

Linear Time-Varying Parameter Estimation: Maximum A Posteriori Approach via Semidefinite Programming

Sasan Vakili, Mohammad Khosravi, Peyman Mohajerin Esfahani and Manuel Mazo Jr.

Abstract—We study the problem of identifying a linear time-varying output map from measurements and linear time-varying system states, which are perturbed with Gaussian observation noise and process uncertainty, respectively. Employing a stochastic model as *prior* knowledge for the parameters of the unknown output map, we reconstruct their estimates from input/output pairs via a Bayesian approach to optimize the *posterior* probability density of the output map parameters. The resulting problem is a non-convex optimization, for which we propose a tractable linear matrix inequalities approximation to warm-start a first-order subsequent method. The efficacy of our algorithm is shown experimentally against classical Expectation Maximization and Dual Kalman Smoother approaches.

I. INTRODUCTION

Bayesian approaches for estimating characteristics of dynamical systems have been a subject of studies for decades and have recently received extensive attention [1], [2]. In systems theory, the significance of the Bayesian approach is highlighted in state estimation of dynamical systems [3], [4], e.g., through the celebrated recursive Kalman *filter*. The Rauch-Tung-Striebel (RTS) Smoother counterparts [4], on the other hand, are (offline) iterative non-causal algorithms incorporating future measurements into the current state estimation.

An alternative to Bayesian estimation, which requires a prior distribution of the parameters of interest, is the min-max estimation approach, assuming instead the knowledge of ambiguity sets. The *least favourable* uncertainty model from this ambiguity set is then used for estimation [5]–[8]. Here, we focus instead on designing a classical smoother for a different problem: system parameters estimation from input/output measurements via Bayesian estimation. This problem arises in, e.g., robot mapping in unknown environments, such as Autonomous Underwater Vehicles operating in the deep sea where global positioning is expensive due to low visibility and lack of radio communications.

Given unknown parameters with random states, applying a Bayesian estimation framework leads to severe non-convexities in the resulting optimization problem. As such, iterative schemes are typically employed to overcome these non-convexities. Assuming the parameters also follow a statistical formulation, two main types of smoother approaches,

This work was supported by the European Union’s Horizon 2020 Research and Innovation Programme through the Marie Skłodowska-Curi Grant under Agreement 956200 and the ERC Starting Grant TRUST-949796. The second author is partially supported by the Swiss National Science Foundation [Post-doc.Mobility Grant 211104].

The authors are with the Delft Center for Systems and Control, Delft University of Technology, Delft, The Netherlands. S.Vakili@tudelft.nl, Mohammad.Khosravi@tudelft.nl, P.MohajerinEsfahani@tudelft.nl, M.Mazo@tudelft.nl

Dual Kalman Smoother (DKS) and *Expectation Maximization* (EM), are available in the literature [9].

Dual Kalman Smoothers (and filters) attempt to maximize the joint probability space of parameters and state (conditioned on input and output observations), iterating between estimating the system states using the last parameters’ estimates followed by estimating the parameters from the currently estimated states. Although DKS is computationally efficient according to its recursive structure, its estimation performance can be significantly suboptimal due to the bilinearity between the parameters and states. Unlike DKS, *Expectation Maximization* learns the parameters of statistical models by maximizing the posterior distribution of the parameters from the observed data and their prior density function when incomplete data or hidden variables exist [10]–[12].

Considering the states of a dynamical system as hidden variables [13]–[15], EM estimates the parameters of a dynamical system in two steps by integrating all possible values of the states in which the model could have generated the observations. The distribution over hidden variables is maximized in the *E-Step* using the parameters estimates from the previous iteration. Subsequently, the *M-Step* maximizes a lower bound of the original cost by fixing the hidden variables distribution to the one optimized in the *E-Step*. A closed-form solution of the *M-Step* is provided in [14] for estimating the parameters of linear time-invariant dynamical systems and in [9, Chapter 6] for estimating the parameters of a Gaussian radial basis function (RBF) approximator. Both solutions consider the maximum likelihood case, where no prior exists for the parameters.

Finding a closed-form expression for the parameters update in the *M-Step* of a Maximum A Posteriori (MAP) smoothing problem when the parameters are time-varying and in the presence of *a priori* knowledge is non-trivial. This challenge leads to a slow convergence of the EM algorithm utilizing computationally demanding approaches to solve the optimization in *M-step*, e.g., first-order methods. The slow convergence of EM is shown experimentally in [16], and further analyses in [17], [18] demonstrate the slow convergence rate of the gradient variant of the EM algorithm for Gaussian Mixture Models.

Alternatives to the iterative schemes can be found in the parameter estimation problem of an elliptically contoured distribution [19, Page 107], employing recent Conic Geometric Optimization methods [20]. These methods, however, require reformulating the MAP problem via techniques, such as those proposed in [21, Section 3], which result in losing

the output map's original structure. Retaining such a structure to leverage the available *a priori* knowledge is essential to our problem.

In this work, we consider systems with known linear time-varying dynamics affected by process and measurement Gaussian noise but with unknown time-varying output maps. We propose a method to estimate the unknown parameters of the output map having *a priori* information a linear stochastic system encoding the evolution of the parameters. We derive an optimization problem applying a fully Bayesian approach, maximizing the exact posterior distribution of the parameters when unfolded over the whole time horizon. A tractable conservative approximation to the resulting optimization problem is derived via semidefinite programming (SDP) using linear matrix inequalities (LMIs) techniques. The solution from this approximation then provides a warm-start for a first-order quasi-Newton algorithm that enjoys a locally superlinear convergence rate. This combination allows us to enjoy both the computational advantage of DKS and outperform the statistical performance of EM. We illustrate the efficacy and performance of our proposed method in comparison with DKS and EM through a Monte Carlo experiment with different signal-to-noise ratios (SNRs) in Section V.

Notation: Throughout this paper, \mathbb{Z}_+ , \mathbb{R} and $\mathbb{R}^{n \times m}$ denote the set of positive integers, real numbers, and n by m real matrices, respectively. We indicate $\text{diag}(A_1, \dots, A_k)$ as a block diagonal matrix with diagonal entries of given matrices A_1, \dots, A_k . The symbol \mathbb{I} denotes the identity matrix, and tr is the trace operator. Given $A \in \mathbb{R}^{m \times n}$, a matrix with columns $a_1, \dots, a_n \in \mathbb{R}^m$, we define $\text{vec}(A)$ as the vector $[a_1^\top, \dots, a_n^\top]^\top \in \mathbb{R}^{mn}$. For a positive symmetric matrix $A \in \mathbb{R}^{n \times n}$, $\Lambda(A) := (\lambda_i(A))_{i=1}^n$ indicates the vector of eigenvalues of A in descending order, i.e., $\lambda_i(A)$ is the i^{th} largest eigenvalue of A . A multivariate normal (Gaussian) distribution with mean μ and covariance matrix Σ is denoted by $\mathcal{N}(\mu, \Sigma)$, and the symbol \sim stands for "distributed according to".

II. PROBLEM DEFINITION

Consider a discrete-time linear time-varying dynamical system described by the process model:

$$x_{k+1} = A_k x_k + B_k u_k + w_k, \quad k \in \mathbb{Z}_+, \quad (1)$$

where k denotes the time index, $x_k \in \mathbb{R}^{n_x}$ is the vector of state variables, $A_k \in \mathbb{R}^{n_x \times n_x}$ is the state transition matrix, $u_k \in \mathbb{R}^{n_u}$ is the vector of inputs, $B_k \in \mathbb{R}^{n_x \times n_u}$ is the input matrix, and $w_k \in \mathbb{R}^{n_x}$ is an independent realization at time k of the process noise with Gaussian distribution $\mathcal{N}(0, \Sigma_{w_k})$. The initial state of system (1), denoted by x_0 , is also assumed to be drawn from a Gaussian distribution $\mathcal{N}(\mu_{x_0}, \Sigma_{x_0})$. For $k \in \mathbb{Z}_+$, the state of the system is observed at time instant k through a perturbed linear time-varying map:

$$y_k = C_k x_k + v_k, \quad k \in \mathbb{Z}_+, \quad (2)$$

where $y_k \in \mathbb{R}^{n_y}$ denotes the output measurements, $C_k \in \mathbb{R}^{n_y \times n_x}$ is an *unknown time-varying* observation matrix, and

$v_k \in \mathbb{R}^{n_y}$ is the vector of measurement noise signals with Gaussian distribution $\mathcal{N}(0, \Sigma_{v_k})$. Let θ_k be the vector of all parameters at each time index k :

$$\theta_k := \text{vec}(C_k^\top), \quad (3)$$

which implies that C_k and θ_k uniquely characterize each other. We introduce the following assumption, providing a form of *a priori* information. This plays a role akin to that of a regularizer in non-Bayesian techniques, such as in *Supervised Learning*, where algorithms without such regularizers are prone to overfitting.

Assumption 1 (Output map dynamics). *The dynamics of the output map are governed by the difference equation*

$$\theta_{k+1} = \theta_k + \eta_k, \quad k \in \mathbb{Z}_+, \quad (4)$$

where k denotes the time index, $\theta_k \in \mathbb{R}^{n_y n_x}$ is the vector of parameters driven by the vector of process noise $\eta_k \in \mathbb{R}^{n_y n_x}$ with Gaussian distribution $\mathcal{N}(\mu_{\eta_k}, \Sigma_{\eta_k})$. Further, assume that the initial parameter of system (4), denoted by θ_0 , is drawn from the normal distribution $\mathcal{N}(\mu_{\theta_0}, \Sigma_{\theta_0})$.

Assumption 1 imposes a Gaussian random walk dynamics on the evolution of the parameters, which is the minimal structure and assumption on the variations of the parameters because of the maximum entropy feature of the Gaussian distributions. This assumption allows us to employ a stochastic belief of a deterministic reality in the *Bayesian* viewpoint.

Let the inputs and outputs of system (1)-(2) be measured for $k = 0, \dots, n_T$, where $(n_T + 1) \in \mathbb{Z}_+$ denotes the length of the measurement data. More precisely, the input-output trajectory data is given by $\mathcal{D} = \{(u_k, y_k) \mid k = 0, \dots, n_T\}$. Additionally, we assume:

Assumption 2 (Noise). *The process, measurement, and output map noise realizations, i.e., w_k , v_k and η_k , respectively, for all $k \in \{0, \dots, n_T\}$, are independent. Furthermore, the means μ_{x_0} , μ_{θ_0} , μ_{η_k} , and covariance matrices Σ_{x_0} , Σ_{w_k} , Σ_{v_k} , Σ_{θ_0} and Σ_{η_k} , for $k \in \{0, \dots, n_T\}$, are known.*

Remark 1 (*A priori* knowledge). While we assume μ_{θ_0} , μ_{η_k} , Σ_{θ_0} , and Σ_{η_k} to be readily known, in practical applications, these parameters can be obtained through various means depending on the context, e.g., employing prior knowledge of the nominal model, empirically from previous experiments' data, or employing a suitable hyperparameter estimation method when $\mu_{\eta_k} = \mu_{\theta_0}$ and $\Sigma_{\eta_k} = \Sigma_{\theta_0}$, for $k \in \mathbb{Z}_+$.

Ultimately, the question is whether the observation model (2) could be estimated. More precisely, we would like to address the following problem:

Problem 1. *Given the process and observation models (1) and (2), input-output measurement data \mathcal{D} , and under Assumptions 1 and 2, estimate the unknown time-varying observation matrices C_k in an efficient and tractable way.*

To address the problem 1, we develop a MAP approach in the next section, followed by a tractable reformulation using LMI techniques in Section IV.

III. MAXIMUM A POSTERIORI ESTIMATION

In this section, we propose a Bayesian method for estimating the unknown observation matrices C_0, \dots, C_{n_τ} . The three main elements in Bayesian estimation methods are a prior density function, an observation model, and a loss function, which we briefly explain for solving our problem with the MAP approach.

A. Lifted Process and Observation Model

Let us represent the process model (1) in the following *lifted matrix form*:

$$\mathbf{x} = \mathbf{A}(\mathbf{u} + \mathbf{w}_x), \quad (5)$$

where $\mathbf{x} = [\mathbf{x}_0^\top, \dots, \mathbf{x}_{n_\tau}^\top]^\top$ includes the system states over the entire horizon up to time n_τ , while the input vector is modified to include the initial state $\mathbf{u} = [\mu_{x_0}^\top, (\mathbf{B}_0 \mathbf{u}_0)^\top, \dots, (\mathbf{B}_{n_\tau-1} \mathbf{u}_{n_\tau-1})^\top]^\top$, and the noise vector consists of the uncertainty of the initial state and process noises $\mathbf{w}_x = [\mathbf{w}_{x_0}^\top, \mathbf{w}_0^\top, \dots, \mathbf{w}_{n_\tau-1}^\top]^\top$ with $\mathbf{w}_{x_0} \sim \mathcal{N}(0, \Sigma_{\mathbf{w}_x})$. Given that the process noises and initial state uncertainty are uncorrelated from Assumption 2, we can specify \mathbf{w}_x in terms of a multivariate normal distribution $\mathcal{N}(0, \Sigma_{\mathbf{w}_x})$ in which $\Sigma_{\mathbf{w}_x} = \text{diag}(\Sigma_{x_0}, \Sigma_{w_0}, \dots, \Sigma_{w_{n_\tau-1}})$. The *lifted* transition matrix \mathbf{A} has the lower triangular form:

$$\mathbf{A} = \begin{bmatrix} \mathbb{I} & & & & & & & \\ & \mathbf{A}_0 & & & & & & \\ & \mathbf{A}_1 \mathbf{A}_0 & & \mathbb{I} & & & & \\ & \vdots & & \vdots & & \mathbb{I} & & \\ & \mathbf{A}_{n_\tau-1} \dots \mathbf{A}_0 & & \mathbf{A}_{n_\tau-1} \dots \mathbf{A}_1 & & \dots & \mathbf{A}_{n_\tau-1} & \mathbb{I} \end{bmatrix}.$$

Similarly, the observation model (2) for the entire trajectory can be expressed as:

$$\mathbf{y} = \mathbf{C}\mathbf{x} + \mathbf{v}, \quad (6)$$

where $\mathbf{y} = [\mathbf{y}_0^\top, \dots, \mathbf{y}_{n_\tau}^\top]^\top$ is the vector of all measurements, $\mathbf{v} \sim \mathcal{N}(0, \Sigma_{\mathbf{v}})$ is the vector of all measurement noise realizations with $\Sigma_{\mathbf{v}} = \text{diag}(\Sigma_{v_0}, \Sigma_{v_1}, \dots, \Sigma_{v_{n_\tau}})$, and \mathbf{C} is the lifted observation matrix:

$$\mathbf{C} = \text{diag}(C_0, C_1, \dots, C_{n_\tau}).$$

Finally, we describe the dynamics of the output map parameters θ for the entire trajectory as:

$$\theta = \mu_\theta + \mathbf{w}_\theta, \quad (7)$$

where $\theta = [\theta_0^\top, \dots, \theta_{n_\tau}^\top]^\top$, and μ_θ results from the summation of the biases of the initial parameter and the noise: $\mu_\theta = [\mu_{\theta_0}^\top, \mu_{\theta_0}^\top + \mu_{\eta_0}^\top, \dots, \mu_{\theta_0}^\top + \sum_{i=0}^{n_\tau-1} \mu_{\eta_i}^\top]^\top$. Similarly, the noise signal \mathbf{w}_θ results from the integration over the entire horizon including the uncertainty of the initial parameter as: $\mathbf{w}_\theta = \mathbf{D}\eta$, where $\eta = [\eta_{\theta_0}^\top, \eta_0^\top, \dots, \eta_{n_\tau-1}^\top]^\top$ and $\eta_{\theta_0} \sim \mathcal{N}(0, \Sigma_{\theta_0})$, with

$$\mathbf{D} = \begin{bmatrix} \mathbb{I} & & & & \\ \mathbb{I} & \mathbb{I} & & & \\ \vdots & \vdots & \ddots & & \\ \mathbb{I} & \mathbb{I} & \dots & \mathbb{I} \end{bmatrix}.$$

Since the parameters are assumed to be independent (Assumption 2), we have $\mathbf{w}_\theta \sim \mathcal{N}(0, \Sigma_{\mathbf{w}_\theta})$, $\Sigma_{\mathbf{w}_\theta} = \mathbf{D}\Sigma_\eta \mathbf{D}^\top$, with $\Sigma_\eta = \text{diag}(\Sigma_{\theta_0}, \Sigma_{\eta_0}, \dots, \Sigma_{\eta_{n_\tau-1}})$. Ultimately, the model (7) is used to specify the prior density function.

In what follows, we represent \mathbf{C} as $\mathbf{C}(\theta)$ to emphasize the dependence of \mathbf{C} on θ according to (3). Consequently, substituting \mathbf{x} in (6) with the expression from (5) results in the observation model, with unknown $\mathbf{C}(\theta)$, describing the measurements \mathbf{y} as a function of the applied inputs \mathbf{u} :

$$\mathbf{y} = \mathbf{C}(\theta)\mathbf{A}\mathbf{u} + \mathbf{w}_y(\theta), \quad (8)$$

where $\mathbf{w}_y(\theta) = \mathbf{C}(\theta)\mathbf{A}\mathbf{w}_x + \mathbf{v}$. Also, from Assumption 2:

$$\mathbf{w}_y(\theta)|\theta \sim \mathcal{N}(0, \Sigma_{\mathbf{w}_y}(\theta)),$$

where $\Sigma_{\mathbf{w}_y}(\theta) = \mathbf{C}(\theta)\mathbf{A}\Sigma_{\mathbf{w}_x}\mathbf{A}^\top\mathbf{C}(\theta)^\top + \Sigma_{\mathbf{v}}$. This model is used later to specify the conditional probability density function of the measurements. Note that $\mathbf{w}_y(\theta)|\theta$ remains Gaussian with the derived covariance since both noise sources, \mathbf{w}_x and \mathbf{v} , are Gaussian and independent.

B. MAP Loss Function

In MAP estimation, one aims to find an estimate $\hat{\theta}$ for the parameters by minimizing the cost function [22]: $\mathbb{E}[1 - \mathbb{1}_{\theta: \|\theta - \hat{\theta}\|_\infty \leq \epsilon}(\theta)]$, where θ is the vector of random variables, $\mathbb{1}(\cdot)$ is an indicator function, and ϵ is a small scalar. Minimizing the expectation of such a loss function implies maximizing the conditional probability density of θ given the vector of observations and inputs, $\hat{\theta} = \underset{\theta}{\text{argmax}} p(\theta|\mathbf{y}, \mathbf{u})$,

where $\hat{\theta}$ is the estimate of the true parameter θ [23, Chapter 4]. The following lemma formalizes this first step to compute the MAP estimation.

Lemma 1 (MAP optimization problem). *Let us define the function $\mathcal{J} : \mathbb{R}^{n_y n_x(n_\tau+1)} \rightarrow \mathbb{R}$ as*

$$\mathcal{J}(\theta) := \log \det(\Sigma_{\mathbf{w}_y}(\theta)) + \|\mathbf{y} - \mathbf{C}(\theta)\mathbf{A}\mathbf{u}\|_{\Sigma_{\mathbf{w}_y}^{-1}(\theta)}^2 + \|\theta - \mu_\theta\|_{\Sigma_{\theta}^{-1}}^2. \quad (9)$$

The MAP estimation is equivalent to

$$\hat{\theta} = \underset{\theta}{\text{argmin}} \mathcal{J}(\theta). \quad (10)$$

Proof. Using Bayes' rule, the MAP estimation can be reformulated as

$$\max_{\theta} p(\theta|\mathbf{y}, \mathbf{u}) = \max_{\theta} \frac{p(\mathbf{y}|\theta, \mathbf{u})p(\theta|\mathbf{u})}{p(\mathbf{y}|\mathbf{u})}. \quad (11)$$

We first note that the denominator of (11) does not depend on θ and, hence, can be neglected without changing the optimizer. Moreover, we note that the dynamics of θ in (4) (or equivalently in the lifted form in (7)) do not depend on the input sequence of \mathbf{u} (i.e., $p(\theta|\mathbf{u}) = p(\theta)$). Next, using a straightforward computation, one can derive the probability density functions $p(\theta)$ and $p(\mathbf{y}|\theta, \mathbf{u})$. Specifically, from (7), we know that the variable θ is Gaussian with the probability density function

$$p(\theta) = \frac{\exp\left(-\frac{1}{2}(\theta - \mu_\theta)^\top \Sigma_{\theta}^{-1}(\theta - \mu_\theta)\right)}{\sqrt{(2\pi)^{(n_\tau+1)n_y n_x} \det(\Sigma_{\theta_0})}}.$$

Similarly, we know from (8) that given θ and the input sequence u , the output sequence y is also Gaussian with the conditional probability density function

$$p(y|\theta, u) = \frac{\exp\left(-\frac{1}{2}(y - C(\theta)Au)^\top \Sigma_{w_y}^{-1}(\theta)(y - C(\theta)Au)\right)}{\sqrt{(2\pi)^{(n_\tau+1)n_y} \det(\Sigma_{w_y}(\theta))}}.$$

Finally, applying the monotonically increasing function \log and observing that all terms in the denominators except $\det(\Sigma_{w_y}(\theta))$ are constant, we arrive at the minimization problem of the function \mathcal{J} defined in (10). ■

In the next section, we propose a tractable conservative approximation using LMI techniques to tackle the non-convex objective function $\mathcal{J}(\theta)$ defined in (9).

Remark 2 (Robust estimation). Alternatively, a robust minimax estimation formulation similar to [8] could be employed. This approach, however, requires finding an ambiguity set to approximate non-Gaussian observation uncertainties due to the multiplication of Gaussian variables in $w_y(\theta)$.

IV. PROPOSED SOLUTION

The optimization problem (10) is non-convex not only because of the weight $\Sigma_{w_y}^{-1}(\theta)$, being quadratic in the parameters θ , in the second term but also because of the log-determinant operator in the first term. A typical approach is to use first-order algorithms to find a solution due to the mentioned non-convexities. These algorithms, however, only guarantee convergence to a local optimum. Therefore, selecting an appropriate initial starting point is crucial to the obtained quality of the solution. We propose to solve the problem in two steps: first, we perform a convex relaxation of (10) into a set of LMIs, which we use to compute an initial approximate minimizer; next, we employ this approximate minimizer to initialize (warm-start) a first-order optimization method, e.g., steepest descent [24] or quasi-Newton algorithms [25], to solve (10) thus refining our initial minimizer estimate.

Theorem 2 (LMI conservative approximation). *Consider the following LMIs:*

$$\begin{aligned} \min_{S, \theta, \gamma, \beta} \quad & \text{tr}(S - \mathbb{I}) + \gamma + \beta \\ \text{s.t.} \quad & \begin{bmatrix} -\Sigma_{w_x}^{-1} & A^\top C(\theta)^\top \\ C(\theta)A & \Sigma_v - S \end{bmatrix} \preceq 0, \\ & \begin{bmatrix} -S & (y - C(\theta)Au) \\ (y - C(\theta)Au)^\top & -\gamma \end{bmatrix} \preceq 0, \\ & \begin{bmatrix} -\Sigma_{w_\theta} & (\theta - \mu_\theta) \\ (\theta - \mu_\theta)^\top & -\beta \end{bmatrix} \preceq 0. \end{aligned} \quad (12)$$

Then, the optimal value of the nonlinear program (10) is upper bounded by $J^* + \|y - C(\theta^*)Au\|_{\Sigma_{w_y}^{-1}(\theta^*) - S^{*-1}}^2$, where J^* and (S^*, θ^*) are the optimal value and the optimizer of (12), respectively.

Proof. Consider a matrix $S \succ 0$ upper bounding the covariance matrix $\Sigma_{w_y}(\theta) \succ 0$ as

$$\Sigma_{w_y}(\theta) = C(\theta)A\Sigma_{w_x}A^\top C(\theta)^\top + \Sigma_v \preceq S. \quad (13)$$

Thus, $\lambda_i(\Sigma_{w_y}(\theta)) \leq \lambda_i(S)$, for $i = 1, \dots, n_y(n_\tau + 1)$, which implies that

$$\log \det(\Sigma_{w_y}(\theta)) \leq \log \det(S).$$

Since $\log \det(S) = \sum_{i=1}^{n_y(n_\tau+1)} \log \lambda_i(S)$ and $\text{tr}(S) = \sum_{i=1}^{n_y(n_\tau+1)} \lambda_i(S)$, we also have

$$\log \det(\Sigma_{w_y}(\theta)) \leq \log \det(S) \leq \text{tr}(S - \mathbb{I}).$$

Using the *Schur complement*, one can see that (13) is equivalent to the following linear matrix inequality

$$\begin{bmatrix} -\Sigma_{w_x}^{-1} & A^\top C(\theta)^\top \\ C(\theta)A & \Sigma_v - S \end{bmatrix} \preceq 0.$$

Similarly, considering $\gamma \geq 0$ and $\beta \geq 0$ such that

$$\begin{aligned} (y - C(\theta)Au)^\top S^{-1}(y - C(\theta)Au) &\leq \\ (y - C(\theta)Au)^\top \Sigma_{w_y}(\theta)^{-1}(y - C(\theta)Au) &\leq \gamma, \end{aligned} \quad (14)$$

and

$$(\theta - \mu_\theta)^\top \Sigma_{w_\theta}^{-1}(\theta - \mu_\theta) \leq \beta, \quad (15)$$

we can apply again the *Schur complement* to the inequalities in (14) and (15) to obtain the last two LMIs in (12). Finally, replacing the terms in the cost function $\mathcal{J}(\theta)$ in (9) with their bounds and including the corresponding LMIs as constraints arrives at the LMIs (12). Note further that by definition, we have

$$\begin{aligned} J^* + \|y - C(\theta^*)Au\|_{\Sigma_{w_y}^{-1}(\theta^*) - S^{*-1}}^2 = \\ \mathcal{J}(\theta^*) + \text{tr}(S^* - \mathbb{I}) - \log \det(\Sigma_{w_y}(\theta^*)) \geq \mathcal{J}(\theta^*), \end{aligned} \quad (16)$$

where the function $\mathcal{J}(\theta^*)$ is defined in (9), and the last inequality follows from (13). ■

The tightness of the inequality in (16) mainly depends on the gap between $\log \det(S)$ and $\text{tr}(S - \mathbb{I})$ since $\log \det(S)$ is bounded from above by $\text{tr}(S - \mathbb{I})$, which is negligible when $\lambda_i(S) \approx 1$, for $i = 1, \dots, n_y(n_\tau + 1)$. One may employ a suitable matrix W to scale the eigenvalues of S , replace $\log \det(S)$ with $\log \det(WSW) - 2 \log \det(W)$ and approximate $\log \det(WSW)$ with $\text{tr}(WSW - \mathbb{I})$. Furthermore, the closeness of J^* and $\mathcal{J}(\theta^*)$ in (16) is proportional to the fitness quality of the measurements and whether S^* is close to the covariance matrix accordingly.

In addition, Theorem 2 provides an approximation of (10), producing an initial near-optimal solution. As indicated earlier, we propose to employ this solution to warm-start a local (non-convex) optimizer. Due to its fast convergence, we propose to employ, as a refining optimizer, the BFGS algorithm [25, Chapter 6], a variant of quasi-Newton methods. The BFGS algorithm approximates the Hessian matrix for its search directions relying on an analytical expression of the gradient $\nabla_\theta \mathcal{J}(\theta)$. The gradient of the cost function (9) with respect to the parameters θ

$$\nabla_\theta \mathcal{J}(\theta) = \left[\frac{\partial \mathcal{J}}{\partial \theta_1}, \dots, \frac{\partial \mathcal{J}}{\partial \theta_{n_y n_x (n_\tau + 1)}} \right]^\top \quad (17)$$

can be easily derived applying the chain rule:

$$\begin{aligned} \frac{\partial \mathcal{J}}{\partial \theta_{ijk}} &= 2\text{tr} \left[\left(A \Sigma_{w_x} A^\top C(\theta)^\top \Sigma_{w_y}(\theta)^{-1} \right) C_k^{ij}(\theta) \right. \\ &\quad - \left(A \Sigma_{w_x} A^\top C(\theta)^\top \Sigma_{w_y}(\theta)^{-1} \times \right. \\ &\quad \quad \left. (y - C(\theta)Au)(y - C(\theta)Au)^\top \Sigma_{w_y}(\theta)^{-1} \right) C_k^{ij}(\theta) \\ &\quad - \left(Au(y - C(\theta)Au)^\top \Sigma_{w_y}(\theta)^{-1} \right) C_k^{ij}(\theta) \\ &\quad \left. + \left((\theta - \mu_\theta)^\top \Sigma_{w_\theta}^{-1} \right) \theta^{ijk} \right], \end{aligned}$$

where $C_k^{ij}(\theta)$ is the single-entry matrix of $C(\theta)$ with the block matrix of $C_k(\theta)$ having 1 at index (i, j) and zero elsewhere, and θ^{ijk} is the single-entry vector of θ with 1 at index ijk and zero elsewhere.

The LMIs (12) initializes the original non-convex problem with a locally optimal solution. Thus, the computational complexity of the proposed method consists of the well-known computational complexity of solving the SDP problems [26], i.e., a one-time solution of (12), and the computation of the gradient (17) per iteration of the first-order method:

$$\mathcal{O} \left(\frac{n_x^3(n_\tau + 1)^3 + n_x^2(n_\tau + 1)^2}{2} + n_y n_x^2(n_\tau + 1)^2 + n_x n_y^2(n_\tau + 1)^3 + (n_y)^3(n_\tau + 1)^3 \right),$$

which is $\mathcal{O}(n_\tau^3)$ when $n_x, n_y \ll n_\tau$.

V. NUMERICAL EXAMPLE

In this section, we provide a numerical example to verify the efficacy and performance of the proposed method: employing the LMIs (12) to warm-start the solution of (10) via the BFGS optimizer. Additionally, we compare the resulting solution with the estimates obtained from EM and DKS algorithms. To have a fair comparison, we also employ the same BFGS optimizer for the *M-Step* of the EM estimation.

We demonstrate our solution on the following system:

$$x_{k+1} = \begin{bmatrix} 0.7 & 0.25 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0.25 & 0.7 \end{bmatrix} x_k + \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} (3.5 + \cos(2k)) + w_k,$$

with $\mu_{x_0} = [1, 0.5, 2]^\top$. The observation is a two-dimensional model, i.e., the number of measurements per time instant is $n_y = 2$. The system has sampled input and measurement pairs in \mathcal{D} every 100 milliseconds for 10 seconds, i.e., $n_\tau = 100$. Thus, the number of parameters to be estimated is $n_y n_x (n_\tau + 1) = 606$. The noise covariance of process, observation and output map dynamics, Σ_{w_k} , Σ_{v_k} and Σ_{η_k} , are assumed to remain constant across the entire horizon. Moreover, the initial state and parameter noise covariance are similar to the noise covariance of process and output map dynamics, respectively (i.e., $\Sigma_{x_0} = \Sigma_{w_k}$ and $\Sigma_{\theta_0} = \Sigma_{\eta_k}$). The output map noise biases μ_{η_k} are generated such that $\mu_{1,\eta_k} = 5 + e^{-0.6k} \cos(0.4k)$, $\mu_{2,\eta_k} = 1.5 + e^{-0.6k} \sin(0.025k)$, $\mu_{3,\eta_k} = 2$, $\mu_{4,\eta_k} = 5 + e^{-0.6k} \cos(0.4k)$, $\mu_{5,\eta_k} = 1.5 + e^{-0.6k} \sin(0.025k)$, and $\mu_{6,\eta_k} = 2$. The initial parameter bias, μ_{θ_0} , is derived from μ_{η_0} by setting $k = 0$,

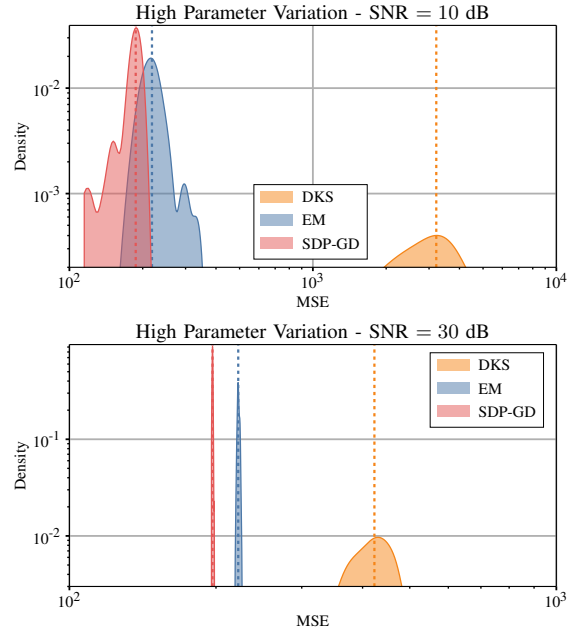


Fig. 1: The Mean Squared Error of the three methods in high noise of Σ_η and two different SNRs for 100 experiments.

and the DKS and EM algorithms' initialization is according to these noise bias values. We examine the performance of our algorithm, SDP-GD, compared with EM and DKS on four scenarios generated by employing High/Low SNRs for the process and observation noise, particularly 30 and 10 decibels (dBs), and *High/Low* parameter variation of:

$$\text{High: } \Sigma_{\eta_k} = \text{diag}(2.17, 0.076, 1.19, 1.38, 0.87, 1.27)$$

$$\text{Low: } \Sigma_{\eta_k} = \text{diag}(6.9, 0.2, 3.8, 4.4, 2.8, 4) \cdot 10^{-2}.$$

Combined results from 100 experiments for each of the four scenarios, keeping the same ground-truth realization in each scenario, are illustrated in Figures 1 and 2. The figures demonstrate the median (vertical dotted lines) and distribution across experiments based on the mean squared error (MSE), i.e. $\text{MSE} = \frac{1}{n_\tau + 1} \sum_{k=0}^{n_\tau} \|\theta_k - \hat{\theta}_k\|_2^2$. One can observe how the DKS underperforms compared to the EM and our SDP-GD solutions in more challenging scenarios where the process-observation model noise is high or in the presence of *High* parameter variation.

The average and standard deviation of the computation time of each method across 100 experiments are reported in Table I. We performed all the experiments on a cluster node with 384G memory and 40 CPU cores (2 Intel Xeon Gold 6148 @ 2.40GHz). The elapsed execution times confirm our hypothesis that EM is computationally more expensive than the other alternatives.

The performance of EM and DKS algorithms highly depends on the initialization, while in contrast, our proposed solution takes advantage of a warm-start initializer obtained from solving a convexified approximation of the original optimization problem. This initialization helps converge to a better local optimum faster than the EM algorithm. In addition, the *M-step* of the EM algorithm, in this problem,

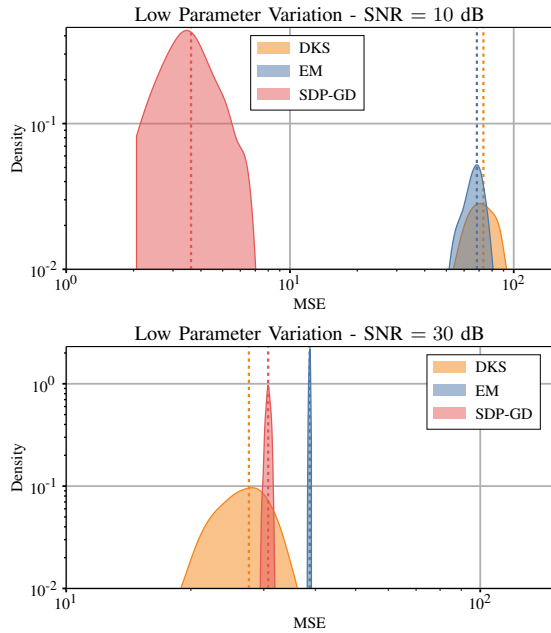


Fig. 2: The Mean Squared Error of the three methods in low noise of Σ_{η} and two different SNRs for 100 experiments.

Elapsed time per seconds (mean \pm std)				
	Experiment Scenario	DKS	EM	SDP-GD
10 dB	Low Parameter Variation	18 \pm 4	14265 \pm 3112	1265 \pm 109
	High Parameter Variation	27 \pm 13	22547 \pm 6768	1542 \pm 178
30 dB	Low Parameter Variation	9 \pm 3	7999 \pm 802	1543 \pm 132
	High Parameter Variation	21 \pm 13	3796 \pm 412	1605 \pm 128

TABLE I: The Average computation performance on all scenarios for 100 experiments.

does not hold a closed-form solution, which results in utilizing a first-order method. This gradient M -Step also plays a part in the general slowness of the EM algorithm. Our algorithm, however, requires a one-time execution of the set of LMIs followed by an iterative quasi-Newton method with a superlinear convergence rate. Hence, it provides the best of both worlds, i.e., better estimations than EM and DKS with less computation time than EM.

VI. CONCLUSION

We have introduced a method for estimating an unknown output map of a linear time-varying system by employing a stochastic characterization of the evolution of the output map parameters, which serves as *a priori* information for MAP estimation. The derived MAP optimization problem is solved by relaxing the optimization as a set of LMIs, whose solution provides a warm-start for a gradient descent algorithm. Compared with standard approaches to solve this problem, namely EM and DKS, we showed experimentally the superiority of our method in estimation performance and lower computational demands compared to EM. Future work will explore the incorporation of other types of *a priori* knowledge on the output map, the development of efficient causal filters following similar approaches, the minimax formulation for robust estimation, considering noise models

with more general structures, and introducing methods for efficient design of the control input.

REFERENCES

- [1] T. D. Barfoot, *State Estimation for Robotics*. Cambridge University Press, 2017.
- [2] R. van de Schoot, S. Depaoli, R. King, B. Kramer, K. Märtens, M. G. Tadesse, M. Vannucci, A. Gelman, D. Veen, J. Willemsen, and C. Yau, "Bayesian statistics and modelling," *Nature Reviews Methods Primers*, vol. 1, no. 1, p. 1, 2021.
- [3] R. E. Kalman and R. S. Bucy, "New results in linear filtering and prediction theory," *Journal of Basic Engineering*, vol. 83, no. 1, pp. 95–108, 1961.
- [4] H. Rauch, "Solutions to the linear smoothing problem," *IEEE Transactions on Automatic Control*, vol. 8, no. 4, pp. 371–372, 1963.
- [5] S. Yi and M. Zorzi, "Robust kalman filtering under model uncertainty: The case of degenerate densities," *IEEE Transactions on Automatic Control*, vol. 67, no. 7, pp. 3458–3471, 2021.
- [6] S. Shafieezadeh Abadeh, V. A. Nguyen, D. Kuhn, and P. M. Mohajerin Esfahani, "Wasserstein distributionally robust kalman filtering," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [7] V. A. Nguyen, S. Shafieezadeh-Abadeh, D. Kuhn, and P. Mohajerin Esfahani, "Bridging bayesian and minimax mean square error estimation via wasserstein distributionally robust optimization," *Mathematics of Operations Research*, vol. 48, no. 1, pp. 1–37, 2023.
- [8] S. Yi and M. Zorzi, "Robust fixed-lag smoothing under model perturbations," *Journal of the Franklin Institute*, vol. 360, no. 1, 2023.
- [9] S. Haykin, *Kalman filtering and Neural Networks*. John Wiley & Sons, 2004.
- [10] L. E. Baum and J. A. Eagon, "An inequality with applications to statistical estimation for probabilistic functions of markov processes and to a model for ecology," *Bulletin of the American Mathematical Society*, vol. 73, no. 3, pp. 360–363, 1967.
- [11] S. Liu, X. Zhang, L. Xu, and F. Ding, "Expectation-maximization algorithm for bilinear systems by using the rauch-tung-striebl smoother," *Automatica*, vol. 142, p. 110365, 2022.
- [12] M. Zheng and Y. Ohta, "Bayesian positive system identification: Truncated Gaussian prior and hyperparameter estimation," *Systems & Control Letters*, vol. 148, p. 104857, 2021.
- [13] T. B. Schön, A. Wills, and B. Ninness, "System identification of nonlinear state-space models," *Automatica*, vol. 47, no. 1, 2011.
- [14] Z. Ghahramani and G. E. Hinton, "Parameter estimation for linear dynamical systems," *Technical Report CRG-TR-96-2*, 1996.
- [15] N. Sammaknejad, Y. Zhao, and B. Huang, "A review of the expectation maximization algorithm in data-driven process identification," *Journal of process control*, vol. 73, pp. 123–136, 2019.
- [16] I. Naim and D. Gildea, "Convergence of the EM algorithm for Gaussian mixtures with unbalanced mixing coefficients," *arXiv preprint arXiv:1206.6427*, 2012.
- [17] S. Balakrishnan, M. J. Wainwright, and B. Yu, "Statistical guarantees for the EM algorithm: From population to sample-based analysis," *The Annals of Statistics*, 2017.
- [18] B. Yan, M. Yin, and P. Sarkar, "Convergence of gradient EM on multi-component mixture of Gaussians," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [19] M. E. Johnson, *Multivariate statistical simulation: A guide to selecting and generating continuous multivariate distributions*. John Wiley & Sons, 1987, vol. 192.
- [20] S. Sra and R. Hosseini, "Conic geometric optimization on the manifold of positive definite matrices," *SIAM Journal on Optimization*, vol. 25, no. 1, pp. 713–739, 2015.
- [21] R. Hosseini and S. Sra, "An alternative to EM for Gaussian mixture models: Batch and stochastic Riemannian optimization," *Mathematical programming*, vol. 181, no. 1, pp. 187–223, 2020.
- [22] Z. Chen, "Bayesian filtering: From kalman filters to particle filters, and beyond," *Statistics*, vol. 182, no. 1, pp. 1–69, 2003.
- [23] B. C. Levy, *Principles of signal detection and parameter estimation*. New York, NY, USA: Springer, 2008.
- [24] D. G. Luenberger, *Optimization by vector space methods*. John Wiley & Sons, 1997.
- [25] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2nd ed. New York, NY, USA: Springer, 2006.
- [26] Y. Nesterov and A. Nemirovskii, *Interior-point polynomial algorithms in convex programming*. SIAM, 1994.