

# Variational System Identification for Nonlinear State-Space Models

Jarrad Courts<sup>\*1</sup>, Adrian Wills<sup>†1</sup>, Thomas B. Schön<sup>‡, §2</sup>, and Brett Ninness<sup>¶1</sup>

<sup>1</sup>*University of Newcastle, School of Engineering, Australia*

<sup>2</sup>*Department of Information Technology, Uppsala University, Uppsala, Sweden*

September 15, 2022

## Abstract

This paper considers parameter estimation for nonlinear state-space models, which is an important but challenging problem. We address this challenge by employing a variational inference (VI) approach, which is a principled method that has deep connections to maximum likelihood estimation. This VI approach ultimately provides estimates of the model as solutions to an optimisation problem, which is deterministic, tractable and can be solved using standard optimisation tools. A specialisation of this approach for systems with additive Gaussian noise is also detailed. The proposed method is examined numerically on a range of simulated and real examples focusing on the robustness to parameter initialisation; additionally, favourable comparisons are performed against state-of-the-art alternatives.

## 1 Introduction

The problem of system identification is long-standing and has significant practical applications with a large body of related literature across many different fields [27, 18]. In many important applications, the underlying system exhibits nonlinear dynamic behaviour, which is assumed to be adequately captured by a time-indexed data record of system inputs and outputs. The system identification problem is then to obtain a suitable model that best explains this data. Such a model can then be employed as a surrogate for the actual system of interest for the purposes of prediction, control, decision making and analysis.

Due to the enormous variety of systems, there is now a wide array of model classes. These range from linear models [27], to block-structured nonlinear models [47], through to highly flexible nonlinear black-box models [30]. This paper focuses on probabilistic nonlinear state-space models, a highly flexible class of models that relies on a so-called state, the time evolution of which models the system dynamics. It is important to highlight that the state is typically unknown and, as such, is sometimes called hidden or latent. It is also well recognised that the system identification problem for these models is important and difficult [30, 35]; the states being unknown is a significant source of this difficulty [42].

In this paper, it is assumed that the structure of the state-space model is known but that the associated model parameters are unknown and are to be estimated from the available data. There are many approaches available for estimating these parameters, but we focus on maximum likelihood (ML) parameter estimation. The ML estimate has desirable statistical properties [2] and is commonly used in system identification [27, 35]. However, computing the ML estimate is generally intractable for nonlinear state-space models. This intractability leads to various approximations of the ML estimate; see, e.g. [43, 27, 35]. These approaches can be broadly grouped into two categories according to the strategy used to handle unknown states; the so-called ‘marginalisation’ and ‘data augmentation’ approaches [42].

---

\*jarrad.courts@uon.edu.au

†Adrian.Wills@newcastle.edu.au

‡thomas.schon@it.uu.se

§This research was financially supported by the Swedish Foundation for Strategic Research (SSF) via the project *ASSEMBLE* (contract number: RIT15-0012), by the Swedish Research Council via the projects *Deep probabilistic regression – new models and learning algorithms* (contract number: 2021-04301) and *NewLEADS - New Directions in Learning Dynamical Systems* (contract number: 621-2016-06079) and by *Kjell och Märta Beijer Foundation*.

¶brett.ninness@newcastle.edu.au

Marginalisation approaches compute the likelihood recursively by marginalising the unknown state and then maximise the resultant likelihood directly over the model parameters [40, 43]. This likelihood calculation is generally intractable due to the required marginalisation step over an unknown state distribution. This fundamental difficulty has led to various approximations; for example, [23] and [43] consider using an unscented Kalman filter and a particle filter, respectively. While promising results using stochastic quasi-Newton optimisation are shown in [48], the marginalisation approach commonly has difficulties regarding undesirable local minima in the resulting optimisation problems leading to issues regarding robustness to the initial parameter estimate [29, 28, 35].

Data augmentation approaches treat the unknown state as an auxiliary variable that is estimated alongside the model parameters. A commonly used data augmentation method is the expectation maximisation (EM) approach [12], which iterates between estimating the state and updating the parameters. Obtaining the state estimate for EM is also generally intractable, and instead, various approximations have been used [8]. Both particle (PSEM) [43] and assumed Gaussian [14, 22] approximations have been used within EM to deliver approximate ML estimates. Stochastic approximation EM (SAEM) [11] offers improved performance over PSEM for nonlinear state-space models [26], primarily due to more efficient iterations that use particles from previous iterations [42].

As an alternative to these methods, the primary contribution of this paper is the presentation of a variational inference (VI) [20, 5, 3, 4] based approach to approximate the ML parameter estimate for nonlinear state-space models. The developed approach is a data augmentation method that provides parameter estimates that *approximate* the maximum likelihood solution. The provided parameter estimate is obtained by solving a single optimisation problem of a standard form with readily available exact first- and second-order derivatives. The presented VI approach jointly estimates the state and parameters and exhibits rapid convergence while remaining robust. This behaviour contrasts with EM methods, where robustness is often observed, but the convergence rate decreases as the parameter estimate approaches the ML estimate.

Variational methods have also been applied to identify models of different classes in several works; see, for example, [3, 4, 39, 25]. However, these methods do not apply to nonlinear state-space models and, along with many other parameter estimation methods that use VI, are all variational-EM approaches [5, 49]. Compared to EM approaches, a *key innovation* of the proposed approach in this paper is the simultaneous optimisation over both the state and the model parameters, which appears to offer *significantly* improved convergence rates.

## 2 Problem Formulation

The problem of system identification for nonlinear state-space models considered in this paper consists of using the input and output measurements  $\mathbf{u} \triangleq u_{1:T} \triangleq \{u_1, \dots, u_T\}$  and  $\mathbf{y} \triangleq y_{1:T} \triangleq \{y_1, \dots, y_T\}$ , respectively, to compute an estimate of the model parameters  $\theta \in \mathcal{R}^{n_\theta}$  for the following model structure

$$x_{k+1} = f_k(x_k, u_k, v_k, \theta), \quad (1a)$$

$$y_k = h_k(x_k, u_k, e_k, \theta), \quad (1b)$$

where  $x_k \in \mathcal{R}^{n_x}$  is the state,  $y_k \in \mathcal{R}^{n_y}$  is the observed measurement,  $u_k \in \mathcal{R}^{n_u}$  is the measured input, and the functions  $f_k(\cdot)$  and  $h_k(\cdot)$  describe the process and measurement models, respectively. The process and measurement noise,  $v_k$  and  $e_k$ , respectively, are random variables assumed to belong to distributions of a known form. Parameters of these distributions are included within  $\theta$ . For ease of exposition, we henceforth drop the explicit dependence on the inputs  $\mathbf{u}$ .

Throughout this paper, we will frequently employ the alternative probabilistic representation of (1), given by

$$x_{k+1} \sim p_\theta(x_{k+1} | x_k), \quad (2a)$$

$$y_k \sim p_\theta(y_k | x_k), \quad (2b)$$

where  $p$  denotes a probability density function.

In this paper, the problem considered is approximating the maximum likelihood parameter estimate, given by

$$\theta_{\text{ML}} = \arg \max_{\theta} \log p_\theta(\mathbf{y}). \quad (3)$$

A well-known fundamental difficulty in solving (3) is that the log-likelihood function cannot be exactly computed, which extends to the computation of gradient information, and therefore optimisation is challenging. In the following section, we address this challenge by employing a so-called variational inference approach.

### 3 Variational Nonlinear System Identification

In this section, the use of variational inference applied to nonlinear state-space systems will be presented in Section 3.1. The relationship between variational inference and expectation maximisation is then examined in Section 3.2.

#### 3.1 Variational Inference

Variational inference is a widely used method where the primary aim is to approximate intractable distributions by a tractable parametric density of an assumed form. The use of the epithet ‘variational’ stems from the idea that these methods rely on optimisation as the primary tool for choosing the member of the parametric assumed density that best matches the intractable distribution of interest. In this paper, we will employ the Kullback-Leibler (KL) divergence [24] as the cost function for determining optimality. For ease of reference, the KL divergence between two densities  $p(z)$  and  $q(z)$  is defined as

$$\text{KL}[p(z) \parallel q(z)] \triangleq \int p(z) \log \frac{p(z)}{q(z)} dz. \quad (4)$$

In the context of this paper, we motivate the use of VI by noting that the log-likelihood can be expressed using KL divergence via

$$\log p_\theta(\mathbf{y}) = \int p_\theta(\mathbf{x} \mid \mathbf{y}) \log p_\theta(\mathbf{y}) d\mathbf{x} \quad (5a)$$

$$= - \int p_\theta(\mathbf{x} \mid \mathbf{y}) \log \frac{p_\theta(\mathbf{x} \mid \mathbf{y})}{p_\theta(\mathbf{x}, \mathbf{y})} d\mathbf{x} \quad (5b)$$

$$= -\text{KL}[p_\theta(\mathbf{x} \mid \mathbf{y}) \parallel p_\theta(\mathbf{x}, \mathbf{y})], \quad (5c)$$

where  $\mathbf{x} \triangleq x_{1:T+1} \triangleq \{x_1, \dots, x_{T+1}\}$  and the first equality comes from the fact that  $\log p_\theta(\mathbf{y})$  does not depend on  $\mathbf{x}$  and is, therefore, invariant to expectation relative to any density in  $\mathbf{x}$ . The second equality stems from the fact that  $-\log p_\theta(\mathbf{y}) = \log p_\theta(\mathbf{x} \mid \mathbf{y}) - \log p_\theta(\mathbf{x}, \mathbf{y})$  and the third equality from the definition of KL divergence (4). Therefore, the maximum likelihood problem (3) can be equivalently stated as

$$\theta_{\text{ML}} = \arg \min_{\theta} \text{KL}[p_\theta(\mathbf{x} \mid \mathbf{y}) \parallel p_\theta(\mathbf{x}, \mathbf{y})]. \quad (6)$$

Unfortunately, the smoothed state distribution  $p_\theta(\mathbf{x} \mid \mathbf{y})$  cannot be expressed in closed form, rendering the associated optimisation problem (6) intractable.

In light of this, we propose to replace the intractable smoothed distribution  $p_\theta(\mathbf{x} \mid \mathbf{y})$  with an  $\eta$ -parameterised distribution  $q_\eta(\mathbf{x})$ , called the assumed density, and solve a similar problem

$$\theta^*, \eta^* = \arg \min_{\theta, \eta} \text{KL}[q_\eta(\mathbf{x}) \parallel p_\theta(\mathbf{x}, \mathbf{y})]. \quad (7)$$

Note that the notation introduced here for an assumed density is general and is used with other parameters than just  $\eta$ , and for distributions over other variables than just  $\mathbf{x}$ , within this paper.

A major benefit of this approach is that  $q_\eta(\mathbf{x})$  can be chosen in a convenient manner such that the resulting problem is tractable. This procedure of approximating an intractable density with an assumed density and minimising the KL divergence is known as variational inference (see, e.g. [5]). The utility of this approach is then highly dependent on the choice of assumed density. In principle, this assumed density is selected to achieve two sometimes competing goals, 1) the assumed density family should be flexible enough to ‘closely’ match  $p_\theta(\mathbf{x} \mid \mathbf{y})$ , and 2) the assumed density  $q_\eta(\mathbf{x})$  should result in a tractable optimisation problem (7).

To further discuss this approach, it will be convenient to introduce a cost function

$$\mathcal{L}(\eta, \theta) = -\text{KL}[q_\eta(\mathbf{x}) \parallel p_\theta(\mathbf{x}, \mathbf{y})], \quad (8)$$

so that (7) becomes  $\theta^*, \eta^* = \arg \max_{\theta, \eta} \mathcal{L}(\eta, \theta)$ . The following lemma establishes that the log-likelihood  $\log p_\theta(\mathbf{y})$  is bounded below by  $\mathcal{L}(\eta, \theta)$ , and the gap between the lower bound and log-likelihood is provided by  $\text{KL}[q_\eta(\mathbf{x}) \parallel p_\theta(\mathbf{x} \mid \mathbf{y})]$ . Therefore, as the KL divergence between  $q_\eta(\mathbf{x})$  and  $p_\theta(\mathbf{x} \mid \mathbf{y})$  diminishes, the lower bound becomes tight.

**Lemma 1.** *The log-likelihood can be expressed as*

$$\log p_\theta(\mathbf{y}) = \mathcal{L}(\eta, \theta) + \text{KL}[q_\eta(\mathbf{x}) \parallel p_\theta(\mathbf{x} \mid \mathbf{y})], \quad (9)$$

and is bounded below according to

$$\log p_\theta(\mathbf{y}) \geq \mathcal{L}(\eta, \theta). \quad (10)$$

*Proof.* Using conditional probability

$$\log p_\theta(\mathbf{y}) = \log p_\theta(\mathbf{x}, \mathbf{y}) - \log p_\theta(\mathbf{x} \mid \mathbf{y}). \quad (11)$$

Through addition and subtraction of  $\log q_\eta(\mathbf{x})$  to the right-hand side of (11), this leads to

$$\log p_\theta(\mathbf{y}) = \log \frac{p_\theta(\mathbf{x}, \mathbf{y})}{q_\eta(\mathbf{x})} + \log \frac{q_\eta(\mathbf{x})}{p_\theta(\mathbf{x} \mid \mathbf{y})}. \quad (12)$$

Taking expectation of both sides relative to  $q_\eta(\mathbf{x})$  delivers

$$\begin{aligned} \log p_\theta(\mathbf{y}) &= - \int q_\eta(\mathbf{x}) \log \frac{q_\eta(\mathbf{x})}{p_\theta(\mathbf{x}, \mathbf{y})} d\mathbf{x} \\ &\quad + \int q_\eta(\mathbf{x}) \log \frac{q_\eta(\mathbf{x})}{p_\theta(\mathbf{x} \mid \mathbf{y})} d\mathbf{x}, \end{aligned} \quad (13)$$

which, from the definition of KL divergence, leads to the expression in (9). Since KL divergence is non-negative, then  $\text{KL}[q_\eta(\mathbf{x}) \parallel p_\theta(\mathbf{x} \mid \mathbf{y})] \geq 0$  and (10) follows immediately from (9).  $\square$

### 3.2 Comparison to Expectation Maximisation

There are close connections between solutions based on EM and VI; see, e.g. [34] and [44]. This section examines some of these similarities and differences in the context of system identification for nonlinear state-space models. Towards this end, the EM method can be summarised as iterating

$$\theta_{k+1} = \arg \min_{\theta} \text{KL}[p_{\theta_k}(\mathbf{x} \mid \mathbf{y}) \parallel p_\theta(\mathbf{x}, \mathbf{y})]. \quad (14)$$

It is interesting to compare this with the ML problem in (6), where the only difference for EM is that the smoothed density  $p_{\theta_k}(\mathbf{x} \mid \mathbf{y})$  is fixed at the  $k$ 'th parameter values. Importantly, EM relies on the smoothed density, which is generally intractable. In theory, the EM approach generates parameters estimates  $\theta_k$  that monotonically increases  $\log p_{\theta_k}(\mathbf{y})$  [7].

The following lemma reveals that EM can be viewed as block-ascent of the VI cost  $\mathcal{L}(\eta, \theta)$  over  $\eta$  and then  $\theta$ , respectively. This requires an assumption on  $q_\eta(\mathbf{x})$ , which essentially means that it is possible to match the smoothed density  $p_{\theta_k}(\mathbf{x} \mid \mathbf{y})$ .

**Lemma 2.** *Assume there exists an  $\eta$  such that*

$$q_\eta(\mathbf{x}) = p_{\theta_k}(\mathbf{x} \mid \mathbf{y}). \quad (15)$$

Then the EM iterations can be expressed as

$$\eta_k = \arg \max_{\eta} \mathcal{L}(\eta, \theta_k), \quad (16a)$$

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}(\eta_k, \theta). \quad (16b)$$

*Proof.* According to (9), we can state (16a) as

$$\eta_k = \arg \max_{\eta} \log p_{\theta_k}(\mathbf{y}) - \text{KL}[q_\eta(\mathbf{x}) \parallel p_{\theta_k}(\mathbf{x} \mid \mathbf{y})].$$

Notice that  $\log p_{\theta_k}(\mathbf{y})$  does not depend on  $\eta$  and, therefore, the above problem becomes

$$\eta_k = \arg \min_{\eta} \text{KL}[q_\eta(\mathbf{x}) \parallel p_{\theta_k}(\mathbf{x} \mid \mathbf{y})]. \quad (17)$$

Under assumption (15),  $\text{KL}[q_\eta(\mathbf{x}) \parallel p_{\theta_k}(\mathbf{x} \mid \mathbf{y})]$  has a global minimum value of zero, which is achieved for any  $\eta_k$  that solves (17). Therefore,  $q_{\eta_k}(\mathbf{x}) = p_{\theta_k}(\mathbf{x} \mid \mathbf{y})$ , and (16b) coincides with (14).  $\square$

It is important to recognise that, generally, it is not possible to select a *tractable*  $q_\eta(\mathbf{x})$  such that  $q_\eta(\mathbf{x}) = p_{\theta_k}(\mathbf{x} | \mathbf{y})$  for some  $\eta$ . Similarly, exactly computing the expectations with respect to the smoothed distribution  $p_{\theta_k}(\mathbf{x} | \mathbf{y})$  is generally not possible for nonlinear state-space models. As such, exactly performing EM on nonlinear state-space models is generally intractable, which removes the guaranteed non-decrease of the log-likelihood. So-called variational EM is obtained from iterating (16a)–(16b), where the assumed density  $q_\eta(\mathbf{x})$  will not generally match  $p_{\theta_k}(\mathbf{x} | \mathbf{y})$  [5, 44]. This approach ensures a monotonic sequence of lower bounds to the log-likelihood.

Relative to these approaches, a major advantage of the proposed VI approach in (7) is that both  $\eta$  and  $\theta$  are jointly optimised; typically, this offers improved convergence rates compared with coordinate descent methods [36]. While the simulations in Section 6 confirm this improved convergence rate, we recognise the limitations of drawing general conclusions from these results.

## 4 Assumed Gaussian Distribution

As mentioned in Section 3, there are two competing goals when choosing the assumed density  $q_\eta(\mathbf{x})$ , namely that it should be flexible enough to closely match  $p_\theta(\mathbf{x} | \mathbf{y})$  and that it should result in a manageable optimisation problem (7).

In light of this, we propose  $q_\eta(\mathbf{x})$  to be a multivariate Normal (MVN) distribution

$$q_\eta(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mu, \Sigma), \quad \eta \triangleq (\mu, \Sigma), \quad (18)$$

where for clarity of exposition, the mean and covariance have the following structure

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_{T+1} \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{1,1} & \Sigma_{1,2} & \cdots & \Sigma_{1,T+1} \\ \Sigma_{1,2}^\top & \Sigma_{2,2} & \cdots & \Sigma_{2,T+1} \\ \vdots & & \ddots & \vdots \\ \Sigma_{1,T+1}^\top & \Sigma_{2,T+1}^\top & \cdots & \Sigma_{T+1,T+1} \end{bmatrix}, \quad (19)$$

where each  $\mu_k \in \mathbb{R}^{n_x}$  and  $\Sigma_{k,j} \in \mathbb{R}^{n_x \times n_x}$ . It is difficult to comment on whether or not this proposal will achieve the first goal since  $p_\theta(\mathbf{x} | \mathbf{y})$  is unknown in general. Nevertheless, this assumption aligns with other approaches where the smoothed distribution is assumed to be MVN [40]. Note that this assumption on the form of the underlying state *may* limit the performance of the overall system identification approach on some systems; examples may include multi-modal systems. Importantly, the MVN assumption leads to a tractable and deterministic optimisation problem (7). This section discusses the relevant details of this assumption and presents a different parameterisation of  $\eta$  with a dimension significantly less than that of  $(\mu, \Sigma)$ .

### 4.1 Assumed Density $q_\eta(\mathbf{x})$ Details

In this section, we show that it is not necessary to compute the full joint state density  $q_\eta(\mathbf{x})$  to compute  $\theta^*$  under the MVN assumption. Towards this, we begin by splitting  $\eta$  into two parts. First,  $\alpha$  containing all the mean values and the block-tridiagonal covariance matrices from (19), and second,  $\gamma$  containing all the remaining non-block-tridiagonal covariance matrices from (19).

$$\alpha \triangleq \left( \mu, (\Sigma_{k,k})_{k=1}^{T+1}, (\Sigma_{k,k+1})_{k=1}^T \right), \quad (20a)$$

$$\gamma \triangleq \left( (\Sigma_{k,j})_{j=k+2}^{T+1} \right)_{k=1}^{T-1}. \quad (20b)$$

The following lemma shows that  $\gamma$  does not affect the optimal  $\alpha$  parameters. This result means that  $\gamma$  can be removed from the problem without impacting  $\theta^*$ . The benefit is that the dimension of  $\alpha$  is significantly smaller than  $\eta$ .

**Lemma 3.** Assume that  $\eta = (\alpha, \gamma)$  according to (20) and let

$$\theta^*, \alpha^* = \arg \max_{\theta, \alpha} \mathcal{L}_R(\alpha, \theta), \quad (21)$$

where  $\mathcal{L}_R(\alpha, \theta)$  is defined as

$$\mathcal{L}_R(\alpha, \theta) = I_1(\alpha) + I_2(\alpha, \theta) - I_3(\alpha), \quad (22)$$

and

$$\begin{aligned}
I_1(\alpha) &= \int q_\alpha(x_1) \log p(x_1) dx_1, \\
I_2(\alpha, \theta) &= \sum_{k=1}^T \int q_\alpha(x_{k:k+1}) \log p_\theta(x_{k+1}, y_k | x_k) dx_{k:k+1}, \\
I_3(\alpha) &= \sum_{k=1}^T \int q_\alpha(x_{k:k+1}) \log q_\alpha(x_{k:k+1}) dx_{k:k+1} \\
&\quad - \sum_{k=2}^T \int q_\alpha(x_k) \log q_\alpha(x_k) dx_k.
\end{aligned}$$

Then

$$\theta^*, \alpha^*, \gamma^* = \arg \max_{\theta, \alpha, \gamma} \mathcal{L}((\alpha, \gamma), \theta). \quad (23)$$

*Proof.* See Appendix A. □

Therefore, since we are ultimately interested in  $\theta^*$ , it suffices to solve (21), which has the major benefit of involving a much-reduced parameter space. To solve (21) using standard gradient-based optimisation methods, it is crucial  $\mathcal{L}_R(\theta, \alpha)$  can be evaluated, along with its gradient and possibly Hessian. For  $I_1$  (with a suitable choice for the prior) and  $I_3$ , these calculations are relatively straightforward since the required expectations have known closed-form solutions. At the same time, these calculations are challenging for  $I_2$  since they cannot be computed in closed-form in general.

To ameliorate this, here we introduce a parameterisation of the joint state  $(x_k, x_{k+1})$  distribution that leads to both computationally efficient approximations and straightforward computation of the first- and second-order derivatives. To this end, let  $\beta_k$  be defined as the set of parameters

$$\beta_k = (\mu_k, \bar{\mu}_k, A_k, B_k, C_k), \quad (24)$$

where  $A_k, B_k, C_k \in \mathcal{R}^{n_x \times n_x}$  and  $A_k, C_k$  are upper triangular. Then we can parameterise the assumed joint state  $(x_k, x_{k+1})$  density as

$$q_{\beta_k}(x_k, x_{k+1}) = \mathcal{N}\left(\begin{bmatrix} x_k \\ x_{k+1} \end{bmatrix}; \begin{bmatrix} \mu_k \\ \bar{\mu}_k \end{bmatrix}, P_k^\top P_k\right), \quad (25)$$

where  $P_k$  is an upper triangular Cholesky factor

$$P_k = \begin{bmatrix} A_k & B_k \\ 0 & C_k \end{bmatrix}. \quad (26)$$

We can collect all the  $\beta_k$ 's into the tuple

$$\beta \triangleq (\beta_1, \beta_2, \dots, \beta_T). \quad (27)$$

Note that this formulation is over-parameterised and it will lead to inconsistent results since the marginal distribution for  $x_k$  can be computed from either  $q_{\beta_{k-1}}(x_{k-1}, x_k)$  or  $q_{\beta_k}(x_k, x_{k+1})$ . The following lemma introduces a constraint set  $\Omega$  that eliminates the unnecessary degrees of freedom in  $\beta$ .

**Lemma 4.** Assume that  $\beta$  is defined by (27). Let

$$\begin{aligned}
\Omega \triangleq \{ \beta \mid & B_k^\top B_k + C_k^\top C_k = A_{k+1}^\top A_{k+1}, \\
& \mu_{k+1} = \bar{\mu}_k, \quad k = 1, \dots, T-1 \}.
\end{aligned}$$

Then, for  $\beta \in \Omega$ , it follows that

$$q_\alpha(x_{k:k+1}) = q_{\beta_k}(x_{k:k+1}), \quad (28)$$

whenever

$$\Sigma_{k,k} = A_k^\top A_k, \quad \Sigma_{k,k+1} = A_k^\top B_k. \quad (29)$$

*Proof.* According to (18), note that  $q_\alpha(x_{k:k+1})$  is given by

$$\mathcal{N}\left(\begin{bmatrix} x_k \\ x_{k+1} \end{bmatrix}; \begin{bmatrix} \mu_k \\ \mu_{k+1} \end{bmatrix}, \begin{bmatrix} \Sigma_{k,k} & \Sigma_{k,k+1} \\ \Sigma_{k+1,k} & \Sigma_{k+1,k+1} \end{bmatrix}\right). \quad (30)$$

From (25), we note that  $q_{\beta_k}(x_{k:k+1})$  is

$$\mathcal{N}\left(\begin{bmatrix} x_k \\ x_{k+1} \end{bmatrix}; \begin{bmatrix} \mu_k \\ \bar{\mu}_k \end{bmatrix}, \begin{bmatrix} A_k^\top A_k & A_k^\top B_k \\ B_k^\top A_k & B_k^\top B_k + C_k^\top C_k \end{bmatrix}\right). \quad (31)$$

If  $\beta \in \Omega$ , then  $\bar{\mu}_k = \mu_{k+1}$  and  $B_k^\top B_k + C_k^\top C_k = A_{k+1}^\top A_{k+1}$ , and therefore (28) follows from (29).  $\square$

The next section makes use of this parameterisation to approximate the required expectation in  $I_2$  in a computationally efficient manner.

## 4.2 VI Approximation

As mentioned in the previous section, it is vital for gradient-based optimisation that the expectation in  $I_2$  can be approximated, along with its first- and second-order derivatives. The key difficulty in calculating  $I_2$  is the computation of  $\int q_\alpha(x_{k:k+1}) \log p_\theta(x_{k+1}, y_k | x_k) dx_{k:k+1}$ . Employing Lemma 4, we can write this integral in terms of  $\beta_k$ , instead of  $\alpha$ , as

$$E_k(\beta_k, \theta) \triangleq \int q_{\beta_k}(x_{k:k+1}) \log p_\theta(x_{k+1}, y_k | x_k) dx_{k:k+1}.$$

Importantly, the choice in parameterisation affords a straightforward approximation of  $E_k(\beta_k, \theta)$ , denoted as  $\hat{E}_k(\beta_k, \theta)$ , using Gaussian quadrature (see, e.g. [21, 46]) via

$$\hat{E}_k(\beta_k, \theta) \triangleq \sum_{j=1}^{n_s} w_j \log p_\theta(\bar{x}_{k+1}^j, y_k | x_k^j). \quad (32)$$

In the above,  $w_j$  are the so-called weights, and  $x_k^j$  and  $\bar{x}_{k+1}^j$  are the so-called sigma points. Importantly, the weights are predefined constants, and the sigma points are defined as

$$\begin{bmatrix} x_k^j \\ \bar{x}_{k+1}^j \end{bmatrix} = \begin{bmatrix} \mu_k & A_k^\top & 0 \\ \bar{\mu}_k & B_k^\top & C_k^\top \end{bmatrix} a_j, \quad (33)$$

where the vectors  $a_j$  are also constant. Both the weights  $w_j$  and vectors  $a_j$  depend on the choice of Gaussian quadrature, but they are nevertheless constant for this choice. For further details on specific quadrature methods and the accuracy of different methods in a similar context, see, for example, [19, 41, 23]. Furthermore, the sigma points being linear combinations of the elements of  $\beta_k$  is critical as it significantly simplifies the calculation of both first- and second-order derivatives used in the optimisation.

Using this approximation, we can define a tractable approximation to problem (21) via

$$\hat{\beta}, \hat{\theta} = \arg \max_{\beta, \theta} \hat{\mathcal{L}}_R(\beta, \theta), \quad \text{s.t. } \beta \in \Omega, \quad (34)$$

where

$$\hat{\mathcal{L}}_R(\beta, \theta) = I_1(\beta) + \hat{I}_2(\beta, \theta) - I_3(\beta), \quad (35)$$

and

$$\hat{I}_2(\beta, \theta) = \sum_{k=1}^T \hat{E}_k(\beta_k, \theta). \quad (36)$$

This constrained optimisation problem is of standard form and can be solved using exact first- and second-order derivatives without any further approximations.

The assumptions contained within (34) are summarised as follows: first, an MVN distribution of the state is assumed. Second,  $\log p_\theta(\bar{x}_{k+1}^j, y_k | x_k^j)$  can be evaluated and it is twice continuously differentiable, and third, the use of quadrature to approximate the intractable expectations.

### 4.3 Additive Gaussian Noise

In this section, we consider a simplification for nonlinear systems with additive Gaussian noise. This simplification enables a reduction in the number of optimisation variables, brings numerical benefits, and enables an effective approximation of the Hessian to be formed utilising only first-order derivatives. We consider a model structure of the form

$$\begin{bmatrix} x_{k+1} \\ y_k \end{bmatrix} = \begin{bmatrix} f(x_k, \phi) \\ h(x_k, \phi) \end{bmatrix} + \begin{bmatrix} v_k \\ e_k \end{bmatrix}, \quad \begin{bmatrix} v_k \\ e_k \end{bmatrix} \sim \mathcal{N}(0, \Pi), \quad (37)$$

where  $\theta = (\phi, \Pi)$  includes model parameters  $\phi$  and noise covariance  $\Pi$ . As detailed in the following Lemma 5, the structure of this system allows for system identification to be performed using a more structured, reduced size optimisation problem.

**Lemma 5.** *For systems in the form of (37), identification can be performed by solving the reduced problem*

$$\hat{\beta}, \hat{\phi} = \arg \max_{\phi, \beta} \hat{\mathcal{L}}_{\text{AG}}(\phi, \beta), \quad \text{s.t. } \beta \in \Omega, \quad (38)$$

where

$$\hat{\mathcal{L}}_{\text{AG}}(\phi, \beta) = I_1(\beta) + \hat{I}_2^{\text{AG}}(\phi, \beta) - I_3(\beta), \quad (39a)$$

$$\hat{I}_2^{\text{AG}}(\phi, \beta) = c + \frac{T}{2} \log |\hat{\Pi}(\phi, \beta)|, \quad (39b)$$

$$\hat{\Pi}(\phi, \beta) = \frac{1}{T} \sum_{k=1}^T \sum_{j=1}^{n_s} w_j \xi_k^j (\xi_k^j)^\top, \quad (39c)$$

$$\xi_k^j = \begin{bmatrix} \bar{x}_{k+1}^j - f(x_k^j, \phi) \\ y_k - h(x_k^j, \phi) \end{bmatrix}, \quad (39d)$$

and  $c$  is a constant that does not depend on  $\phi$  or  $\beta$ .

*Proof.* Using (36) and (37), we can evaluate  $\hat{I}_2$  as

$$\hat{I}_2 = \frac{T}{2} \log |2\pi\Pi| + \frac{1}{2} \sum_{k=1}^T \sum_{j=1}^{n_s} w_j (\xi_k^j)^\top \Pi^{-1} (\xi_k^j). \quad (40)$$

We note that  $(\xi_k^j)^\top \Pi^{-1} (\xi_k^j) = \text{tr}\{(\xi_k^j)^\top \Pi^{-1} (\xi_k^j)\} = \text{tr}\{\Pi^{-1} (\xi_k^j) (\xi_k^j)^\top\}$  since the trace operator is invariant to cyclic permutations. Due to the fact that the trace is also a linear operator, we have that

$$\hat{I}_2 = \frac{T}{2} \log |2\pi\Pi| + \text{tr} \left\{ \Pi^{-1} \frac{1}{2} \sum_{k=1}^T \sum_{j=1}^{n_s} w_j (\xi_k^j)^\top (\xi_k^j) \right\}.$$

First-order necessary conditions of optimality require that  $\partial \hat{I}_2(\beta, \theta) / \partial \Pi = 0$ , which occurs when

$\Pi = \frac{1}{T} \sum_{k=1}^T \sum_{j=1}^{n_s} w_j \xi_k^j (\xi_k^j)^\top$ , hence (39c). Substituting (39c) into (40) yields (39b), which completes the proof.  $\square$

Note that the final estimate of  $\Pi$  is subsequently given by evaluating  $\hat{\Pi}(\hat{\phi}, \hat{\beta})$  using (39c). Further, decoupled noise can be accommodated by constraining the off-diagonal block of  $\Pi$  to zero.

This specialisation possesses several benefits compared with a general approach. First, the number of variables required to be optimised is reduced while still addressing the original underlying problem. Second, using only the Jacobian of  $f(x_k, \phi)$  and  $h(x_k, \phi)$ , the exact gradient and a good approximation for the Hessian of  $\hat{\mathcal{L}}_{\text{AG}}(\phi, \beta)$  can be obtained; this is discussed further in Section 5.2, which focuses on the resulting optimisation problem.

## 5 Implementation Details

In this section, several key points regarding implementation details and the resulting optimisation problems are considered. Section 5.1 examines the initialisation of the resultant optimisation problems, and Section 5.2 provides some details regarding the optimisations required.

## 5.1 Initialisation

Similar to all identification approaches for nonlinear systems, an initial estimate of  $\theta$  must be provided and influences both the run-time required, and potentially, the parameter estimate produced. However, compared to EM-based approaches, the proposed approach also requires an initial estimate of each pairwise joint distribution. Similar to EM identification approaches, a smoother can provide this initialisation.

Alternatively, other approaches to initialise the state distributions can be used. This initialisation does not need to satisfy the constraints and introduces an added level of flexibility to exploit. Generally, using the distributions from a filtering pass has proven both effective and straightforward. The initial distribution can also be selected using problem-specific knowledge; an example is when state estimates can be readily calculated from the measurements. This alternative initialisation is particularly beneficial for the additive noise specialisation (Section 4.3), as it removes the requirement to provide initial estimates for  $\Pi$ , which is not well known a priori.

Due to the general nature of the optimisation problems, it is not possible to describe the ‘best’ initialisation scheme or how to ensure that undesirable local minima are avoided. As such, the numerical examples in Section 6 focus on robustness with respect to the initialisation of parameters and use all of the initialisation schemes for the state distributions discussed.

## 5.2 Optimisation

As previously stated, the resulting optimisation problems are of standard form and can be solved using standard solvers. In the numeric examples, the solvers `fmincon` [32] and `KNITRO` [6] are both used.

To perform the optimisation effectively, the exact first- and second-order derivatives are used for the general form of the proposed approach. Due to the assumed Gaussian distributions, closed-form expressions, and first- and second-order derivatives exist for all terms, except for  $\hat{I}_2(\beta, \theta)$ . Due to the parametrisation of the optimisation problem, the exact gradient and Hessian of  $\hat{I}_2(\beta, \theta)$  can be efficiently obtained using automatic differentiation [16]; we have used `CasADi` [1] for this purpose.

To further structure the optimisation problem a copy of  $\theta$  for each time step, denoted  $\theta_k$ , and the constraints  $\theta_1 = \theta_2 = \dots = \theta_T$ , are introduced to result in a sparse block-diagonal Hessian. This structure allows the Hessian to be both efficiently formed and used within optimisation routines with a computational complexity that grows linearly in data length, i.e.  $\mathcal{O}(T)$ .

For the additive Gaussian noise specialisation, the derivatives for  $\hat{I}_2^{AG}(\phi, \beta)$  are calculated using Lemma 6.

**Lemma 6.** *The gradient and Hessian of  $\hat{I}_2^{AG}(\phi, \beta)$  with respect to  $x$  are given by*

$$-\nabla_x \hat{I}_2^{AG}(\phi, \beta) = T J_n^T \text{vec} \left( \bar{\Pi}^{-\frac{1}{2}} e(x) \right), \quad (41a)$$

$$-\nabla_{x,x}^2 \hat{I}_2^{AG}(\phi, \beta) = T J_n^T J_n - \frac{T}{2} V^T V + S, \quad (41b)$$

where

$$J_n = \left( I_{Tn_s} \otimes \bar{\Pi}^{-\frac{1}{2}} \right) \frac{\partial \text{vec}(e(x))}{\partial x^T}, \quad (41c)$$

$$V = \tilde{V} + P_{n_x+n_y} \tilde{V}, \quad (41d)$$

$$\tilde{V} = \left( \bar{\Pi}^{-\frac{1}{2}} e(x) \otimes I_{n_x+n_y} \right) J_n, \quad (41e)$$

$$\bar{\Pi} = e(x) e(x)^T, \quad (41f)$$

$$e(x) = [e(x)_1, \dots, e(x)_T], \quad (41g)$$

$$e(x)_k = [\sqrt{w_1} \xi_k^1, \dots, \sqrt{w_{n_s}} \xi_k^{n_s}]. \quad (41h)$$

The matrix  $P_{n_x+n_y}$  is a square vec-permutation matrix [17],  $S$  contains second-order derivative terms,  $x$  is the variables being optimised,  $J_n^T J_n$  is block diagonal,  $V^T V$  is dense,  $V \in \mathcal{R}^{(n_x+n_y)^2 \times Tn_b}$ ,  $n_b$  is the dimension of the diagonal blocks, and, as  $e(x)$  is a linear function of the sigma points, both  $J_n$  and  $V$  can be obtained using only the Jacobian of  $f(x_k, \phi)$  and  $h(x_k, \phi)$ .

*Proof.* These equations can be verified by applying the approach to matrix calculus given in [31], properties of Kronecker products, and extensive algebraic manipulations.  $\square$

The Hessian of  $-\hat{J}_2^{\text{AG}}(\phi, \beta)$  is approximated as

$$H_2 \approx T J_n^T J_n - \frac{T}{2} V^T V. \quad (42)$$

For large  $T$ , the dense term renders both forming and factorising  $H_2$  computationally intractable. To avoid this issue, we have developed a custom trust-region optimisation routine that, by exploiting the structure of the Hessian approximation, performs an equivalent factorisation with a computational complexity linear in time with respect to the number of measurements. The developed routine is based on the approaches in [36, 9, 33]; the full details are beyond the scope of this paper.

## 6 Examples

In this section, we present three numerical examples. First, a stochastic volatility model (Section 6.1) that does not possess additive Gaussian noise. Second, a simulated example of a differential drive robot (Section 6.2) and third, an inverted pendulum using real data (Section 6.3). These latter two examples have additive Gaussian noise and structures arising from first-principles modelling, representative of practical applications. The pendulum example differs in that it compares with a particle, rather than an assumed Gaussian method, and that the system is unstable, a property that can be challenging [38].

For all examples, unless specified otherwise, the proposed approach refers to the general form given by (34). A third-order unscented transform [21, 46] has been used for all Gaussian quadrature approximations. All numerical examples were conducted on a laptop with an i7-7820HK processor and 32GB of memory.

### 6.1 Stochastic Volatility Model

In this example, we consider the estimation of  $\theta = [a, b, c]$  for the following stochastic volatility model

$$x_{k+1} = a + bx_k + \sqrt{c}v_k, \quad y_k = \sqrt{e^{x_k}}e_k, \quad (43)$$

where  $v_k \sim \mathcal{N}(0, 1)$  and  $e_k \sim \mathcal{N}(0, 1)$  is considered using 726 simulated measurements. Results obtained using stochastic approximation EM (SAEM) with a conditional particle filter (CPF) [26], referred to as CPF-SAEM, which is an asymptotically convergent method, are used to compare against.

For the proposed method, each joint state distribution is initialised with a mean of  $[2, 2]^T$  and a diagonal covariance with standard deviations of 0.1. An initial parameter estimate of  $\theta = [0, 0.5, 1]^T$  is used for both approaches, and 50 particles were used for the CPF-SAEM approach.

The proposed method converged in 19 iterations and 3.18s using `fmincon` [32] to perform the optimisation. In contrast, CPF-SAEM did not converge and is limited to 1000 iterations, which required 326s to complete. Fig. 1 shows the parameter trajectory of both methods and the true values. This result highlights that both methods provide similar estimates that are close to the true parameter values. CPF-SAEM, however, requires a significantly larger quantity of iterations and has only approached the parameter values obtained using the proposed method.

To examine if, at least for this example, the parameter estimates produced using the developed method converge as expected for a maximum likelihood method, we have additionally estimated  $\theta$  for 50 different data realisations for differing numbers of simulated measurements. Table 1 shows the results of this experiment and illustrates an asymptotic trend of decreasing bias and variance of the estimated parameters as the number of measurement samples increases.

The robustness to initial estimates of the proposed method has also been examined using 100 random initial parameter estimates sampled from

$$a \sim \mathcal{U}(-0.5, 0.5), \quad b \sim \mathcal{U}(0, 1.5), \quad c \sim \mathcal{U}(0.25, 2).$$

Parameter trajectories for each trial are plotted in Fig. 2 and shows that each initialisation has converged to the same parameter estimate.

In this section, the proposed system identification approach is applied to a system without additive Gaussian noise and compared favourably with CPF-SAEM. The lack of additive Gaussian noise required no alterations to the proposed approach, which proved to be both effective and robust to initial parameter estimates.

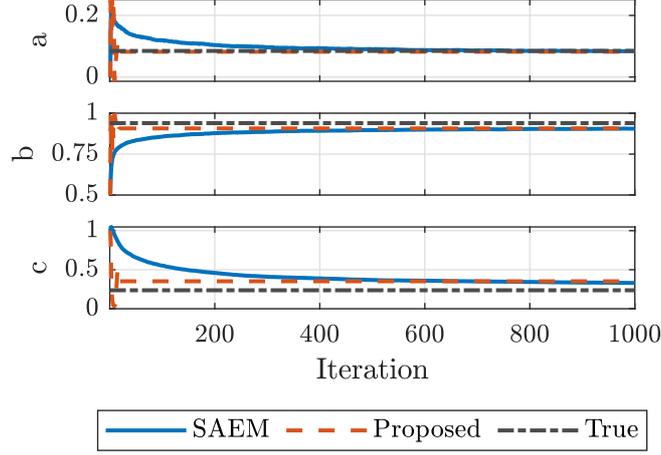


Figure 1: Parameter trajectory vs. iteration count. The lines for the proposed method have been extended beyond 19 iterations for clarity.

Table 1: Mean and standard deviation of the estimated parameters from the true values over 50 different realisations for differing numbers of measurement samples.

Samples	Error of the Estimated Parameter		
	$a$ (true = 0.0848)	$b$ (true = 0.9393)	$c$ (true = 0.2369)
100	$3.19 \times 10^{-1} \pm 5.98 \times 10^{-1}$	$-3.80 \times 10^{-1} \pm 4.28 \times 10^{-1}$	$4.76 \times 10^{-1} \pm 6.51 \times 10^{-1}$
500	$2.28 \times 10^{-2} \pm 5.07 \times 10^{-2}$	$-1.93 \times 10^{-2} \pm 2.75 \times 10^{-2}$	$4.85 \times 10^{-2} \pm 8.11 \times 10^{-2}$
1000	$2.31 \times 10^{-2} \pm 4.12 \times 10^{-2}$	$-1.40 \times 10^{-2} \pm 2.26 \times 10^{-2}$	$3.79 \times 10^{-2} \pm 6.10 \times 10^{-2}$
5000	$8.63 \times 10^{-3} \pm 1.28 \times 10^{-2}$	$-6.70 \times 10^{-3} \pm 7.68 \times 10^{-3}$	$3.05 \times 10^{-2} \pm 2.36 \times 10^{-2}$
10 000	$7.59 \times 10^{-3} \pm 8.66 \times 10^{-3}$	$-5.34 \times 10^{-3} \pm 5.04 \times 10^{-3}$	$2.86 \times 10^{-2} \pm 1.74 \times 10^{-2}$
25 000	$5.65 \times 10^{-3} \pm 5.45 \times 10^{-3}$	$-4.87 \times 10^{-3} \pm 3.29 \times 10^{-3}$	$2.50 \times 10^{-2} \pm 1.10 \times 10^{-2}$

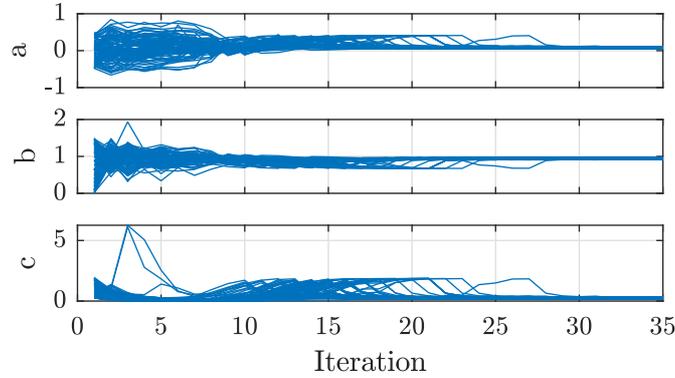


Figure 2: Parameter trajectory vs. iteration count for 100 random initial estimates using the proposed method.

## 6.2 Differential Drive Robot

In this example, we consider a continuous-time model of a differential drive robot given by

$$\begin{bmatrix} \dot{q}_1(t) \\ \dot{q}_2(t) \\ \dot{q}_3(t) \\ \dot{p}_1(t) \\ \dot{p}_2(t) \end{bmatrix} = \begin{bmatrix} \frac{\cos(q_3(t))p_1(t)}{m} \\ \frac{\sin(q_3(t))p_1(t)}{m} \\ \frac{p_2(t)}{J+ml^2} \\ \frac{-r_1 p_1(t)}{m} - \frac{mlp_2^2(t)}{(J+ml^2)^2} + u_1(t) + u_2(t) \\ \frac{(lp_1(t)-r_2)p_2(t)}{J+ml^2} + au_1(t) - au_2(t) \end{bmatrix},$$

where  $r_1 = 1$ ,  $r_2 = 1$ ,  $a = 0.5$ ,  $m = 5$ ,  $J = 0.2$ ,  $l = 0.15$ ,  $u_1(t)$  is the force applied to the left wheel,  $u_2(t)$  is the force applied to the right wheel, and the state vector  $x(t) = [q_1(t), q_2(t), q_3(t), p_1(t), p_2(t)]^T$  con-

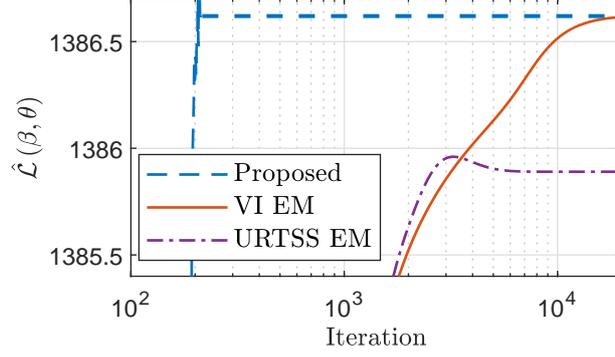


Figure 3: Comparison of proposed, VI-EM, and URTSS-EM approaches to identification on a simulated robot example. Proposed extended past the 194 iterations required for convergence for clarity.

sists of x-position, y-position, heading, linear momentum, and angular momentum states, respectively. A 50 s simulated trajectory is generated using an ODE solver disturbed by noise sampled from  $\mathcal{N}(0, Q)$  where  $Q = \text{diag}(0.001, 0.001, 1.745 \times 10^{-3}, 0.001, 0.001)$  at 0.1 s intervals. Measurements at each interval are obtained according to

$$y_k = [q_1(t), q_2(t), q_3(t)]^T + e_k, \quad e_k \sim \mathcal{N}(0, R),$$

where  $R = \text{diag}(0.1^2, 0.1^2, 0.0349^2)$ . As the system identification methods considered in this paper are all for discrete-time models, a Euler discretisation of the continuous-time dynamics over each 0.1 s interval was used to obtain a model in the form of (37).

### 6.2.1 Convergence

This section examines the convergence when estimating  $m$ ,  $l$ ,  $r$ , and a decoupled additive Gaussian covariance. For the proposed method, KNITRO [6] is used to perform the optimisation. The results are compared against two assumed Gaussian EM approaches that approximate the smoothing step using the URTSS and the VI smoother in [10]; these are denoted URTSS-EM and VI-EM, respectively and differ in how the assumed Gaussian smoothed density is obtained. From Section 3.2, the VI-EM corresponds to variational EM. As such, VI-EM is expected to converge towards the same values as the proposed method, albeit at a slower rate. For each method an initial parameter estimate of  $m = 10$ ,  $J = 4$ ,  $l = 0.3$ , and

$$\Pi = \text{diag}(0.01^2, 0.01^2, 0.0035^2, 0.01^2, 0.01^2, 0.1^2, 0.1^2, 0.0349^2),$$

is used. For the proposed method, the state distributions have been initialised using the filter developed in [10].

The proposed method converged to a locally optimal solution in 194 iterations in 297 s; neither the URTSS EM nor the VI EM approaches converged within the  $20 \times 10^3$  iteration limit, which required 1.07 h and 5.81 h, respectively, to reach. Fig. 3 shows the cost obtained using each method as a function of iteration count and illustrates the 194 iterations of the proposed method outperformed the  $20 \times 10^3$  iterations of both the EM approaches. As expected, VI-EM asymptotically approaches the cost achieved by the proposed method. Contrarily, URTSS-EM does not approach this cost; neither does it maintain the desired monotonic behaviour of EM. This highlights the benefit of following the more principled approach regarding any approximations introduced.

### 6.2.2 Robustness to Initialisation

This section examines the robustness of the proposed method to the initial parameter estimate. For this, estimation of  $m$ ,  $l$ ,  $r$ , and coupled noise term  $\Pi$  is considered from 75 differing parameter initialisations given by

$$m \sim \mathcal{U}(0.5, 15), \quad J \sim \mathcal{U}(0.01, 10), \quad l \sim \mathcal{U}(0.01, 0.5), \\ \Pi = \text{diag}(0.1^2, 0.1^2, 0.0349^2, 0.1^2, 0.1^2, 0.5^2, 0.5^2, 0.1745^2).$$

The initial state distributions for each optimisation were obtained by running the URTSS from the initial parameter estimates.

Table 2 shows a subset of the estimated parameters for this experiment, where  $\Pi_{ii}$  indicates the  $i^{\text{th}}$  diagonal element of  $\Pi$ . In the third column, the mean parameter estimate for each of the 75 trials is shown and is an effective estimate

Table 2: Parameter estimate for 75 differing initial parameter estimates to illustrate robustness.

Parameter	True	Estimated	Maximum difference
$m$	5	5.07	$6.15 \times 10^{-9}$
$J$	2	2.24	$1.20 \times 10^{-8}$
$l$	0.15	0.14	$6.64 \times 10^{-10}$
$\Pi_{11}$	$1.00 \times 10^{-3}$	$1.12 \times 10^{-3}$	$5.75 \times 10^{-9}$
$\Pi_{22}$	$1.00 \times 10^{-3}$	$1.41 \times 10^{-3}$	$3.04 \times 10^{-8}$
$\Pi_{33}$	$1.74 \times 10^{-3}$	$1.85 \times 10^{-3}$	$3.76 \times 10^{-9}$
$\Pi_{44}$	$1.00 \times 10^{-3}$	$4.90 \times 10^{-3}$	$1.52 \times 10^{-7}$
$\Pi_{55}$	$1.00 \times 10^{-3}$	$1.40 \times 10^{-3}$	$1.15 \times 10^{-8}$
$\Pi_{66}$	$1.00 \times 10^{-2}$	$1.00 \times 10^{-2}$	$4.69 \times 10^{-7}$
$\Pi_{77}$	$1.00 \times 10^{-2}$	$8.25 \times 10^{-3}$	$4.14 \times 10^{-7}$
$\Pi_{88}$	$1.22 \times 10^{-3}$	$1.38 \times 10^{-3}$	$1.33 \times 10^{-7}$

of the true parameters. The fourth column shows the maximum difference between this mean and each trial and indicates that each of the 75 trials converged to the same parameter estimate, subject to a small tolerance.

These results indicate the robustness of the proposed approach and highlights that, at least for this example, the proposed method does not suffer from a potentially large quantity of undesirable local minima. The estimates for the off-diagonal elements of  $\Pi$  performed similarly as the parameters shown in Table 2. Due to the multi-variable nature of  $\Pi$ , these numerical values are excluded for brevity.

### 6.3 Inverted Pendulum

This section considers a rotational inverted pendulum, or Furata pendulum [13], using data collected from a Quanser QUBE-Servo 2. Letting the state vector used to model the Furata pendulum be  $x = [\vartheta \ \alpha \ \dot{\vartheta} \ \dot{\alpha}]^\top$ , where  $\vartheta$  and  $\alpha$  are the base arm and pendulum angles, respectively, and the controllable input to the system is the motor voltage  $V_m$ . Then, the continuous time dynamics are then given by

$$\begin{aligned}
 M(\alpha) \begin{bmatrix} \ddot{\vartheta} \\ \ddot{\alpha} \end{bmatrix} + \nu(\dot{\vartheta}, \dot{\alpha}) \begin{bmatrix} \dot{\vartheta} \\ \dot{\alpha} \end{bmatrix} &= \begin{bmatrix} \frac{k_m(V_m - k_m \dot{\vartheta})}{R_m} - D_r \dot{\vartheta} \\ -\frac{1}{2} m_p L_p g \sin(\alpha) - D_p \dot{\alpha} \end{bmatrix}, \\
 M(\alpha) &= \begin{bmatrix} m_p L_r^2 + \frac{1}{4} m_p L_p^2 (1 - \cos(\alpha)^2) + J_r & \frac{1}{2} m_p L_p L_r \cos(\alpha) \\ \frac{1}{2} m_p L_p L_r \cos(\alpha) & J_p + \frac{1}{4} m_p L_p^2 \end{bmatrix}, \\
 \nu(\dot{\vartheta}, \dot{\alpha}) &= \begin{bmatrix} \frac{1}{2} m_p L_p^2 \sin(\alpha) \cos(\alpha) \dot{\alpha} & -\frac{1}{2} m_p L_p L_r \sin(\alpha) \dot{\alpha} \\ -\frac{1}{4} m_p L_p^2 \cos(\alpha) \sin(\alpha) \dot{\vartheta} & 0 \end{bmatrix},
 \end{aligned}$$

where  $m_p$  is the pendulum mass,  $L_r$ ,  $L_p$  and the rod and pendulum lengths,  $J_r$ ,  $J_p$  are the rod and pendulum inertias,  $R_m$  and  $k_m$  are the motor resistance and constant,  $D_p$  and  $D_r$  are the pendulum and arm damping constants. The considered process model is a two-step Euler discretisation of these continuous-time dynamics over an 8 ms sampling time subsequently disturbed by noise  $v_k$ . The available measurements are from encoders on the rod and pendulum angle, and the motor current. The resulting measurement model is

$$y_k = \left[ \vartheta \ \alpha \ \frac{V_m - k_m \dot{\vartheta}}{R_m} \right]^\top + e_k, \quad \begin{bmatrix} v_k^\top & e_k^\top \end{bmatrix}^\top \sim \mathcal{N}(0, \Pi),$$

which, together with the discrete process model obtained using the two-step Euler integration, results in a model in the form of (37). We are interested in estimating the parameters  $\theta = (D_r, D_p, J_r, J_p, k_m, R_m, \Pi)$ .

This section compares the effectiveness and robustness of the additive noise specialisation of the proposed method with CPF-SAEM from 50 differing parameter initialisations. These initialisations sampled from

$$\begin{aligned}
 J_r &\sim \mathcal{U}(1 \times 10^{-7}, 0.01), & D_r &\sim \mathcal{U}(1 \times 10^{-7}, 0.01), \\
 J_p &\sim \mathcal{U}(1 \times 10^{-7}, 0.01), & D_p &\sim \mathcal{U}(1 \times 10^{-7}, 0.01), \\
 k_m &\sim \mathcal{U}(0.001, 1), & R_m &\sim \mathcal{U}(4, 15),
 \end{aligned}$$

with the noise covariance term initialised as

$$\Pi = \text{diag}\left(7.6154 \times 10^{-5}, 7.6154 \times 10^{-5}, 0.0012, 0.0012, 3.0462 \times 10^{-4}, 3.0462 \times 10^{-4}, 0.01\right),$$

Table 3: Mean and standard deviation of the parameter estimates for each successful trial from 50 differing initial values. Twelve of the CPF-SAEM runs were unsuccessful and have been censored.

Parameter	Proposed	CPF-SAEM
$J_r$	$1.07 \times 10^{-4} \pm 8.17 \times 10^{-11}$	$1.08 \times 10^{-4} \pm 3.27 \times 10^{-7}$
$J_p$	$2.92 \times 10^{-5} \pm 3.32 \times 10^{-11}$	$2.90 \times 10^{-5} \pm 7.80 \times 10^{-8}$
$k_m$	$4.57 \times 10^{-2} \pm 2.73 \times 10^{-8}$	$4.99 \times 10^{-2} \pm 4.60 \times 10^{-4}$
$R_m$	$9.59 \pm 4.06 \times 10^{-7}$	$1.02 \times 10^1 \pm 7.42 \times 10^{-2}$
$D_p$	$4.68 \times 10^{-5} \pm 7.66 \times 10^{-11}$	$5.03 \times 10^{-5} \pm 7.24 \times 10^{-7}$
$D_r$	$2.81 \times 10^{-4} \pm 1.87 \times 10^{-9}$	$2.15 \times 10^{-4} \pm 4.32 \times 10^{-6}$
$\Pi_{11}$	$3.94 \times 10^{-6} \pm 1.58 \times 10^{-10}$	$6.12 \times 10^{-6} \pm 3.41 \times 10^{-7}$
$\Pi_{22}$	$1.99 \times 10^{-6} \pm 2.16 \times 10^{-10}$	$4.63 \times 10^{-6} \pm 2.98 \times 10^{-7}$
$\Pi_{33}$	$2.72 \times 10^{-2} \pm 2.71 \times 10^{-7}$	$2.58 \times 10^{-2} \pm 1.10 \times 10^{-3}$
$\Pi_{44}$	$3.10 \times 10^{-2} \pm 2.17 \times 10^{-7}$	$3.08 \times 10^{-2} \pm 1.61 \times 10^{-3}$
$\Pi_{55}$	$2.42 \times 10^{-6} \pm 5.21 \times 10^{-10}$	$1.31 \times 10^{-6} \pm 9.42 \times 10^{-8}$
$\Pi_{66}$	$1.75 \times 10^{-6} \pm 2.03 \times 10^{-10}$	$1.28 \times 10^{-6} \pm 6.68 \times 10^{-8}$
$\Pi_{77}$	$1.91 \times 10^{-2} \pm 1.02 \times 10^{-9}$	$1.96 \times 10^{-2} \pm 1.18 \times 10^{-4}$

on a single data set consisting of 375 measurements, during which the pendulum undergoes full rotations. Note that the noise terms are only utilised for the CPF-SAEM approach; in accordance with Section 5.1 they are not required for the additive noise version of the proposed approach.

For these trials, the initial position and velocity states used the respective measurements and a finite-difference approximation for the means, and standard deviations of  $1^\circ$  and  $10^\circ \text{s}^{-1}$ , respectively. For the proposed method, a relative function tolerance of  $1 \times 10^{-8}$  on  $\hat{\mathcal{L}}_{\text{AG}}(\phi, \beta)$  was the termination condition, which required a median of 170 iterations and 112 s to achieve. Importantly, the proposed method proved to be robust to the differing initialisations, with all trials converging to the same value subject to a small numeric tolerance.

The CPF-SAEM approach was performed using a fully adapted auxiliary particle filter [37] with 100 particles, 50 particles for the backward simulation smoother [15], and a burn-in of 100 iterations before applying the stochastic approximation. The iterations terminated when the change in all parameter estimates between successive iterations fell below  $1 \times 10^{-3}$  for five of the last ten iterations, which required a median of 156 iterations and 308 s to achieve. In contrast to the proposed method, CPF-SAEM failed on 12 of the 50 trials.

Table 3 shows the mean and standard deviation of a subset of the estimated parameter values of each successful trial, where  $\Pi_{ii}$  indicates the  $i^{\text{th}}$  diagonal element of  $\Pi$ . The closeness of the estimates for the proposed method to the successful CPF-SAEM trials illustrates that, despite the approximations introduced, similar parameter estimates are obtained. While, for brevity, the numeric values of off-diagonal elements of  $\Pi$  are omitted here, they were similarity estimated.

This section has demonstrated the applicability of the developed system identification approach on a real unstable system of practical interest. In particular, this example has shown the robustness to initialisation, computational efficiency, and ability to provide parameter estimates close to those of particle methods.

## 7 Conclusion

The contribution of this paper is to present a VI based approach to system identification for nonlinear state-space models. The resulting system identification approach consists of a single, deterministic, optimisation problem. Due to the assumed density and constrained parametrisation chosen, this optimisation problem is of a standard form and is efficiently solved using readily available exact first- and second-order derivatives. A specialisation for systems with additive Gaussian noise was also presented. The proposed method has been numerically examined on a range of simulated and real examples to illustrate the robustness and effectiveness and is compared favourably to state-of-the-art alternatives.

## References

- [1] Joel A. E. Andersson, Joris Gillis, Greg Horn, James B. Rawlings, and Moritz Diehl. CasADi: a software framework for nonlinear optimization and optimal control. *Mathematical Programming Computation*, 11(1):1–36, July 2018.

- [2] K. J. Åström. Maximum likelihood and prediction error methods. *Automatica*, 16(5):551–574, September 1980.
- [3] Matthew Beal and Zoubin Ghahramani. The Variational Bayesian EM Algorithm for Incomplete Data: with Application to Scoring Graphical Model Structures. *Bayesian Statistics*, 7, 2003.
- [4] Matthew J. Beal. *Variational Algorithms for Approximate Bayesian Inference*. phdthesis, Gatsby Computational Neuroscience Unit, University College London, 2003.
- [5] D. Blei, A. Kucukelbir, and J. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- [6] Richard H. Byrd, Jorge Nocedal, and Richard A. Waltz. Knitro: An Integrated Package for Nonlinear Optimization. In *Nonconvex Optimization and Its Applications*, pages 35–59. Springer US, 2006.
- [7] Olivier Cappe, Eric Moulines, and Tobias Ryden. *Inference in Hidden Markov Models*. Springer, 2007.
- [8] S. B. Chitralakha, J. Prakash, H. Raghavan, R. B. Gopaluni, and S. L. Shah. Comparison of Expectation-Maximization based parameter estimation using Particle Filter, Unscented and Extended Kalman Filtering techniques. *IFAC Proceedings Volumes*, 42(10):804–809, 2009.
- [9] Andrew R. Conn, Nicholas I. M. Gould, and Philippe L. Toint. *Trust Region Methods*. Society for Industrial and Applied Mathematics, January 2000.
- [10] Jarrad Courts, Adrian Wills, and Thomas B. Schön. Gaussian Variational State Estimation for Nonlinear State-Space Models. *IEEE Transactions on Signal Processing*, 69:5979–5993, 2021.
- [11] Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the EM algorithm. *The Annals of Statistics*, 27(1), March 1999.
- [12] 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*, 39(1):1–38, 1977.
- [13] K. Furuta, M. Yamakita, and S. Kobayashi. Swing-up Control of Inverted Pendulum Using Pseudo-State Feedback. *Proceedings of the Institution of Mechanical Engineers, Part I: Journal of Systems and Control Engineering*, 206(4):263–269, November 1992.
- [14] Matej Gašperin and Dani Juričić. Application of Unscented Transformation in Nonlinear System Identification. *IFAC Proceedings Volumes*, 44(1):4428–4433, January 2011.
- [15] Simon J. Godsill, Arnaud Doucet, and Mike West. Monte Carlo Smoothing for Nonlinear Time Series. *Journal of the American Statistical Association*, 99(465):156–168, March 2004.
- [16] Andreas Griewank and Andrea Walther. *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*. Society for Industrial and Applied Mathematics, January 2008.
- [17] Harold V. Henderson and S. R. Searle. Vec and vech operators for matrices, with some uses in jacobians and multivariate statistics. *Canadian Journal of Statistics*, 7(1):65–81, 1979.
- [18] Rolf Isermann and Marco Münchhof. *Identification of Dynamic Systems*. Springer Berlin Heidelberg, 2011.
- [19] Bin Jia, Ming Xin, and Yang Cheng. High-degree cubature Kalman filter. *Automatica*, 49(2):510–518, February 2013.
- [20] Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An Introduction to Variational Methods for Graphical Models. *Machine Learning*, 37(2):183–233, November 1999.
- [21] Simon J. Julier and Jeffrey K. Uhlmann. New extension of the Kalman filter to nonlinear systems. In Ivan Kadar, editor, *Signal Processing, Sensor Fusion, and Target Recognition VI*. SPIE, July 1997.
- [22] Juho Kokkala, Arno Solin, and Simo Särkkä. Expectation Maximization Based Parameter Estimation by Sigma-Point and Particle Smoothing. In *17<sup>th</sup> International Conference on Information Fusion (FUSION)*, July 2014.
- [23] Juho Kokkala, Arno Solin, and Simo Särkkä. Sigma-Point Filtering and Smoothing Based Parameter Estimation in Nonlinear Dynamic Systems. *Journal of Advances in Information Fusion*, 11(1):15–30, 2016.
- [24] S. Kullback and R. A. Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86, March 1951.

- [25] Martin Lindfors and Tianshi Chen. Regularized LTI system identification in the presence of outliers: A variational EM approach. *Automatica*, 121:109152, November 2020.
- [26] Fredrik Lindsten. An efficient stochastic approximation EM algorithm using conditional particle filters. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, May 2013.
- [27] L. Ljung. *System Identification: Theory for the User*. Prentice Hall information and system sciences series. Prentice Hall PTR, 1999.
- [28] L. Ljung. On Convexification of System Identification Criteria. *Automation and Remote Control*, 80(9):1591–1606, September 2019.
- [29] Lennart Ljung. Some aspects on nonlinear system identification. *IFAC Proceedings Volumes*, 39(1):553–564, 2006.
- [30] Lennart Ljung. Perspectives on system identification. *Annual Reviews in Control*, 34(1):1–12, April 2010.
- [31] Jan R. Magnus and Jan R. Magnus. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Wiley, February 2019.
- [32] MATLAB. *Optimization Toolbox Release 2018b*. The MathWorks, Inc., Natick, Massachusetts, United States, 2018.
- [33] Jorge J. Moré and D. C. Sorensen. Computing a Trust Region Step. *SIAM Journal on Scientific and Statistical Computing*, 4(3):553–572, September 1983.
- [34] Radford M. Neal and Geoffrey E. Hinton. A View of the EM Algorithm that Justifies Incremental, Sparse, and other Variants. In *Learning in Graphical Models*, pages 355–368. Springer Netherlands, 1998.
- [35] Brett Ninness. Some System Identification Challenges and Approaches. *IFAC Proceedings Volumes*, 42(10):1–20, 2009.
- [36] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer New York, second edition, 2006.
- [37] Michael K. Pitt and Neil Shephard. Filtering via Simulation: Auxiliary Particle Filters. *Journal of the American Statistical Association*, 94(446):590–599, June 1999.
- [38] Antônio H. Ribeiro, Koen Tiels, Jack Umenberger, Thomas B. Schön, and Luis A. Aguirre. On the smoothness of nonlinear system identification. *Automatica*, 121:109158, November 2020.
- [39] Riccardo S. Risuleo, Giulio Bottegal, and Håkan Hjalmarsson. Variational Bayes identification of acyclic dynamic networks. *IFAC-PapersOnLine*, 50(1):10556–10561, July 2017.
- [40] Simo Särkkä. *Bayesian Filtering and Smoothing*. Institute of Mathematical Statistics Textbooks. Cambridge University Press, 2013.
- [41] Simo Särkkä, Jouni Hartikainen, Lennart Svensson, and Fredrik Sandblom. Gaussian process quadratures in nonlinear sigma-point filtering and smoothing. In *17<sup>th</sup> International Conference on Information Fusion (FUSION)*, pages 1–8, 2014.
- [42] Thomas B. Schön, Fredrik Lindsten, Johan Dahlin, Johan Wågberg, Christian A. Naesseth, Andreas Svensson, and Liang Dai. Sequential Monte Carlo Methods for System Identification. *IFAC-PapersOnLine*, 48(28):775–786, 2015.
- [43] Thomas B. Schön, Adrian Wills, and Brett Ninness. System identification of nonlinear state-space models. *Automatica*, 47(1):39–49, 2011.
- [44] Dimitris G. Tzikas, Aristidis C. Likas, and Nikolaos P. Galatsanos. The variational approximation for Bayesian inference. *IEEE Signal Processing Magazine*, 25(6):131–146, November 2008.
- [45] Michail D. Vrettas, Yuan Shen, and Dan Cornford. Derivations of Variational Gaussian Process Approximation Framework. Technical report, Aston University, March 2008.
- [46] E. A. Wan and R. Van Der Merwe. The unscented Kalman filter for nonlinear estimation. In *Proceedings of the IEEE Adaptive Systems for Signal Processing, Communications, and Control Symposium*. IEEE, October 2000.
- [47] Adrian Wills, Thomas Schön, Lennart Ljung, and Brett Ninness. Identification of Hammerstein–Wiener Models. *Automatica*, 49(1):70–81, January 2013.

- [48] Adrian G. Wills and Thomas B. Schön. Stochastic quasi-Newton with line-search regularisation. *Automatica*, 127:109503, May 2021.
- [49] Jeremy Nathan Wong, David Juny Yoon, Angela P. Schoellig, and Timothy D. Barfoot. Variational Inference With Parameter Learning Applied to Vehicle Trajectory Estimation. *IEEE Robotics and Automation Letters*, 5(4):5291–5298, October 2020.

## A Proof of Lemma 3

*Proof.* The cost function  $\mathcal{L}((\alpha, \gamma), \theta)$  is given by

$$\mathcal{L}((\alpha, \gamma), \theta) = \int q_\eta(\mathbf{x}) \log \frac{p_\theta(\mathbf{x}, \mathbf{y})}{q_\eta(\mathbf{x})} d\mathbf{x} \quad (44)$$

$$\begin{aligned} &= \int q_\eta(\mathbf{x}) \log p_\theta(\mathbf{x}, \mathbf{y}) d\mathbf{x} \\ &\quad - \int q_\eta(\mathbf{x}) \log q_\eta(\mathbf{x}) d\mathbf{x}. \end{aligned} \quad (45)$$

Due to the Markovian nature of state-space models, conditional probability, and that  $q_\eta(\mathbf{x})$  is a probability distribution, similar to [43, 8, 40, 45, 10], we can write

$$\int q_\eta(\mathbf{x}) \log p_\theta(\mathbf{x}, \mathbf{y}) d\mathbf{x} = I_1(\alpha) + I_2(\alpha, \theta). \quad (46)$$

This is as only the pairwise joint distributions are required in the calculation of the first integral of (45), i.e. it is independent of  $\gamma$ .

As shown in [10], the optimal assumed density  $q_{\eta^*}(\mathbf{x})$  factors according to

$$q_{\eta^*}(\mathbf{x}) = q_{\eta^*}(x_1) \prod_{k=1}^T q_{\eta^*}(x_{k+1} | x_k), \quad (47)$$

and that  $\gamma^* = g(\alpha^*)$ , where  $g(\alpha)$  is a function (defined by the process detailed in Appendix C of [10]) that provides a  $\gamma$  such that the above factorisation is satisfied. As such, by substituting  $\gamma = g(\alpha)$ , the second integral in (45) is given by

$$\int q_\eta(\mathbf{x}) \log q_\eta(\mathbf{x}) d\mathbf{x} = I_3(\alpha). \quad (48)$$

Therefore, we can equivalently find  $\theta^*$  and  $\alpha^*$  from either (21) or (23).  $\square$