

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

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

### 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 (4) via a set of *transition probability measures*  $G_j = \sum_{k=1}^K \pi_{jk} \delta_{\theta_k}$ , where  $\delta_{\theta}$  is a mass concentrated at  $\theta$ . We then operate directly in the parameter space  $\Theta$  and transition between emission parameters with probabilities given by  $\{G_j\}$ . That is

$$\theta'_t \mid \theta'_{t-1} \sim G_{j:\theta'_{t-1}=\theta_j} \\ \mathbf{y}_t \mid \theta'_t \sim F(\theta'_t). \quad (6)$$

Here,  $\theta'_t \in \{\theta_1, \dots, \theta_K\}$  and is equivalent to  $\theta_{z_t}$  of (4). A Bayesian nonparametric HMM takes  $G_j$  to be *random*<sup>1</sup> 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  $\Theta$  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  $\beta_k$ : the  $k^{\text{th}}$  weight is a random proportion  $\nu_k$  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  $\theta'_i$  taking identical values within the set  $\{\theta_k\}$ , with  $\theta_k$  defined as in (7). For each value  $\theta'_i$ , let  $z_i$  be an indicator random variable that picks out the unique value  $\theta_k$  such that  $\theta'_i = \theta_{z_i}$ . Blackwell and MacQueen [39] introduced a Pólya urn representation of the  $\theta'_i$

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

<sup>1</sup>Formally, a random measure on a measurable space  $\Theta$  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  $A \in \mathcal{A}$ .

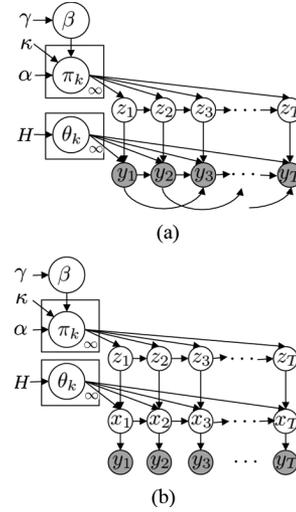


Fig. 1. Sticky HDP-HMM prior on (a) switching VAR(2) and (b) SLDS processes with the mode evolving as  $z_{t+1} \mid \{\pi_k\}_{k=1}^{\infty}, z_t \sim \pi_{z_t}$  for  $\pi_k \mid \alpha, \kappa, \beta \sim \text{DP}(\alpha + \kappa, (\alpha\beta + \kappa\delta_k) / (\alpha + \kappa))$ . Here,  $\beta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

Here,  $n_k$  is the number of observations  $\theta'_i$  taking the value  $\theta_k$ . From (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  $G_j$  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 DP  $(\gamma, H)$ . Through this construction, one can show that the probability measures are described as

$$G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k} \quad \beta \mid \gamma \sim \text{GEM}(\gamma) \\ G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\theta_k} \quad \pi_j \mid \alpha, \beta \sim \text{DP}(\alpha, \beta) \\ \theta_k \mid H \sim H. \quad (10)$$

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:

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

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