampling $z_{1:T}$ (HDP-SLDS only)*: For the HDP-SLDS, iterating between the previous sampling stages can lead to slow mixing rates since the mode sequence is sampled conditioned on a sample of the state sequence. For high-dimensional state spaces \mathbb{R}^n , this problem is exacerbated. Instead, one can analytically marginalize the state sequence and sequentially sample the mode sequence from $p(z_t | z_{\setminus t}, \mathbf{y}_{1:T}, \pi, \theta)$. This marginalization is accomplished by once again

harnessing the fact that conditioned on the mode sequence, our model reduces to a time-varying linear dynamical system. When sampling z_t and conditioning on the mode sequence at all *other* time steps, we can run a forward Kalman filter to marginalize the state sequence $\mathbf{x}_{1:t-2}$ producing $p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}, z_{1:t-1}, \theta)$, and a backward filter to marginalize $\mathbf{x}_{t+1:T}$ producing $p(\mathbf{y}_{t+1:T} | \mathbf{x}_t, z_{t+1:T}, \theta)$. Then, for each possible value of z_t , we combine these forward and backward messages with the local likelihood $p(\mathbf{y}_t | \mathbf{x}_t)$ and local dynamic $p(\mathbf{x}_t | \mathbf{x}_{t-1}, \theta, z_t = k)$ and marginalize over \mathbf{x}_t and \mathbf{x}_{t-1} resulting in the likelihood of the observation sequence $\mathbf{y}_{1:T}$ as a function of z_t . This likelihood is combined with the prior probability of transitioning from z_{t-1} to $z_t = k$ and from $z_t = k$ to z_{t+1} . The resulting distribution is given by:

$$\begin{aligned} p(z_t = k | z_{\setminus t}, \mathbf{y}_{1:T}, \pi, \theta) &\propto \pi_{z_{t-1}, k} \pi_{k, z_{t+1}} \\ &\quad \frac{|\Lambda_t^{(k)}|^{1/2}}{|\Lambda_t^{(k)} + \Lambda_{t|t}^b|^{1/2}} \exp \left(-\frac{1}{2} \vartheta_t^{(k)^T} \Lambda_t^{-(k)} \vartheta_t^{(k)} + \right. \\ &\quad \left. \frac{1}{2} (\vartheta_t^{(k)} + \vartheta_{t|t}^b)^T (\Lambda_t^{(k)} + \Lambda_{t|t}^b)^{-1} (\vartheta_t^{(k)} + \vartheta_{t|t}^b) \right) \end{aligned} \quad (34)$$

with

$$\begin{aligned} \Lambda_t^{(k)} &= (\Sigma^{(k)} + \mathbf{A}^{(z_t)} \Lambda_{t-1|t-1}^{-f} \mathbf{A}^{(z_t)^T})^{-1} \\ \vartheta_t^{(k)} &= (\Sigma^{(k)} + \mathbf{A}^{(z_t)} \Lambda_{t-1|t-1}^{-f} \mathbf{A}^{(z_t)^T})^{-1} \\ &\quad \mathbf{A}^{(z_t)} \Lambda_{t-1|t-1}^{-f} \vartheta_{t-1|t-1}^f. \end{aligned} \quad (35)$$

See [42] for full derivations. Here, $\vartheta_{t|t}^f$ and $\Lambda_{t|t}^f$ are the updated information parameters for a forward running Kalman filter, defined recursively as

$$\begin{aligned} \Lambda_{t|t}^f &= C^T R^{-1} C + \Sigma^{-(z_t)} - \Sigma^{-(z_t)} \mathbf{A}^{(z_t)} \\ &\quad \cdot (\mathbf{A}^{(z_t)^T} \Sigma^{-(z_t)} \mathbf{A}^{(z_t)} + \Lambda_{t-1|t-1}^f)^{-1} \mathbf{A}^{(z_t)^T} \Sigma^{-(z_t)} \\ \vartheta_{t|t}^f &= C^T R^{-1} y_t + \Sigma^{-(z_t)} \mathbf{A}^{(z_t)} \\ &\quad \cdot (\mathbf{A}^{(z_t)^T} \Sigma^{-(z_t)} \mathbf{A}^{(z_t)} + \Lambda_{t-1|t-1}^f)^{-1} \vartheta_{t-1|t-1}^f. \end{aligned} \quad (36)$$

Note that a sequential node ordering for this sampling step allows for efficient updates to the recursively defined filter parameters. However, this sequential sampling is still computationally intensive, so our Gibbs sampler iterates between blocked sampling of the state and mode sequences many times before interleaving a sequential mode sequence sampling step.

The resulting Gibbs sampler is outlined in Algorithm 1.

Given a previous set of mode-specific transition probabilities $\pi^{(n-1)}$, the global transition distribution $\beta^{(n-1)}$, and dynamic parameters $\theta^{(n-1)}$:

1) Set $\pi = \pi^{(n-1)}$, $\beta = \beta^{(n-1)}$, and $\theta = \theta^{(n-1)}$.

2) If HDP-SLDS,

a) For each $t \in \{1, \dots, T\}$, compute $\{\vartheta_{t|t}^f, \Lambda_{t|t}^f\}$ as in Eq. (36).

b) For each $t \in \{T, \dots, 1\}$,

i) Compute $\{\vartheta_{t|t}^b, \Lambda_{t|t}^b\}$ as in Eq. (33).

ii) For each $k \in \{1, \dots, L\}$, compute $\{\vartheta_t^{(k)}, \Lambda_t^{(k)}\}$ as in Eq. (35) and set

$$f_k(\mathbf{y}_{1:T}) = |\Lambda_t^{(k)}|^{1/2} |\Lambda_{t|t}^{(k)} + \Lambda_{t|t}^b|^{-1/2} \exp \left(-\frac{1}{2} \vartheta_t^{(k)T} \Lambda_t^{-(k)} \vartheta_t^{(k)} + \frac{1}{2} (\vartheta_t^{(k)} + \vartheta_{t|t}^b)^T (\Lambda_t^{(k)} + \Lambda_{t|t}^b)^{-1} (\vartheta_t^{(k)} + \vartheta_{t|t}^b) \right).$$

iii) Sample a mode assignment

$$z_t \sim \sum_{k=1}^L \pi_{z_{t-1}, k} \pi_{k, z_{t+1}} f_k(\mathbf{y}_{1:T}) \delta(z_t, k).$$

c) Working sequentially forward in time sample

$$\mathbf{x}_t \sim \mathcal{N}((\Sigma^{-(z_t)} + \Lambda_{t|t}^b)^{-1} (\Sigma^{-(z_t)} A^{(z_t)} \mathbf{x}_{t-1} + \vartheta_{t|t}^b), (\Sigma^{-(z_t)} + \Lambda_{t|t}^b)^{-1}).$$

d) Set pseudo-observations $\psi_{1:T} = \mathbf{x}_{1:T}$.

3) If HDP-AR-HMM, set pseudo-observations $\psi_{1:T} = \mathbf{y}_{1:T}$.

4) Block sample $z_{1:T}$ given transition distributions π , dynamic parameters θ , and pseudo-observations $\psi_{1:T}$ as in Algorithm 2.

5) Update the global transition distribution β (utilizing auxiliary variables \mathbf{m} , \mathbf{w} , and $\bar{\mathbf{m}}$), mode-specific transition distributions π_k , and hyperparameters α , γ , and κ as in [21].

6) For each $k \in \{1, \dots, L\}$, sample dynamic parameters $\{\mathbf{A}^{(k)}, \Sigma^{(k)}\}$ given the pseudo-observations $\psi_{1:T}$ and mode sequence $z_{1:T}$ as in Algorithm 3 for the MNIW prior and Algorithm 4 for the ARD prior.

7) If HDP-SLDS, also sample the measurement noise covariance

$$R \sim \text{IW} \left(T + r_0, \sum_{t=1}^T (\mathbf{y}_t - C \mathbf{x}_t)(\mathbf{y}_t - C \mathbf{x}_t)^T + R_0 \right).$$

8) Fix $\pi^{(n)} = \pi$, $\beta^{(n)} = \beta$, and $\theta^{(n)} = \theta$.

Algorithm 1: HDP-SLDS and HDP-AR-HMM Gibbs sampler.

Given mode-specific transition probabilities π , dynamic parameters θ , and pseudo-observations $\psi_{1:T}$:

1) Calculate messages $m_{t,t-1}(k)$, initialized to $m_{T+1,T}(k) = 1$, and the sample mode sequence $z_{1:T}$:

a) For each $t \in \{T, \dots, 1\}$ and $k \in \{1, \dots, L\}$, compute

$$m_{t,t-1}(k) = \sum_{j=1}^L \pi_{kj} \mathcal{N} \left(\psi_t; \sum_{i=1}^r A_i^{(j)} \psi_{t-i}, \Sigma^{(j)} \right) m_{t+1,t}(j)$$

b) Working sequentially forward in time, starting with transitions counts $n_{jk} = 0$:

i) For each $k \in \{1, \dots, L\}$, compute the probability

$$f_k(\psi_t) = \mathcal{N} \left(\psi_t; \sum_{i=1}^r A_i^{(k)} \psi_{t-i}, \Sigma^{(k)} \right) m_{t+1,t}(k)$$

ii) Sample a mode assignment z_t as follows and increment $n_{z_{t-1} z_t}$:

$$z_t \sim \sum_{k=1}^L \pi_{z_{t-1}, k} f_k(\psi_t) \delta(z_t, k)$$

Note that the likelihoods can be precomputed for each $k \in \{1, \dots, L\}$.

Algorithm 2: Blocked mode-sequence sampler for HDP-AR-HMM or HDP-SLDS.

Given pseudo-observations $\psi_{1:T}$ and mode sequence $z_{1:T}$, for each $k \in \{1, \dots, K\}$:

- 1) Construct $\Psi^{(k)}$ and $\bar{\Psi}^{(k)}$ as in Eq. (15).
- 2) Compute sufficient statistics using pseudo-observations ψ_t associated with $z_t = k$:

$$\mathbf{S}_{\bar{\psi}\bar{\psi}}^{(k)} = \bar{\Psi}^{(k)} \bar{\Psi}^{(k)T} + K \quad \mathbf{S}_{\psi\bar{\psi}}^{(k)} = \Psi^{(k)} \bar{\Psi}^{(k)T} + MK \quad \mathbf{S}_{\psi\psi}^{(k)} = \Psi^{(k)} \Psi^{(k)T} + MKM^T.$$

- 3) Sample dynamic parameters:

$$\Sigma^{(k)} \sim \text{IW} \left(n_k + n_0, \mathbf{S}_{\bar{\psi}\bar{\psi}}^{(k)} + S_0 \right) \quad \mathbf{A}^{(k)} | \Sigma^{(k)} \sim \mathcal{MN} \left(\mathbf{S}_{\psi\bar{\psi}}^{(k)} \mathbf{S}_{\bar{\psi}\bar{\psi}}^{-(k)}, \Sigma^{(k)}, \mathbf{S}_{\bar{\psi}\bar{\psi}}^{(k)} \right).$$

Algorithm 3: Parameter sampling using MNIW prior.

Given pseudo-observations $\psi_{1:T}$, mode sequence $z_{1:T}$, and a previous set of dynamic parameters $(\mathbf{A}^{(k)}, \Sigma^{(k)}, \alpha^{(k)})$, for each $k \in \{1, \dots, K\}$:

- 1) Construct $\tilde{\Psi}_t$ as in Eq. (24).
- 2) Iterate multiple times between the following steps:
 - a) Construct $\Sigma_0^{(k)}$ given $\alpha^{(k)}$ as in Eq. (23) and sample the dynamic matrix:

$$\text{vec}(\mathbf{A}^{(k)}) | \Sigma^{(k)}, \alpha^{(k)} \sim \mathcal{N}^{-1} \left(\sum_{t|z_t=k} \tilde{\Psi}_{t-1}^T \Sigma^{-(k)} \psi_t, \Sigma_0^{-(k)} + \sum_{t|z_t=k} \tilde{\Psi}_{t-1}^T \Sigma^{-(k)} \tilde{\Psi}_{t-1} \right).$$

- b) For each $\ell \in \{1, \dots, m\}$, with $m = n$ for the SLDS and $m = r$ for the switching VAR, sample ARD precision parameters:

$$\alpha_\ell^{(k)} | \mathbf{A}^{(k)} \sim \text{Gamma} \left(a + \frac{|\mathcal{S}_\ell|}{2}, b + \frac{\sum_{(i,j) \in \mathcal{S}_\ell} a_{ij}^{(k)^2}}{2} \right).$$

- c) Compute sufficient statistic:

$$\mathbf{S}_{\psi\bar{\psi}}^{(k)} = \sum_{t|z_t=k} (\psi_t - \mathbf{A}^{(k)} \bar{\psi}_{t-1})(\psi_t - \mathbf{A}^{(k)} \bar{\psi}_{t-1})^T$$

and sample process noise covariance:

$$\Sigma^{(k)} | \mathbf{A}^{(k)} \sim \text{IW} \left(n_k + n_0, \mathbf{S}_{\psi\bar{\psi}}^{(k)} + S_0 \right).$$

Algorithm 4: Parameter sampling using ARD prior.

IV. RESULTS

A. MNIW prior

We begin by examining a set of three synthetic datasets displayed in Fig. 2(a) in order to analyze the relative modeling power of the HDP-VAR(1)-HMM, HDP-VAR(2)-HMM, and HDP-SLDS using the MNIW prior. Here, we use the notation HDP-VAR(r)-HMM to explicitly denote an order r HDP-AR-HMM with vector observations. We compare to a baseline sticky HDP-HMM using first difference observations, imitating a HDP-VAR(1)-HMM with $A^{(k)} = I$ for all k . In Fig. 2(b)-(e) we display Hamming distance errors that are calculated by choosing the optimal mapping of indices maximizing overlap between the true and estimated mode sequences.

We place a $\text{Gamma}(a, b)$ prior on the sticky HDP-

HMM concentration parameters $\alpha + \kappa$ and γ , and a $\text{Beta}(c, d)$ prior on the self-transition proportion parameter $\rho = \kappa / (\alpha + \kappa)$. We choose the weakly informative setting of $a = 1$, $b = 0.01$, $c = 10$, and $d = 1$. The details on setting the MNIW hyperparameters from statistics of the data are discussed in the Appendix.

For the first scenario (Fig. 2 (top)), the data were generated from a five-mode switching VAR(1) process with a 0.98 probability of self-transition and equally likely transitions to the other modes. The same mode-transition structure was used in the subsequent two scenarios, as well. The three switching linear dynamical models provide comparable performance since both the HDP-VAR(2)-HMM and HDP-SLDS with $C = I_3$ contain the class of HDP-VAR(1)-HMMs. In the second scenario (Fig. 2 (middle)), the data were generated from a 3-mode switching AR(2) process. The HDP-AR(2)-HMM has

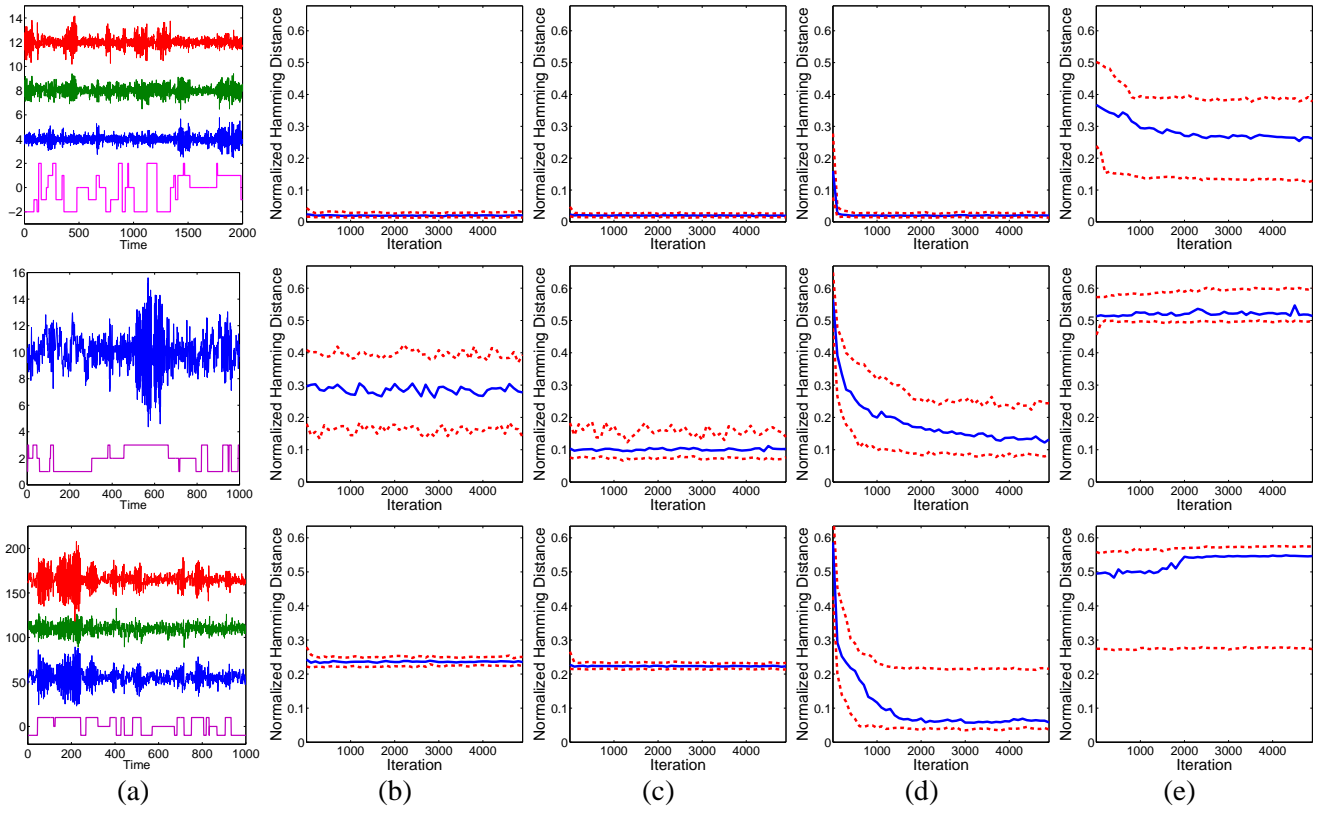


Fig. 2. (a) Observation sequence (blue, green, red) and associated mode sequence (magenta) for a 5-mode switching VAR(1) process (top), 3-mode switching AR(2) process (middle), and 3-mode SLDS (bottom). The components of the observation vector are offset for clarity. The associated 10th, 50th, and 90th Hamming distance quantiles over 100 trials are shown for the (b) HDP-AR(1)-HMM, (c) HDP-AR(2)-HMM, (d) HDP-SLDS with $C = I$ (top and bottom) and $C = [1 \ 0]$ (middle), and (e) sticky HDP-HMM using first difference observations.

significantly better performance than the HDP-AR(1)-HMM while the performance of the HDP-SLDS with $C = [1 \ 0]$ performs similarly, but has greater posterior variability because the HDP-AR(2)-HMM model family is smaller. Note that the HDP-SLDS sampler is slower to mix since the hidden, continuous state is also sampled. The data in the third scenario (Fig. 2 (bottom)) were generated from a three-mode SLDS model with $C = I_3$. Here, we clearly see that neither the HDP-AR(1)-HMM nor HDP-AR(2)-HMM is equivalent to the HDP-SLDS. Note that all of the switching models yielded significant improvements relative to the baseline sticky HDP-HMM. This input representation is more effective than using raw observations for HDP-HMM learning, but still much less effective than richer models which switch among learned LDS. Together, these results demonstrate both the differences between our models as well as the models' ability to learn switching processes with varying numbers of modes.

B. ARD prior

We now compare the utility of the ARD prior to the MNIW prior using the HDP-SLDS model when the true underlying dynamical modes have sparse dependencies

relative to the assumed model order. That is, the HDP-SLDS may have dynamical regimes reliant on lower state dimensions, or the HDP-AR-HMM may have modes described by lower order VAR processes. We generated data from a two-mode SLDS with 0.98 probability of self-transition and

$$\mathbf{A}^{(1)} = \begin{bmatrix} 0.8 & -0.2 & 0 \\ -0.2 & 0.8 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad \mathbf{A}^{(2)} = \begin{bmatrix} -0.2 & 0 & 0.8 \\ 0.8 & 0 & -0.2 \\ 0 & 0 & 0 \end{bmatrix},$$

with $C = [I_2 \ 0]$, $\Sigma^{(1)} = \Sigma^{(2)} = I_3$, and $R = I_2$. The first dynamical process can be equivalently described by just the first and second state components since the third component is simply white noise that does not contribute to the state dynamics and is not directly (or indirectly) observed. For the second dynamical process, the third state component is once again a white noise process, but *does* contribute to the dynamics of the first and second state components. However, we can equivalently represent the dynamics of this mode as

$$\begin{aligned} x_{1,t} &= -0.2x_{1,t-1} + \tilde{e}_{1,t} \\ x_{2,t} &= 0.8x_{1,t-1} + \tilde{e}_{2,t} \end{aligned} \quad \tilde{\mathbf{A}}^{(2)} = \begin{bmatrix} -0.2 & 0 & 0 \\ 0.8 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

where \tilde{e}_t is a white noise term defined by the original process noise combined with $x_{3,t}$, and $\tilde{\mathbf{A}}^{(2)}$ is

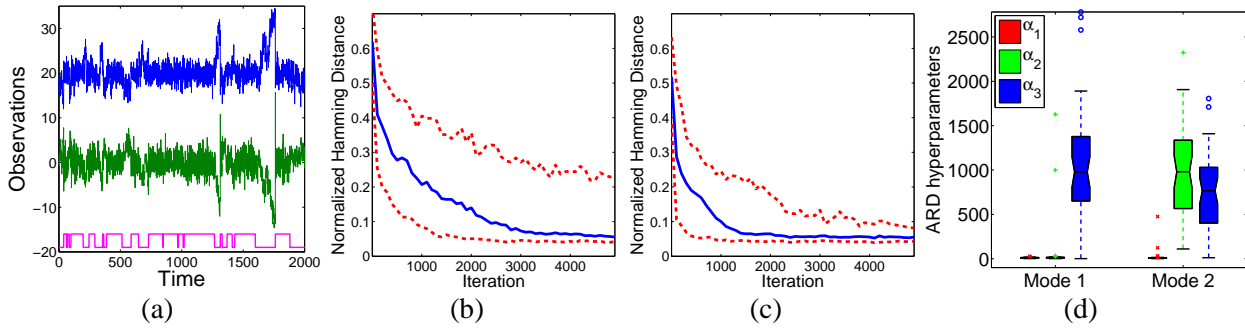


Fig. 3. (a) Observation sequence (green, blue) and mode sequence (magenta) of a 2-mode SLDS, where the first mode can be realized by the first two state components and the second mode solely by the first. The associated 10th, 50th, and 90th Hamming distance quantiles over 100 trials are shown for the (b) MNIW and (c) ARD prior. (d) Box plots of inferred ARD precisions associated with the first and second dynamical modes at the 5000th Gibbs iteration. The center line indicates the median, edges the 25th and 75th quantiles, and whiskers the range of data excluding outliers which are separately marked. Larger ARD precision values correspond to non-dynamical components.

the dynamical matrix associated with this equivalent representation of the second dynamical mode. Notice that this SLDS does not satisfy Criterion 3.1 since the second column of $\mathbf{A}^{(2)}$ is zero while the second column of \mathbf{C} is not. Nevertheless, because the realization is in our canonical form with $\mathbf{C} = [\mathbf{I}_2 \ 0]$, we still expect to recover the $\mathbf{a}_2^{(2)} = \mathbf{a}_3^{(2)} = 0$ sparsity structure. We set the parameters of the Gamma(a, b) prior on the ARD precisions as $a = |\mathcal{S}_\ell|$ and $b = a/1000$, where we recall the definition of \mathcal{S}_ℓ from Eq. (26). This specification fixes the mean of the prior to 1000 while aiming to provide a prior that is roughly equally informative for various choices of model order (i.e., sizes $|\mathcal{S}_\ell|$).

In Fig. 3, we see that even in this low-dimensional example, the ARD provides superior mode-sequence estimates, as well as a mechanism for identifying non-dynamical state components. The box plots of the inferred $\alpha^{(k)}$ are shown in Fig. 3(d). From the clear separation between the sampled dynamic range of $\alpha_3^{(1)}$ and $(\alpha_1^{(1)}, \alpha_2^{(1)})$, and between that of $(\alpha_2^{(2)}, \alpha_3^{(2)})$ and $\alpha_1^{(2)}$, we see that we are able to correctly identify dynamical systems with $\mathbf{a}_3^{(1)} = 0$ and $\mathbf{a}_2^{(2)} = \mathbf{a}_3^{(2)} = 0$.

C. Dancing Honey Bees

Honey bees perform a set of dances within the beehive in order to communicate the location of food sources. Specifically, they switch between a set of *waggle*, *turn-right*, and *turn-left* dances. During the waggle dance, the bee walks roughly in a straight line while rapidly shaking its body from left to right. The turning dances simply involve the bee turning in a clockwise or counterclockwise direction. We display six such sequences of honey bee dances in Fig. 4. The data consist of measurements $\mathbf{y}_t = [\cos(\theta_t) \ \sin(\theta_t) \ x_t \ y_t]^T$, where (x_t, y_t) denotes the 2D coordinates of the bee's body and θ_t its head angle. Both Oh et. al. [10] and Xuan and Murphy [15] used switching dynamical models to analyze these honey bee dances. We wish to analyze the

performance of our Bayesian nonparametric variants of these models in segmenting the six sequences into the dance labels displayed in Fig. 4.

MNIW Prior — Unsupervised: We start by testing the HDP-VAR(1)-HMM using a MNIW prior. (Note that we did not see performance gains by considering the HDP-SLDS, so we omit showing results for that architecture.) We set the prior distributions on the dynamic parameters and hyperparameters as in Sec. IV-A for the synthetic data examples, with the MNIW prior based on a pre-processed observation sequence. The pre-processing involves centering the position observations around 0 and scaling each component of \mathbf{y}_t to be within the same dynamic range. We compare our results to those of Xuan and Murphy [15], who used a change-point detection technique for inference on this dataset. As shown in Fig. 5(a)-(b), our model achieves a superior segmentation compared to the change-point formulation in almost all cases, while also identifying modes which reoccur over time. Example segmentations are shown in Fig. 6. Oh et. al. [10] also presented an analysis of the honey bee data, using an SLDS with a fixed number of modes. Unfortunately, that analysis is not directly comparable to ours, because Oh et. al. [10] used their SLDS in a supervised formulation in which the ground truth labels for all but one of the sequences are employed in the inference of the labels for the remaining held-out sequence, and in which the kernels used in the MCMC procedure depend on the ground truth labels. (The authors also considered a “parameterized segmental SLDS (PS-SLDS),” which makes use of domain knowledge specific to honey bee dancing and requires additional supervision during the learning process.) Nonetheless, in Table III we report the performance of these methods as well as the median performance (over 100 trials) of the unsupervised HDP-VAR(1)-HMM in order to provide a sense of the level of performance achievable without detailed, manual supervision. As seen in Table III, the

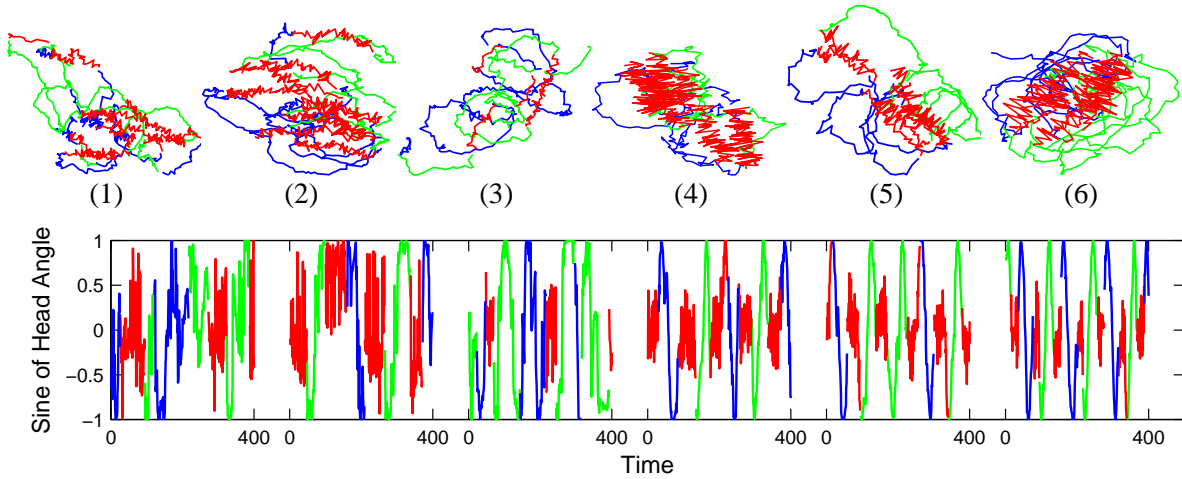


Fig. 4. *Top*: Trajectories of the dancing honey bees for sequences 1 to 6, colored by *waggle* (red), *turn right* (blue), and *turn left* (green) dances. *Bottom*: Sine of the bee's head angle measurements colored by ground truth labels for 400 frames of each sequence. The data are available at http://www.cc.gatech.edu/~borg/ijcv_psslds/.

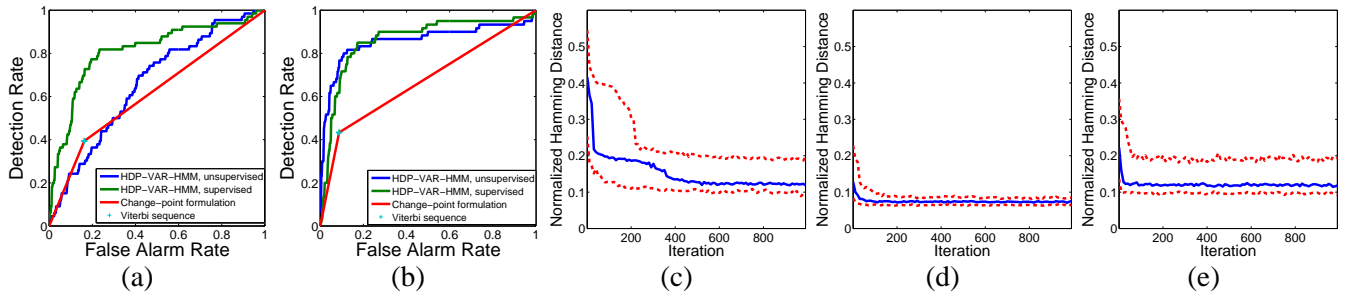


Fig. 5. (a)-(b) ROC curves for the unsupervised HDP-VAR-HMM, partially supervised HDP-VAR-HMM, and change-point formulation of [15] using the Viterbi sequence for segmenting datasets 1-3 and 4-6, respectively. (c)-(e) The 10th, 50th, and 90th Hamming distance quantiles over 100 trials are shown for sequences 4, 5, and 6, respectively.

HDP-VAR(1)-HMM yields very good performance on sequences 4 to 6 in terms of the learned segmentation and number of modes (see Fig. 6); the performance approaches that of the supervised method.

For sequences 1 to 3—which are much less regular than sequences 4 to 6—the performance of the unsupervised procedure is substantially worse. In Fig. 4, we see the extreme variation in head angle during the waggle dances of sequences 1 to 3.³ As noted by Oh, the tracking results based on the vision-based tracker are noisier for these sequences and the patterns of switching between dance modes is more irregular. This dramatically affects our performance since we do not use domain-specific information. For sequence 2 in particular, our learned segmentations often create new, sequence-specific waggle dance modes contributing to our calculated Hamming distance errors on this sequence. Overall, however, we are able to achieve reasonably good segmentations without having to manually input domain-specific knowledge.

³From Fig. 4, we also see that even in sequences 4 to 6, the ground truth labeling appear to be inaccurate at times. Specifically, certain time steps are labeled as waggle dances (red) that look more typical of a turning dance (green, blue).

MNIW Prior — Partially Supervised: The discrepancy in performance between our results and the supervised approach of Oh et. al. [10] motivated us to also consider a partially supervised variant of the HDP-VAR(1)-HMM in which we fix the ground truth mode sequences for five out of six of the sequences, and jointly infer both a combined set of dynamic parameters and the left-out mode sequence. This is equivalent to informing the prior distributions with the data from the five fixed sequences, and using these updated posterior distributions as the prior distributions for the held-out sequence. As we see in Table III and the segmentations of Fig. 6, this partially supervised approach considerably improves performance for these three sequences, especially sequences 2 and 3. In this analysis, we hand-aligned sequences so that the waggle dances tended to have head angle measurements centered about $\pi/2$ radians. Aligning the waggle dances is possible by looking at the high frequency portions of the head angle measurements. Additionally, the pre-processing of the unsupervised approach is not appropriate here as the scalings and shiftings are dance-specific, and such transformations modify the associated switching VAR(1) model. Instead, to account for the varying frames of reference (i.e., point of origin for each

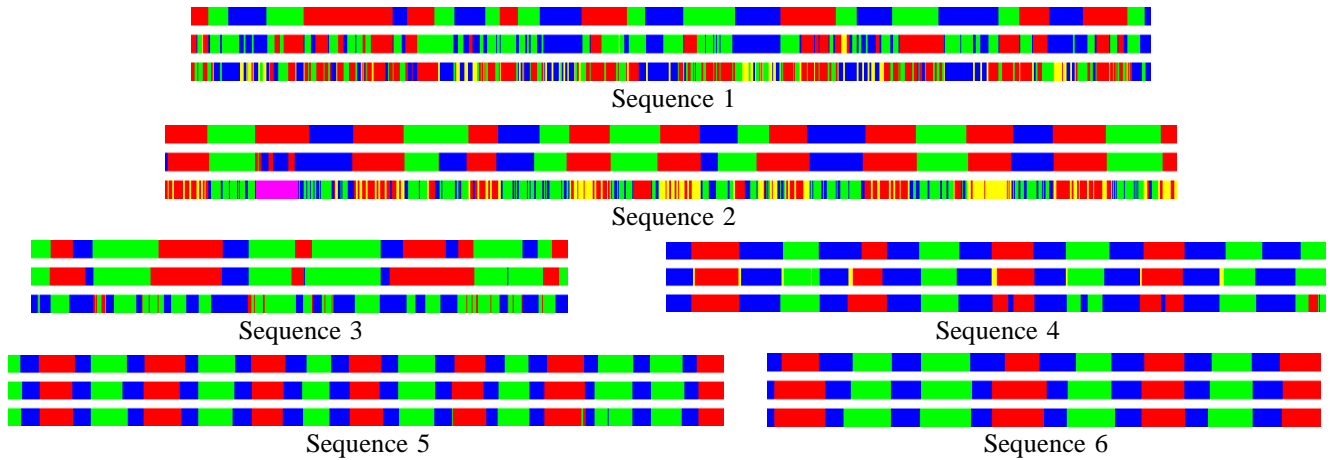


Fig. 6. Estimated mode sequences at the 1000th Gibbs iteration representing the median error over 100 trials. For each sequence we plot the true labels (top), labels from the partially supervised HDP-VAR-HMM (middle), and unsupervised HDP-VAR-HMM (bottom). Colors are as in Fig. 4.

| Sequence | 1 | 2 | 3 | 4 | 5 | 6 |
|-------------------------------------|------|------|------|------|------|------|
| HDP-VAR(1)-HMM unsupervised | 45.0 | 42.7 | 47.3 | 88.1 | 92.5 | 88.2 |
| HDP-VAR(1)-HMM partially supervised | 55.0 | 86.3 | 81.7 | 89.0 | 92.4 | 89.6 |
| SLDS DD-MCMC | 74.0 | 86.1 | 81.3 | 93.4 | 90.4 | 90.2 |
| PS-SLDS approx. Viterbi | 75.9 | 92.4 | 83.1 | 93.4 | 91.0 | 90.4 |

TABLE III

Median LABEL ACCURACY OF THE HDP-VAR(1)-HMM USING UNSUPERVISED AND PARTIALLY SUPERVISED GIBBS SAMPLING, COMPARED TO ACCURACY OF THE SUPERVISED SLDS DATA-DRIVEN MCMC (DD-MCMC) MAP SEGMENTATION AND PS-SLDS APPROXIMATE *Viterbi* SEGMENTATION PROCEDURES OF OH ET. AL. [10].

bee body) we allowed for a mean $\mu^{(k)}$ on the process noise, and placed an independent $\mathcal{N}(0, \Sigma_0)$ prior on this parameter. See the Appendix for details on how the hyperparameters of these prior distributions are set.

ARD Prior: Using the cleaner sequences 4 to 6, we investigate the affects of the sparsity-inducing ARD prior by assuming a higher order switching VAR model and computing the likelihood of the second half of each dance sequence based on parameters inferred from Gibbs sampling using the data from the first half of each sequence. In Fig. 7, we specifically compare the performance of an HDP-VAR(r)-HMM with a conjugate MNIW prior for $r = 1, 2, 7$ to that of an HDP-VAR(7)-HMM with an ARD prior. We use the same approach to setting the hyperparameters as in Sec. IV-B. We see that assuming a higher order model improves the predictive likelihood performance, but only when combined with a regularizing prior (e.g., the ARD) that avoids overfitting in the presence of limited data. Although not depicted here (see instead [42]), the ARD prior also informs us of the variable-order nature of this switching dynamical process. When considering an HDP-VAR(2)-HMM with an ARD prior, the posterior distribution of the ARD hyperparameters for the first and second order lag components associated with each of the three dominant inferred dances clearly indicates that two of the dances (turning dances) simply rely on the first lag

component while the other dance (waggle dance) relies on both lag components. To verify these results, we provided the data and ground truth labels to MATLAB's `lpc` implementation of Levinson's algorithm, which indicated that the turning dances are well approximated by an order 1 process, while the waggle dance relies on an order 2 model. Thus, our learned orders for the three dances match what is indicated by Levinson's algorithm on ground-truth segmented data.

V. MODEL VARIANTS

There are many variants of the general SLDS and switching VAR models that are pervasive in the literature. One important example is when the dynamic matrix is shared between modes; here, the dynamics are instead distinguished based on a switching mean, such as the Markov switching stochastic volatility (MSSV) model. In the maneuvering target tracking community, it is often further assumed that the dynamic matrix is shared and *known* (due to the understood physics of the target). We explore both of these variants in the following sections.

A. Shared Dynamic Matrix, Switching Driving Noise

In many applications, the dynamics of the switching process can be described by a shared linear dynamical system matrix A ; the dynamics within a given mode are

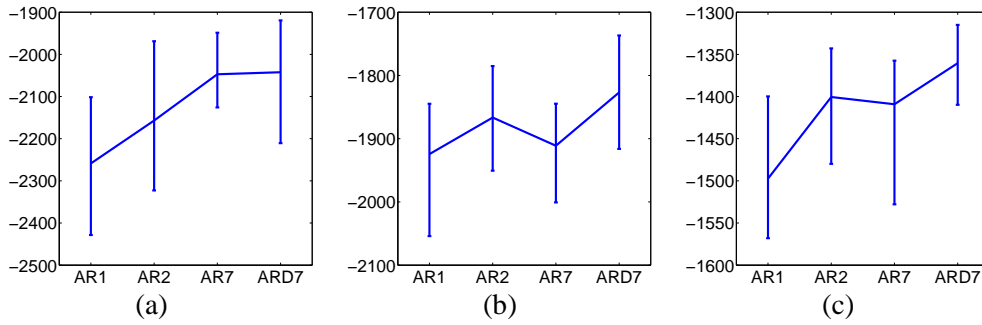


Fig. 7. For an order 1, 2, and 7 HDP-AR-HMM with a MNIW prior and an order 7 HDP-AR-HMM with an ARD prior, we plot the shortest intervals containing 95% of the held-out log-likelihoods calculated based on a set of Gibbs samples taken at iteration 1000 from 100 chains. (a) Log-likelihood of the second half of honey bee dance sequence 4 based on model parameters inferred from the first half of the sequence. (b)-(c) Similarly for sequences 5 and 6, respectively.

then determined by some external force acting upon this LDS, and it is how this force is exerted that is mode-specific. The general form for such an SLDS is given by

$$z_t | z_{t-1} \sim \pi_{z_{t-1}} \quad (37)$$

$$x_t = Ax_{t-1} + e_t^{(z_t)} \quad y_t = Cx_t + w_t,$$

with process and measurement noise $e_t^{(k)} \sim \mathcal{N}(\mu^{(k)}, \Sigma^{(k)})$ and $w_t \sim \mathcal{N}(0, R)$, respectively. In this scenario, the data are generated from one dynamic matrix, A , and multiple process noise covariance matrices, $\Sigma^{(k)}$. Thus, one cannot place a MNIW prior jointly on these parameters (conditioned on $\mu^{(k)}$) due to the coupling of the parameters in this prior. We instead consider independent priors on A , $\Sigma^{(k)}$, and $\mu^{(k)}$. We will refer to the choice of a normal prior on A , inverse-Wishart prior on $\Sigma^{(k)}$, and normal prior on $\mu^{(k)}$ as the *N-IW-N* prior. See [42] for details on deriving the resulting posterior distributions given these independent priors.

Stochastic Volatility: An example of an SLDS in a similar form to that of Eq. (37) is the Markov switching stochastic volatility (MSSV) model [5], [6], [44]. The MSSV assumes that the log-volatilities follow an AR(1) process with a Markov switching mean. This underlying process is observed via conditionally independent and normally distributed daily returns. Specifically, let y_t represent, for example, the daily returns of a stock index. The state x_t is then given the interpretation of log-volatilities and the resulting state space model is given by [7]

$$z_t | z_{t-1} \sim \pi_{z_{t-1}} \quad (38)$$

$$x_t = ax_{t-1} + e_t^{(z_t)} \quad y_t = u_t(x_t),$$

with $e_t^{(k)} \sim \mathcal{N}(\mu^{(k)}, \sigma^2)$ and $u_t(x_t) \sim \mathcal{N}(0, \exp(x_t))$. Here, only the mean of the process noise is mode-specific. Note, however, that the measurement equation is non-linear in the state x_t . Carvalho and Lopes [7] employ a particle filtering approach to cope with these

non-linearities. In [6], the MSSV is instead modeled in the log-squared-daily-returns domain such that

$$\log(y_t^2) = x_t + w_t, \quad (39)$$

where w_t is additive, non-Gaussian noise. This noise is sometimes approximated by a moment-matched Gaussian [45], while So et. al. [6] use a mixture of Gaussians approximation. The MSSV is then typically bestowed a fixed set of two or three regimes of volatility.

We examine the IBOVESPA stock index (Sao Paulo Stock Exchange) over the period of 01/03/1997 to 01/16/2001, during which ten key world events are cited in [7] as affecting the emerging Brazilian market. The key world events are summarized in Table IV and shown in the plots of Fig. 8. Use of this dataset was motivated by the work of Carvalho and Lopes [7], in which a two-mode MSSV model is assumed. We consider a variant of the HDP-SLDS to match the MSSV model of Eq. (38). Specifically we examine log-squared daily returns, as in Eq. (39), and use a DP mixture of Gaussians to model the measurement noise:

$$e_t^{(k)} \sim \mathcal{N}(\mu^{(k)}, \Sigma^{(k)})$$

$$w_t \sim \sum_{\ell=1}^{\infty} \omega_{\ell} \mathcal{N}(0, R_{\ell}) \quad \omega \sim \text{GEM}(\sigma_r), \quad (40)$$

$$R_{\ell} \sim \text{IW}(n_r, S_r).$$

We truncate the measurement noise DP mixture to 10 components. For the HDP concentration hyperparameters, α , γ , and κ , we use the same prior distributions as in Sec. IV-A-IV-C. For the dynamic parameters, we rely on the N-IW-N prior described in Sec. V-A and once again set the hyperparameters of this prior from statistics of the data as described in the Appendix. Since we allow for a mean on the process noise and examine log-squared daily returns, we do not preprocess the data.

The posterior probability of an HDP-SLDS inferred change point is shown in Fig. 8(a), and in Fig. 8(b) we display the corresponding plot for a non-sticky variant (i.e., with $\kappa = 0$ so that there is no bias towards mode self-transitions.) The HDP-SLDS is able to infer very

| Date | Event |
|------------|--|
| 07/02/1997 | Thailand devalues the Baht by as much as 20% |
| 08/11/1997 | IMF and Thailand set a rescue agreement |
| 10/23/1997 | Hong Kongs stock index falls 10.4%. South Korea won starts to weaken |
| 12/02/1997 | IMF and South Korea set a bailout agreement |
| 06/01/1998 | Russias stock market crashes |
| 06/20/1998 | IMF gives final approval to a loan package to Russia |
| 08/19/1998 | Russia officially falls into default |
| 10/09/1998 | IMF and World Bank joint meeting to discuss global economic crisis. The Fed cuts interest rates |
| 01/15/1999 | The Brazilian government allows its currency, the Real, to float freely by lifting exchange controls |
| 02/02/1999 | Arminio Fraga is named President of Brazils Central Bank |

TABLE IV

TABLE OF 10 KEY WORLD EVENTS AFFECTING THE IBOVESPA STOCK INDEX (SAO PAULO STOCK EXCHANGE) OVER THE PERIOD OF 01/03/1997 TO 01/16/2001, AS CITED BY CARVALHO AND LOPES [7].

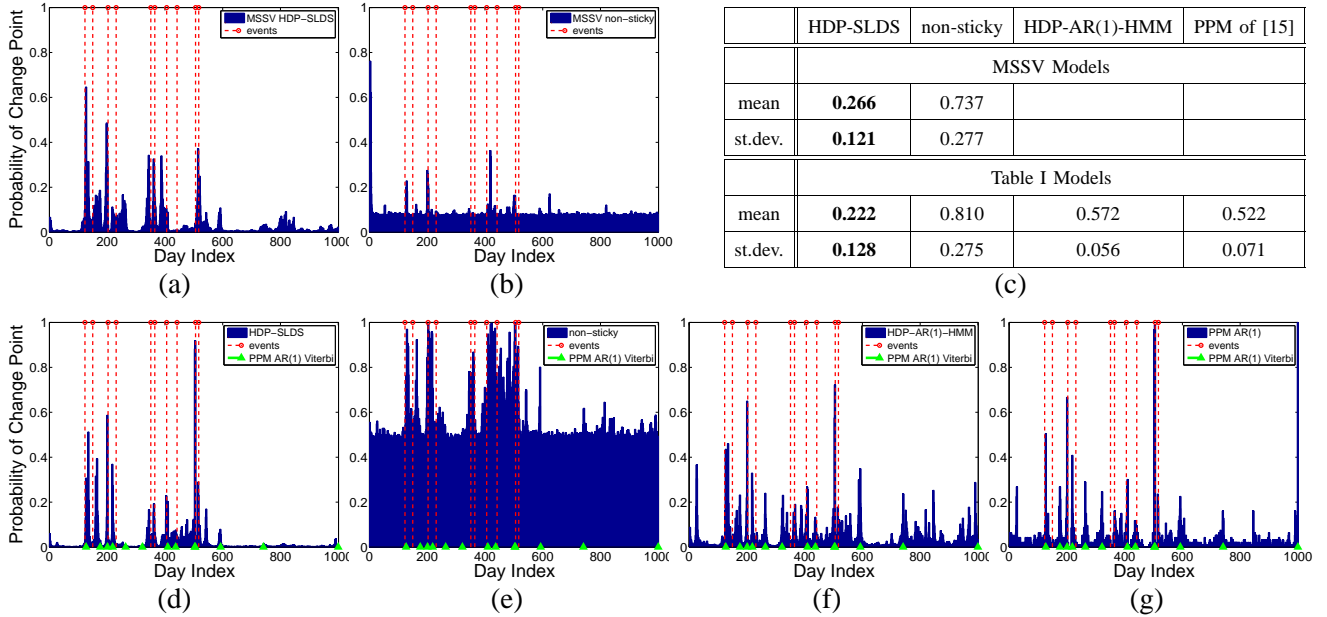


Fig. 8. (a) Plot of the estimated probability of a change point on each day using 3,000 Gibbs samples for a MSSV variant of the HDP-SLDS using a shared dynamic matrix and allowing a mean on the mode-specific process noise and a mixture of Gaussian measurement noise model. The observations are log-squared daily return measurements, and the 10 key events are indicated with red lines. (b) Similar plot for the *non-sticky* HDP-SLDS with no bias towards self-transitions. (d)-(g) Analogous plots for the HDP-SLDS of Table I, a non-sticky variant, an HDP-AR(1)-HMM, and the switching AR(1) product partition model (PPM) of Xuan and Murphy [15], all using raw daily return measurements. The Viterbi change points provided by the formulation of [15] are shown with green triangles. (c) For each of the compared models, mean and standard deviation of the normalized Hamming distance between a label sequence associated with the true event dates and that formed for each of the Gibbs sampled change points.

similar change points to those presented in [7]. Without the sticky extension, the non-sticky model variant over-segments the data and rapidly switches between redundant states leading to many inferred change points that do not align with any world event. As a quantitative comparison of the inferred change points, we compute a Hamming distance metric as follows. For the cited event dates, we form a “true” label sequence with labels that increment at each event. Then, for each inferred set of change points we form a separate label sequence in an analogous manner (i.e., with incrementing label numbers at each inferred change point.) We then compute the Hamming distance between the true and estimated label sequences after an optimal mapping between these labels. The resulting performances are summarized in the

table of Fig. 8(c).

We also analyzed the performance of an HDP-SLDS as defined in Table I. We used raw daily-return observations, and first pre-processed the data in the same manner as the honey bee data by centering the observations around 0 and scaling the data to be roughly within a $[-10, 10]$ dynamic range. We then took a MNIW prior on the dynamic parameters, as outlined in the Appendix. Overall, although the state of this HDP-SLDS does not have the interpretation of log-volatilities, we see are still able to capture regime-changes in the dynamics of this stock index and find change points that align better with the true world events than in the MSSV HDP-SLDS model. See Fig. 8(e)-(h), which also provides a comparison with the change points inferred by an HDP-

AR(1)-HMM⁴ and a switching AR(1) product partition model (PPM) of Xuan and Murphy [15]. The PPM inferred change points align well with those of the HDP-AR(1)-HMM—we expect this similar performance in such low-dimensional, long time series where the penalty incurred (in terms of quality of parameter estimates) by not revisiting modes is minimal.

B. Fixed Dynamic Matrix, Switching Driving Noise

There are some cases in which the dynamical model is well-defined through knowledge of the physics of the system being observed, such as simple kinematic motion. More complicated motions can typically be modeled using the same fixed dynamical model, but with a more complex description of the driving force. A generic LDS driven by an unknown control input \mathbf{u}_t can be represented as

$$\mathbf{x}_t = A\mathbf{x}_{t-1} + B\mathbf{u}_t + \mathbf{v}_t \quad \mathbf{y}_t = C\mathbf{x}_t + D\mathbf{u}_t + \mathbf{w}_t, \quad (41)$$

where $\mathbf{v}_t \sim \mathcal{N}(0, Q)$ and $\mathbf{w}_t \sim \mathcal{N}(0, R)$. It is often appropriate to assume $D = 0$, as we do herein.

Maneuvering Target Tracking: Target tracking provides an application domain in which one often assumes that the dynamical model is known. One method of describing a maneuvering target is to consider the control input as a random process [46]. For example, a *jump-mean* Markov process [47] yields dynamics described as

$$\begin{aligned} z_t \mid z_{t-1} &\sim \pi_{z_{t-1}} \\ \mathbf{x}_t &= A\mathbf{x}_{t-1} + B\mathbf{u}_t^{(z_t)} + \mathbf{v}_t \quad \mathbf{y}_t = C\mathbf{x}_t + \mathbf{w}_t \\ \mathbf{u}_t^{(k)} &\sim \mathcal{N}(\boldsymbol{\mu}^{(k)}, \Sigma^{(k)}) \quad \mathbf{v}_t \sim \mathcal{N}(0, Q) \quad \mathbf{w}_t \sim \mathcal{N}(0, R). \end{aligned} \quad (42)$$

Classical approaches rely on defining a fixed set of dynamical modes and associated transition distributions. The state dynamics of Eq. (42) can be equivalently described as

$$\mathbf{x}_t = A\mathbf{x}_{t-1} + \mathbf{e}_t^{(z_t)} \quad (43)$$

$$\mathbf{e}_t^{(k)} \sim \mathcal{N}(B\boldsymbol{\mu}^{(k)}, B\Sigma^{(k)}B^T + Q). \quad (44)$$

This model can be captured by our HDP-SLDS formulation of Eq. (37) with a fixed dynamic matrix (e.g., constant velocity or constant acceleration models [46]) and mode-specific, non-zero mean process noise. Such a formulation was explored in [9] along with experiments that compare the performance to that of standard multiple model techniques, demonstrating the flexibility of the Bayesian nonparametric approach. Fox et.

al. [9] also present an alternative sampling scheme that harnesses the fact that the control input may be much lower-dimensional than the state and sequentially block-samples (z_t, \mathbf{u}_t) analytically marginalizing over the state sequence $\mathbf{x}_{1:T}$. Note that this variant of the HDP-SLDS can be viewed as an extension of the work by Caron et. al. [20] in which the exogenous input is modeled as an independent noise process (i.e., no Markov structure on z_t) generated from a DP mixture model.

VI. CONCLUSION

In this paper, we have addressed the problem of learning switching linear dynamical models with an unknown number of modes for describing complex dynamical phenomena. We presented a Bayesian nonparametric approach and demonstrated both the utility and versatility of the developed HDP-SLDS and HDP-AR-HMM on real applications. Using the same parameter settings, although different model choices, in one case we are able to learn changes in the volatility of the IBOVESPA stock exchange while in another case we learn segmentations of data into *waggle*, *turn-right*, and *turn-left* honey bee dances. We also described a method of applying automatic relevance determination (ARD) as a sparsity-inducing prior, leading to flexible and scalable dynamical models that allow for identification of variable order structure. We concluded by considering adaptations of the HDP-SLDS to specific forms often examined in the literature such as the Markov switching stochastic volatility model and a standard multiple model target tracking formulation.

The batch processing of the Gibbs samplers derived herein may be impractical and offline-training online-tracking infeasible for certain applications. Due both to the nonlinear dynamics and uncertainty in model parameters, exact recursive estimation is infeasible. One could leverage the *conditionally linear* dynamics and use *Rao-Blackwellized particle filtering* (RBPF) [48]. However, one challenge is that such particle filters can suffer from a progressively impoverished particle representation. A possible direction of future research is to consider building on the recent work of [49] and embedding a RBPF within an MCMC algorithm. Another interesting avenue of research is to analyze high-dimensional time series. Although there is nothing fundamentally different in considering such datasets, based on experiments in related models [21] we expect to run into mixing rate issues with the Gibbs sampler since the parameter associated with each new considered dynamical mode is a sample from the (high-dimensional) prior. Developing split-merge algorithms similar to those

⁴We do not compare to an HDP-AR(1)-HMM for the MSSV formulation since there is no adequate way to capture the complex MSSV observation model with an autoregressive process.

developed in [50] for the DP mixture model could be useful in ameliorating these issues.

Overall, the formulation we developed herein represents a flexible, Bayesian nonparametric model for describing complex dynamical phenomena and discovering simple underlying temporal structures.

APPENDIX

a) MNIW General Method: For the experiments of Sec. IV-A, we set $M = \mathbf{0}$ and $K = I_m$. This choice centers the mass of the prior around stable dynamic matrices while allowing for considerable variability. The inverse-Wishart portion is given $n_0 = m + 2$ degrees of freedom. For the HDP-AR-HMM, the scale matrix $S_0 = 0.75\bar{\Sigma}$, where $\bar{\Sigma} = \frac{1}{T} \sum (\mathbf{y}_t - \bar{\mathbf{y}})(\mathbf{y}_t - \bar{\mathbf{y}})^T$. Setting the prior directly from the data can help move the mass of the distribution to reasonable values of the parameter space. Since each new considered dynamical mode is associated with a set of parameters sampled from the prior distribution, and this dynamical mode is compared against others that have already been informed by the data, setting the base measure in this manner can improve mixing rates over a non-informative setting. For an HDP-SLDS with $\mathbf{x}_t \in \mathbb{R}^n$ and $\mathbf{y}_t \in \mathbb{R}^d$ and $n = d$, we set $S_0 = 0.675\bar{\Sigma}$. We then set the inverse-Wishart prior on the measurement noise, R , to have $r_0 = d + 2$ and $R_0 = 0.075\bar{\Sigma}$. For $n > d$, see [42].

b) Partially Supervised Honey Bee Experiments: For the partially supervised experiments of Sec. IV-C, we set $\Sigma_0 = 0.75S_0$. Since we are not shifting and scaling the observations, we set S_0 to 0.75 times the empirical covariance of the *first difference* observations. We also use $n_0 = 10$, making the distribution tighter than in the unsupervised case. Examining first differences is appropriate since the bee's dynamics are better approximated as a random walk than as i.i.d. observations. Using raw observations in the unsupervised approach creates a larger expected covariance matrix making the prior on the dynamic matrix less informative, which is useful in the absence of other labeled data.

c) IBOVESPA Stock Index Experiments: For the HDP-SLDS variant of the MSSV model of Eq. (38), we rely on the N-IW-N prior described in Sec. V-A. For the dynamic parameter a and process noise mean $\mu^{(k)}$, we use $\mathcal{N}(0, 0.75\bar{\Sigma})$ priors. The IW prior on $\Sigma^{(k)}$ was given 3 degrees of freedom and an expected value of $0.75\bar{\Sigma}$. Finally, each component of the mixture-of-Gaussian measurement noise was given an IW prior with 3 degrees of freedom and an expected value of $5 * \pi^2$, which matches with the moment-matching technique of Harvey et. al. [45].

For the HDP-SLDS comparison using the model of Table I, we use a MNIW prior with $M = 0$, $K = 1$, $n_0 = 3$, and $S_0 = 0.75\bar{\Sigma}$. The IW prior on R was given $r_0 = 100$ and an expected covariance of 25. Our sampler initializes parameters from the prior, and we found it useful to set the prior around large values of R in order to avoid initial samples chattering between dynamical regimes caused by the state sequence having to account for the noise in the observations. After accounting for the residuals of the data in the posterior distribution, we typically learned $R \approx 10$.

REFERENCES

- [1] E. Fox, E. Sudderth, M. Jordan, and A. Willsky, "Nonparametric Bayesian learning of switching dynamical systems," in *Advances in Neural Information Processing Systems*, vol. 21, 2009, pp. 457–464.
- [2] —, "Nonparametric Bayesian identification of jump systems with sparse dependencies," in *Proc. 15th IFAC Symposium on System Identification*, July 2009.
- [3] V. Pavlović, J. Rehg, and J. MacCormick, "Learning switching linear models of human motion," in *Advances in Neural Information Processing Systems*, vol. 13, 2001, pp. 981–987.
- [4] L. Ren, A. Patrick, A. Efros, J. Hodgins, and J. Rehg, "A data-driven approach to quantifying natural human motion," in *SIGGRAPH*, August 2005.
- [5] C.-J. Kim, "Dynamic linear models with Markov-switching," *Journal of Econometrics*, vol. 60, pp. 1–22, 1994.
- [6] M. So, K. Lam, and W. Li, "A stochastic volatility model with Markov switching," *Journal of Business & Economic Statistics*, vol. 16, no. 2, pp. 244–253, 1998.
- [7] C. Carvalho and H. Lopes, "Simulation-based sequential analysis of Markov switching stochastic volatility models," *Computational Statistics & Data Analysis*, vol. 51, pp. 4526–4542, 9 2007.
- [8] X. Rong Li and V. Jilkov, "Survey of maneuvering target tracking. Part V: Multiple-model methods," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 41, no. 4, pp. 1255–1321, 2005.
- [9] E. Fox, E. Sudderth, and A. Willsky, "Hierarchical Dirichlet processes for tracking maneuvering targets," in *Proc. International Conference on Information Fusion*, July 2007.
- [10] S. Oh, J. Rehg, T. Balch, and F. Dellaert, "Learning and inferring motion patterns using parametric segmental switching linear dynamic systems," *International Journal of Computer Vision*, vol. 77, no. 1–3, pp. 103–124, 2008.
- [11] L. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [12] V. Jilkov and X. Rong Li, "Online Bayesian estimation of transition probabilities for Markovian jump systems,"

- IEEE Transactions on Signal Processing*, vol. 52, no. 6, pp. 1620–1630, June 2004.
- [13] C. Li and S. Andersen, “Efficient blind system identification of non-Gaussian autoregressive models with HMM modeling of the excitation,” *IEEE Transactions on Signal Processing*, vol. 55, no. 6, pp. 2432–2445, June 2007.
- [14] J. Chiang, Z. Wang, and M. McKeown, “A hidden Markov, multivariate autoregressive (HMM-mAR) network framework for analysis of surface emg (sEMG) data,” *IEEE Transactions on Signal Processing*, vol. 56, no. 8, pp. 4069–4081, August 2008.
- [15] X. Xuan and K. Murphy, “Modeling changing dependency structure in multivariate time series,” in *Proc. International Conference on Machine Learning*, June 2007.
- [16] Y. Teh, M. Jordan, M. Beal, and D. Blei, “Hierarchical Dirichlet processes,” *Journal of the American Statistical Association*, vol. 101, no. 476, pp. 1566–1581, 2006.
- [17] M. Beal, Z. Ghahramani, and C. Rasmussen, “The infinite hidden Markov model,” in *Advances in Neural Information Processing Systems*, vol. 14, 2002, pp. 577–584.
- [18] J. Paisley and L. Carin, “Hidden Markov models with stick-breaking priors,” *IEEE Transactions on Signal Processing*, vol. 57, no. 10, pp. 3905–3917, October 2009.
- [19] Y. Qi, J. Paisley, and L. Carin, “Music analysis using hidden Markov mixture models,” *IEEE Transactions on Signal Processing*, vol. 55, no. 11, pp. 5209–5224, November 2007.
- [20] F. Caron, M. Davy, A. Doucet, E. Duflos, and P. Vanheeghe, “Bayesian inference for dynamic models with Dirichlet process mixtures,” *IEEE Transactions on Signal Processing*, vol. 56, no. 1, pp. 71–84, January 2008.
- [21] E. Fox, E. Sudderth, M. Jordan, and A. Willsky, “A sticky HDP-HMM with application to speaker diarization,” *To appear in Annals of Applied Statistics*, 2010.
- [22] D. MacKay, *Bayesian methods for backprop networks*, ser. Models of Neural Networks, III. Springer, 1994, ch. 6, pp. 211–254.
- [23] R. Neal, Ed., *Bayesian Learning for Neural Networks*, ser. Lecture Notes in Statistics. Springer, 1996, vol. 118.
- [24] M. Beal, “Variational algorithms for approximate Bayesian inference,” Ph.D. Thesis, University College London, London, UK, 2003.
- [25] S. Paoletti, A. Juloski, G. Ferrari-Trecate, and R. Vidal, “Identification of hybrid systems: A tutorial,” *European Journal of Control*, vol. 2–3, pp. 242–260, 2007.
- [26] R. Vidal, S. Soatto, Y. Ma, and S. Sastry, “An algebraic geometric approach to the identification of a class of linear hybrid systems,” in *Proc. IEEE Conference on Decision and Control*, December 2003.
- [27] Z. Psaradakis and N. Spagnolo, “Joint determination of the state dimension and autoregressive order for models with Markov regime switching,” *Journal of Time Series Analysis*, vol. 27, pp. 753–766, 2006.
- [28] K. Huang, A. Wagner, and Y. Ma, “Identification of hybrid linear time-invariant systems via subspace embedding and segmentation SES,” in *Proc. IEEE Conference on Decision and Control*, December 2004.
- [29] R. Vidal, Y. Ma, and S. Sastry, “Generalized principal component analysis (GPCA): Subspace clustering by polynomial factorization, differentiation, and division,” *UC Berkeley, Technical Report UCB/ERL*, August 2003.
- [30] G. Kotsalis, A. Megretski, and M. Dahleh, “Model reduction of discrete-time Markov jump linear systems,” in *Proc. American Control Conference*, June 2006.
- [31] B. Anderson, “The realization problem for hidden Markov models,” *Mathematics of Control, Signals, and Systems*, vol. 12, pp. 80–120, 1999.
- [32] M. Petreczky and R. Vidal, “Realization theory of stochastic jump-Markov linear systems,” in *Proc. IEEE Conference on Decision and Control*, December 2007.
- [33] Z. Ghahramani and G. Hinton, “Variational learning for switching state-space models,” *Neural Computation*, vol. 12, no. 4, pp. 831–864, 2000.
- [34] R. Vidal, A. Chiuso, and S. Soatto, “Observability and identifiability of jump linear systems,” in *Proc. IEEE Conference on Decision and Control*, December 2002.
- [35] H. Akaike, “A new look at the statistical model identification,” *IEEE Transactions on Automatic Control*, vol. 19, no. 6, pp. 716–723, 1974.
- [36] G. Schwarz, “Estimating the dimension of a model,” *The Annals of Statistics*, pp. 461–464, 1978.
- [37] M. Aoki and A. Havenner, “State space modeling of multiple time series,” *Econometric Reviews*, vol. 10, no. 1, pp. 1–59, 1991.
- [38] J. Sethuraman, “A constructive definition of Dirichlet priors,” *Statistica Sinica*, vol. 4, pp. 639–650, 1994.
- [39] D. Blackwell and J. MacQueen, “Ferguson distributions via Polya urn schemes,” *The Annals of Statistics*, vol. 1, no. 2, pp. 353–355, 1973.
- [40] M. West and J. Harrison, *Bayesian Forecasting and Dynamic Models*. Springer, 1997.
- [41] O. Costa, M. Fragoso, and R. Marques, *Discrete-Time Markov Jump Linear Systems*. Springer, 2005.
- [42] E. Fox, “Bayesian nonparametric learning of complex dynamical phenomena,” Ph.D. dissertation, MIT, July 2009.
- [43] C. Carter and R. Kohn, “Markov chain Monte Carlo in conditionally Gaussian state space models,” *Biometrika*, vol. 83, pp. 589–601, 3 1996.
- [44] J. Hamilton, “A new approach to the economic analysis of nonstationary time series and the business cycle,” *Econometrica*, vol. 57, no. 2, pp. 357–384, 1989.
- [45] A. Harvey, E. Ruiz, and N. Shephard, “Multivariate stochastic variance models,” *Review of Economic Studies*, vol. 61, pp. 247–264, 1994.
- [46] X. Rong Li and V. Jilkov, “Survey of maneuvering target tracking. Part I: Dynamic models,” *IEEE Transactions*

- on Aerospace and Electronic Systems*, vol. 39, no. 4, pp. 1333–1364, 2003.
- [47] R. Moose, H. VanLandingham, and D. McCabe, “Modeling and estimation of tracking maneuvering targets,” *IEEE Transactions on Aerospace and Electronic Systems*, vol. 15, no. 3, pp. 448–456, 1979.
- [48] A. Doucet, N. de Freitas, K. Murphy, and S. Russell, “Rao-Blackwellised particle filtering for dynamic Bayesian networks,” in *Proc. Conference on Uncertainty in Artificial Intelligence*, August 2000, pp. 176–183.
- [49] C. Andrieu, A. Doucet, and R. Holenstein, “Particle Markov chain Monte carlo methods,” *Journal of the Royal Statistical Society, Series B*, vol. 72, no. 3, pp. 269–342, 2010.
- [50] S. Jain and R. Neal, “A split-merge Markov chain Monte Carlo procedure for the Dirichlet process mixture model,” *Journal of Computational and Graphical Statistics*, vol. 13, pp. 158–182, 2004.