# System Identification for Linear Dynamics with Bilinear Observation Models: An Expectation–Maximization Approach

Diyou Liu and Mohammad Khosravi

Abstract—In this paper, we study the system identification problem for linear time-invariant dynamics with bilinear observation models. Accordingly, we consider a suitable parametric description for the system model and formulate the identification problem as estimating the parameters characterizing the mathematical representation of the system through input-output measurement data. To this end, we employ a probabilistic framework aiming to obtain the maximum likelihood estimates of the parameters. Accordingly, we propose utilizing the expectationmaximization approach to improve the tractability of the identification procedure. Through the numerical experiments, we verify the efficacy of the proposed scheme and demonstrate its performance.

## I. INTRODUCTION

System identification, as initially introduced in [1], is an active area of research [2] focusing on the theory and methods for data-driven modeling of dynamical systems [3]. Due to the ubiquitous nature of dynamical systems across various science and technology domains and considering the significance of models for prediction and control, system identification has received extensive attention [4]. As a result, numerous methodologies are developed for diverse categories of systems, spanning from linear dynamics to nonlinear dynamical systems [5–8]. The complexity of real systems poses a significant challenge in accurately capturing their behaviors, necessitating the utilization of various tools introduced in statistics and optimization theory, such as different parametric and nonparametric estimation techniques [9–14].

Parametric system identification is a common method that considers a rich class of models with specific structures, and subsequently, the parameters characterizing the model are estimated using measurement data. In this context, various classes of parametric models and a broad range of systems are considered, and different parameter estimation methods are employed. For example, the prediction error method (PEM) [11] is widely used in identification of linear time-invariant systems, which adjusts the parameters iteratively to minimize the norm of the losses, evaluated based on the difference between the model's predicted output and the actual output from the real system. For estimating the parameters, probabilistic frameworks are also commonly utilized. For instance, the Bayesian approach [15] treats parameters as random variables with probability distributions of specific forms. With sufficient data, the Bayesian approach obtains

accurate probability distributions of the parameters of interest. As an example, maximum likelihood (ML) estimation [16] is commonly used in system identification, aiming to find parameter values that maximize the likelihood function (i.e., the probability of observing the given measurement data given the parameter values). Another widely used method is maximum a posteriori (MAP) estimation [17]. The MAP estimation method is similar to the ML estimation method, but instead of maximizing the likelihood function, it finds the parameter values that maximize the posterior distribution. MAP estimation is more robust than ML estimation when informative prior is available.

In addition to linear system identification, there is a growing interest in methods for identifying nonlinear systems [18, 19]. For example, in [20], the authors proposed an MLbased algorithm for Hammerstein-Wiener models. Also, in [21], a Bayesian approach-based method is introduced to identify Wiener-Hammerstein models. Besides the study of general nonlinear systems, bilinear systems, as a special and simple case of nonlinear systems, are also widely studied due to their technical tractability and their relevance in various fields [22]. The characteristics of bilinear systems can be exploited in designing effective identification methods. For instance, in [23, 24], subspace techniques are used to identify bilinear state space systems. Additionally, in [25] and [26], a bilinear system in the observability canonical form is considered. In the former paper, the authors proposed an approach using the Kalman filter to estimate the system states and a gradient-based iterative algorithm to identify system parameters. In the latter one, the authors use the Rauch-Tung-Striebel smoother (RTS) to estimate the state variables and the Expectation-Maximization (EM) algorithm to identify parameters. These papers consider systems with bilinear dynamics and linear observation models. Inspired by the Wiener-Hammerstein models and Hammerstein-Wiener models, in the current paper, we consider bilinear systems with linear dynamics and bilinear observation models.

In this paper, we propose a scheme for identifying systems with linear dynamics and bilinear observation models, based on the RTS smoother [27] and the EM algorithm [28]. The system matrices, together with the mean and covariance of noise distributions, are treated as parameters. The scheme consists of two steps. In the first step, state estimates are computed, and a log-likelihood function of the parameters is defined. To estimate the states, we employ a Kalman filter [29], which recursively updates the state estimates

Delft Center for Systems and Control, Delft University of Technology, Delft, Netherlands. Email: {d.liu-9, mohammad.khosravi}@tudelft.nl

by combining the information from previous measurements and new measurements. There are several extensions of Kalman filter to enhance estimation accuracy[30, 31]. In this paper, the RTS smoother is considered. Rather than utilizing only the previous and current measurement data, the RTS smoother also employs estimates of future measurements to smooth the estimates, resulting in more accurate estimation results compared to using the Kalman filter alone. After obtaining these estimates, the expected value of the loglikelihood function of the parameters can be defined. In the second step, with the estimation of states, we compute an optimal parameter that maximizes the mentioned expected value based on the partial derivatives. These two steps are executed iteratively until convergence. The contributions of this paper are summarized as follows.

- This paper formulates a tractable approach to identify systems with linear dynamics and bilinear observation models. Additionally, using RTS smoother and EM algorithm, a scheme is proposed to solve the problem.
- Suitable examples are provided to demonstrate the performance of the proposed scheme. Moreover, simulations are presented to study the performance of the proposed scheme under different SNR levels.

The rest of this paper is structured as follows. In Section II, the main notations used in this paper are listed. In Section III, a tractable identification problem is formulated. In Section IV, we present the proposed scheme and derive the EMbased approach. Finally, in Section V, numerical examples are provided to verify the performance of proposed scheme.

#### II. NOTATION

In this paper, we use  $\mathbb{Z}$ ,  $\mathbb{Z}_+$ ,  $\mathbb{R}$ , and  $\mathbb{R}^{n \times m}$  to denote, the set of integers, the set of positive integers, the set of real numbers, and the set of n by m matrices with real value respectively. For a matrix  $A \in \mathbb{R}^{m \times n}$ , vec(A) is denoted as a column vector in  $\mathbb{R}^{mn}$  obtained from stacking the columns of the matrix A. Additionally, we denote  $\otimes$  as the Kronecker product. The vector 2-norm of a vector  $\mathbf{x} \in \mathbb{R}^n$  is denoted as  $\|\mathbf{x}\|$ . Finally, the conditional probability of A given B is denoted as p(A|B).

## III. IDENTIFICATION OF LINEAR DYNAMICS WITH BILINEAR OBSERVATION MODELS

Consider an *unknown* time-invariant random dynamical system S with linear dynamics and a bilinear observation model. More precisely, let the process model describing the dynamics of S be as

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{w}_k, \qquad \forall k \in \mathbb{Z}_+, \tag{1}$$

where  $\mathbf{x}_k \in \mathbb{R}^{n_x}$ ,  $\mathbf{u}_k \in \mathbb{R}^{n_u}$ , and  $\mathbf{w}_k \in \mathbb{R}^{n_x}$  are respectively the vectors of state variables, input, and process noise, at time instant  $k \in \mathbb{Z}_+$ , and,  $\mathbf{A} \in \mathbb{R}^{n_x \times n_x}$  and  $\mathbf{B} \in \mathbb{R}^{n_x \times n_u}$ are *unknown* matrices characterizing the dynamics of system. Also, let the observation model of the system have a bilinear form as

$$\mathbf{y}_{k+1} = \left(\mathbf{C}_0 + \sum_{i=1}^{n_u} \mathbf{C}_i u_{k,i}\right) \mathbf{x}_k + \mathbf{D} \mathbf{u}_k + \mathbf{v}_k, \qquad \forall k \in \mathbb{Z}_+, \ (2)$$

where  $u_{k,i}$  denotes the  $i^{\text{th}}$  entry of  $u_k$ ,  $\forall i = 1, \ldots, n_u$ ,  $k \in \mathbb{Z}_+$ ,  $y_k \in \mathbb{R}^{n_y}$  and  $v_k \in \mathbb{R}^{n_y}$  are respectively the vectors of output observations and measurement noise, and,  $C_0, C_1, \ldots, C_{n_u} \in \mathbb{R}^{n_y \times n_x}$  and  $D \in \mathbb{R}^{n_y \times n_u}$  are *unknown* matrices describing the observation model of the system. Suppose the initial state  $x_0$ , measurement noise  $(v_k)_{k \in \mathbb{Z}_+}$ , and process noise  $(w_k)_{k \in \mathbb{Z}_+}$  are mutually independent Gaussian random variables as

$$\mathbf{x}_0 \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}_0}, \mathbf{S}_{\mathbf{x}_0}), \tag{3}$$

$$\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{S}_{\mathbf{v}}), \qquad \forall k \in \mathbb{Z}_+,$$
 (4)

$$\mathbf{w}_k \sim \mathcal{N}(0, \mathbf{S}_{\mathbf{w}}), \qquad \forall k \in \mathbb{Z}_+,$$
 (5)

where  $\mu_{x_0} \in \mathbb{R}^{n_x}$  is an *unknown* vector and  $S_{x_0} \in \mathbb{R}^{n_x \times n_x}$ is an *unknown* positive definite matrix respectively denoting the mean and covariance of  $x_0$ , and,  $S_w \in \mathbb{R}^{n_x \times n_x}$  and  $S_v \in \mathbb{R}^{n_y \times n_y}$  are *unknown* positive definite matrices representing respectively the covariance of vector  $w_k$  and the covariance of vector  $v_k$ , for any  $k \in \mathbb{Z}_+$ .

Assume a measurement dataset of  $n_{\mathcal{D}} \in \mathbb{N}$  input-output pairs is given as

$$\mathcal{D} := \{ (\mathbf{u}_k, \mathbf{y}_k) \, | \, k \in [0, n_{\mathcal{D}} - 1] \}.$$
(6)

Accordingly, we present the main problem as identifying system S through estimating the unknown vector and matrices mentioned above using the dataset D. For the ease of discussion, assume D = 0 throughout this paper.

**Problem (Identification Problem for Linear Dynamics with Bilinear Observation Models).** Given the measurement set of data  $\mathcal{D}$ , estimate A, B, C<sub>0</sub>, C<sub>1</sub>, ..., C<sub>n<sub>u</sub></sub>,  $\mu_{x_0}$ , S<sub>x0</sub>, S<sub>w</sub>, and S<sub>v</sub>.

To address this problem, we employ an expectationmaximization (EM) approach [32] and obtain a tractable procedure. More details are discussed in the next section.

## IV. AN EXPECTATION-MAXIMIZATION APPROACH

This section presents an expectation-maximization (EM) algorithm for identifying linear time-invariant dynamics with bilinear observation models. The EM algorithm is a powerful iterative estimation scheme, particularly beneficial when deriving the solution of maximum likelihood estimation is computationally intractable. The EM approach has an iterative procedure where each of its iterations consists of two main steps: the expectation (E) step and the maximization (M) step. In the expectation (E) step, given measurement data and an initial estimation for the parameters of the system, a probability distribution for the trajectory of state variables is obtained, e.g., through Rauch-Tung-Striebel smoother [27]. Subsequently, in the maximization (M) step, using these state estimates and the resulting distributions, the system parameter estimates are updated by maximizing the expected loglikelihood function obtained from the distributions estimated in the E step.

### A. Rauch–Tung–Striebel smoother

The RTS smoother is widely used in states variables estimation, which adds an additional smoothing part using the estimated states distributions computed by the Kalman filter [29].

For the ease of discussion, we define vector of parameters, denoted by  $\theta$ , as

$$\theta := [\operatorname{vec}(\mathbf{A})^{\mathsf{T}}, \operatorname{vec}(\mathbf{B})^{\mathsf{T}}, \operatorname{vec}(\mathbf{C}_{0})^{\mathsf{T}}, \dots, \operatorname{vec}(\mathbf{C}_{n_{u}})^{\mathsf{T}}, \\ \operatorname{vec}(\mu_{\mathbf{x}_{0}})^{\mathsf{T}}, \operatorname{vec}(\mathbf{S}_{\mathbf{x}_{0}})^{\mathsf{T}}, \operatorname{vec}(\mathbf{S}_{\mathbf{w}})^{\mathsf{T}}, \operatorname{vec}(\mathbf{S}_{\mathbf{v}})^{\mathsf{T}}].$$

$$(7)$$

Given measurement data  $\mathcal{D}$ , we want to obtain the estimates of  $\theta$ , which is denoted as  $\hat{\theta}$ . From the Kalman filter, we can compute the estimated mean

$$\hat{\mathbf{x}}_{t|t} = \mathbb{E}[\mathbf{x}_t | \mathbf{y}_0, ..., \mathbf{y}_t, \hat{\theta}], \tag{8}$$

and, the estimated covariance

$$P_{t|t} = \mathbb{E}[(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t})(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t})^{\mathsf{T}} | \mathbf{y}_0, ..., \mathbf{y}_t, \hat{\theta}]$$
(9)

using following iterative scheme

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t (\mathbf{y}_t - \ddot{\Xi}_t \hat{\mathbf{x}}_{t|t-1}), 
\mathbf{P}_{t|t} = (\mathbb{I}_{n_{\mathbf{x}}} - \mathbf{K}_t \hat{\Xi}_t) \mathbf{P}_{t|t-1}, 
\mathbf{K}_t = \mathbf{P}_{t|t-1} \hat{\Xi}_t^{\mathsf{T}} [\hat{\Xi}_t \mathbf{P}_{t|t-1} \hat{\Xi}_t^{\mathsf{T}} + \hat{\mathbf{S}}_{\mathbf{v}}]^{-1}, \quad (10) 
\hat{\mathbf{x}}_{t+1|t} = \hat{\mathbf{A}} \hat{\mathbf{x}}_{t|t} + \hat{\mathbf{B}} \mathbf{u}_t,$$

$$\mathbf{P}_{t+1|t} = \hat{\mathbf{A}} \mathbf{P}_{t|t} \hat{\mathbf{A}}^{\mathsf{T}} + \hat{\mathbf{S}}_{\mathbf{w}},$$

where  $\hat{\Xi}_t$  is defined as

$$\hat{\Xi}_t = \hat{C}_0 + \sum_{i=1}^{n_u} \hat{C}_i u_{t,i}.$$
(11)

Define  $X := \{x_0, x_1, ..., x_{n_D}\}, Y := \{y_0, y_1, ..., y_{n_D-1}\}$ and  $U := \{u_0, u_1, ..., u_{n_D-1}\}$ . Using the RTS smoother, the estimated mean  $\hat{x}_{t|n_D} = \mathbb{E}[x_t|Y, \hat{\theta}]$  and estimated covariance  $P_{t|n_D} = \mathbb{E}[(x_t - \hat{x}_{t|n_D})(x_t - \hat{x}_{t|n_D})^{\mathsf{T}}|Y, \hat{\theta}]$  can be recursively computed as

$$\hat{\mathbf{x}}_{t|n_{\mathcal{D}}} = \hat{\mathbf{x}}_{t|t} + \mathbf{H}_{t}(\hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} - \hat{\mathbf{x}}_{t+1|t}), 
\mathbf{P}_{t|n_{\mathcal{D}}} = \mathbf{P}_{t|t} + \mathbf{H}_{t}(\mathbf{P}_{t+1|n_{\mathcal{D}}} - \mathbf{P}_{t+1|t})\mathbf{H}_{t}^{\mathsf{T}}, 
\mathbf{H}_{t} = \mathbf{P}_{t|t}\hat{\mathbf{A}}^{\mathsf{T}}\mathbf{P}_{t+1|t}^{-1},$$
(12)

where the initial conditions  $\hat{x}_{n_{\mathcal{D}}|n_{\mathcal{D}}}$  and  $P_{n_{\mathcal{D}}|n_{\mathcal{D}}}$  can be obtained from the results of Kalman filter. Moreover, we can further derive the following estimations [33]

$$\mathbb{E}[\mathbf{x}_{t}\mathbf{x}_{t}^{\mathsf{T}}|\mathbf{Y},\hat{\theta}] = \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}\hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t|n_{\mathcal{D}}}, \\
\mathbb{E}[\mathbf{x}_{t+1}\mathbf{x}_{t}^{\mathsf{T}}|\mathbf{Y},\hat{\theta}] = \hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}}\hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t+1|n_{\mathcal{D}}}\mathbf{H}_{t}^{\mathsf{T}},$$
(13)

which will be used in EM algorithm as discussed in the remainder of this section.

#### B. Parameters Estimation using EM algorithm

The EM contains two basic steps, which are discussed in the sequel. The expectation step computes a log likelihood function of parameters from the current estimates and given data. The maximization step finds the parameters that maximize the log likelihood function. In this work, we use EM algorithm with RTS smoother to estimate parameters of dynamical system S.

Denote  $\ddot{\theta}_k$  as the estimated parameters at iteration k and define

$$Q(\theta|\hat{\theta}_k) = \mathbb{E}_{p(\mathbf{X}|\mathbf{Y},\hat{\theta}_k)} \Big[ \ln p(\mathbf{X},\mathbf{Y}|\theta) \Big], \tag{14}$$

Instead of directly optimizing the log-likelihood  $\ln p(Y|\theta)$ , EM algorithm iteratively improves  $Q(\theta|\hat{\theta}_k)$  at each step [9], which guarantees an improvement in  $\ln p(Y|\theta)$  at least as much [12]. Consider the likelihood function

$$p(\mathbf{X}, \mathbf{Y}|\theta) = p(\mathbf{X}|\theta)p(\mathbf{Y}|\mathbf{X}, \theta).$$
(15)

Using Markov property,  $p(\mathbf{X}|\boldsymbol{\theta})$  and  $p(\mathbf{Y}|\mathbf{X},\boldsymbol{\theta})$  can be decomposed as

$$p(\mathbf{X}|\theta) = p(\mathbf{x}_0|\theta) \prod_{t=1}^{n_{\mathcal{D}}} p(\mathbf{x}_t|\mathbf{x}_{t-1},\theta).$$
(16)

and

$$p(\mathbf{Y}|\mathbf{X},\theta) = \prod_{t=0}^{n_{\mathcal{D}}-1} p(\mathbf{y}_t|\mathbf{X},\theta) = \prod_{t=0}^{n_{\mathcal{D}}-1} p(\mathbf{y}_t|\mathbf{x}_t,\theta), \quad (17)$$

Thus, the log-likelihood function of (15) can be derived as

$$\ln p(\mathbf{X}, \mathbf{Y}|\theta) = \ln p(\mathbf{x}_0|\theta) + \sum_{t=0}^{n_{\mathcal{D}}-1} \ln p(\mathbf{y}_t|\mathbf{x}_t, \theta) + \sum_{t=1}^{n_{\mathcal{D}}} \ln p(\mathbf{x}_t|\mathbf{x}_{t-1}, \theta).$$
(18)

From the dynamics function in (2),  $p(y_t|x_t, \theta)$  is subject to a Gaussian distribution  $\mathcal{N}(\mu_{y_t}, S_y)$ , where

$$\mu_{\mathbf{y}_t} = \Xi_t \mathbf{x}_t,$$
  
$$\Xi_t = \mathbf{C}_0 + \sum_{i=1}^{n_u} \mathbf{C}_i \mathbf{u}_{t,i}.$$
 (19)

Analogously from (1),  $p(\mathbf{x}_t|\theta)$  is subject to a Gaussian distribution  $\mathcal{N}(\mu_{\mathbf{x}_t}, \mathbf{S}_{\mathbf{w}})$ , where  $\mu_{\mathbf{x}_t} = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{B}\mathbf{u}_{t-1}, \forall t \ge 1$ . Therefore, using these Gaussian distributions,  $Q(\theta|\hat{\theta}_k)$  can be rewritten as

$$Q(\theta|\hat{\theta}_{k}) = \mathbb{E}_{p(X|Y,\hat{\theta}_{k})} \Big\{ -\frac{n_{\mathcal{D}}}{2} \ln \det S_{v} - \frac{1}{2} \ln \det S_{x_{0}} \\ -\frac{n_{\mathcal{D}}}{2} \ln \det S_{w} - \frac{1}{2} \sum_{t=0}^{n_{\mathcal{D}}-1} (y_{t} - \Xi_{t} x_{t})^{\mathsf{T}} S_{v}^{-1} (y_{t} - \Xi_{t} x_{t}) \\ -\frac{1}{2} \sum_{t=1}^{n_{\mathcal{D}}} (x_{t+1} - A x_{t} - B u_{t})^{\mathsf{T}} S_{w}^{-1} (x_{t+1} - A x_{t} - B u_{t}) \\ -\frac{1}{2} (x_{0} - \mu_{x_{0}})^{\mathsf{T}} S_{x_{0}}^{-1} (x_{0} - \mu_{x_{0}}) \Big\}.$$

$$(20)$$

For the convenience of notation, we define  $M = [A, B], C = [C_0, C_1, ..., C_{n_u}]$ , and  $z_t = \begin{bmatrix} x_t \\ u_t \end{bmatrix}$ , for all  $t = 0, ..., n_D - 1$ .

Using the property that  $\mathbf{x}^{\mathsf{T}}\mathbf{M}\mathbf{x} = \mathrm{tr}(\mathbf{M}\mathbf{x}\mathbf{x}^{\mathsf{T}}), \forall \mathbf{M} \in \mathbb{R}^{n_{x} \times n_{x}},$ we can further simplify  $Q(\theta|\hat{\theta}_{k})$  as

$$Q(\theta|\hat{\theta}_{k}) = \mathbb{E}_{p(\mathbf{X}|\mathbf{Y},\hat{\theta}_{k})} \left\{ -\frac{n_{\mathcal{D}}}{2} \ln \det \mathbf{S}_{\mathbf{v}} - \frac{1}{2} \ln \det \mathbf{S}_{\mathbf{x}_{0}} - \frac{n_{\mathcal{D}}}{2} \ln \det \mathbf{S}_{\mathbf{w}} - \frac{1}{2} \sum_{t=0}^{n_{\mathcal{D}}-1} \operatorname{tr} \{\mathbf{S}_{\mathbf{v}}^{-1}[\mathbf{y}_{t} - \mathbf{C}(\begin{bmatrix} 1\\\mathbf{u}_{t} \end{bmatrix} \otimes \mathbf{x}_{t})] \right]$$

$$[\mathbf{y}_{t} - \mathbf{C}(\begin{bmatrix} 1\\\mathbf{u}_{t} \end{bmatrix} \otimes \mathbf{x}_{t})]^{\mathsf{T}} - \frac{1}{2} \sum_{t=1}^{n_{\mathcal{D}}} \operatorname{tr} \{\mathbf{S}_{\mathbf{w}}^{-1}(\mathbf{x}_{t+1} - \operatorname{Mz}_{t}) - (\mathbf{x}_{t+1} - \operatorname{Mz}_{t})^{\mathsf{T}} \} - \frac{1}{2} \operatorname{tr} \{\mathbf{S}_{\mathbf{x}_{0}}^{-1}(\mathbf{x}_{0} - \mu_{\mathbf{x}_{0}})(\mathbf{x}_{0} - \mu_{\mathbf{x}_{0}})^{\mathsf{T}} \} \right\}.$$

$$(21)$$

To obtain the optimal  $\theta$  which maximizes  $Q(\theta|\hat{\theta}_k)$ , we take the partial derivative of  $Q(\theta|\hat{\theta}_k)$  with respect to each parameter and set it equal to zero. Based on the calculation in Appendix A, the optimal parameters can be found as

$$\begin{split} \hat{\mathbf{C}}_{k+1} &= \sum_{t=0}^{n_{\mathcal{D}}-1} \mathbf{y}_{t} (\begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix} \otimes \hat{\mathbf{x}}_{t|n_{\mathcal{D}}})^{\mathsf{T}} \\ & \left( \sum_{t=0}^{n_{\mathcal{D}}-1} \begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix} \begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix}^{\mathsf{T}} \otimes (\hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t|n_{\mathcal{D}}}) \right)^{-1}, \\ \hat{\mathbf{S}}_{\mathbf{v},k+1} &= \frac{1}{n_{\mathcal{D}}} \sum_{t=0}^{n_{\mathcal{D}}-1} \mathbf{y}_{t} \mathbf{y}_{t}^{\mathsf{T}} - \frac{1}{n_{\mathcal{D}}} \sum_{t=0}^{n_{\mathcal{D}}-1} \mathbf{y}_{t} (\begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix} \otimes \hat{\mathbf{x}}_{t|n_{\mathcal{D}}})^{\mathsf{T}} \\ & \left( \sum_{t=0}^{n_{\mathcal{D}}-1} \begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix} \begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix}^{\mathsf{T}} \otimes (\hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t|n_{\mathcal{D}}}) \right)^{-1} \\ & \sum_{t=0}^{n_{\mathcal{D}}-1} \left( \begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix} \otimes \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \right) \mathbf{y}_{t}^{\mathsf{T}}, \\ \hat{\mathbf{M}}_{k+1} &= \sum_{t=1}^{n_{\mathcal{D}}} \left( \left[ \hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t+1|n_{\mathcal{D}}} \mathbf{H}_{t}^{\mathsf{T}} - \hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} \mathbf{u}_{t}^{\mathsf{T}} \right] \\ & \left[ \begin{bmatrix} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t|n_{\mathcal{D}}} - \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \mathbf{u}_{t}^{\mathsf{T}} \right]^{-1} \\ & \frac{1}{n_{\mathcal{D}}} \sum_{t=0}^{n_{\mathcal{D}}-1} (\hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} + \mathbf{P}_{t+1|n_{\mathcal{D}}} \mathbf{H}_{t}^{\mathsf{T}} - \hat{\mathbf{x}}_{t+1|N} \mathbf{u}_{t}^{\mathsf{T}} \right] \\ & - \frac{1}{n_{\mathcal{D}}} \sum_{t=0}^{n_{\mathcal{D}}-1} \begin{bmatrix} \hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t+1|n_{\mathcal{D}}} \mathbf{H}_{t}^{\mathsf{T}} - \hat{\mathbf{x}}_{t+1|N} \mathbf{u}_{t}^{\mathsf{T}} \right] \\ & \left( \sum_{t=0}^{n_{\mathcal{D}}-1} \begin{bmatrix} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t+1|n_{\mathcal{D}}} \mathbf{H}_{t}^{\mathsf{T}} - \hat{\mathbf{x}}_{t+1|N} \mathbf{u}_{t}^{\mathsf{T}} \right] \right)^{-1} \\ & \sum_{t=0}^{n_{\mathcal{D}}-1} \begin{bmatrix} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t+1|n_{\mathcal{D}}} \mathbf{H}_{t}^{\mathsf{T}} - \hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} \mathbf{u}_{t}^{\mathsf{T}} \right] \right)^{-1} \\ & \hat{\mathbf{x}}_{t,n_{\mathcal{D}}, \mathbf{x}_{t,n_{\mathcal{D}}}, \\ & \hat{\mathbf{x}}_{n,k+1} = \hat{\mathbf{x}}_{0|n_{\mathcal{D}}, \\ & \hat{\mathbf{x}}_{n,k+1} = \mathbf{P}_{0|n_{\mathcal{D}}. \end{aligned} \right$$

**Remark 1.** Equation (22) implies there is one optimal parameter  $\theta$  which let the partial derivatives to be zero. Intuitively, if we only consider  $S_w$  and M and assume  $n_x = 1$ ,  $Q(\theta|\hat{\theta}_k)$  can be simplified as  $\mathbb{E}_{p(X|Y,\hat{\theta}_k)} - \frac{n_{\mathcal{D}}}{2} \ln S_w - \frac{1}{2} \sum_{t=1}^{n_{\mathcal{D}}} \frac{(x_{t+1} - Mz_t)^2}{S_w}$ . As  $f(x) = a \ln x + \frac{b}{x}$ ,  $\forall a, b \in \mathbb{R}^-$  only has one maximum Algorithm 1 An EM Approach for Identification of Linear Dynamics with Bilinear Observation Models

Dynamics with Diffiear Observation Models	
	Input: D.
	<b>Output:</b> <i>θ</i> .
1:	Initial guess: $\hat{\theta}_0$
2:	$k \leftarrow 0$
3:	while 1 do
4:	Current parameters estimates: $\hat{\theta} \leftarrow \hat{\theta}_k$
5:	for $t \leftarrow 0$ to $n_{\mathcal{D}}$ do
6:	Kalman Filter: (10)
7:	for $t \leftarrow n_{\mathcal{D}}$ to 0 do
8:	RTS smoother: (12)
9:	EM approach: compute (22) to find a new parameters
	estimates. $\hat{\theta}_{k+1} \leftarrow \arg \max_{\theta} Q(\theta   \hat{\theta}_k)$ .
10:	if $\ \hat{ heta}_{k+1} - \hat{ heta}_k\  < \epsilon$ then
11:	break
12:	else
13:	$k \leftarrow k+1$
14:	$ heta \leftarrow \hat{ heta}_k$

at  $x = \frac{b}{a}$ , there is only one optimal  $S_w$ . In addition, for any given  $S_w$ , there is also only one optimal M. For other parameters, using the similar method, there is one local maximum which is also global maximum for  $Q(\theta|\hat{\theta}_k)$ .

Algorithm 1 summarizes the introduced EM based approach for identification of linear dynamics with bilinear observation models.

## V. NUMERICAL EXPERIMENTS

In this section, we evaluate the performance of Algorithm 1 through numerical examples. Accordingly, we consider a time-invariant system as introduced in (1)-(2) with matrices A, B,  $C_0$  and  $C_1$  as

$$A = \begin{bmatrix} 0.85 & -0.4 \\ 0.35 & 0.65 \end{bmatrix}, \qquad B = \begin{bmatrix} 0.5 \\ -0.5 \end{bmatrix}, \qquad (23)$$
$$C_0 = \begin{bmatrix} 0.5 & -0.2 \end{bmatrix}, \qquad C_1 = \begin{bmatrix} 0.15 & 0.1 \end{bmatrix}.$$

Additionally, for process and measurement noise covariance, and prior of the initial state, we assume that

$$S_{\rm w} = \begin{bmatrix} 0.05^2 & 0\\ 0 & 0.05^2 \end{bmatrix}, \quad S_{\rm x_0} = \begin{bmatrix} 0.05^2 & 0\\ 0 & 0.05^2 \end{bmatrix}, \quad (24)$$
$$S_{\rm v} = 0.05^2, \qquad \qquad \mu_{\rm x_0} = \begin{bmatrix} 1\\ 1 \end{bmatrix}.$$

A random binary signal is generated as the control input U for the system, with a length of  $n_D = 1000$ . Following the problem settings introduced in Section III, the process noise and measurement noise are generated from Gaussian distributions  $\mathcal{N}(\mathbf{0}, S_w)$  and  $\mathcal{N}(\mathbf{0}, S_v)$ , respectively. In this experiment, we set a stopping condition as  $\|\hat{\theta}_{k+1} - \hat{\theta}_k\| < \epsilon$ , where  $\epsilon = 10^{-4}$ .

Figure 1 illustrates the performance of Algorithm 1. In this figure, we plot the trajectories of relative errors  $\|\hat{C} - C\| / \|C\|$  and  $\|\hat{M} - M\| / \|M\|$ , where  $C = \begin{bmatrix} C_0 & C_1 \end{bmatrix}$ and  $M = \begin{bmatrix} A & B \end{bmatrix}$  as defined in (21). It can be seen that,



Figure 1. The relative error of system matrices estimates with respect to iteration steps.



Figure 2. Comparison of real system outputs and identified system outputs.

initially, the relative error is large due to the gap between the initial guess and the true parameters. However, the relative error diminishes with each iteration. The convergence rate is fast in the initial iterations, however, it gets slower later. In Figure 1, for the sake of clarity, we have shown the first 200 iterations. Nonetheless, for  $\epsilon = 10^{-4}$ , the algorithm stops after about 650 EM steps. Figure 2 shows the performance of the proposed algorithm. To this end, another random input sequence with a length T = 100 is generated to validate the identified system. It can be seen that the predicted trajectory of outputs for the identified system is close to that of the real system. The relatively error, computed as  $\sum_{t=0}^{T} \frac{\|y_t - \hat{y}_t\|}{\|y_t\|}$ , is approximately 0.02.

**Remark 2.** As shown in Figure 1, the convergence speed for EM algorithm is relatively slow. As discussed after (14), with each iteration step, improving  $Q(\theta|\hat{\theta}_k)$  will improve the log-likelihood  $\ln p(Y|\theta)$ . However, the rate of improvement can be slow. Furthermore, there is no assurance that the loglikelihood will reach the global optimum. More precisely, the EM algorithm may converge to a local optimum. Consequently, in simulations, it may be necessary to restart the algorithm with different initial guesses to obtain accurate



Figure 3. Relative error of system matrices under four different SNRs.

parameter estimates.

To further investigate the performance of Algorithm 1, we analyze it under four different SNR levels and perform a Monte Carlo experiment with respect to each of these levels. Given the same system (23), prior knowledge (24) and the same input sequences, 100 different noise sequences are generated under four different SNR levels: 5 dB, 10 dB, 15 dB, and 20 dB. For each realization, the proposed algorithm is used to identify the system. Figure 3 demonstrates the performance for estimating system matrices. As shown in Figure 3, in general, the mean relative errors of both C and M are lower for scenarios with higher SNRs, and consequently, the accuracy of identified parameters is higher. When SNR = 5 dB, the estimated system matrices are the worst. In this case, the performance depends considerably on the initial guess. Indeed, for some realizations, the relative errors are low; meanwhile, for other initial guesses, the relative errors may be higher, resulting in higher covariance. Moreover, for SNR = 15 dB or 20 dB, the covariance of error is relatively small, indicating that the algorithm is less dependent on the initial guess.

### VI. CONCLUSION

In this paper, we have studied the identification problem for unknown dynamical systems with linear dynamics and bilinear observation models, and proposed a tractable identification procedure based on EM procedure. In the proposed scheme, the RTS smoother is used to find the estimates of the states, while the EM-based approach is introduced to estimate the unknown system parameters. The performance of the introduced scheme is evaluated through a numerical example. Furthermore, we also compare the performance of the proposed scheme under different SNRs through a Monte Carlo numerical experiment.

#### APPENDIX

To obtain the optimal solution, we use first order necessary condition. Accordingly, we compute the partial derivatives with respect to each parameter. Additionally, from Kalman filter and RTS smoother in (12) and (13), we can obtain  $\mathbb{E}_{p(X|Y,\hat{\theta}_k)}[x_t|Y,\hat{\theta}_k]$ ,  $\mathbb{E}_{p(X|Y,\hat{\theta}_k)}[x_tx_t^{\mathsf{T}}|Y\hat{\theta}_k]$ , and  $\mathbb{E}_{p(X|Y,\hat{\theta}_k)}[x_{t+1}x_t^{\mathsf{T}}|Y,\hat{\theta}_k]$ . Thus, the partial derivatives are

$$\begin{split} \frac{\partial Q(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}_{k})}{\partial \mathbf{C}} = & \mathbf{S}_{\mathbf{v}}^{-1} \Big\{ \sum_{t=0}^{n_{\mathcal{D}}-1} \mathbf{y}_{t} (\begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix} \otimes \hat{\mathbf{x}}_{t|n_{\mathcal{D}}})^{\mathsf{T}} - \mathbf{C} \\ & \Big( \sum_{t=0}^{n_{\mathcal{D}}-1} (\begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix} \begin{bmatrix} 1\\ \mathbf{u}_{t} \end{bmatrix}^{\mathsf{T}}) \otimes (\hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t|n_{\mathcal{D}}}) \Big) \Big\}, \end{split}$$

$$\begin{aligned} \frac{\partial Q(\theta|\theta_k)}{\partial \mathbf{S}_{\mathbf{v}}} = & \frac{1}{2} \sum_{t=0}^{n_{\mathcal{D}}-1} \left\{ -n_{\mathcal{D}} \mathbf{S}_{\mathbf{v}}^{-1} + \mathbf{S}_{\mathbf{v}}^{-1} \left[ \mathbf{y}_t \mathbf{y}_t^\mathsf{T} - \mathbf{C}(\begin{bmatrix} 1\\ \mathbf{u}_t \end{bmatrix} \otimes \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}) \mathbf{y}_t^\mathsf{T} - \mathbf{y}_t (\begin{bmatrix} 1\\ \mathbf{u}_t \end{bmatrix} \otimes \hat{\mathbf{x}}_{t|n_{\mathcal{D}}})^\mathsf{T} \mathbf{C}^\mathsf{T} + \mathbf{C} \Big( \sum_{t=0}^{n_{\mathcal{D}}-1} \begin{bmatrix} 1\\ \mathbf{u}_t \end{bmatrix} \begin{bmatrix} 1\\ \mathbf{u}_t \end{bmatrix}^\mathsf{T} \otimes (\hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^\mathsf{T} + \mathbf{P}_{t|n_{\mathcal{D}}}) \Big) \end{aligned}$$

$$= \begin{bmatrix} a_{t} \end{bmatrix} \begin{bmatrix} a_{t} \end{bmatrix}$$

$$C^{\mathsf{T}} S_{\mathsf{v}}^{-1} \Big\}, [$$

$$\begin{split} \frac{\partial Q(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}_{k})}{\partial \mathbf{M}} = & \mathbf{S}_{\mathbf{w}}^{-1} \Big[ \sum_{t=1}^{n_{\mathcal{D}}} \Big[ \hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t+1|n_{\mathcal{D}}} \mathbf{H}_{t}^{\mathsf{T}} \quad \hat{\mathbf{x}}_{t+1|n_{\mathcal{D}}} \mathbf{u}_{t}^{\mathsf{T}} \Big] \\ & - \mathbf{M} \begin{bmatrix} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t|n_{\mathcal{D}}} \quad \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \mathbf{u}_{t}^{\mathsf{T}} \\ \mathbf{u}_{t} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} \quad \mathbf{u}_{t} \mathbf{u}_{t}^{\mathsf{T}} \end{bmatrix} \Big], \end{split}$$

$$\frac{\partial Q(\theta|\hat{\theta}_k)}{\partial S_w} = \frac{1}{2} \sum_{t=1}^{n_{\mathcal{D}}} \left\{ -n_{\mathcal{D}} S_w^{-1} + S_w^{-1} \left( \left( \hat{x}_{t+1|n_{\mathcal{D}}} \hat{x}_{t+1|n_{\mathcal{D}}}^{\mathsf{T}} + P_{t+1|n_{\mathcal{D}}} \right) - \left[ \hat{x}_{t+1|n_{\mathcal{D}}} \hat{x}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + P_{t+1|n_{\mathcal{D}}} H_t^{\mathsf{T}} - \hat{x}_{t+1|n_{\mathcal{D}}} u_t^{\mathsf{T}} \right] M^{\mathsf{T}} - M \left[ \hat{x}_{t+1|n_{\mathcal{D}}} \hat{x}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + P_{t+1|n_{\mathcal{D}}} H_t^{\mathsf{T}} - \hat{x}_{t+1|N} u_t^{\mathsf{T}} \right]^{\mathsf{T}}$$

$$M \begin{bmatrix} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \hat{\mathbf{x}}_{t|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{t|n_{\mathcal{D}}} & \hat{\mathbf{x}}_{t|n_{\mathcal{D}}} \mathbf{u}_{t}^{\mathsf{T}} \end{bmatrix} \mathbf{M}^{\mathsf{T}} \mathbf{S}_{\mathsf{w}}^{-1},$$
$$\frac{\partial Q(\theta|\hat{\theta}_{k})}{\partial \mathbf{x}_{k}} = \mathbf{S}^{-1} \{ (\hat{\mathbf{x}}_{0|n_{\mathcal{D}}} - u_{\mathbf{x}_{k}}) \boldsymbol{\mu}^{\mathsf{T}} \},$$

$$\frac{\partial \mu_{\mathbf{x}_{0}}}{\partial \mathbf{S}_{\mathbf{x}_{0}}} = \frac{1}{2} \{ -\mathbf{S}_{\mathbf{x}_{0}}^{-1} + \mathbf{S}_{\mathbf{x}_{0}}^{-1} (\hat{\mathbf{x}}_{0|n_{\mathcal{D}}} \hat{\mathbf{x}}_{0|n_{\mathcal{D}}}^{\mathsf{T}} + \mathbf{P}_{0|n_{\mathcal{D}}} \\ - \hat{\mathbf{x}}_{0|n_{\mathcal{D}}} \mu_{\mathbf{x}_{0}}^{\mathsf{T}} - \mu_{\mathbf{x}_{0}} \hat{\mathbf{x}}_{0|n_{\mathcal{D}}}^{\mathsf{T}} + \mu_{\mathbf{x}_{0}} \mu_{\mathbf{x}_{0}}^{\mathsf{T}} ) \mathbf{S}_{\mathbf{x}_{0}}^{-1} \}.$$

We set the partial derivatives to be zero, and obtain (22).

#### REFERENCES

- L. Zadeh, "On the identification problem," *IRE Transactions on Circuit Theory*, vol. 3, no. 4, pp. 277–281, 1956.
- [2] L. Ljung, "Perspectives on system identification," Annual Reviews in Control, vol. 34, no. 1, pp. 1–12, 2010.
- [3] D. G. Luenberger, Dynamic Systems. J. Wiley Sons, 1979.
- [4] L. Ljung, *System identification: Theory for the user*. Prentice Hall, 1999.
- [5] J. Schoukens and L. Ljung, "Nonlinear system identification: A useroriented road map," *IEEE Control Systems Magazine*, vol. 39, no. 6, pp. 28–99, 2019.
- [6] M. Khosravi, "Representer theorem for learning Koopman operators," IEEE Transactions on Automatic Control, 2023.
- [7] M. Khosravi and R. S. Smith, "Kernel-based identification with frequency domain side-information," *Automatica*, vol. 150, p. 110813, 2023.
- [8] D. Liu and M. Khosravi, "Learning stable evolutionary PDE dynamics: A scalable system identification approach," in *IEEE Conference on Control Technology and Applications*, 2024, pp. 79–84.

- [9] R. J. Little and D. B. Rubin, Statistical analysis with missing data. John Wiley & Sons, 2019, vol. 793.
- [10] M. Khosravi and R. S. Smith, "The existence and uniqueness of solutions for kernel-based system identification," *Automatica*, vol. 148, p. 110728, 2023.
- [11] K. Åstrom, "Maximum likelihood and prediction error methods," *IFAC Proceedings Volumes*, vol. 12, no. 8, pp. 551–574, 1979.
- [12] P. Brémaud, An introduction to probabilistic modeling. Springer Science & Business Media, 2012.
- [13] M. Khosravi and R. S. Smith, "Kernel-based impulse response identification with side-information on steady-state gain," *IEEE Transactions* on Automatic Control, 2023.
- [14] M. Zorzi, "Nonparametric identification of kronecker networks," Automatica, vol. 145, p. 110518, 2022.
- [15] J. M. Bernardo and A. F. Smith, *Bayesian theory*. John Wiley & Sons, 2009, vol. 405.
- [16] R. J. Rossi, Mathematical statistics: an introduction to likelihood based inference. John Wiley & Sons, 2018.
- [17] J.-L. Gauvain and C.-H. Lee, "Maximum a posteriori estimation for multivariate gaussian mixture observations of markov chains," *IEEE transactions on speech and audio processing*, vol. 2, no. 2, pp. 291– 298, 1994.
- [18] M. Khosravi and R. S. Smith, "Nonlinear system identification with prior knowledge on the region of attraction," *IEEE Control Systems Letters*, vol. 5, no. 3, pp. 1091–1096, 2021.
  - [19] M. Khosravi and R. S. Smith, "Convex nonparametric formulation for identification of gradient flows," *IEEE Control Systems Letters*, vol. 5, no. 3, pp. 1097–1102, 2021.
  - [20] A. Wills, T. B. Schön, L. Ljung, and B. Ninness, "Identification of hammerstein-wiener models," *Automatica*, vol. 49, no. 1, pp. 70–81, 2013.
  - [21] Q. Liu, X. Tang, J. Li, J. Zeng, K. Zhang, and Y. Chai, "Identification of wiener–hammerstein models based on variational bayesian approach in the presence of process noise," *Journal of the Franklin Institute*, vol. 358, no. 10, pp. 5623–5638, 2021.
  - [22] P. M. Pardalos and V. A. Yatsenko, *Optimization and control of bilinear systems: theory, algorithms, and applications.* Springer Science & Business Media, 2010, vol. 11.
  - [23] V. Verdult and M. Verhaegen, "Identification of multivariable bilinear state space systems based on subspace techniques and separable least squares optimization," *International Journal of Control*, vol. 74, no. 18, pp. 1824–1836, 2001.
  - [24] V. Verdult and M. Verhaegen, "Kernel methods for subspace identification of multivariable lpv and bilinear systems," *Automatica*, vol. 41, no. 9, pp. 1557–1565, 2005.
  - [25] S. Liu, F. Ding, and T. Hayat, "Moving data window gradient-based iterative algorithm of combined parameter and state estimation for bilinear systems," *International Journal of Robust and Nonlinear Control*, vol. 30, no. 6, pp. 2413–2429, 2020.
  - [26] S. Liu, X. Zhang, L. Xu, and F. Ding, "Expectation-maximization algorithm for bilinear systems by using the rauch-tung-striebel smoother," *Automatica*, vol. 142, p. 110365, 2022.
  - [27] H. E. Rauch, F. Tung, and C. T. Striebel, "Maximum likelihood estimates of linear dynamic systems," *AIAA journal*, vol. 3, no. 8, pp. 1445–1450, 1965.
  - [28] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the royal statistical society: series B (methodological)*, vol. 39, no. 1, pp. 1– 22, 1977.
  - [29] R. E. Kalman, "A new approach to linear filtering and prediction problems," 1960.
  - [30] S. J. Julier and J. K. Uhlmann, "Unscented filtering and nonlinear estimation," *Proceedings of the IEEE*, vol. 92, no. 3, pp. 401–422, 2004.
  - [31] S. J. Julier, J. K. Uhlmann, and H. F. Durrant-Whyte, "A new approach for filtering nonlinear systems," in *Proceedings of 1995 American Control Conference-ACC'95*, vol. 3. IEEE, 1995, pp. 1628–1632.
  - [32] S. Theodoridis, Pattern recognition. Academic press, 2006.
  - [33] S. Särkkä and L. Svensson, *Bayesian filtering and smoothing*. Cambridge university press, 2023, vol. 17.