

# LEARNING LINEAR DYNAMICAL SYSTEMS WITH SPARSE SYSTEM MATRICES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Due to the tractable analysis and control, linear dynamical systems (LDSs) provide a fundamental mathematical tool for time-series data modeling in various disciplines. Particularly, many LDSs have sparse system matrices because interactions among variables are limited or only a few significant relationships exist. However, available learning algorithms for LDSs lack the ability to learn system matrices with the sparsity constraint. To address this issue, we impose sparsity-promoting priors on system matrices and explore the expectation-maximization (EM) algorithm to give a maximum a posteriori (MAP) estimate of both hidden states and system matrices from noisy observations. In addition, we find that many learning algorithms based on the gradient descent method use an inappropriate derivative rule, because they neglect the inherent symmetry of noise covariance matrices. Here, we consider the derivative rule of structured matrices during the optimization process to guarantee their symmetry. Experimental results on simulation and real-world problems illustrate that the proposed algorithm significantly improves learning accuracy over classical ones.

## 1 INTRODUCTION

Linear dynamical systems (LDSs) are fundamental mathematical models for analyzing time-series data with application in robotics (Mamakoukas et al., 2019; 2020), systems biology (Jin et al., 2020b; Pillonetto & Ljung, 2023), and natural language processing (Smith et al., 1999; Belanger & Kakade, 2015). Basically, LDSs consider a linear transformation between finite-dimensional hidden states disturbed by input signals and noise to capture the time evolution of systems (Hazan et al., 2017). In addition, LDSs are also widely used to approximate complex nonlinear systems in industrial processes given their relative simplicity (Yuan et al., 2017; Lusch et al., 2018). Due to a complete rigorous theory available on LDSs, learning LDSs from noisy observations can enable us to make tractable analysis and control of systems (Chen & Poor, 2022; Bakshi et al., 2023).

In this paper, we focus on learning LDSs with sparse system matrices for two important reasons. First, the learned LDSs should include the minimally required parameters to explain time-series data following the *Occam's razor* principle. Additionally, many real-world systems have sparse topology because each state or measurement variable only depends on a few other state variables and inputs (Efroni et al., 2022). For example, a gene only regulates the expression of a limited number of other genes in gene regulatory networks (He et al., 2024). In industry, communication systems usually have sparse topology to reduce energy consumption (Jin et al., 2020a;b). However, available learning algorithms lack the ability to learn LDSs with the sparsity constraint on system matrices. In addition, many algorithms neglect the inherent symmetry of noise covariance matrices during the optimization process.

To learn the LDSs with sparse system matrices, we impose the sparsity-promoting prior on them to balance model complexity and modeling error in this paper (Wang et al., 2024). Subsequently, we can combine the likelihood and prior functions to derive the posterior distribution of system matrices following the Bayes' rule. However, directly maximizing such a posterior distribution to estimate system matrices is intractable because the states of LDSs are unknown. To address this issue, we explore the expectation-maximization (EM) algorithm to give an alternate maximum a posteriori (MAP) estimate of hidden states and system matrices. In the expectation step, we use the Rauch-Tung-Striebel (RTS) smoother, also known as the Kalman smoother, to give a closed-

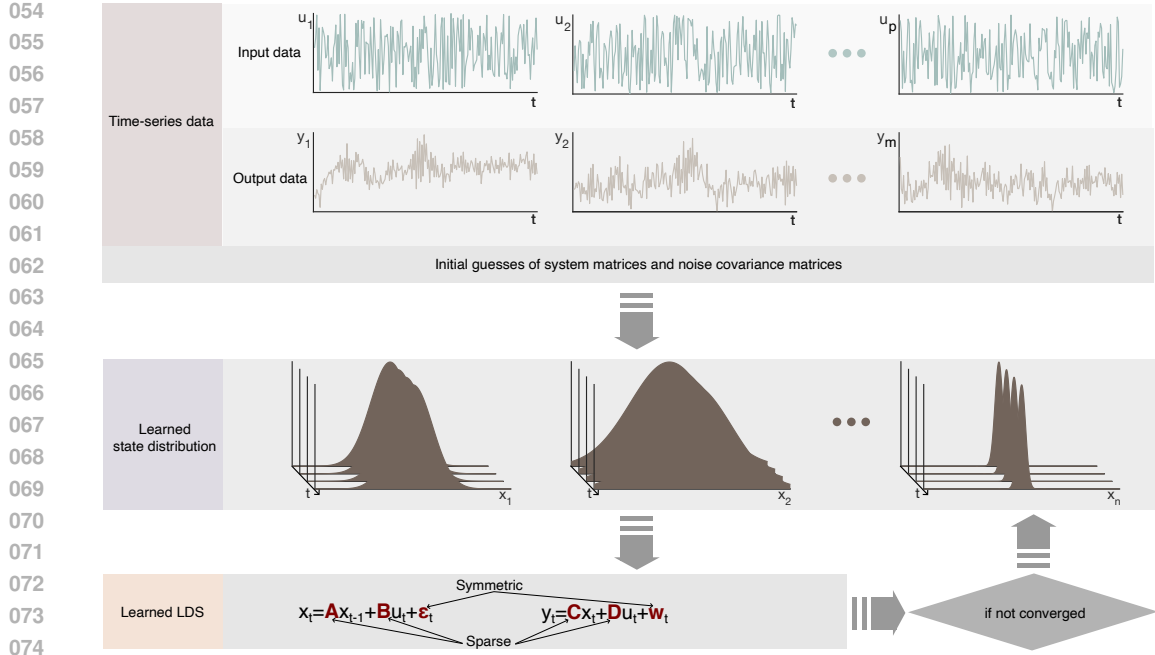


Figure 1: The pipeline of the proposed algorithm. Given time-series data and initial guesses of system matrices and noise covariance matrices, the proposed algorithm alternately learns the hidden state distribution and unknown LDS until it converges.

form update rule for the hidden states. In the maximization step, we leverage the block gradient descent method to optimize the system matrices in turn because they are highly coupled in the objective function. Given the inherent symmetry of noise variance, we consider the derivative rule of structured matrices during the optimization process. By alternately performing the expectation and maximization steps until convergence, the proposed algorithm can determine the sparse system matrices of LDSs from noisy observations.

**Contributions.** The contributions of this paper are summarized as follows:

- Leveraging sparsity-promoting techniques, we propose an algorithm to learn LDSs with sparse system matrices from noisy observations. Particularly, the proposed algorithm gives the MAP estimate of both hidden states and system matrices by exploring the EM algorithm.
- The proposed algorithm utilizes the derivative rule of structured matrices to ensure the symmetry of noise covariance matrices during the optimization process. While many EM-based learning algorithms provide the same update rule for noise covariance matrices, we argue that they use an inappropriate derivative rule.
- Experimental results on simulation and real-world datasets demonstrate that the proposed algorithm outperforms the classical ones on learning LDSs with sparse system matrices.

## 2 RELATED WORK

**Prediction error minimization.** Prediction error minimization (PEM) learns LDSs by minimizing one-step prediction error objective via gradient-based optimization methods (Ljung, 2002; Katayama et al., 2005). Expanding LDSs into linear AutoRegressive Moving Average with exogenous excitation (ARMAX) or AutoRegressive Moving Average (ARMA) models, Li & Zhang (2006) leverage the PEM-based method to estimate the second-order structural parameters of linear structural systems. Given a symmetric transition matrix, Hazan et al. (2017) present an efficient method for the online prediction of discrete-time LDSs by formulating system identification as an online PEM problem. In particular, Abdalmoaty & Hjalmarsson (2019) extends PEM for learning

108 stochastic nonlinear models recently. However, PEM is found to be sensitive to initial values and  
 109 cannot characterize the sparsity of system matrices (Martens, 2010).

110  
 111 **Subspace state-space system identification.** Basically, subspace state-space system identification  
 112 (4SID) algorithms project data Hankel matrices onto certain subspaces to estimate the extended  
 113 observability matrix and hidden states using linear algebra tools (Larimore, 1990; Verhaegen &  
 114 Dewilde, 1992; Van Overschee & De Moor, 1994). Leveraging the least squares method, system  
 115 matrices can thus be recovered from either the extended observability matrix or hidden states (Fa-  
 116 voreel et al., 2000). Based on principal component analysis, Wang & Qin (2002) present a new 4SID  
 117 algorithm to learn LDSs under the errors-in-variables situation. By choosing different weighting  
 118 matrices to perform the singular value decomposition, Van Overschee & De Moor (2012) provide  
 119 a geometric framework to unify almost all classical 4SID methods. Further, Huang et al. (2016)  
 120 presents the Weight-Least-Square method to learn stable LDSs by multiplying the unstable compo-  
 121 nent with a weight matrix. However, 4SID algorithms learn system matrices via the least squares  
 122 method and thus cannot produce sparse system matrices (Tibshirani, 1996). In particular, it is widely  
 123 recognized that such algorithms generally cannot obtain accurate system matrices as required (Qin,  
 2006; Martens, 2010).

124 **Maximum likelihood estimation.** Because the joint likelihood function of LDSs involves hidden  
 125 states, the EM algorithm is employed to give the maximum likelihood estimation (MLE) of system  
 126 matrices (Shumway & Stoffer, 1982; Ghahramani & Hinton, 1996). Leveraging the EM algorithm,  
 127 the distribution of hidden states can be explicitly derived using the Kalman smoother based on the  
 128 current estimate of system matrices. Subsequently, it updates system matrices by maximizing the  
 129 expected log-likelihood with respect to the hidden states. Gibson & Ninness (2005) present a robust  
 130 MLE of LDSs by implementing the expectation and maximization steps via the LR and Cholesky  
 131 factorisation respectively. To increase the efficiency of EM for learning LDSs, Martens (2010) pro-  
 132 poses an approximate second-order statistics (ASOS) scheme to approximate the expectation step.  
 133 Combining EM and Lagrangian relaxation, Umenberger et al. (2018) use semidefinite program-  
 134 ming to optimize the tight bounds on the likelihood to learn LDSs with model stability constraints.  
 135 However, such learning algorithms lack the ability to deal with sparse system matrices. Particularly,  
 136 an inappropriate derivative rule is used to take the derivatives of the likelihood function with respect  
 137 to noise covariance matrices due to the neglect of the inherent symmetry of them.

### 138 3 PROBLEM FORMULATION

139 Generally, LDSs describe time-series data  $\{(\mathbf{u}_t, \mathbf{y}_t)\}_{t=1}^T$  through the following stochastic difference  
 140 equation (Shumway et al., 2000):

$$141 \mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{B}\mathbf{u}_t + \boldsymbol{\varepsilon}_t, \quad (1)$$

$$142 \mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{D}\mathbf{u}_t + \boldsymbol{\omega}_t, \quad (2)$$

143 where  $\mathbf{u}_t \in \mathbb{R}^p$  is the input signal,  $\mathbf{y}_t \in \mathbb{R}^m$  is the noisy observation,  $\mathbf{x}_t \in \mathbb{R}^n$  is the hidden  
 144 state,  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{C} \in \mathbb{R}^{m \times n}$ , and  $\mathbf{D} \in \mathbb{R}^{m \times p}$  are the system matrices, and  $\boldsymbol{\varepsilon}_t \sim$   
 145  $\mathcal{N}(0, \mathbf{R})$  and  $\boldsymbol{\omega}_t \sim \mathcal{N}(0, \mathbf{Q})$  are the process and measurement noise, respectively. In fact, LDSs are  
 146 widely used to model complex systems and have been received successful applications in industrial  
 147 processes (Favoreel et al., 2000). To make predictions about future outputs and unknown states and  
 148 realize the control of systems, it is necessary to propose an algorithm to learn the model parameters  
 149  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ ,  $\mathbf{D}$ ,  $\mathbf{R}$ , and  $\mathbf{Q}$  from time-series data  $\{(\mathbf{u}_t, \mathbf{y}_t)\}_{t=1}^T$ . In particular, many systems have  
 150 sparse topology to enable efficient working mechanisms (Jin et al., 2020a). As a result, the proposed  
 151 algorithm needs to integrate such priori information into the learning process.

152 **Sparsity-promoting prior.** Sparsity-promoting priors can enforce the sparsity of model parameters  
 153 by balancing model complexity and modeling error (Wang et al., 2023; Tripura & Chakraborty,  
 154 2023). Because the likelihood function of LDSs is Gaussian distributed, the Student’s  $t$ -distribution  
 155 prior severing as its conjugate prior can be imposed on each component of the unknown system  
 156 matrices to promote their sparsity. Generally, the Student’s  $t$ -distribution prior is implemented in a  
 157 hierarchical way (Tipping, 2001). It imposes a Gaussian prior on the system matrices and then adopt  
 158 an Inverse-Gamma hyperprior on the unknown variance of the Gaussian distribution. For example,  
 159 we can impose the Student’s  $t$ -distribution prior on the system matrix  $\mathbf{A}$  to promote its sparsity as  
 160  
 161

162 follows:

$$163 p(\mathbf{A} \mid \Gamma_a) = \prod_{i=1}^n \prod_{j=1}^n p(\mathbf{A}_{ij} \mid \Gamma_{a,ij}) = \prod_{i=1}^n \prod_{j=1}^n \frac{1}{\sqrt{2\pi\Gamma_{a,ij}}} \exp\left(-\frac{\mathbf{A}_{ij}^2}{2\Gamma_{a,ij}}\right), \quad (3)$$

$$164 p(\Gamma_a) = \prod_{i=1}^n \prod_{j=1}^n \frac{a_0^{b_0}}{\Gamma(a_0)} \Gamma_{a,ij}^{-a_0-1} \exp\left(-\frac{b_0}{\Gamma_{a,ij}}\right), \quad (4)$$

165 where  $\mathbf{A}_{ij}$  and  $\Gamma_{a,ij}$  are the  $ij$ -th components of  $\mathbf{A}$  and  $\Gamma_a$ , respectively. To generate non-  
166 informative hyperprior on  $\Gamma_{a,ij}$ ,  $a_0$  and  $b_0$  are typically set to very small values (e.g.,  $10^{-6}$ ). In  
167 addition,  $\Gamma_b$ ,  $\Gamma_c$ ,  $\Gamma_d$ ,  $\Gamma_{b,ij}$ ,  $\Gamma_{c,ij}$ , and  $\Gamma_{d,ij}$  are defined in a similar manner (see Appendix A).

168 **Loss function.** Following the Bayes' rule, we can combine the marginal likelihood function and  
169 sparsity-promoting prior to estimate the model parameters:

$$170 p(\Theta \mid \mathbf{Y}) \propto \underbrace{p(\mathbf{Y} \mid \Theta)}_{\text{Marginal likelihood}} \times \underbrace{p(\Theta)}_{\text{Prior}}, \quad (5)$$

171 where  $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T]$  and  $\Theta = \{\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{R}, \mathbf{Q}, \Gamma_a, \Gamma_b, \Gamma_c, \Gamma_d\}$ . Note that directly  
172 maximizing equation 5 is generally intractable because  $p(\mathbf{Y} \mid \Theta)$  is hard to be explicitly computed.  
173 However, the EM algorithm provides an iterative optimization framework to address such a problem.  
174 Instead of maximizing equation 5, the EM algorithm focuses on iteratively improving the expected  
175 value of the log posterior function of  $\Theta$  with respect to the hidden state vector  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T]$   
176 as follows:

$$177 H(\Theta \mid \Theta^k) = \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} \mid \mathbf{Y}, \Theta^k)}[\log p(\mathbf{Y}, \mathbf{X} \mid \Theta)p(\Theta)]. \quad (6)$$

178 Notably, improving equation 6 is equivalent to improving equation 5 at each iteration (Little &  
179 Rubin, 2019).

### 180 3.1 RAUCH–TUNG–STRIEBEL SMOOTHER

181 To explicitly compute equation 6, we first need to derive the conditional distribution of  $\mathbf{x}_t$  given the  
182 noisy observation  $\mathbf{Y}$  and current  $\Theta^k = \{\mathbf{A}^k, \mathbf{B}^k, \mathbf{C}^k, \mathbf{D}^k, \mathbf{R}^k, \mathbf{Q}^k, \Gamma_a^k, \Gamma_b^k, \Gamma_c^k, \Gamma_d^k\}$ , which can be  
183 formulated as a classical smoothing problem. For LDSs, the RTS smoother provides a closed-form  
184 smoothing solution for  $p(\mathbf{x}_t \mid \mathbf{Y}, \Theta^k)$ .

185 **Lemma 1.** (RTS smoother (Särkkä & Svensson, 2023)) For LDSs, the RTS smoother states that

$$186 p(\mathbf{x}_t \mid \mathbf{Y}, \Theta^k) = \mathcal{N}(\mathbf{x}_t \mid \mathbf{m}_t^k, \mathbf{P}_t^k), \quad (7)$$

187 where  $t = 0, \dots, T$ . Here,  $\mathbf{m}_t^k$  and  $\mathbf{P}_t^k$  are derived via the reverse-time recursions as follows:

$$188 \mathbf{m}_t^k = \boldsymbol{\mu}_t^k + \mathbf{G}_t^k (\mathbf{m}_{t+1}^k - \bar{\boldsymbol{\mu}}_{t+1}^k), \quad (8)$$

$$189 \mathbf{P}_t^k = \boldsymbol{\Sigma}_t^k + \mathbf{G}_t^k (\mathbf{P}_{t+1}^k - \bar{\boldsymbol{\Sigma}}_{t+1}^k) (\mathbf{G}_t^k)', \quad (9)$$

190 with  $\mathbf{G}_t^k = \boldsymbol{\Sigma}_t^k (\mathbf{A}^k)' (\bar{\boldsymbol{\Sigma}}_{t+1}^k)^{-1}$ . The quantities  $\boldsymbol{\mu}_t^k$ ,  $\bar{\boldsymbol{\mu}}_t^k$ ,  $\boldsymbol{\Sigma}_t^k$ , and  $\bar{\boldsymbol{\Sigma}}_t^k$  coupled in equation 8  
191 and equation 9 are pre-computed using the Kalman filter as follows:

$$192 \bar{\boldsymbol{\mu}}_t^k = \mathbf{A}^k \boldsymbol{\mu}_{t-1}^k + \mathbf{B}^k \mathbf{u}_t, \quad (10)$$

$$193 \bar{\boldsymbol{\Sigma}}_t^k = \mathbf{A}^k \boldsymbol{\Sigma}_{t-1}^k (\mathbf{A}^k)' + \mathbf{R}^k, \quad (11)$$

$$194 \mathbf{K}_t^k = \bar{\boldsymbol{\Sigma}}_t^k (\mathbf{C}^k)' (\mathbf{C}^k \bar{\boldsymbol{\Sigma}}_t^k (\mathbf{C}^k)' + \mathbf{Q}^k)^{-1}, \quad (12)$$

$$195 \boldsymbol{\mu}_t^k = \bar{\boldsymbol{\mu}}_t^k + \mathbf{K}_t^k (\mathbf{Y}_t - \mathbf{C}^k \bar{\boldsymbol{\mu}}_t^k - \mathbf{D}^k \mathbf{u}_t), \quad (13)$$

$$196 \boldsymbol{\Sigma}_t^k = (\mathbf{I}_n - \mathbf{K}_t^k \mathbf{C}^k) \bar{\boldsymbol{\Sigma}}_t^k, \quad (14)$$

197 where  $\mathbf{I}_n$  is an identity matrix of dimension  $n$ . Note that the reverse-time recursions of equation 8  
198 and equation 9 start from the initial conditions  $\mathbf{m}_T^k = \boldsymbol{\mu}_T^k$  and  $\mathbf{P}_T^k = \boldsymbol{\Sigma}_T^k$ , and the recursions  
199 of equation 10–equation 14 start from the mean  $\boldsymbol{\mu}_0^k$  and covariance  $\boldsymbol{\Sigma}_0^k$  of the initial state  $\mathbf{x}_0$ .

Besides  $p(\mathbf{x}_t | \mathbf{Y}, \Theta^k)$ , we also need to derive the covariance matrix between the adjacent states  $\mathbf{x}_t$  and  $\mathbf{x}_{t-1}$  given  $\mathbf{Y}$  and  $\Theta^k$  to compute equation 6. To address this issue, the following lemma gives necessary recursions.

**Lemma 2.** (The lag-one covariance smoother (Särkkä & Svensson, 2023)) For the LDSs, the covariance matrix  $\mathbf{P}_{t,t-1}^k$  between the adjacent states  $\mathbf{x}_t$  and  $\mathbf{x}_{t-1}$  given  $\mathbf{Y}$  and  $\Theta^k$  can be recursively derived as follows:

$$\mathbf{P}_{t,t-1}^k = \Sigma_t^k (\mathbf{G}_{t-1}^k)' + \mathbf{G}_t^k (\mathbf{P}_{t+1,t}^k - \mathbf{A}^k \Sigma_t^k) (\mathbf{G}_{t-1}^k)' \quad (15)$$

with the initial condition  $\mathbf{P}_{T,T-1}^k = (\mathbf{I}_n - \mathbf{K}_T^k \mathbf{C}^k) \mathbf{A}^k \Sigma_{T-1}^k$ .

Based on Lemmas 1 and 2, we are able to calculate the loss function in equation 6 as follows:

$$H(\Theta | \Theta^k) = H_1(\mathbf{A}, \mathbf{B}, \mathbf{R}) + H_2(\mathbf{C}, \mathbf{D}, \mathbf{Q}) + H_3(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \Gamma_a, \Gamma_b, \Gamma_c, \Gamma_d), \quad (16)$$

where

$$H_1(\mathbf{A}, \mathbf{B}, \mathbf{R}) = \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{X} | \mathbf{A}, \mathbf{B}, \mathbf{R})], \quad (17)$$

$$H_2(\mathbf{C}, \mathbf{D}, \mathbf{Q}) = \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{Y} | \mathbf{X}, \mathbf{C}, \mathbf{D}, \mathbf{Q})], \quad (18)$$

$$H_3(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \Gamma_a, \Gamma_b, \Gamma_c, \Gamma_d) = \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{A}, \Gamma_a) p(\mathbf{B}, \Gamma_b) p(\mathbf{C}, \Gamma_c) p(\mathbf{D}, \Gamma_d)]. \quad (19)$$

Due to the limited space, the detailed derivation of equation 16 and explicit mathematical expressions of  $H_1(\mathbf{A}, \mathbf{B}, \mathbf{R})$ ,  $H_2(\mathbf{C}, \mathbf{D}, \mathbf{Q})$ , and  $H_3(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \Gamma_a, \Gamma_b, \Gamma_c, \Gamma_d)$  are given in Appendix B.1.

### 3.2 PARAMETER AND HYPERPARAMETER LEARNING

As  $H(\Theta | \Theta^k)$  is a non-convex function and unknown parameters are highly coupled, it is difficult to obtain an efficient algorithm with theoretical guarantees for solving such a problem. A heuristic method is to leverage the block gradient descent method to iteratively optimize the model parameters.

**Update procedures of  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$ .** For MLE, leveraging the EM algorithm can give a closed-form solution to update  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$  (Ghahramani & Hinton, 1996; Gibson & Ninness, 2005). However, it is intractable to obtain a similar update procedure in this case due to the introduction of the sparsity-promoting prior. For example, we can calculate the derivative of  $H(\Theta | \Theta^k)$  with respect to  $\mathbf{A}$  at the  $k$ th iteration as follows:

$$\begin{aligned} & \frac{\partial H_1(\mathbf{A}, \mathbf{B}^k, \mathbf{R}^k)}{\partial \mathbf{A}} + \frac{\partial H_3(\mathbf{A}, \mathbf{B}^k, \mathbf{C}^k, \mathbf{D}^k, \Gamma_a^k, \Gamma_b^k, \Gamma_c^k, \Gamma_d^k)}{\partial \mathbf{A}} \\ &= \sum_{t=1}^T (\mathbf{R}^k)^{-1} (\mathbf{P}_{t,t-1}^k + (\mathbf{m}_t^k - \mathbf{A} \mathbf{m}_{t-1}^k - \mathbf{B}^k \mathbf{u}_t) (\mathbf{m}_{t-1}^k)' - \mathbf{A} \mathbf{P}_{t-1}^k) - \mathbf{A} \odot \bar{\Gamma}_a^k, \end{aligned} \quad (20)$$

where the  $ij$ th component of  $\bar{\Gamma}_a^k$  is  $1/\Gamma_{a,ij}^k$  and  $\odot$  is the Hadamard product. Obviously, setting equation 20 to zero and solving for  $\mathbf{A}$  cannot give a closed-form solution. To address this issue, we approximate  $\mathbf{R}^k$  using the diagonal matrix formed by its diagonal components to facilitate the optimization process. As such, we can calculate the derivative of  $H(\Theta | \Theta^k)$  with respect to the  $r$ th row of  $\mathbf{A}$ , denoted as  $\mathbf{A}_r$ , at the  $k$ th iteration as follows:

$$\begin{aligned} & \frac{\partial H_1(\mathbf{A}, \mathbf{B}^k, \mathbf{R}^k)}{\partial \mathbf{A}_r} + \frac{\partial H_3(\mathbf{A}, \mathbf{B}^k, \mathbf{C}^k, \mathbf{D}^k, \Gamma_a^k, \Gamma_b^k, \Gamma_c^k, \Gamma_d^k)}{\partial \mathbf{A}_r} \\ &= \sum_{t=1}^T (\mathbf{R}_{rr}^k)^{-1} (\mathbf{P}_{t,t-1,r}^k + (\mathbf{m}_{t,r}^k - \mathbf{A}_r \mathbf{m}_{t-1}^k - \mathbf{B}_r^k \mathbf{u}_t) (\mathbf{m}_{t-1}^k)' - \mathbf{A}_r \mathbf{P}_{t-1}^k) - \mathbf{A}_r \bar{\Gamma}_{a,r}^{kd}, \end{aligned} \quad (21)$$

where  $\mathbf{R}_{rr}^k$  is the  $rr$ th component of  $\mathbf{R}^k$ ,  $\mathbf{P}_{t,t-1,r}^k$ ,  $\mathbf{m}_{t,r}^k$ , and  $\mathbf{B}_r^k$  are the  $r$ th rows of  $\mathbf{P}_{t,t-1}^k$ ,  $\mathbf{m}_t^k$ , and  $\mathbf{B}^k$ , respectively. In particular,  $\bar{\Gamma}_{a,r}^{kd} = \text{diag}[\bar{\Gamma}_{a,r}^k]$  with  $\bar{\Gamma}_{a,r}^k$  being the  $r$ th row of  $\bar{\Gamma}_a^k$ . Set-

ting equation 21 to zero leads to

$$\begin{aligned} \mathbf{A}_r^{k+1} &= \left( \sum_{t=1}^T ((\mathbf{m}_{t,r}^k - \mathbf{B}_r^k \mathbf{u}_t)(\mathbf{m}_{t-1}^k)' + \mathbf{P}_{t,t-1,r}^k) \right) \\ &\quad \times \left( \sum_{t=1}^T (\mathbf{P}_{t-1}^k + \mathbf{m}_{t-1}^k (\mathbf{m}_{t-1}^k)') + \mathbf{R}_{rr}^k \bar{\Gamma}_{a,r}^{kd} \right)^{-1}. \end{aligned} \quad (22)$$

The detailed derivation of equation 22 can be found in Appendix B.2. Similarly, we can update the  $r$ th row of  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$  as follows:

$$\mathbf{B}_r^{k+1} = \left( \sum_{t=1}^T (\mathbf{m}_{t,r}^k - \mathbf{A}_r^{k+1} \mathbf{m}_{t-1}^k) \mathbf{u}_t' \right) \left( \sum_{t=1}^T \mathbf{u}_t \mathbf{u}_t' + \mathbf{R}_{rr}^k \bar{\Gamma}_{b,r}^{kd} \right)^{-1}, \quad (23)$$

$$\mathbf{C}_r^{k+1} = \left( \sum_{t=1}^T (\mathbf{y}_{t,r} - \mathbf{D}_r^k \mathbf{u}_t) (\mathbf{m}_t^k)' \right) \left( \sum_{t=1}^T (\mathbf{P}_t^k + \mathbf{m}_t^k (\mathbf{m}_t^k)') + \mathbf{Q}_{rr}^k \bar{\Gamma}_{c,r}^{kd} \right)^{-1}, \quad (24)$$

$$\mathbf{D}_r^{k+1} = \left( \sum_{t=1}^T (\mathbf{y}_{t,r} - \mathbf{C}_r^{k+1} \mathbf{m}_t^k) \mathbf{u}_t' \right) \left( \sum_{t=1}^T \mathbf{u}_t \mathbf{u}_t' + \mathbf{Q}_{rr}^k \bar{\Gamma}_{d,r}^{kd} \right)^{-1}, \quad (25)$$

where  $\mathbf{y}_{t,r}$  is the  $r$ th component of  $\mathbf{y}_t$ ,  $\mathbf{Q}_{rr}^k$  is the  $rr$ th component of  $\mathbf{Q}^k$ , and  $\bar{\Gamma}_{b,r}^{kd}$ ,  $\bar{\Gamma}_{c,r}^{kd}$ , and  $\bar{\Gamma}_{d,r}^{kd}$  are defined as that of  $\bar{\Gamma}_{a,r}^{kd}$ .

**Update procedures of  $\mathbf{R}$  and  $\mathbf{Q}$ .** Because many learning algorithms do not consider the inherent symmetry of noise covariance matrices, we argue that they use an inappropriate derivative rule to calculate the derivatives of the loss function with respect to  $\mathbf{R}$  and  $\mathbf{Q}$  (Gibson & Ninness, 2005; Umenberger et al., 2018). Based on the derivative rule of structured matrices (Petersen et al., 2008), we can calculate the derivative of  $H(\Theta | \Theta^k)$  with respect to  $\mathbf{R}$  at the  $k$ th iteration as follows:

$$\frac{\partial H_1(\mathbf{A}^{k+1}, \mathbf{B}^{k+1}, \mathbf{R})}{\partial \mathbf{R}} = \frac{2\mathbf{L}(\mathbf{R}) - \mathbf{L}(\mathbf{R}) \odot \mathbf{I}_n}{2}, \quad (26)$$

where

$$\begin{aligned} \mathbf{L}(\mathbf{R}) &= \sum_{t=1}^T \mathbf{R}^{-1} (\mathbf{m}_t^k - \mathbf{A}^{k+1} \mathbf{m}_{t-1}^k - \mathbf{B}^{k+1} \mathbf{u}_t) (\mathbf{m}_t^k - \mathbf{A}^{k+1} \mathbf{m}_{t-1}^k - \mathbf{B}^{k+1} \mathbf{u}_t)' \mathbf{R}^{-1} \\ &\quad + \sum_{t=1}^T \mathbf{R}^{-1} (\mathbf{P}_t^k - \mathbf{A}^{k+1} \mathbf{P}_{t,t-1}^k - \mathbf{P}_{t,t-1}^k (\mathbf{A}^{k+1})' + \mathbf{A}^{k+1} \mathbf{P}_{t-1}^k (\mathbf{A}^{k+1})') \mathbf{R}^{-1} - T\mathbf{R}^{-1}. \end{aligned} \quad (27)$$

To update  $\mathbf{R}$ , we first introduce the following lemma to simplify the derivation.

**Lemma 3.** For a square matrix  $\mathbf{H} \in \mathbb{R}^{n \times n}$ , if  $2\mathbf{H} - \mathbf{H} \odot \mathbf{I}_n = \mathbf{0}$ , we have  $\mathbf{H} = \mathbf{0}$ .

The proof of Lemma 3 is straightforward and thus is omitted here. Based on Lemma 3, setting  $\mathbf{L}(\mathbf{R})$  to zero yields

$$\begin{aligned} \mathbf{R}^{k+1} &= \frac{\sum_{t=1}^T (\mathbf{m}_t^k - \mathbf{A}^{k+1} \mathbf{m}_{t-1}^k - \mathbf{B}^{k+1} \mathbf{u}_t) (\mathbf{m}_t^k - \mathbf{A}^{k+1} \mathbf{m}_{t-1}^k - \mathbf{B}^{k+1} \mathbf{u}_t)'}{T} \\ &\quad + \frac{\sum_{t=1}^T (\mathbf{P}_t^k - \mathbf{A}^{k+1} \mathbf{P}_{t,t-1}^k - \mathbf{P}_{t,t-1}^k (\mathbf{A}^{k+1})' + \mathbf{A}^{k+1} \mathbf{P}_{t-1}^k (\mathbf{A}^{k+1})')}{T}. \end{aligned} \quad (28)$$

**Remark 1.** Without considering the symmetry of  $\mathbf{R}$ , the derivative of  $H(\Theta | \Theta^k)$  with respect to  $\mathbf{R}$  is equal to  $\mathbf{L}(\mathbf{R})$ . Hence, many learning algorithms give the same update procedure of  $\mathbf{R}$  as ours. However, we argue that they use an inappropriate derivative rule during the optimization process.

Similarly, we can update  $\mathbf{Q}$  as follows:

$$\mathbf{Q}^{k+1} = \frac{\sum_{t=1}^T ((\mathbf{y}_t - \mathbf{C}^{k+1} \mathbf{m}_t^k - \mathbf{D}^{k+1} \mathbf{u}_t)(\mathbf{y}_t - \mathbf{C}^{k+1} \mathbf{m}_t^k - \mathbf{D}^{k+1} \mathbf{u}_t)' + \mathbf{C}^{k+1} \mathbf{P}_t^k (\mathbf{C}^{k+1})')}{T}. \quad (29)$$

**Update procedures of  $\Gamma_a$ ,  $\Gamma_b$ ,  $\Gamma_c$ , and  $\Gamma_d$ .** Because each component of  $\Gamma_a$ ,  $\Gamma_b$ ,  $\Gamma_c$ , and  $\Gamma_d$  is independent, we can update them individually. For example, we can calculate the derivative of  $H(\Theta | \Theta^k)$  with respect to  $\Gamma_{a,ij}$  at the  $k$ th iteration as follows:

$$\frac{H_3(\mathbf{A}^{k+1}, \mathbf{B}^{k+1}, \mathbf{C}^{k+1}, \mathbf{D}^{k+1}, \Gamma_a, \Gamma_b^k, \Gamma_c^k, \Gamma_d^k)}{\partial \Gamma_{a,ij}} = -\frac{2a_0 + 3}{2\Gamma_{a,ij}} + \frac{(\mathbf{A}_{ij}^{k+1})^2 + 2b_0}{2\Gamma_{a,ij}^2}. \quad (30)$$

Setting equation 30 to zero and solving for  $\Gamma_{a,ij}$  leads to:

$$\Gamma_{a,ij}^{k+1} = \frac{(\mathbf{A}_{ij}^{k+1})^2 + 2b_0}{2a_0 + 3}. \quad (31)$$

Similarly, we can update each component of  $\Gamma_b$ ,  $\Gamma_c$ , and  $\Gamma_d$  as follows:

$$\Gamma_{b,ij}^{k+1} = \frac{(\mathbf{B}_{ij}^{k+1})^2 + 2b_0}{2a_0 + 3}, \quad (32)$$

$$\Gamma_{c,ij}^{k+1} = \frac{(\mathbf{C}_{ij}^{k+1})^2 + 2b_0}{2a_0 + 3}, \quad (33)$$

$$\Gamma_{d,ij}^{k+1} = \frac{(\mathbf{D}_{ij}^{k+1})^2 + 2b_0}{2a_0 + 3}. \quad (34)$$

Based on the block gradient descent method, we derive an analytical update procedure for learning  $\Theta$ . During the optimization process of the system matrices, we use diagonal matrices to approximate  $\mathbf{R}$  and  $\mathbf{Q}$  to give a closed-form update rule for them. During the optimization process of  $\mathbf{R}$  and  $\mathbf{Q}$ , we employ the derivative rule of structured matrices to ensure their symmetry. Experimental results demonstrate that such a learning algorithm can learn LDSs with sparse system matrices accurately. Finally, Algorithm 1 summarizes the procedure for learning LDSs with sparse system matrices.

**Remark 2.** For LDSs without input signals, the proposed method can also learn LDSs with sparse system matrices from  $\{\mathbf{y}_t\}_{t=1}^T$  in an unsupervised manner by simply removing  $\mathbf{B}$ ,  $\mathbf{D}$ , and  $\mathbf{u}_t$  from the related update procedures or directly setting them to zero in the optimization process.

---

**Algorithm 1:** The proposed learning algorithm for LDSs

---

**Input:** Time-series data  $\{(\mathbf{u}_t, \mathbf{y}_t)\}_{t=1}^T$ , initial guess of  $\Theta$  and maximum number of iterations

$k_{max}$

**Output:** MAP estimate of  $\Theta$  and  $\{\mathbf{x}_t\}_{t=1}^T$

```

1 for  $k = 1, \dots, k_{max}$  do
2   // MAP estimate of  $\{\mathbf{x}_t\}_{t=1}^T$ 
3   for  $t = 1, \dots, T$  do
4     Update the mean  $\boldsymbol{\mu}_t^k$  of  $\mathbf{x}_t$  via equation 13 ;
5     Update the variance  $\boldsymbol{\Sigma}_t^k$  of  $\mathbf{x}_t$  via equation 14 ;
6     Update the covariance  $\mathbf{P}_{t,t-1}^k$  between  $\mathbf{x}_t$  and  $\mathbf{x}_{t-1}$  via equation 15 ;
7   end
8   // MAP estimate of  $\Theta$ 
9   Update system matrices  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$  via equation 22– equation 25, respectively;
10  Update noise covariance matrices  $\mathbf{R}$  and  $\mathbf{Q}$  via equation 28 and equation 29, respectively;
11  Update hyperparameter matrices  $\Gamma_a$ ,  $\Gamma_b$ ,  $\Gamma_c$ , and  $\Gamma_d$  via equation 31– equation 34,
    respectively;
12  if a stopping criterion is satisfied then
13    | Break;
14  end
15 end
```

---

## 4 SIMILARITY TRANSFORMATION OF LDSS

For LDSs, the similarity transformation is an important mathematical operation to transform them into different coordinate systems, making it easier to analyze system properties like controllability,

observability, and stability. Specifically, we can transform the state vector  $\mathbf{x}_t$  into a new state vector  $\bar{\mathbf{x}}_t$  through the relation:

$$\bar{\mathbf{x}}_t = \mathbf{P}\mathbf{x}_t, \quad (35)$$

where  $\mathbf{P} \in \mathbb{R}^{n \times n}$  is a nonsingular matrix. As such, we can derive an equivalent realization of the original LDSs as follows (see Appendix C):

$$\bar{\mathbf{x}}_t = \bar{\mathbf{A}}\bar{\mathbf{x}}_{t-1} + \bar{\mathbf{B}}\mathbf{u}_t + \bar{\boldsymbol{\varepsilon}}_t, \quad (36)$$

$$\mathbf{y}_t = \bar{\mathbf{C}}\bar{\mathbf{x}}_t + \mathbf{D}\mathbf{u}_t + \boldsymbol{\omega}_t, \quad (37)$$

where  $\bar{\mathbf{A}} = \mathbf{P}\mathbf{A}\mathbf{P}^{-1}$ ,  $\bar{\mathbf{B}} = \mathbf{P}\mathbf{B}$ ,  $\bar{\mathbf{C}} = \mathbf{C}\mathbf{P}^{-1}$ , and  $\bar{\boldsymbol{\varepsilon}}_t \sim \mathcal{N}(0, \mathbf{P}\mathbf{R}\mathbf{P}')$ . However, the similarity transformation makes it particularly difficult to accurately learn system matrices. Given the input signals  $\{\mathbf{u}_t\}_{t=1}^T$ , the transformed LDSs can produce the same output data  $\{\mathbf{y}_t\}_{t=1}^T$  as that of the original LDSs. Hence, classical learning algorithms for LDSs only learn the system matrices up to a similar transformation (Viberg, 1994). For LDSs with sparse system matrices, such a transformation changes not only the values but, more importantly, the topological structure of the system matrices, resulting in misinterpretation of intrinsic working mechanisms.

#### 4.1 BENEFIT OF SPARSE-PROMOTING PRIORS

Unlike classical learning algorithms, the proposed algorithm learns LDSs with sparse system matrices by adopting a sparsity-promoting prior to balance model complexity and modeling error. Given the sparsity constraint of system matrices, the similarity transformation cannot be applied using any arbitrary nonsingular matrix. For the LDSs with sparse system matrices following the *Occam's razor* principle, the nonsingular matrix is typically restricted to be a generalized permutation matrix; otherwise, the transformed LDSs will include redundant parameters to describe the systems. For example, if we consider the LDSs with sparse system matrices as follows:

$$\mathbf{A} = \begin{bmatrix} 0 & 0.9 \\ 0.9 & 0 \end{bmatrix}, \mathbf{B} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \mathbf{C} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \mathbf{D} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}, \quad (38)$$

it is easy to verify that such a system follows the *Occam's razor* principle because the rank of system matrices is equal to the number of nonzero components. Hence, we can derive the nonsingular matrix  $\mathbf{P}$  must satisfies

$$\mathbf{P} = \begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix} \quad \text{or} \quad \mathbf{P} = \begin{bmatrix} 0 & a \\ b & 0 \end{bmatrix}, \quad (39)$$

where  $a$  and  $b$  are arbitrary constants. As such, the transformed system matrices do not introduce additional parameters to increase model complexity.

Note that applying the similarity transformation with a generalized permutation matrix to the original state variables will scale their magnitudes and reorder them. However, it will scale the nonzero components and permute the rows or columns of system matrices accordingly. Hence, an additional advantage of the sparse-promoting prior is its ability to maximally preserve the inherent topological structure among the variables. While the learned system matrices differ from the true ones in scale, such a difference is only caused by the scaled definition of state variables. Hence, the learned LDS has the same topological structure and dynamic behavior as the real one.

## 5 EXPERIMENT

In this section, we validate the proposed algorithm on simulation and real-world datasets. In addition, we compare the proposed algorithm with classical ones mentioned previously to demonstrate its superior performance, including PEM, 4SID, and MLE. Here, we use the built-in functions **n4sid** and **pem** of Matlab to implement the 4SID and PEM algorithms, respectively. To implement MLE, we remove the sparsity-promoting priors from our derivation and revise the code accordingly. In all experiments, the dataset is split into training and testing sets with a 2:1 ratio, where 66.7% of the data is used for training and 33.3% for testing. Here, we use the mean relative error (MRE) to evaluate the performance of all the algorithms defined as follows:

$$\text{MRE} = \sum_{t=1}^T \frac{\|\mathbf{y}_t - \hat{\mathbf{y}}_t\|_2^2}{T\|\mathbf{y}_t\|_2^2}, \quad (40)$$



Table 1: Learned result of all the algorithms on the simulation system

| Method | Ours   | PEM  | 4SID  | MLE  |
|--------|--|--|---|--|
| $A$    | $\begin{bmatrix} 0 & 0.900 \\ 0.897 & 0 \end{bmatrix}$         | $\begin{bmatrix} 0 & 0.877 \\ 0.917 & -0.002 \end{bmatrix}$          | $\begin{bmatrix} 0.897 & 0.023 \\ 0.038 & -0.898 \end{bmatrix}$         | $\begin{bmatrix} 0.007 & 0.889 \\ 0.905 & -0.009 \end{bmatrix}$  |
| $B$    | $\begin{bmatrix} 4.004 & 0 \\ 0 & 4.024 \end{bmatrix}$         | $\begin{bmatrix} -14.284 & 16.858 \\ 16.979 & -14.985 \end{bmatrix}$ | $\begin{bmatrix} 0.005 & 0.004 \\ -0.001 & 0 \end{bmatrix}$             | $\begin{bmatrix} 2.887 & 1.235 \\ 1.308 & 2.848 \end{bmatrix}$   |
| $C$    | $\begin{bmatrix} 1.001 & 0 \\ 0 & 1.011 \end{bmatrix}$         | $\begin{bmatrix} 0.880 & 0.745 \\ 0.766 & 0.861 \end{bmatrix}$       | $\begin{bmatrix} 422.168 & -170.034 \\ 415.973 & 179.648 \end{bmatrix}$ | $\begin{bmatrix} 1.777 & -0.778 \\ -0.760 & 1.763 \end{bmatrix}$ |
| $D$    | $\begin{bmatrix} 1.473 & 0 \\ 0 & 1.480 \end{bmatrix}$         | $\begin{bmatrix} 5.521 & 0.099 \\ 0.017 & 5.550 \end{bmatrix}$       | $\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$                          | $\begin{bmatrix} 1.404 & 0.113 \\ -0.092 & 1.468 \end{bmatrix}$  |
| $R$    | $\begin{bmatrix} 1.879 & 0.986 \\ 0.986 & 1.877 \end{bmatrix}$ | $\begin{bmatrix} - & - \\ - & - \end{bmatrix}$                       | $\begin{bmatrix} - & - \\ - & - \end{bmatrix}$                          | $\begin{bmatrix} 1.543 & 1.402 \\ 1.402 & 1.546 \end{bmatrix}$   |
| $Q$    | $\begin{bmatrix} 0.473 & 0.199 \\ 0.199 & 0.481 \end{bmatrix}$ | $\begin{bmatrix} - & - \\ - & - \end{bmatrix}$                       | $\begin{bmatrix} - & - \\ - & - \end{bmatrix}$                          | $\begin{bmatrix} 0.455 & 0.191 \\ 0.191 & 0.469 \end{bmatrix}$   |
| MRE    | <b>7.35%</b>   | 16.37%   | 17.25%  | 7.38%  |

where  $\{\hat{y}_t\}_{t=1}^T$  is the sequence of data points generated by the learned systems in response to the same input signals. Experimental results illustrate that the proposed algorithm outperforms classical ones on learning LDSs with sparse system matrices.

### 5.1 A SYNTHETIC SYSTEM FOLLOWING THE OCCAM’S RAZOR PRINCIPLE

First, we consider a synthetic system to facilitate the comparison between the proposed algorithm and classical ones as follows:

$$\mathbf{x}_t = \begin{bmatrix} 0 & 0.9 \\ 0.9 & 0 \end{bmatrix} \mathbf{x}_{t-1} + \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{u}_t + \boldsymbol{\varepsilon}_t, \quad \boldsymbol{\varepsilon}_t \sim \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.49 & 0.25 \\ 0.25 & 0.49 \end{bmatrix} \right), \quad (41)$$

$$\mathbf{y}_t = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x}_t + \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \mathbf{u}_t + \mathbf{w}_t, \quad \mathbf{w}_t \sim \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.49 & 0.25 \\ 0.25 & 0.49 \end{bmatrix} \right). \quad (42)$$

To generate data points, the initial values of  $\mathbf{x}_0$  are drawn from the Gaussian distribution with mean  $[1, 1]'$  and an identity matrix as the covariance, and the input signal  $\mathbf{u}_t$  is drawn from the uniform distribution on  $[0, 2]$ . As for algorithm implementation, we collect 2000 data points and set the initial value of  $A$ ,  $B$ ,  $C$ ,  $D$ ,  $R$ , and  $Q$  both to be an identity matrix. The learned parameters less than the threshold 0.001 are removed from the result. In Table 1, we give the learned system matrices of all the algorithms and the corresponding MRE.

Because **n4sid** and **pem** consider the innovation representation of the LDSs (Qin, 2006), they focus on learning the innovation covariance matrix instead of  $R$  and  $Q$ , which are thus omitted here. Obviously, the learned system matrices of classical algorithms are completely different with the original ones, making it difficult for us to understand the system. However, sparse-promoting priors will restrict the nonsingular matrix  $P$  of the similarity transformation to be a generalized permutation matrix for this system. Comparing the learned  $B$  and  $C$  of the proposed algorithm with the real ones, we can derive  $P \approx 2I_2$ , which is indeed consistent with the theoretical analysis. Hence, the learned LDS of the proposed algorithm preserves the topological structure of the system, differing only in the scale of the parameters. Consequently, the learned LDS can enable us to explore the working mechanisms of the system.

### 5.2 INDUSTRIAL PROCESS SYSTEMS

Next, we validate the proposed algorithm on the real-world datasets obtained from the Database for the Identification of Systems, which are standard datasets used for learning LDSs (Zhu et al., 1994; Martens, 2010).

Table 2: Learned result of all the algorithms on the industrial process systems

| Dataset | Industrial evaporator |        |        |        | Glass furnace |        |        |        |
|---------|-----------------------|--------|--------|--------|---------------|--------|--------|--------|
| Method  | Ours                  | PEM    | 4SID   | MLE    | Ours          | PEM    | 4SID   | MLE    |
| MRE     | <b>13.74%</b>         | 17.90% | 43.77% | 18.00% | <b>18.74%</b> | 62.47% | 24.32% | 30.27% |

**Industrial evaporation systems.** In industry, multiple-stage evaporators are widely used to reduce the water content of a product such as milk. The dataset is composed of 3-dimensional time-series with a length of 6305. The inputs consist of the feed flow, vapor flow to the first evaporator stage, and cooling water flow to the condenser, while the outputs include the dry matter content, flow rate, and temperature of the product.

**Glass furnaces.** The second dataset comes from the Philips glass furnace, which is used to melt raw materials into glass. The glass furnace has two burners and one ventilator. Hence, the dataset includes two heating inputs and one cooling input with a length of 1247. In addition, we collect three outputs from temperature sensors in a cross section of the furnace.

Table 2 displays the MRE between the predicted outputs of all the learned LDSs and real ones. Due to the lack of the ground truth, the learned system matrices of all the algorithms are not depicted for comparison. Note that the proposed algorithm obtains minimum MRE on both datasets, demonstrating its superiority over classical algorithms.

## 6 DISCUSSION

To learn the LDSs with sparse system matrices, we impose sparsity-promoting priors on system matrices to balance model complexity and modeling error in this paper. Following the MAP principle, we then learn system matrices by exploring the EM algorithm to maximize the loss function composed of the priors and likelihood function. During the optimization process, we use the derivative rule of structured matrices to ensure the symmetry of noise covariance matrices. In addition, we find that the sparsity-promoting prior is capable of retaining the topological structure of the LDSs, as the nonsingular matrix of the similarity transformation is typically limited to be a generalized permutation matrix. Hence, the proposed algorithm is more useful for us to explore the interacting laws of the LDSs compared to the classical ones.

There still remains some potential limitations for the proposed algorithm. First, it cannot determine the order  $n$  of the system from data directly. While we can compare the performance of the learned systems across different orders to select the best one, such a method is quite exhaustive. The other limitation is that the similarity transformation may shrink many parameters to very small values, potentially leading to numerical errors. However, we believe that the proposed algorithm sheds a light on the learning of LDSs with sparse system matrices. In our future work, we hope to explore how to exactly learn LDSs with additional constraints.

## REFERENCES

- Mohamed Rasheed-Hilmy Abdalmoaty and Håkan Hjalmarsson. Linear prediction error methods for stochastic nonlinear models. *Automatica*, 105:49–63, 2019.
- Ainesh Bakshi, Allen Liu, Ankur Moitra, and Morris Yau. Tensor decompositions meet control theory: learning general mixtures of linear dynamical systems. In *International Conference on Machine Learning*, pp. 1549–1563. PMLR, 2023.
- David Belanger and Sham Kakade. A linear dynamical system model for text. In *International Conference on Machine Learning*, pp. 833–842. PMLR, 2015.
- Yanxi Chen and H Vincent Poor. Learning mixtures of linear dynamical systems. In *International Conference on Machine Learning*, pp. 3507–3557. PMLR, 2022.
- Yonathan Efroni, Sham Kakade, Akshay Krishnamurthy, and Cyril Zhang. Sparsity in partially controllable linear systems. In *International Conference on Machine Learning*, pp. 5851–5860. PMLR, 2022.

- 540 Wouter Favoreel, Bart De Moor, and Peter Van Overschee. Subspace state space system identifica-  
541 tion for industrial processes. *Journal of Process Control*, 10(2-3):149–155, 2000.
- 542 Zoubin Ghahramani and Geoffrey E Hinton. Parameter estimation for linear dynamical systems.  
543 1996.
- 544 Stuart Gibson and Brett Ninness. Robust maximum-likelihood estimation of multivariable dynamic  
545 systems. *Automatica*, 41(10):1667–1682, 2005.
- 546 Elad Hazan, Singh Karan, and Zhang Cyril. Learning linear dynamical systems via spectral filtering.  
547 In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- 548 Xin He, Yasen Wang, and Junyang Jin. Bayesian inference and optimisation of stochastic dynamical  
549 networks. *International Journal of Systems Science*, pp. 1–15, 2024.
- 550 Wenbing Huang, Lele Cao, Fuchun Sun, Deli Zhao, Huaping Liu, and Shanshan Yu. Learning  
551 stable linear dynamical systems with the weighted least square method. In *International Joint*  
552 *Conferences on Artificial Intelligence*, pp. 1599–1605, 2016.
- 553 Junyang Jin, Ye Yuan, and Jorge Gonalves. A full Bayesian approach to sparse network inference  
554 using heterogeneous datasets. *IEEE Transactions on Automatic Control*, 66(7):3282–3288, 2020a.
- 555 Junyang Jin, Ye Yuan, and Jorge Gonalves. High precision variational Bayesian inference of sparse  
556 linear networks. *Automatica*, 118:109017, 2020b.
- 557 Tohru Katayama et al. *Subspace methods for system identification*, volume 1. Springer, 2005.
- 558 Wallace E Larimore. Canonical variate analysis in identification, filtering, and adaptive control. In  
559 *IEEE Conference on Decision and Control*, pp. 596–604. IEEE, 1990.
- 560 Jian Li and Yunfeng Zhang. Prediction error method-based second-order structural identification  
561 algorithm in stochastic state space formulation. *Earthquake Engineering & Structural Dynamics*,  
562 35(6):761–779, 2006.
- 563 Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*, volume 793. John  
564 Wiley & Sons, 2019.
- 565 Lennart Ljung. Prediction error estimation methods. *Circuits, Systems and Signal Processing*, 21:  
566 11–21, 2002.
- 567 Bethany Lusch, J Nathan Kutz, and Steven L Brunton. Deep learning for universal linear embeddings  
568 of nonlinear dynamics. *Nature Communications*, 9(1):4950, 2018.
- 569 Giorgos Mamakoukas, Maria Castano, Xiaobo Tan, and Todd Murphey. Local Koopman operators  
570 for data-driven control of robotic systems. In *Robotics: Science and Systems*, 2019.
- 571 Giorgos Mamakoukas, Orest Xherija, and Todd Murphey. Memory-efficient learning of stable linear  
572 dynamical systems for prediction and control. In *Advances in Neural Information Processing*  
573 *Systems*, pp. 13527–13538, 2020.
- 574 James Martens. Learning the linear dynamical system with ASOS. In *International Conference on*  
575 *Machine Learning*, pp. 743–750. Citeseer, 2010.
- 576 Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University*  
577 *of Denmark*, 7(15):510, 2008.
- 578 G Pillonetto and Lennart Ljung. Full Bayesian identification of linear dynamic systems using stable  
579 kernels. *Proceedings of the National Academy of Sciences*, 120(18):e2218197120, 2023.
- 580 S Joe Qin. An overview of subspace identification. *Computers & Chemical Engineering*, 30:1502–  
581 1513, 2006.
- 582 Simo Särkkä and Lennart Svensson. *Bayesian filtering and smoothing*, volume 17. Cambridge  
583 University Press, 2023.

- 594 Robert H Shumway and David S Stoffer. An approach to time series smoothing and forecasting  
595 using the EM algorithm. *Journal of Time Series Analysis*, 3(4):253–264, 1982.  
596
- 597 Robert H Shumway, David S Stoffer, and David S Stoffer. *Time series analysis and its applications*,  
598 volume 3. Springer, 2000.  
599
- 600 Gavin Smith, João de Freitas, Tony Robinson, and Mahesan Niranjan. Speech modelling using  
601 subspace and EM techniques. In *Advances in Neural Information Processing Systems*, volume 12,  
602 1999.  
603
- 604 Robert Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical*  
605 *Society Series B: Statistical Methodology*, 58(1):267–288, 1996.  
606
- 607 Michael E Tipping. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine*  
608 *Learning Research*, 1:211–244, 2001.  
609
- 610 Tapas Tripura and Souvik Chakraborty. A sparse Bayesian framework for discovering interpretable  
611 nonlinear stochastic dynamical systems with Gaussian white noise. *Mechanical Systems and*  
612 *Signal Processing*, 187:109939, 2023.  
613
- 614 Jack Umenberger, Johan Wågberg, Ian R Manchester, and Thomas B Schön. Maximum likelihood  
615 identification of stable linear dynamical systems. *Automatica*, 96:280–292, 2018.  
616
- 617 Peter Van Overschee and Bart De Moor. N4sid: Subspace algorithms for the identification of com-  
618 bined deterministic-stochastic systems. *Automatica*, 30(1):75–93, 1994.  
619
- 620 Peter Van Overschee and BL De Moor. *Subspace identification for linear systems: The-*  
621 *ory—Implementation—Applications*. Springer Science & Business Media, 2012.  
622
- 623 M Verahegen and Patrick Dewilde. Subspace model identification. part i: The output-error state-  
624 space model identification class of algorithm. *International Journal of Control*, 56:1187–1210,  
625 1992.  
626
- 627 Mats Viberg. Subspace methods in system identification. *IFAC Proceedings Volumes*, 27(8):1–12,  
628 1994.  
629
- 630 Jin Wang and S Joe Qin. A new subspace identification approach based on principal component  
631 analysis. *Journal of Process Control*, 12(8):841–855, 2002.  
632
- 633 Wenlong Wang, Feifei Qi, David Wipf, Chang Cai, Tianyou Yu, Yuanqing Li, Zhuliang Yu, and  
634 Wei Wu. Sparse Bayesian learning for end-to-end EEG decoding. *IEEE Transactions on Pattern*  
635 *Analysis and Machine Intelligence*, 2023.  
636
- 637 Yasen Wang, Junlin Li, Zuogong Yue, and Ye Yuan. An iterative min-min optimization method for  
638 sparse Bayesian learning. In *International Conference on Machine Learning*, volume 235, pp.  
639 50859–50873. PMLR, 2024.  
640
- 641 Xiaofeng Yuan, Yalin Wang, Chunhua Yang, Zhiqiang Ge, Zhihuan Song, and Weihua Gui.  
642 Weighted linear dynamic system for feature representation and soft sensor application in nonlin-  
643 ear dynamic industrial processes. *IEEE Transactions on Industrial Electronics*, 65(2):1508–1517,  
644 2017.  
645
- 646 Yucai Zhu, Peter van Overschee, Ban de Moor, and Lennan Ljung. Comparison of three classes of  
647 identification methods. *IFAC Proceedings Volumes*, 27(8):169–174, 1994.

## APPENDIX

## A SPARSITY-PROMOTING PRIOR

Besides  $\mathbf{A}$ , we also impose the sparsity-promoting priors on  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$  as follows:

$$p(\mathbf{B} | \Gamma_b) = \prod_{i=1}^n \prod_{j=1}^p p(\mathbf{B}_{ij} | \Gamma_{b,ij}) = \prod_{i=1}^n \prod_{j=1}^p \frac{1}{\sqrt{2\pi\Gamma_{b,ij}}} \exp\left(-\frac{\mathbf{B}_{ij}^2}{2\Gamma_{b,ij}}\right), \quad (43)$$

$$p(\mathbf{C} | \Gamma_c) = \prod_{i=1}^m \prod_{j=1}^n p(\mathbf{C}_{ij} | \Gamma_{c,ij}) = \prod_{i=1}^m \prod_{j=1}^n \frac{1}{\sqrt{2\pi\Gamma_{c,ij}}} \exp\left(-\frac{\mathbf{C}_{ij}^2}{2\Gamma_{c,ij}}\right), \quad (44)$$

$$p(\mathbf{D} | \Gamma_d) = \prod_{i=1}^m \prod_{j=1}^p p(\mathbf{D}_{ij} | \Gamma_{d,ij}) = \prod_{i=1}^m \prod_{j=1}^p \frac{1}{\sqrt{2\pi\Gamma_{d,ij}}} \exp\left(-\frac{\mathbf{D}_{ij}^2}{2\Gamma_{d,ij}}\right), \quad (45)$$

where  $\Gamma_{b,ij}$ ,  $\Gamma_{c,ij}$ , and  $\Gamma_{d,ij}$  are the  $ij$ th component of  $\Gamma_b$ ,  $\Gamma_c$ , and  $\Gamma_d$ , respectively. To complete the hierarchy, the Inverse-Gamma distribution prior is imposed on each component of  $\Gamma_b$ ,  $\Gamma_c$ , and  $\Gamma_d$  as follows:

$$p(\Gamma_b) = \prod_{i=1}^n \prod_{j=1}^p \frac{a_0^{b_0}}{\Gamma(a_0)} \Gamma_{b,ij}^{-a_0-1} \exp\left(-\frac{b_0}{\Gamma_{b,ij}}\right), \quad (46)$$

$$p(\Gamma_c) = \prod_{i=1}^m \prod_{j=1}^n \frac{a_0^{b_0}}{\Gamma(a_0)} \Gamma_{c,ij}^{-a_0-1} \exp\left(-\frac{b_0}{\Gamma_{c,ij}}\right), \quad (47)$$

$$p(\Gamma_d) = \prod_{i=1}^m \prod_{j=1}^p \frac{a_0^{b_0}}{\Gamma(a_0)} \Gamma_{d,ij}^{-a_0-1} \exp\left(-\frac{b_0}{\Gamma_{d,ij}}\right). \quad (48)$$

## B DETAILED MATHEMATICAL DERIVATION

## B.1 DERIVATION OF EQUATION 16

Given the conditional independence between the variables, we can derive

$$\begin{aligned} & H(\Theta | \Theta^k) \\ &= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{Y}, \mathbf{X} | \Theta) p(\Theta)] \\ &= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{Y} | \mathbf{X}, \Theta) p(\mathbf{X} | \Theta) p(\Theta)] \\ &= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{Y} | \mathbf{X}, \mathbf{C}, \mathbf{D}, \mathbf{Q}) p(\mathbf{X} | \mathbf{A}, \mathbf{B}, \mathbf{R}) p(\mathbf{A}, \Gamma_a) p(\mathbf{B}, \Gamma_b) p(\mathbf{C}, \Gamma_c) p(\mathbf{D}, \Gamma_d)] \\ &= \underbrace{\mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{X} | \mathbf{A}, \mathbf{B}, \mathbf{R})]}_{H_1(\mathbf{A}, \mathbf{B}, \mathbf{R})} + \underbrace{\mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{Y} | \mathbf{X}, \mathbf{C}, \mathbf{D}, \mathbf{Q})]}_{H_2(\mathbf{C}, \mathbf{D}, \mathbf{Q})} \\ &+ \underbrace{\mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{A}, \Gamma_a) p(\mathbf{B}, \Gamma_b) p(\mathbf{C}, \Gamma_c) p(\mathbf{D}, \Gamma_d)]}_{H_3(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \Gamma_a, \Gamma_b, \Gamma_c, \Gamma_d)}. \end{aligned} \quad (49)$$

**Explicit mathematical expression of  $H_1(\mathbf{A}, \mathbf{B}, \mathbf{R})$ .** Based on equation 1 and the chain rule in probability, we can derive

$$\begin{aligned} & p(\mathbf{X} | \mathbf{A}, \mathbf{B}, \mathbf{R}) \\ &= p(\mathbf{x}_0) \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}, \mathbf{B}, \mathbf{R}) \end{aligned} \quad (50)$$

$$\propto \prod_{t=1}^T |\mathbf{R}|^{-\frac{1}{2}} \exp\left(-\frac{(\mathbf{x}_t - \mathbf{A}\mathbf{x}_{t-1} - \mathbf{B}\mathbf{u}_t)' \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{A}\mathbf{x}_{t-1} - \mathbf{B}\mathbf{u}_t)}{2}\right). \quad (51)$$

Hence,

$$\begin{aligned}
H_1(\mathbf{A}, \mathbf{B}, \mathbf{R}) &= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{X} | \mathbf{A}, \mathbf{B}, \mathbf{R})] \\
&= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} \left[ -\frac{T \log |\mathbf{R}| + \sum_{t=1}^T (\mathbf{x}_t - \mathbf{A}\mathbf{x}_{t-1} - \mathbf{B}\mathbf{u}_t)' \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{A}\mathbf{x}_{t-1} - \mathbf{B}\mathbf{u}_t)}{2} \right] \\
&= -\frac{T \log |\mathbf{R}| + \sum_{t=1}^T \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} (\mathbf{x}_t - \mathbf{A}\mathbf{x}_{t-1} - \mathbf{B}\mathbf{u}_t)' \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{A}\mathbf{x}_{t-1} - \mathbf{B}\mathbf{u}_t)}{2} \\
&= -\frac{T \log |\mathbf{R}| + \sum_{t=1}^T (\mathbf{m}_t^k - \mathbf{A}\mathbf{m}_{t-1}^k - \mathbf{B}\mathbf{u}_t)' \mathbf{R}^{-1} (\mathbf{m}_t^k - \mathbf{A}\mathbf{m}_{t-1}^k - \mathbf{B}\mathbf{u}_t)}{2} \\
&= -\frac{\sum_{t=1}^T (\text{Tr}(\mathbf{R}^{-1} \mathbf{P}_t^k) - \text{Tr}(\mathbf{R}^{-1} \mathbf{A} \mathbf{P}_{t,t-1}^k) - \text{Tr}(\mathbf{A}' \mathbf{R}^{-1} \mathbf{P}_{t,t-1}^k) + \text{Tr}(\mathbf{A}' \mathbf{R}^{-1} \mathbf{A} \mathbf{P}_{t-1}^k))}{2}. \quad (52)
\end{aligned}$$

**Explicit mathematical expression of  $H_2(\mathbf{C}, \mathbf{D}, \mathbf{Q})$ .** Based on equation 2, we can derive

$$\begin{aligned}
p(\mathbf{Y} | \mathbf{X}, \mathbf{C}, \mathbf{D}, \mathbf{Q}) &= \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{C}, \mathbf{D}) \\
&\propto \prod_{t=1}^T |\mathbf{Q}|^{-\frac{1}{2}} \exp \left( -\frac{(\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{D}\mathbf{u}_t)' \mathbf{Q}^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{D}\mathbf{u}_t)}{2} \right). \quad (53)
\end{aligned}$$

Hence,

$$\begin{aligned}
H_2(\mathbf{C}, \mathbf{D}, \mathbf{Q}) &= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{Y} | \mathbf{X}, \mathbf{C}, \mathbf{D}, \mathbf{Q})] \\
&= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} \left[ -\frac{T \log |\mathbf{Q}| + \sum_{t=1}^T (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{D}\mathbf{u}_t)' \mathbf{Q}^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{D}\mathbf{u}_t)}{2} \right] \\
&= -\frac{T \log |\mathbf{Q}| + \sum_{t=1}^T \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{D}\mathbf{u}_t)' \mathbf{Q}^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{D}\mathbf{u}_t)}{2} \\
&= -\frac{T \log |\mathbf{Q}| + \sum_{t=1}^T ((\mathbf{y}_t - \mathbf{C}\mathbf{m}_t^k - \mathbf{D}\mathbf{u}_t)' \mathbf{Q}^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{m}_t^k - \mathbf{D}\mathbf{u}_t) + \text{Tr}(\mathbf{C}' \mathbf{Q}^{-1} \mathbf{C} \mathbf{P}_t^k))}{2}. \quad (54)
\end{aligned}$$

**Explicit mathematical expression of  $H_3(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \Gamma_a, \Gamma_b, \Gamma_c, \Gamma_d)$ .** Based on the priors imposed on the system matrices and corresponding hyperparameters, we can derive:

$$\begin{aligned}
&p(\mathbf{A}, \Gamma_a) p(\mathbf{B}, \Gamma_b) p(\mathbf{C}, \Gamma_c) p(\mathbf{D}, \Gamma_d) \\
&= p(\mathbf{A} | \Gamma_a) p(\Gamma_a) p(\mathbf{B} | \Gamma_b) p(\Gamma_b) p(\mathbf{C} | \Gamma_c) p(\Gamma_c) p(\mathbf{D} | \Gamma_d) p(\Gamma_d) \\
&\propto \prod_{i=1}^n \prod_{j=1}^n \Gamma_{a,ij}^{-\frac{2a_0+3}{2}} \exp \left( -\frac{\mathbf{A}_{ij}^2 + 2b_0}{2\Gamma_{a,ij}} \right) \times \prod_{i=1}^n \prod_{j=1}^p \Gamma_{b,ij}^{-\frac{2a_0+3}{2}} \exp \left( -\frac{\mathbf{B}_{ij}^2 + 2b_0}{2\Gamma_{b,ij}} \right) \\
&\times \prod_{i=1}^m \prod_{j=1}^n \Gamma_{c,ij}^{-\frac{2a_0+3}{2}} \exp \left( -\frac{\mathbf{C}_{ij}^2 + 2b_0}{2\Gamma_{c,ij}} \right) \times \prod_{i=1}^m \prod_{j=1}^p \Gamma_{d,ij}^{-\frac{2a_0+3}{2}} \exp \left( -\frac{\mathbf{D}_{ij}^2 + 2b_0}{2\Gamma_{d,ij}} \right). \quad (55)
\end{aligned}$$

Hence, we have

$$\begin{aligned}
& H_3(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \Gamma_a, \Gamma_b, \Gamma_c, \Gamma_d) \\
&= \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \Theta^k)} [\log p(\mathbf{A}, \Gamma_a) p(\mathbf{B}, \Gamma_b) p(\mathbf{C}, \Gamma_c) p(\mathbf{D}, \Gamma_d)]. \\
&= - \sum_{i=1}^n \sum_{j=1}^n \left( \frac{(2a_0 + 3) \log |\Gamma_{a,ij}|}{2} + \frac{\mathbf{A}_{ij}^2 + 2b_0}{2\Gamma_{a,ij}} \right) \\
&\quad - \sum_{i=1}^n \sum_{j=1}^p \left( \frac{(2a_0 + 3) \log |\Gamma_{b,ij}|}{2} + \frac{\mathbf{B}_{ij}^2 + 2b_0}{2\Gamma_{b,ij}} \right) \\
&\quad - \sum_{i=1}^m \sum_{j=1}^n \left( \frac{(2a_0 + 3) \log |\Gamma_{c,ij}|}{2} + \frac{\mathbf{C}_{ij}^2 + 2b_0}{2\Gamma_{c,ij}} \right) \\
&\quad - \sum_{i=1}^m \sum_{j=1}^p \left( \frac{(2a_0 + 3) \log |\Gamma_{d,ij}|}{2} + \frac{\mathbf{D}_{ij}^2 + 2b_0}{2\Gamma_{d,ij}} \right). \tag{56}
\end{aligned}$$

## B.2 DERIVATION OF EQUATION 22

To provide an efficient closed-form update rule for  $\mathbf{A}$ , we only keep the diagonal components of  $\mathbf{R}^k$  and set the others to zero during the optimization process. As such, we have

$$\begin{aligned}
& H_1(\mathbf{A}, \mathbf{B}^k, \mathbf{R}^k) + H_3(\mathbf{A}, \mathbf{B}^k, \mathbf{C}^k, \mathbf{D}^k, \Gamma_a^k, \Gamma_b^k, \Gamma_c^k, \Gamma_d^k) \\
&= - \frac{\sum_{t=1}^T \sum_{r=1}^n (\mathbf{R}_{rr}^k)^{-1} (\mathbf{m}_{t,r}^k - \mathbf{A}_r \mathbf{m}_{t-1}^k - \mathbf{B}_r^k \mathbf{u}_t)^2}{2} - \sum_{r=1}^n \frac{\mathbf{A}_r \bar{\Gamma}_{a,r}^{kd} \mathbf{A}_r'}{2} \\
&\quad - \frac{\sum_{t=1}^T \sum_{r=1}^n \left( \text{Tr} \left( \mathbf{A}_r' (\mathbf{R}_{rr}^k)^{-1} \mathbf{A}_r \mathbf{P}_{t-1}^k \right) - 2 (\mathbf{R}_{rr}^k)^{-1} \mathbf{A}_r (\mathbf{P}_{t,t-1,r}^k)' \right)}{2} + c, \tag{57}
\end{aligned}$$

where  $c$  is the term unrelated to  $\mathbf{A}$ . Hence, we can calculate the derivative of  $H(\Theta | \Theta^k)$  with respect to  $\mathbf{A}_r$  at the  $k$ th iteration as follows:

$$\begin{aligned}
& \frac{\partial H_1(\mathbf{A}, \mathbf{B}^k, \mathbf{R}^k)}{\partial \mathbf{A}_r} + \frac{\partial H_3(\mathbf{A}, \mathbf{B}^k, \mathbf{C}^k, \mathbf{D}^k, \Gamma_a^k, \Gamma_b^k, \Gamma_c^k, \Gamma_d^k)}{\partial \mathbf{A}_r} \\
&= \sum_{t=1}^T (\mathbf{R}_{rr}^k)^{-1} (\mathbf{m}_{t,r}^k - \mathbf{A}_r \mathbf{m}_{t-1}^k - \mathbf{B}_r^k \mathbf{u}_t) (\mathbf{m}_{t-1}^k)' - \mathbf{A}_r \bar{\Gamma}_{a,r}^{kd} \\
&\quad - \sum_{t=1}^T (\mathbf{R}_{rr}^k)^{-1} (\mathbf{A}_r \mathbf{P}_{t-1}^k - \mathbf{P}_{t,t-1,r}^k) \\
&= \sum_{t=1}^T (\mathbf{R}_{rr}^k)^{-1} \left( \mathbf{P}_{t,t-1,r}^k + (\mathbf{m}_{t,r}^k - \mathbf{A}_r \mathbf{m}_{t-1}^k - \mathbf{B}_r^k \mathbf{u}_t) (\mathbf{m}_{t-1}^k)' - \mathbf{A}_r \mathbf{P}_{t-1}^k \right) - \mathbf{A}_r \bar{\Gamma}_{a,r}^{kd}. \tag{58}
\end{aligned}$$

Setting equation 58 to zero leads to

$$\begin{aligned}
& \sum_{t=1}^T (\mathbf{R}_{rr}^k)^{-1} \left( \mathbf{A}_r \mathbf{m}_{t-1}^k (\mathbf{m}_{t-1}^k)' + \mathbf{A}_r \mathbf{P}_{t-1}^k \right) + \mathbf{A}_r \bar{\Gamma}_{a,r}^{kd} \\
&= \sum_{t=1}^T (\mathbf{R}_{rr}^k)^{-1} \left( \mathbf{P}_{t,t-1,r}^k + (\mathbf{m}_{t,r}^k - \mathbf{B}_r^k \mathbf{u}_t) (\mathbf{m}_{t-1}^k)' \right), \tag{59}
\end{aligned}$$

Hence, we can update  $\mathbf{A}$  at the  $k$ th iteration as follows:

$$\begin{aligned}
\mathbf{A}_r^{k+1} &= \left( \sum_{t=1}^T ((\mathbf{m}_{t,r}^k - \mathbf{B}_r^k \mathbf{u}_t) (\mathbf{m}_{t-1}^k)' + \mathbf{P}_{t,t-1,r}^k) \right) \\
&\quad \times \left( \sum_{t=1}^T (\mathbf{P}_{t-1}^k + \mathbf{m}_{t-1}^k (\mathbf{m}_{t-1}^k)') + \mathbf{R}_{rr}^k \bar{\Gamma}_{a,r}^{kd} \right)^{-1}. \tag{60}
\end{aligned}$$

## C EQUIVALENT REALIZATION OF LDSs

Based on the transformed coordinates, we can derive

$$\bar{\mathbf{x}}_t = \mathbf{P}\mathbf{x}_t = \mathbf{P}\mathbf{A}\mathbf{x}_{t-1} + \mathbf{P}\mathbf{B}\mathbf{u}_t + \mathbf{P}\boldsymbol{\varepsilon}_t = (\mathbf{P}\mathbf{A}\mathbf{P}^{-1})\bar{\mathbf{x}}_{t-1} + (\mathbf{P}\mathbf{B})\mathbf{u}_t + \mathbf{P}\boldsymbol{\varepsilon}_t, \quad (61)$$

$$\mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{D}\mathbf{u}_t + \boldsymbol{\omega}_t = (\mathbf{C}\mathbf{P}^{-1})\bar{\mathbf{x}}_t + \mathbf{D}\mathbf{u}_t + \boldsymbol{\omega}_t. \quad (62)$$

Hence, an equivalent realization of the original LDSs is as follows:

$$\bar{\mathbf{x}}_t = \bar{\mathbf{A}}\bar{\mathbf{x}}_{t-1} + \bar{\mathbf{B}}\mathbf{u}_t + \bar{\boldsymbol{\varepsilon}}_t, \quad (63)$$

$$\mathbf{y}_t = \bar{\mathbf{C}}\bar{\mathbf{x}}_t + \mathbf{D}\mathbf{u}_t + \boldsymbol{\omega}_t, \quad (64)$$

where  $\bar{\mathbf{A}} = \mathbf{P}\mathbf{A}\mathbf{P}^{-1}$ ,  $\bar{\mathbf{B}} = \mathbf{P}\mathbf{B}$ ,  $\bar{\mathbf{C}} = \mathbf{C}\mathbf{P}^{-1}$ , and  $\bar{\boldsymbol{\varepsilon}}_t = \mathbf{P}\boldsymbol{\varepsilon}_t$ .