

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

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

form a matrix  $\boldsymbol{\Psi}^{(k)}$  with  $n_k$  columns consisting of the  $\boldsymbol{\psi}_t$  with  $z_t = k$ . Then

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

where  $\bar{\boldsymbol{\Psi}}^{(k)}$  is a matrix of the associated  $\bar{\boldsymbol{\psi}}_{t-1}$  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 (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 (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  $\mathbf{D}^{(k)} = \{\boldsymbol{\Psi}^{(k)}, \bar{\boldsymbol{\Psi}}^{(k)}\}$ . The posterior distribution of the dynamic parameters for the  $k^{\text{th}}$  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])

$$p(\mathbf{A}^{(k)} \mid \Sigma^{(k)}, \mathbf{D}^{(k)}) = \mathcal{MN}(\mathbf{S}_{\boldsymbol{\psi}|\bar{\boldsymbol{\psi}}}^{(k)} \mathbf{S}_{\bar{\boldsymbol{\psi}}|\boldsymbol{\psi}}^{-1}, \Sigma^{(k)}, \mathbf{S}_{\bar{\boldsymbol{\psi}}|\boldsymbol{\psi}}^{(k)}) \quad (17)$$

with  $\mathbf{B}^{-1}$  denoting  $(\mathbf{B}^{(k)})^{-1}$  for a given matrix  $\mathbf{B}$ ,  $\mathcal{MN}(M, V, K)$  denoting a matrix-normal prior<sup>2</sup> for  $\mathbf{A}^{(k)}$  with mean matrix  $M$  and left and right covariances  $K^{-1}$  and  $V$  and

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

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

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

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

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

by data terms  $\mathbf{S}_{\boldsymbol{\psi}|\bar{\boldsymbol{\psi}}}^{(k)} = \mathbf{S}_{\boldsymbol{\psi}|\boldsymbol{\psi}}^{(k)} - \mathbf{S}_{\bar{\boldsymbol{\psi}}|\boldsymbol{\psi}}^{(k)} \mathbf{S}_{\bar{\boldsymbol{\psi}}|\boldsymbol{\psi}}^{-1} \mathbf{S}_{\boldsymbol{\psi}|\bar{\boldsymbol{\psi}}}^{(k)T}$  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  $\mathbf{A}^{(k)}$

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

Each precision parameter  $\alpha_j^{(k)}$  is given a  $\text{Gamma}(a, b)$  prior. The zero-mean Gaussian prior penalizes nonzero 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  $j^{\text{th}}$  state component does not contribute to the dynamics of the  $k^{\text{th}}$  mode. Thus, examining the set of large  $\alpha_j^{(k)}$  provides insight into the order of that mode. Looking at the  $k^{\text{th}}$  dynamical mode alone, having  $\mathbf{a}_j^{(k)} = 0$  implies that the realization of *that mode* is not minimal since the associated Hankel matrix

$$\mathcal{H} = [C^T \quad CA^T \quad \dots \quad (CA^{d-1})^T]^T \times [G \quad AG \quad \dots \quad A^{d-1}G] \equiv \mathcal{OR} \quad (21)$$

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.

3) *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  $j^{\text{th}}$  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  $C = [I_d \quad 0]$  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