



Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

# Scalable Variational Inference for Stochastic Differential Equations

Master Thesis

Ruifeng Xu

September 1, 2017

Supervisor: Prof. Dr. Joachim M. Buhmann

Advisors: Stefan Bauer & Nico S. Gorbach

Department of Computer Science, ETH Zürich



---

## Abstract

State and parameter estimation in dynamical systems based on sparse, discrete observations is a topical yet challenging problem. Traditional methods suffer from extremely high computational costs due to the need to carry out explicit numerical integration after parameter adaptation. In contrast, the recently proposed gradient matching with Gaussian process model is a promising tool. It is a grid-free inference technique that also eliminates the dependency on numerical integration. However, due to the intractability of the posterior, approximate inference techniques must be used. Sampling-based solutions fall short in this case since most real-world dynamical systems are high-dimensional. On the other hand, variational approaches have shown their potential both in terms of prediction accuracy and runtime performance, even in situations when the system is only partially observed. Extending the state-of-art variational gradient matching framework, this work further improves the flexibility of the inference algorithm by relaxing the structural assumption on the dynamical systems. Support for positivity constraints on the state and parameters that are common in many biochemical and physical applications is also introduced. Finally, a highly efficient parallel solution is devised to address problems involving stochastic differential equations.



---

# Contents

---

<b>Contents</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Challenges . . . . .	3
1.3 Contributions . . . . .	3
1.4 Organization . . . . .	3
<b>2 Inference in Dynamical Systems</b>	<b>5</b>
2.1 Deterministic dynamical systems . . . . .	5
2.2 Random dynamical systems . . . . .	7
2.3 Related work . . . . .	8
<b>3 Gradient Matching with Gaussian Processes</b>	<b>11</b>
3.1 Preliminary . . . . .	11
3.1.1 Gaussian process regression . . . . .	11
3.1.2 Kullback-Leibler divergence . . . . .	12
3.1.3 Mean-field variational inference . . . . .	13
3.2 Sampling-based gradient matching with Gaussian processes .	16
3.3 Variational gradient matching with Gaussian processes . . . .	20
<b>4 Laplace Mean-Field Approximation</b>	<b>25</b>
4.1 Laplace approximation . . . . .	26
4.2 Laplace mean-field approximation . . . . .	27
4.3 Derivation for the gradients and Hessians . . . . .	29
4.4 Positivity constraints . . . . .	31
<b>5 Extension to Random Dynamical Systems</b>	<b>33</b>
5.1 Random ordinary differential equations . . . . .	33
5.2 Doss-Sussmann/Imkeller-Schmalfuss correspondence . . . .	34
5.3 Laplace mean-field for random dynamical systems . . . . .	35

## CONTENTS

---

<b>6</b>	<b>Experiments</b>	<b>37</b>
6.1	Implementation . . . . .	37
6.2	Lotka-Volterra model . . . . .	38
6.3	Protein signalling transduction pathway . . . . .	42
6.4	Lorenz 96 model . . . . .	45
6.5	Lorenz 63 model . . . . .	54
<b>7</b>	<b>Conclusion</b>	<b>61</b>
	<b>Bibliography</b>	<b>63</b>

## Chapter 1

---

# Introduction

---

Modelling complicated natural phenomena using mathematical abstractions is a common practice in modern scientific research and real-world applications. A successful model of a phenomenon helps us interpret its key features while omits extraneous details (Babtie et al., 2014). For various dynamic processes underlying a broad range of fields including chemistry, physics, biology, economics, meteorology, etc., one convenient and useful abstraction is to describe them as a set of *ordinary differential equations (ODEs)* or *stochastic (ordinary) differential equations (SDEs)* (Ellner and Guckenheimer, 2011; Gardiner, 2009), which are also called *deterministic dynamical systems* and *random dynamical systems* respectively.

Using the term *dynamical system* to refer to the above two systems in general, this work addresses the topical yet challenging problem of statistical inference of the states and parameters of dynamical systems given noisy, sparse or even incomplete states observation. By extending the state-of-art approaches that combine the techniques of Gaussian process regression, gradient matching, and variational inference, the aim is to build an inference pipeline that operates efficiently, predicts accurately and scales to large systems.

### 1.1 Motivation

To understand the importance of the problem, it is possible to find numerous alternative models that describe the observations unless there is enough domain expertise to designate a specific model (Babtie et al., 2014). In order to select a good candidate, each of the alternatives must be trained using the observations, and then the results will be evaluated according to certain criteria. Even if the general underlying process is known, the parameters that control the dynamics of the system still need to be inferred from the observations. It is therefore critical for the data fitting process to be accu-

rate and robust. Meanwhile, the procedure should also be performant so that complex models can be trained within reasonable time constraints. For more concrete illustration of modelling using dynamical systems, consider the following two examples.

**Deterministic dynamical systems** One essential task in system biology is to compare and select an appropriate model to characterize a biochemical system, e.g. protein signalling transduction pathway (Vysheirsky and Girolami, 2007). Typically, the structure of the system can be viewed as a network of biochemical reactions, which can be formally described by a group of nonlinear ODEs; the transitions among the network components, e.g. protein species, are determined by a set of kinetic parameters (Macdonald and Husmeier, 2015). During experiments, only concentrations on the species over a time span are observed. Therefore, inference of the true parameter values or even the *a posteriori* distributions for different models given the experimental data is of crucial importance for model selection purposes.

**Random dynamical systems** In reality, it is often that the mathematical model is not able to capture all the characteristics of the phenomenon. If we allow some randomness inside a deterministic dynamical system, we would get a random dynamical system (Øksendal, 2013). An example where random dynamical systems have a long history of application is weather forecasting, where the continuous evolution of the atmosphere is described by discretized quantities like pressure, temperature, wind speed, etc., measured at fixed intervals (Archambeau et al., 2007). Therefore, it is reasonable to model the rest of the unknown dynamics as the noise process within the SDEs. For such systems, useful prediction depends not only on the realistic modelling of the atmospheric environment, but also on the precise estimation of the initial conditions (Kalnay, 2003), due to the high sensitivity of future states on the initial conditions that is also known as the famous *butterfly effect* (Lorenz, 2000).

To summarize, the relative importance of state versus parameter estimation varies, depending on the specific application. State estimation is more important for short-term predictions such as weather forecasting, while parameter estimation is more important for long-term targets such as climate pattern modelling (Vrettas et al., 2015).

This work places more emphasis on the inference of system parameters since they shed light on the internal mechanisms of a dynamical system. Nevertheless, the states are also simultaneously estimated as a consequence of the design of the algorithm, which will be shown in later chapters.



## 1.2 Challenges

Although solving statistical inference problems involving ODEs and SDEs is useful in practice, but from a technical point of view, such problems involve many technical challenges as below.

First, since closed-form solutions do not exist for most ODEs and not at all for SDEs, conventional methods based on explicit numerical integrations are computationally expensive, which renders them impractical even for small-sized applications.

Second, the likelihood surfaces in the parameter space are likely multimodal due to nonlinearity within the dynamical systems (Calderhead et al., 2009), and may exhibit many local maxima, which makes the optimization challenging. From a Bayesian perspective, the intractability of the marginalization term is conventionally approximated using *Markov chain Monte Carlo* (MCMC) sampling schemes, which are in general very flexible but come at the cost of high computational intensity and onerous convergence analysis. Hence, the applicability of sampling-based approaches is largely subject to the dimensionality of the system under inference (Vrettas et al., 2015).

Last, in most of the realistic scenarios, only corrupted, and sometimes even only sparse and partial observations are available. Devising inference algorithms that are robust in such situations is itself a difficult topic across all machine learning paradigms.

## 1.3 Contributions

The contributions of this work are as follows. The Laplace mean-field approximation is proposed to relax the structural assumption on the dynamical systems. Through reparameterization on the optimization objectives, positivity constraints on the states and parameters are supported, which are essential for many real-world application. By utilizing the Doss-Sussmann/Imkeller-Schmalfuss correspondence, the Laplace mean-field approximation is further extended to devise a distributed inference method to infer states and parameters of the SDEs. A brand new Python based solution is also implemented.

## 1.4 Organization

This thesis is organized as follows. Chapter 2 gives a general introduction to deterministic and random dynamical systems and provides a brief review of other related work on inference in dynamical systems. Chapter 3 introduces in detail the gradient matching with Gaussian process framework

(Calderhead et al., 2009; Dondelinger et al., 2013) and its recent improvement (Gorbach et al., 2016, 2017) based on variational inference, which are the foundation of this work. In Chapter 4, the Laplace approximation technique is applied to the variational gradient matching model to derive a new solution that relaxes the structural assumption about the ODEs, and further introduces positivity constraints on the states and parameters through reparameterization. In Chapter 5, an ensemble-like inference solution for SDEs is derived by using the Doss-Sussmann/Imkeller-Schmalzfuss correspondence. The examination and comparison of the accuracy, performance and scalability of the solutions proposed in this work are conducted in Chapter 6. Lastly, Chapter 7 draws the conclusion.

## Chapter 2

---

# Inference in Dynamical Systems

---

This chapter reviews the necessary basics about ODEs (Section 2.1) and SDEs (Section 2.2) with an emphasis on the task of inferring states and parameters of dynamical systems given noisy, sparse or even incomplete state observations. Comprehensive discussions about ODEs and SDEs together with their numerical solutions can be found in textbooks such as Butcher (2016) and Øksendal (2013). Furthermore, the notations used to describe the dynamical systems, and the model under noisy observations are also introduced, which will be used throughout this work. In Section 2.3, a survey of related work is provided. The details about the inference technique using variational gradient matching with Gaussian processes, which are the core foundations of this work, are examined in Chapter 3.

### 2.1 Deterministic dynamical systems

A  $K$ -dimensional deterministic dynamical system, i.e. a system with  $K$  states, can be described by a set of ODEs as follows:

$$\dot{\mathbf{x}}(t) = \frac{d\mathbf{x}(t)}{dt} = \mathbf{f}(\mathbf{x}(t), \boldsymbol{\theta}) \quad (2.1)$$

where  $\mathbf{x}(t) = [x_1(t), \dots, x_K(t)]^\top \in \mathbb{R}^K$  is a vector containing the  $K$  states of the system at time point  $t$  with their time derivatives collectively denoted by  $\dot{\mathbf{x}}(t) = \frac{d\mathbf{x}(t)}{dt} = [f_1(\mathbf{x}(t), \boldsymbol{\theta}), \dots, f_K(\mathbf{x}(t), \boldsymbol{\theta})]^\top \in \mathbb{R}^K$ , and  $\mathbf{f} : \mathbb{R}^K \mapsto \mathbb{R}^K$  encodes the functional relationship between the states and their derivatives over time that is in turn governed by a vector of parameters  $\boldsymbol{\theta} = [\theta_1, \dots, \theta_M]^\top \in \mathbb{R}^M$ . Note that in general,  $\mathbf{f}$  may have direct dependency on time, which is suppressed here to simplify the notations.

From Butcher (2016), the specification of the ODEs alone is not interesting since it generally does not guarantee a unique solution. But if the *initial condition*  $\mathbf{x}(t_0) = [x_1(t_0), \dots, x_K(t_0)]^\top$  is given, then together with (Eq. 2.1), they

define a problem known as the *initial value problem*, where the goal is to solve the differential equations to approximate the future states of the dynamical system. Three important aspects about the initial value problem are the existence of a solution, the uniqueness of the solution, and the sensitivity of the solution due to small perturbations to the initial condition.

Within the context of this work, the initial problem can be therefore specified as the estimation of the states and the parameters of the ODEs based on state observations that are usually contaminated by noise. In light of this, the following paragraphs introduce the notations and the probabilistic model of a dynamical system under noisy observations, which will be used later in this work.

### Noisy observation model

Suppose for the  $K$ -dimensional deterministic dynamical system given by (Eq. 2.1), we have a sequence of noisy observations  $Y = [\mathbf{y}(t_1), \dots, \mathbf{y}(t_N)] \in \mathbb{R}^{K \times N}$  over  $N$  time points, whose the corresponding true states values are  $X = [\mathbf{x}(t_1), \dots, \mathbf{x}(t_N)] \in \mathbb{R}^{K \times N}$  such that  $\mathbf{y}(t_n) = \mathbf{x}(t_n) + \boldsymbol{\varepsilon}(t_n)$  for  $n = 1, \dots, N$ , where  $\boldsymbol{\varepsilon}(t_n) \in \mathbb{R}^K$  denotes the observation noises for the  $K$  states  $\mathbf{x}(t_n)$  at time point  $t_n$ . The above description can be succinctly written in matrix notation as

$$Y = X + E$$

where  $Y$  and  $X$  are defined before, and  $E = [\boldsymbol{\varepsilon}(t_1), \dots, \boldsymbol{\varepsilon}(t_N)] \in \mathbb{R}^{K \times N}$ .

For simplicity, we assume that the observation noises  $\boldsymbol{\varepsilon}(t_n)$  at each time point are additive and state specific. Moreover, they follow an *independent and identically distributed (i.i.d.)* multivariate Gaussian distribution with zero mean and a diagonal covariance matrix across all time points, i.e.  $\boldsymbol{\varepsilon}_{(\cdot)} \sim \mathcal{N}(\mathbf{0}, \mathbf{D})$ , where  $D_{ik} = \delta(i, k)\sigma_k^2$  for  $i, k = 1, \dots, K$ ,  $\delta$  is the *Kronecker delta* function, and  $\sigma_k^2$  is the variance of the observation noise for state  $k$ .

Let  $y_k(t_n) \in \mathbb{R}$  and  $x_k(t_n) \in \mathbb{R}$  be the observation and true value for the  $k$ -th state at time point  $t_n$  respectively. We can then collectively use  $\mathbf{y}_k = [y_k(t_1), \dots, y_k(t_N)]^\top \in \mathbb{R}^N$  and  $\mathbf{x}_k = [x_k(t_1), \dots, x_k(t_N)]^\top \in \mathbb{R}^N$  to denote the sequence of observations and the true values for the  $k$ -th state over the  $N$  time points. With the above noise assumption and notations, we have

$$p(Y|X, \sigma) = \prod_k \mathcal{N}(\mathbf{y}_k | \mathbf{x}_k, \sigma_k^2 \mathbf{I}) \quad (2.2)$$

Note that the discussion within this work can be generalized to cases where only combinations of states are observed such that  $Y = \mathbf{H}X + E$ , where  $\mathbf{H}$  describes the relationship between the states and the observations. In order

to simplify the notation, we assume that  $\mathbf{H}$  is equal to the identity matrix  $\mathbf{I}$ . Furthermore, handling of the cases where some states are not observed is also possible and will be presented in the relevant sections later.

## 2.2 Random dynamical systems

As another important family of dynamical systems, random dynamical systems, also referred to as *diffusion processes*, have been widely applied to various domains because of their capability to incorporate unknown processes as internal noise processes (Vrettas et al., 2011). On a high level, a random dynamical system is a continuous time *Markov process* consisting of a deterministic part and a stochastic noise driven component (Riesinger et al., 2016). Such system is described by a set of SDEs and requires the special *stochastic calculus*.

Based on Øksendal (2013) and Vrettas et al. (2015), a  $K$ -dimensional SDE system defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  is represented in the *Itô* form as

$$dx(t) = f(x(t), \theta)dt + g(x(t), \rho)dW_t \quad (2.3)$$

where  $f : \mathbb{R}^K \mapsto \mathbb{R}^K$  is the deterministic *drift function* with *drift parameter* vector  $\theta = [\theta_1, \dots, \theta_M]^\top \in \mathbb{R}^M$ ,  $g : \mathbb{R}^K \mapsto \mathbb{R}^{K \times W}$  is the coefficient function for the noise process with *diffusion parameter* vector  $\rho = [\rho_1, \dots, \rho_V]^\top \in \mathbb{R}^V$ , and  $dW_t = [dW_t^1, \dots, dW_t^W]^\top \in \mathbb{R}^W$  is the differential of a standard *Wiener process*  $W_t$  of dimension  $W$ , i.e.  $dW_t \sim \mathcal{N}(\mathbf{0}, dt\mathbf{I})$ . In its integral form, the above equation is equivalent to

$$x(t_T) = x(t_0) + \int_0^T f(x(t_s), \theta)ds + \sum_i^W \int_0^T g_i(x(t_s), \rho)dW_s^i \quad (2.4)$$

where  $x(t_0)$  denotes the initial condition and  $g_i(x(t_s), \rho)$  is the  $i$ -th column of the corresponding noise coefficient matrix. As  $W_t$  is a stochastic process, each time we solve the above integration, we would obtain mostly likely a different *sample path*. Note that similar to the ODEs in (Eq. 2.1), the dependency of  $f$  and  $g$  on time  $t$  is suppressed here to unclutter the notations.

The above equations define a linear diffusion process with multiplicative noises. For simplicity, we consider only stochastic systems with state-specific, additive white noise in this work. A class of multiplicative noise models can be mapped to this model through reparametrization (Vrettas et al., 2011). This assumption means that we can simplify the SDEs to

$$dx(t) = f(x(t), \theta)dt + \Sigma^{\frac{1}{2}}dW_t \quad (2.5)$$

where  $\Sigma$  is the diagonal noise covariance matrix, i.e.  $\Sigma_{ik} = \delta(i, k)\rho_k^2$  for  $i, k = 1, \dots, K$ , and  $W_t$  becomes a standard  $K$ -dimensional Wiener process.

Skipping over the details of stochastic calculus, here we give an intuitive understanding of the evolution of the system by using the *Euler-Maruyama* (Higham, 2001) representation as follows:

$$x(t_{n+1}) - x(t_n) = f(x(t_n), \theta)\Delta t + \sqrt{\Delta t}\Sigma\epsilon_{t_n} \quad (2.6)$$

where  $\Delta t$  denotes the time increment and  $\epsilon_{t_n}$  is a standard multivariate Gaussian random vectors, i.e.  $\epsilon_{t_n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . The SDEs in (Eq. 2.5) can be considered as the limit of the process described by (Eq. 2.6) (Archambeau et al., 2007).

Without redundancy, the noisy observation model and the notations introduced previously in Section 2.1 can be directly applied to the SDE models in the rest of this text.

### 2.3 Related work

For state and parameter inference in ODEs, *gradient matching* has established itself in recent years as a promising tool. The main idea behind gradient matching is straightforward. It first interpolates the states  $X$  from the observations  $Y$  using a smoothing technique, and then minimizes the discrepancy between slopes of the interpolants and the derivatives obtained from the ODEs  $f(x(t), \theta)$  (Macdonald and Husmeier, 2015). This dramatically improves the inference efficiency since the ODEs never need to be solved explicitly in contrast to traditional techniques.

Early applications of gradient matching trace back to the spline based method proposed by Ramsay et al. (2007). Unfortunately, such a solution requires full observability of the system and is difficult to extend. More recently, gradient matching with Gaussian processes (GMGP) (Calderhead et al., 2009; Dondelinger et al., 2013) model has been proposed. The introduction of Gaussian processes makes the inference technique grid free and is even capable of handling partially observed systems. In the original GMGP model, sampling is used to estimate the intractable posterior. A more performant and scalable solution called variational GMGP (VGMGP) is proposed by Gorbach et al. (2017), which uses variational technique. As the foundation of this work, details of the original GMGP method and the VGMGP method are discussed in Chapter 3.

For inference problems involving random dynamical systems, classical approaches typically resort to Kalman filters and MCMC. The MCMC based solutions scale poorly in practice. As another approximate inference technique, variational methods have gained popularity in recent years. For example, the variational Gaussian process approximation (VGPA) proposed by

Archambeau et al. (2007) uses a linear approximation strategy to infer the states and the parameters of diffusions processes. It consists of two steps. In the forward step, the mean and covariance of the approximate process are estimated. Then in the backward step, the time evolution of the Lagrange multipliers are calculated. The Lagrange multipliers ensure consistency for the mean and covariance. These two steps are then iteratively executed to improved the result. Further improvements based on this method can be found in Archambeau et al. (2008) and Vrettas et al. (2011, 2015).





---

## Gradient Matching with Gaussian Processes

---

### 3.1 Preliminary

This section briefly introduces Gaussian process regression (Section 3.1.1), Kullback-Leibler (KL) divergence (Section 3.1.2), and mean-field variational inference (Section 3.1.3) as the essential machine learning topics for the rest of the work. Formal treatment can be found by following the references listed for each topic.

#### 3.1.1 Gaussian process regression

As a nonparametric, kernel-based Bayesian inference technique, Gaussian process regression is a powerful tool for many non-linear regression tasks. In terms of application, a Gaussian process prior needs first to be specified on the regression function  $f(\mathbf{x}): \mathbb{R}^D \rightarrow \mathbb{R}$ . After observing data, the prior can be turned into a posterior distribution to make predictions (Murphy, 2012). The definition of Gaussian process given by Rasmussen and Williams (2006) is listed as follow.

**Definition 3.1** *A Gaussian process is a collection of random variables such that any finite subset of it forms a multivariate Gaussian distribution.*

The Gaussian process prior on  $f$  is denoted as

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (3.1)$$

where

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] \quad (3.2)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \quad (3.3)$$

The above equations simply mean that for any finite collection of points  $X = \{x_i \in \mathbb{R}^D | i = 1, \dots, N\}$ , we have  $f|X \sim \mathcal{N}(m, K(X, X))$ , where  $f_i = f(x_i)$ ,  $m_i = m(x_i)$ , and  $K_{ij} = k(x_i, x_j)$  for  $i, j = 1, \dots, N$ . This shows that a Gaussian process is completely specified by its *mean function* (Eq. 3.2) and *covariance function* (Eq. 3.3). Furthermore, to form a multivariate Gaussian distribution, the covariance function must generate a covariance matrix that is positive definite for any set of points.

We often only have access to noisy observations denoted by  $y = \{y_i | i = 1, \dots, N\}$  such that  $y_i = f(x_i) + \varepsilon_i$  for  $i = 1, \dots, N$ . We assume i.i.d. additive noise  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ . For any finite collection of test points  $X^* = \{x_i^* \in \mathbb{R}^D | i = 1, \dots, M\}$ , we have

$$\begin{bmatrix} y \\ f^* \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} m \\ m^* \end{bmatrix}, \begin{bmatrix} K(X, X) + \sigma^2 I & K(X, X^*) \\ K(X^*, X) & K(X^*, X^*) \end{bmatrix}\right)$$

The posterior on  $f^*$  can then be easily found in closed-form using the standard conditional Gaussian distribution formula to obtain

$$f^*|X, y, \sigma, X^* \sim \mathcal{N}(m^* + K(X^*, X)[K(X, X) + \sigma^2 I]^{-1}(y - m), \\ K(X^*, X^*) - K(X^*, X)[K(X, X) + \sigma^2 I]^{-1}K(X, X^*)) \quad (3.4)$$

Note that in practice, the mean function is usually assumed to be zero, i.e.  $m(x) = 0$  (Rasmussen and Williams, 2006).

### 3.1.2 Kullback-Leibler divergence

Suppose there are two continuous<sup>1</sup> probability distributions  $p$  and  $q$ , and the task is to approximate the unknown distribution  $p$  using any distribution  $q$ . In information theory, a frequently used dissimilarity measure between  $p$  and  $q$  is the *forward KL divergence* or *relative entropy* (Kullback and Leibler, 1951), which is defined as

$$\begin{aligned} KL(p||q) &= \int p(x) \ln \frac{p(x)}{q(x)} dx \\ &= \int p(x) \ln p(x) dx - \int p(x) \ln q(x) dx \\ &= -H(p) + H(p, q) \end{aligned} \quad (3.5)$$

where  $H(p)$  is the *entropy* of the distribution  $p$ , and  $H(p, q)$  is the *cross entropy* between  $p$  and  $q$ .

---

<sup>1</sup>Given the context of this work, the discussion here focuses more on continuous probability distributions. However similar reasoning and results can be applied to discrete probability distributions as well.

Similarly, the *reverse KL divergence*  $KL(q||p)$  can be defined as

$$\begin{aligned} KL(q||p) &= \int q(x) \ln \frac{q(x)}{p(x)} dx \\ &= \int q(x) \ln q(x) dx - \int q(x) \ln p(x) dx \\ &= -H(q) + H(q, p) \end{aligned} \tag{3.6}$$

To see why the KL divergence can be used as a dissimilarity measure, one can use *Jensen's inequality* (Jensen, 1906) to prove the following stated as a theorem. The proof (Bishop, 2006, Section 1.6) is omitted here. However, it is not a distance measure since  $KL(p||q) \neq KL(q||p)$  in general, i.e. it is not symmetric (Goodfellow et al., 2016).

**Theorem 3.2**  $KL(p||q) \geq 0$  with  $KL(p||q) = 0$  if and only if  $p(x) = q(x)$  almost everywhere. Similar argument also applies to  $KL(q||p)$ .

From the above, the proxy distribution  $q$  can therefore be found by minimizing either  $KL(p||q)$  (Eq. 3.5) or  $KL(q||p)$  (Eq. 3.6). However, the effects of the two minimization strategies are different. The intuition behind is summarized here from the discussions in Bishop (2006) and Murphy (2012).

- The minimization of  $KL(p||q)$  can be characterized as *zero avoiding* for  $q$  because there is a large contribution to the integral at places where  $p(x) > 0$  and  $q(x) = 0$ . As a result, we need to ensure that  $q(x) > 0$  whenever  $p(x) > 0$ , which generally leads to an overestimation of the support for  $p$ .
- Conversely, the minimization of  $KL(q||p)$  can be viewed as *zero forcing* for  $q$  due to the large contribution to the integral at place where  $p(x) = 0$  and  $q(x) > 0$ . Consequently, we need to force  $q(x) = 0$  whenever  $p(x) = 0$ , which tends to underestimate the support for  $p$ .

Visual illustrations of approximating both a single 2-D Gaussian distribution and mixture of two 2-D Gaussian distributions using another 2-D Gaussian distribution can be found in Section 10.1 of Bishop (2006) and Section 21.2 of Murphy (2012).

### 3.1.3 Mean-field variational inference

The evaluation of the posterior distribution given the observations plays a central role in Bayesian statistics since it is essential to make point estimates, form predictive distributions, etc. (Bishop, 2006). Unfortunately, solutions to many of such problems are not analytically tractable, e.g. the Bayesian mixture of Gaussians example in Section 2.1 of Blei et al. (2017), and hence, we have to resort to approximation algorithms. Here, we review *variational*

*inference* (Jordan et al., 1999; Wainwright et al., 2008), which is a deterministic approximation scheme that approximates probability densities through optimization. Compared to the stochastic MCMC methods, variational inference tends to be faster and has been shown empirically to be capable of scaling to large datasets, for example topic modeling using 1.8 million New York Times articles (Hoffman et al., 2013), traffic pattern analysis using 1.7 million taxi rides (Kucukelbir et al., 2017). The core concepts behind variational inference are summarized below based on Bishop (2006) and Blei et al. (2017).

### Evidence lower bound

Suppose given the model specification  $p(\mathbf{x}, \mathbf{z})$ , we want to approximate the unknown posterior  $p(\mathbf{z}|\mathbf{x})$ , where  $\mathbf{z}$  denotes the *latent (hidden) variables*<sup>2</sup>, and  $\mathbf{x}$  denotes the observed variables. In variational inference, we first posit a family of distributions  $\mathcal{Q}$  over  $\mathbf{z}$  parameterized by a set of free *variational parameters*<sup>3</sup>. The “closest” approximation  $q(\mathbf{z})$  to  $p(\mathbf{z}|\mathbf{x})$  can then be found by minimizing the KL divergence  $KL(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}))$ .

$$q^*(\mathbf{z}) \in \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} KL(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x})) \quad (3.7)$$

However, directly solving the above optimization problem is infeasible since expanding the objective reveals its dependency on the *evidence*  $p(\mathbf{x})$ , and we assume that  $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z})d\mathbf{z}$  is intractable.

$$\begin{aligned} KL(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x})) &= \int q(\mathbf{z}) \ln q(\mathbf{z})d\mathbf{z} - \int q(\mathbf{z}) \ln p(\mathbf{z}|\mathbf{x})d\mathbf{z} \\ &= \int q(\mathbf{z}) \ln q(\mathbf{z})d\mathbf{z} - \int q(\mathbf{z}) \ln p(\mathbf{x}, \mathbf{z})d\mathbf{z} + \ln p(\mathbf{x}) \end{aligned} \quad (3.8)$$

Instead, we can rearrange the terms in (Eq. 3.8) to obtain

$$\ln p(\mathbf{x}) = KL(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x})) + \mathcal{L}(q(\mathbf{z})) \quad (3.9)$$

where

$$\mathcal{L}(q(\mathbf{z})) = - \int q(\mathbf{z}) \ln q(\mathbf{z})d\mathbf{z} + \int q(\mathbf{z}) \ln p(\mathbf{x}, \mathbf{z})d\mathbf{z} \quad (3.10)$$

Because  $KL(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x})) \geq 0$  (Theorem 3.2),  $\mathcal{L}(q(\mathbf{z}))$  can be treated as a lower bound on  $\ln p(\mathbf{x})$ , i.e.  $\ln p(\mathbf{x}) \geq \mathcal{L}(q(\mathbf{z}))$ , and hence, it is also called

---

<sup>2</sup>In full-Bayesian treatment, other model parameters, when they exist, are absorbed into  $\mathbf{z}$  to simplify the notation.

<sup>3</sup>The dependency on the variational parameters is left implicit in the notations used in the following text.

the *evidence lower bound*. Furthermore, since  $\ln p(\mathbf{x})$  is a constant, we can state that *maximizing  $\mathcal{L}(q(\mathbf{z}))$  is equivalent to minimizing  $KL(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}))$* .

Rewriting (Eq. 3.10) as (Eq. 3.11) below also provides us with some intuitions about the form of the optimal proxy distribution. The first term in (Eq. 3.11) encourages  $q(\mathbf{z})$  to put weights on  $\mathbf{z}$  that explains  $\mathbf{x}$  through  $p(\mathbf{x}|\mathbf{z})$ , while the second term encourages  $q(\mathbf{z})$  to be close to the prior  $p(\mathbf{z})$ . This reflects the usual balance between likelihood and prior in Bayesian statistics.

$$\begin{aligned}\mathcal{L}(q(\mathbf{z})) &= - \int q(\mathbf{z}) \ln q(\mathbf{z}) d\mathbf{z} + \int q(\mathbf{z}) \ln [p(\mathbf{x}|\mathbf{z})p(\mathbf{z})] d\mathbf{z} \\ &= \mathbb{E}_q[\ln p(\mathbf{x}|\mathbf{z})] - KL(q(\mathbf{z})\|p(\mathbf{z}))\end{aligned}\tag{3.11}$$

### Mean-field variational family

Having specified the optimization objective in terms of  $\mathcal{L}(q(\mathbf{z}))$ , we now need to decide on the family of distributions  $\mathcal{Q}$ . Specifically, we focus on the *mean-field variational family*, where the latent variables and model parameters are factorized mutually independent groups, and each group is controlled by its own variational parameters, i.e.  $q(\mathbf{z}) = \prod_i q_i(\mathbf{z}_i)$ . Note that since the complexity of the optimization is determined by the reach of  $\mathcal{Q}$ , one should be aware of the trade-off between model flexibility and runtime efficiency when choosing the reach of  $\mathcal{Q}$ .

By dissecting out the terms that are dependent on  $q_i(\mathbf{z}_i)$ , and absorbing the terms that are independent of it into constant from  $\mathcal{L}(q(\mathbf{z}))$ , the optimal solution  $q_i^*(\mathbf{z}_i)$  can be expressed as

$$q_i^*(\mathbf{z}_i) \propto \exp \{ \mathbb{E}_{q_{j \neq i}} [\ln p(\mathbf{x}, \mathbf{z})] \} \tag{3.12}$$

where  $\mathbb{E}_{q_{j \neq i}} [\ln p(\mathbf{x}, \mathbf{z})]$  denotes the expectation of the log of the joint distribution  $p(\mathbf{x}, \mathbf{z})$  with respect to all the other factors  $q_j$  except  $q_i$ . The derivation (Bishop, 2006, Section 10.1) is omitted here.

The advantage of the factorization assumption is that we can fix other factors and optimize only one factor at a time in an iterative manner until convergence, or the preconfigured maximum number of iterations has been reached. This approach is called the *coordinate ascent mean-field variational inference*. Detailed pseudocode for the algorithm, handling for numerical stability, test for convergence, etc. are discussed in Blei et al. (2017).

As a consequence of the optimization objective and the factorization assumption, the optimal solution tends to underestimate the variance of the posterior to produce a too compact distribution. Discussion about this issue can be found in Section 3.1.2 of this work and Section 10.1 of Bishop (2006). Lastly, in cases where the objective is non-convex, convergence is only guaranteed to a local optimum subject to the initialization (Blei et al., 2017).

## 3.2 Sampling-based gradient matching with Gaussian processes

This section reviews a recently proposed framework to infer parameters and states of ODEs using gradient matching with Gaussian processes (GMGP), which originally appeared in Calderhead et al. (2009) and was extended by Dondelinger et al. (2013).

Adopting the notations from Section 2.1, for a  $K$ -dimensional dynamical system described by (Eq. 2.1) as:

$$\dot{\mathbf{x}}(t) = \frac{d\mathbf{x}(t)}{dt} = \mathbf{f}(\mathbf{x}(t), \boldsymbol{\theta}) \quad (3.13)$$

Calderhead et al. (2009) introduce independent Gaussian process priors on each state  $k$  for  $k = 1, \dots, K$  such that

$$p(\mathbf{X}|\boldsymbol{\varphi}) = \prod_k \mathcal{N}(\mathbf{x}_k | \mathbf{0}, \mathbf{C}_{\boldsymbol{\varphi}_k}) \quad (3.14)$$

where  $\mathbf{C}_{\boldsymbol{\varphi}_k}$  is the covariance matrix induced by the corresponding kernel function  $\mathcal{K}_{\boldsymbol{\varphi}_k}$  with the hyperparameter vector  $\boldsymbol{\varphi}_k$  over the  $N$  time points.

### States interpolation

If we assume the observation on the states is corrupted by noise described by (Eq. 2.2) as:

$$p(\mathbf{Y}|\mathbf{X}, \sigma) = \prod_k \mathcal{N}(\mathbf{y}_k | \mathbf{x}_k, \sigma_k^2 \mathbf{I}) \quad (3.15)$$

the posterior distribution on the states  $\mathbf{X}$  is obtained as

$$\begin{aligned} p(\mathbf{X}|\mathbf{Y}, \boldsymbol{\varphi}, \sigma) &= \frac{p(\mathbf{X}|\boldsymbol{\varphi})p(\mathbf{Y}|\mathbf{X}, \sigma)}{\int p(\mathbf{X}|\boldsymbol{\varphi})p(\mathbf{Y}|\mathbf{X}, \sigma)d\mathbf{X}} \\ &= \prod_k \mathcal{N}(\mathbf{x}_k | \boldsymbol{\mu}_k(\mathbf{y}_k), \boldsymbol{\Sigma}_k) \end{aligned} \quad (3.16)$$

where<sup>4</sup>  $\boldsymbol{\mu}_k(\mathbf{y}_k) = \mathbf{C}_{\boldsymbol{\varphi}_k}(\mathbf{C}_{\boldsymbol{\varphi}_k} + \sigma_k^2 \mathbf{I})^{-1} \mathbf{y}_k$  and  $\boldsymbol{\Sigma}_k = \sigma_k^2 \mathbf{C}_{\boldsymbol{\varphi}_k}(\mathbf{C}_{\boldsymbol{\varphi}_k} + \sigma_k^2 \mathbf{I})^{-1}$ .

### Gaussian process response model

Because differentiation is a linear operation, the derivative of a Gaussian process is also a Gaussian process (Rasmussen and Williams, 2006, Section 9.4). Hence, the joint distribution of the states and their derivatives within

---

<sup>4</sup>Another equivalent derivation, as stated in Gorbach et al. (2017), is to define  $\boldsymbol{\mu}_k(\mathbf{y}_k) = \sigma_k^{-2} \boldsymbol{\Sigma}_k^{-1} \mathbf{y}_k$ , and  $\boldsymbol{\Sigma}_k = (\sigma_k^{-2} \mathbf{I} + \mathbf{C}_{\boldsymbol{\varphi}_k}^{-1})^{-1}$ .

### 3.2. Sampling-based gradient matching with Gaussian processes

a finite amount of time points is Gaussian. To illustrate, we write the joint distribution of  $\mathbf{x}_k$  and  $\dot{\mathbf{x}}_k$  as follows:

$$\begin{bmatrix} \mathbf{x}_k \\ \dot{\mathbf{x}}_k \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{C}_{\varphi_k} & \mathbf{C}'_{\varphi_k} \\ {}'\mathbf{C}_{\varphi_k} & \mathbf{C}''_{\varphi_k} \end{bmatrix}\right)$$

where

$$\begin{aligned} C_{\varphi_k i,j} &= \mathcal{K}_{\varphi_k}(t_i, t_j) \\ C'_{\varphi_k i,j} &= \frac{\partial \mathcal{K}_{\varphi_k}(t_i, t_j)}{\partial t_j} \\ {}'C_{\varphi_k i,j} &= \frac{\partial \mathcal{K}_{\varphi_k}(t_i, t_j)}{\partial t_i} \\ C''_{\varphi_k i,j} &= \frac{\partial^2 \mathcal{K}_{\varphi_k}(t_i, t_j)}{\partial t_i \partial t_j} \end{aligned}$$

In other words,  $\mathbf{C}_{\varphi_k}$  and  $\mathbf{C}''_{\varphi_k}$  are the *auto-covariance* matrices for  $\mathbf{x}_k$  and  $\dot{\mathbf{x}}_k$  respectively, and  $\mathbf{C}'_{\varphi_k}$  and  ${}'\mathbf{C}_{\varphi_k}$  are the *cross-covariance* matrices between  $\mathbf{x}_k$  and  $\dot{\mathbf{x}}_k$  and between  $\dot{\mathbf{x}}_k$  and  $\mathbf{x}_k$  respectively. Therefore, the conditional distribution over the state derivatives  $\dot{\mathbf{X}}$  is given by

$$p(\dot{\mathbf{X}}|\mathbf{X}, \boldsymbol{\varphi}) = \prod_k \mathcal{N}(\dot{\mathbf{x}}_k | \mathbf{m}_k, \mathbf{A}_k) \quad (3.17)$$

where<sup>5</sup>  $\mathbf{m}_k = {}'\mathbf{C}_{\varphi_k} \mathbf{C}_{\varphi_k}^{-1} \mathbf{x}_k$  and  $\mathbf{A}_k = \mathbf{C}''_{\varphi_k} - {}'\mathbf{C}_{\varphi_k} \mathbf{C}_{\varphi_k}^{-1} \mathbf{C}'_{\varphi_k}$ .

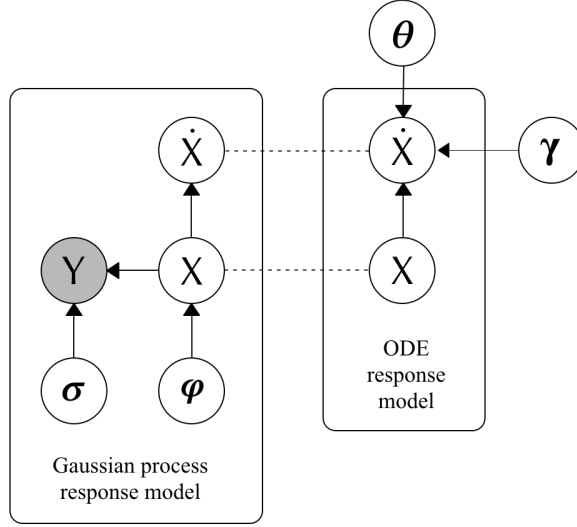
#### ODE response model

Since given certain states and parameters, the ODEs (Eq. 3.13) are the natural information source of the state derivatives, if we assume state specific, additive and normally distributed errors between the true state derivatives and response from the ODEs using the states and parameters estimation results as input, we have

$$p(\dot{\mathbf{X}}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\gamma}) = \prod_k \mathcal{N}(\dot{\mathbf{x}}_k | \mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta}), \gamma_k \mathbf{I}) \quad (3.18)$$

where the vector  $\boldsymbol{\gamma}$  contains the error variances for each state, i.e.  $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_K]^T \in \mathbb{R}^K$ .

<sup>5</sup>The derivation from Dondelinger et al. (2013) is used here since there is a minor inconsistency in the original proposal from Calderhead et al. (2009).



**Figure 3.1:** Graphical representation of the gradient matching with Gaussian processes model proposed by Calderhead et al. (2009). The dashed lines indicate information association from two response models using the product of experts (Eq. 3.19). Details of the model are in Section 3.2.

### Product of experts

Next, the *product of experts* (Hinton, 1999) approach is employed to combine the Gaussian process response model (Eq. 3.17) and the ODE response model (Eq. 3.18) to obtain

$$p(\dot{X}|X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\gamma}) \propto p(\dot{X}|X, \boldsymbol{\varphi})p(\dot{X}|X, \boldsymbol{\theta}, \boldsymbol{\gamma}) \quad (3.19)$$

This product form implies that  $p(\dot{X}|X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\gamma})$  attains high densities at locations where both  $p(\dot{X}|X, \boldsymbol{\varphi})$  and  $p(\dot{X}|X, \boldsymbol{\theta}, \boldsymbol{\gamma})$  have strong support.

### Full model distribution

To summarize, the full joint distribution of the GMGP model is graphically represented in Figure 3.1, which corresponds to the following:

$$\begin{aligned} p(Y, X, \dot{X}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \sigma, \boldsymbol{\gamma}) &= p(Y|X, \sigma)p(X|\boldsymbol{\varphi})p(\dot{X}|X, \boldsymbol{\theta}, \boldsymbol{\varphi}, \boldsymbol{\gamma})p(\boldsymbol{\varphi})p(\boldsymbol{\theta})p(\sigma)p(\boldsymbol{\gamma}) \\ &\propto p(Y|X, \sigma)p(X|\boldsymbol{\varphi})p(\dot{X}|X, \boldsymbol{\varphi})p(\dot{X}|X, \boldsymbol{\theta}, \boldsymbol{\gamma})p(\boldsymbol{\varphi})p(\boldsymbol{\theta})p(\sigma)p(\boldsymbol{\gamma}) \end{aligned}$$

where  $p(\boldsymbol{\varphi})$ ,  $p(\boldsymbol{\theta})$ ,  $p(\sigma)$ , and  $p(\boldsymbol{\gamma})$  denote the corresponding priors.



### Original sampling scheme

Calderhead et al. (2009) propose to integrate over the latent state derivatives  $\dot{X}$  to obtain

$$\begin{aligned} p(\boldsymbol{\theta}, \gamma | \mathbf{X}, \boldsymbol{\varphi}) &= p(\boldsymbol{\theta}) p(\gamma) \int p(\dot{X} | \mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\varphi}, \gamma) d\dot{X} \\ &\propto \frac{p(\boldsymbol{\theta}) p(\gamma)}{\mathcal{Z}(\gamma)} \exp \left[ -\frac{1}{2} \sum_k (\mathbf{f}_k - \mathbf{m}_k)^T \boldsymbol{\Lambda}_k (\mathbf{f}_k - \mathbf{m}_k) \right] \end{aligned} \quad (3.20)$$

where  $\boldsymbol{\Lambda}_k^{-1} = \mathbf{A}_k + \gamma_k \mathbf{I}$ ,  $\mathcal{Z}(\gamma) = \prod_k |2\pi(\mathbf{A}_k + \gamma_k \mathbf{I})|^{\frac{1}{2}}$  and  $\mathbf{f}_k = \mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta})$ .

Then the sampling procedure follows the scheme below

$$\boldsymbol{\varphi}, \sigma \sim p(\boldsymbol{\varphi}, \sigma | \mathbf{Y}) \quad (3.21)$$

$$\mathbf{X} \sim p(\mathbf{X} | \mathbf{Y}, \boldsymbol{\varphi}, \sigma) \quad (3.22)$$

$$\boldsymbol{\theta}, \gamma \sim p(\boldsymbol{\theta}, \gamma | \mathbf{X}, \boldsymbol{\varphi}) \quad (3.23)$$

where

$$\begin{aligned} p(\boldsymbol{\varphi}, \sigma | \mathbf{Y}) &\propto p(\boldsymbol{\varphi}) p(\sigma) p(\mathbf{Y} | \boldsymbol{\varphi}, \sigma) \\ &\propto p(\boldsymbol{\varphi}) p(\sigma) \int p(\mathbf{X} | \boldsymbol{\varphi}) p(\mathbf{Y} | \mathbf{X}, \sigma) d\mathbf{X} \\ &= p(\boldsymbol{\varphi}) p(\sigma) \prod_k \mathcal{N}(\mathbf{y}_k | \mathbf{0}, \mathbf{C}_{\boldsymbol{\varphi}_k} + \sigma_k^2 \mathbf{I}) \end{aligned}$$

The sampling procedure requires two MCMC samplings for (Eq. 3.21) and (Eq. 3.23) respectively, and a direct sampling from the multivariate Gaussian distribution in (Eq. 3.22). The sampling steps are repeated until convergence is reached.

### Adaptive sampling scheme

One fundamental weakness of the previous sampling strategy is that the results on  $\boldsymbol{\theta}$  and  $\gamma$  from (Eq. 3.23) are never propagated back to (Eq. 3.21) and (Eq. 3.22), and hence have no influence on the inference of states  $\mathbf{X}$ . To close the feedback loop, Dondelinger et al. (2013) proposed an improved sampling scheme called *adaptive gradient matching*.

Dondelinger et al. (2013) consider the joint distribution  $p(\dot{X}, \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma)$ , and show that its marginalization over the state derivatives  $\dot{X}$  is tractable:

$$\begin{aligned} p(\mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma) &= \int p(\dot{X}, \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma) d\dot{X} \\ &= p(\mathbf{X} | \boldsymbol{\varphi}) p(\boldsymbol{\varphi}) p(\boldsymbol{\theta}) p(\gamma) \int p(\dot{X} | \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma) d\dot{X} \\ &\propto \exp \left[ -\frac{1}{2} \sum_k (\mathbf{x}_k^T \mathbf{C}_{\boldsymbol{\varphi}_k}^{-1} \mathbf{x}_k + (\mathbf{f}_k - \mathbf{m}_k)^T \boldsymbol{\Lambda}_k (\mathbf{f}_k - \mathbf{m}_k)) \right] \end{aligned} \quad (3.24)$$

The full joint distribution, after integrating over the latent state derivatives  $\dot{X}$ , is then given by

$$\begin{aligned} p(Y, X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \sigma, \gamma) &= p(Y|X, \sigma)p(X|\boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma)p(\boldsymbol{\varphi})p(\boldsymbol{\theta})p(\sigma)p(\gamma) \\ &= p(Y|X, \sigma)p(X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma)p(\sigma) \end{aligned} \quad (3.25)$$

Based on the above joint distribution, Dondelinger et al. (2013) devised an improved Metropolis-Hastings sampling scheme, which allows  $\boldsymbol{\theta}$  to exert an influence on  $X$ .

### Summary

To conclude this section, given the implicit solutions from the Gaussian processes, the states and parameters of the ODEs are inferred simultaneously without explicitly solving the ODEs anymore, which has led to significant speed-up. Lastly, Calderhead et al. (2009) also discuss the handling of partial observations, i.e. when some states are not observed, by utilizing to the prior we have imposed on the states (Eq. 3.14).

## 3.3 Variational gradient matching with Gaussian processes

Instead of using sampling methods, Gorbach et al. (2016, 2017) proposed a solution to the previous problem based on variational inference, which has been shown, for specific types of ODEs (Eq. 3.28) as discussed below, to be much more efficient than the sampling-based solutions, and scales well to large dynamical systems such the deterministic Lorenz 96 model with 1000 states in their experiments. This section gives an outline of the method by following the description from Gorbach et al. (2017) and the method is abbreviated as VGMGP in the rest of this work.

### Maximum a posteriori estimation

By combining the prior  $p(\boldsymbol{\theta})$ , the state posterior (Eq. 3.16), the product of experts result (Eq. 3.19), and then integrating out the latent state derivatives  $\dot{X}$  as in Calderhead et al. (2009), the joint posterior  $p(X, \boldsymbol{\theta}|Y, \boldsymbol{\varphi}, \sigma, \gamma)$  is given as follows:

$$\begin{aligned} p(X, \boldsymbol{\theta}|Y, \boldsymbol{\varphi}, \sigma, \gamma) &= p(\boldsymbol{\theta}) \int p(X|Y, \boldsymbol{\varphi}, \sigma)p(\dot{X}|X, \boldsymbol{\theta}, \boldsymbol{\varphi}, \gamma)d\dot{X} \\ &\propto p(\boldsymbol{\theta}) \prod_k [\mathcal{N}(x_k|\boldsymbol{\mu}_k(\mathbf{y}_k), \boldsymbol{\Sigma}_k)\mathcal{N}(\mathbf{f}_k(X, \boldsymbol{\theta})|\mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1})] \end{aligned} \quad (3.26)$$

Using the inference of the parameters  $\boldsymbol{\theta}$  as an example, the “best” parameters  $\boldsymbol{\theta}^*$  could ideally be estimated based on the *maximum a posteriori* (MAP) of the

posterior  $p(\theta|Y, \boldsymbol{\varphi}, \sigma, \gamma)$ , which is equivalent to

$$\theta^* = \arg \max_{\theta} \ln p(\theta|Y, \boldsymbol{\varphi}, \sigma, \gamma) \quad (3.27)$$

where the posterior on  $\theta$  is obtained by marginalizing (Eq. 3.26) over the states  $X$ , namely

$$p(\theta|Y, \boldsymbol{\varphi}, \sigma, \gamma) = \int p(X, \theta|Y, \boldsymbol{\varphi}, \sigma, \gamma) dX$$

However, the above equation is intractable due to strong non-linear couplings of the states induced by the ODEs (Eq. 3.13). As a workaround, Gorbach et al. (2017) designed a proxy distribution on the states and parameters denoted as  $Q(X, \theta)$ , and derived analytically tractable variational lower bounds based on mean-field variational inference.

### Structural assumption

First of all, the ODEs primarily considered by Gorbach et al. (2017) are state-wise of the following structure:

$$f_k(\mathbf{x}(t), \boldsymbol{\theta}) = \sum_m \theta_m \prod_{i \in \mathcal{M}_{km}} x_i(t) + C \quad (3.28)$$

for  $k = 1, \dots, K$ , where  $\mathcal{M}_{km}$  contains the states in the  $k$ -th equation of the ODEs that are controlled by the  $m$ -th parameter  $\theta_m$ , and  $C$  denotes other terms that are independent of the parameters. The first term of (Eq. 3.28) can be viewed as a linear combination of the parameters and the product of the monomials of the states. For the second term, the only requirement on terms in  $C$  is that each state can only appear as monomial inside a term, but a term in  $C$  can contain the product of the monomials of the states.

The structural assumption imposed on the ODEs leads to an important conclusion that the following conditional distributions  $p(\theta|Y, X, \boldsymbol{\varphi}, \gamma)$  and  $p(\mathbf{x}_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \theta, \sigma, \gamma)$  for  $u = 1, \dots, K$  are Gaussian distributed, where  $X_{/\{u\}} = \{\mathbf{x}_o | o = 1, \dots, K \text{ and } o \neq u\}$ , i.e.  $X_{/\{u\}}$  is the set of all states except  $\mathbf{x}_u$ . But before writing out the distributions, we first need to introduce several more notations.

One implication of the above assumption is that each ODE can be expressed as a linear combination of the parameters plus a term that is independent of them. We therefore use  $\mathbf{B}_{\theta k}$  and  $\mathbf{b}_{\theta k}$  to denote the corresponding coefficients such that  $\mathbf{B}_{\theta k} \boldsymbol{\theta} + \mathbf{b}_{\theta k} = \mathbf{f}_k(X, \boldsymbol{\theta})$  for  $k = 1, \dots, K$ . Similarly, due to the monomial assumption about the states, each ODE can also be transformed into a linear combination of the  $u$ -th state for  $u = 1, \dots, K$  plus an independent term. Here we adopt the terms  $\mathbf{B}_{uk}$  and  $\mathbf{b}_{uk}$  to denote the respective coefficients such that  $\mathbf{B}_{uk} \mathbf{x}_u + \mathbf{b}_{uk} = \mathbf{f}_k(X, \boldsymbol{\theta})$  for  $k, u = 1, \dots, K$ .

With these notations, we now state the conditional distributions  $p(\theta|Y, X, \boldsymbol{\varphi}, \gamma)$  and  $p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \theta, \sigma, \gamma)$  as follows:

$$p(\theta|Y, X, \boldsymbol{\varphi}, \gamma) = \mathcal{N}(\theta|\mathbf{r}_\theta, \boldsymbol{\Omega}_\theta) \quad (3.29)$$

$$p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \theta, \sigma, \gamma) = \mathcal{N}(x_u|\mathbf{r}_u, \boldsymbol{\Omega}_u) \quad (3.30)$$

where

$$\mathbf{r}_\theta = \boldsymbol{\Omega}_\theta \sum_k \mathbf{B}_{\theta k}^T \boldsymbol{\Lambda}_k (\mathbf{m}_k - \mathbf{b}_{\theta k}) \quad (3.31)$$

$$\boldsymbol{\Omega}_\theta^{-1} = \sum_k \mathbf{B}_{\theta k}^T \boldsymbol{\Lambda}_k \mathbf{B}_{\theta k} \quad (3.32)$$

$$\mathbf{r}_u = \boldsymbol{\Omega}_u [\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\mu}_u(y_u) + \sum_k \mathbf{B}_{uk}^T \boldsymbol{\Lambda}_k (\mathbf{m}_k - \mathbf{b}_{uk})] \quad (3.33)$$

$$\boldsymbol{\Omega}_u^{-1} = \boldsymbol{\Sigma}_u^{-1} + \sum_k \mathbf{B}_{uk}^T \boldsymbol{\Lambda}_k \mathbf{B}_{uk} \quad (3.34)$$

The proofs are omitted here and can be found in the supporting materials of Gorbach et al. (2017). They basically utilizes the formula to calculate the product of Gaussians (Petersen and Pedersen, 2012, Section 8.1), and the well-known closure property of Gaussian distribution under linear transformations, to obtain a closed-form solution. As a side note, this conclusion will also be used later in Chapter 4 to derive another approximate inference solution based on the Laplace approximation technique.

To simplify the notations, in the following, we will switch to the canonical form of *exponential family* (Murphy, 2012, Section 9.2) to describe the distributions (Eq. 3.29) and (Eq. 3.30) as:

$$p(\theta|Y, X, \boldsymbol{\varphi}, \gamma) = h_\theta(\theta) \exp(\boldsymbol{\eta}_\theta(Y, X, \boldsymbol{\varphi}, \gamma)^T \mathbf{T}_\theta(\theta) - A_\theta(\boldsymbol{\eta}_\theta)) \quad (3.35)$$

$$p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \theta, \sigma, \gamma) = h_u(x_u) \exp(\boldsymbol{\eta}_u(Y, X_{/\{u\}}, \boldsymbol{\varphi}, \theta, \sigma, \gamma)^T \mathbf{T}_u(x_u) - A_u(\boldsymbol{\eta}_u)) \quad (3.36)$$

where  $\boldsymbol{\eta}_{(\cdot)}(\cdot)$ ,  $\mathbf{T}_{(\cdot)}(\cdot)$ ,  $h_{(\cdot)}(\cdot)$ ,  $A_{(\cdot)}(\cdot)$  are the *natural parameters*, *sufficient statistics*, *base measure* and *log partition function* respectively.

### Variational lower bound

Introducing  $\lambda$  and  $\boldsymbol{\psi}_u$  for  $u = 1, \dots, K$  as the variational parameters, the proxy distribution  $Q(X, \theta)$  is constrained to the family  $\mathcal{Q}$  that factorizes over the parameters and the states, and each factor is in the same exponential family as its corresponding counterparts defined in (Eq. 3.35) and (Eq. 3.36) such that

$$\mathcal{Q} = \{Q(X, \theta) | Q(X, \theta) = q(\theta|\lambda) \prod_u q(x_u|\boldsymbol{\psi}_u)\} \quad (3.37)$$

where

$$\begin{aligned} q(\boldsymbol{\theta}|\boldsymbol{\lambda}) &= h_{q\boldsymbol{\theta}}(\boldsymbol{\theta}) \exp(\boldsymbol{\lambda}^T \mathbf{T}_{q\boldsymbol{\theta}}(\boldsymbol{\theta}) - A_{q\boldsymbol{\theta}}(\boldsymbol{\lambda})) \\ q(\mathbf{x}_u|\boldsymbol{\psi}_u) &= h_{qu}(\mathbf{x}_u) \exp(\boldsymbol{\psi}_u^T \mathbf{T}_{qu}(\boldsymbol{\psi}_u) - A_{qu}(\boldsymbol{\psi}_u)) \end{aligned}$$

Using standard mean-field variation technique, the optimal  $Q^*$  is given by

$$\begin{aligned} Q^* &= \arg \min_{Q(\mathbf{X}, \boldsymbol{\theta}) \in \mathcal{Q}} KL(Q(\mathbf{X}, \boldsymbol{\theta}) \| p(\mathbf{X}, \boldsymbol{\theta} | \mathbf{Y}, \boldsymbol{\varphi}, \boldsymbol{\sigma}, \boldsymbol{\gamma})) \\ &= \arg \max_{Q(\mathbf{X}, \boldsymbol{\theta}) \in \mathcal{Q}} \mathcal{L}_Q(\boldsymbol{\lambda}, \boldsymbol{\psi}) \end{aligned} \quad (3.38)$$

where  $\mathcal{L}_Q(\boldsymbol{\lambda}, \boldsymbol{\psi})$  denotes the evidence lower bound.

Gorbach et al. (2017) showed that maximizing  $\mathcal{L}_Q(\boldsymbol{\lambda}, \boldsymbol{\psi})$  with respect to  $\boldsymbol{\theta}$  is equivalent to maximizing

$$\begin{aligned} \mathcal{L}_{\boldsymbol{\theta}}(\boldsymbol{\lambda}) &= \mathbb{E}_Q[\ln p(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\gamma})] - \mathbb{E}_Q[\ln q(\boldsymbol{\theta} | \boldsymbol{\lambda})] \\ &= \mathbb{E}_Q[\boldsymbol{\eta}_{\boldsymbol{\theta}}(\mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\gamma})^T \nabla_{\boldsymbol{\lambda}} A_{\boldsymbol{\theta}}(\boldsymbol{\eta}_{\boldsymbol{\theta}})] - \boldsymbol{\lambda}^T \nabla_{\boldsymbol{\lambda}} A_{q\boldsymbol{\theta}}(\boldsymbol{\lambda}) \end{aligned} \quad (3.39)$$

and maximizing  $\mathcal{L}_Q(\boldsymbol{\lambda}, \boldsymbol{\psi})$  with respect to  $\mathbf{x}_u$  is equivalent to maximizing

$$\begin{aligned} \mathcal{L}_{\boldsymbol{\psi}_u} &= \mathbb{E}_Q[\ln p(\mathbf{x}_u | \mathbf{Y}, \mathbf{X}_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})] - \mathbb{E}_Q[\ln q(\mathbf{x}_u | \boldsymbol{\psi}_u)] \\ &= \mathbb{E}_Q[\boldsymbol{\eta}_u(\mathbf{Y}, \mathbf{X}_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})^T \nabla_u A_u(\boldsymbol{\eta}_u)] - \boldsymbol{\psi}_u^T \nabla_u A_{qu}(\boldsymbol{\psi}_u) \end{aligned} \quad (3.40)$$

where  $\nabla_{(\cdot)}$  is the gradient operator.

Given the assumption about the conditional distributions (Eq. 3.35) and (Eq. 3.36), plus the design of the proxy distribution family (Eq. 3.37), they further showed that the lower bounds (Eq. 3.39) and (Eq. 3.40) can be optimized analytically by setting their gradients with respect to their variational parameters to zero. The optimal variational parameters are thus given by

$$\boldsymbol{\lambda}^* = \mathbb{E}_Q[\boldsymbol{\eta}_{\boldsymbol{\theta}}(\mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\gamma})] \quad (3.41)$$

$$\boldsymbol{\psi}_u^* = \mathbb{E}_Q[\boldsymbol{\eta}_u(\mathbf{Y}, \mathbf{X}_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})] \quad (3.42)$$

where the natural parameters are defined as  $\boldsymbol{\eta}_{\boldsymbol{\theta}}(\mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\gamma}) = [\boldsymbol{\Omega}_{\boldsymbol{\theta}}^{-1} \mathbf{r}_{\boldsymbol{\theta}}, -\frac{1}{2} \boldsymbol{\Omega}_{\boldsymbol{\theta}}^{-1}]^T$  and  $\boldsymbol{\eta}_u(\mathbf{Y}, \mathbf{X}_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma}) = [\boldsymbol{\Omega}_u^{-1} \mathbf{r}_u, -\frac{1}{2} \boldsymbol{\Omega}_u^{-1}]^T$ .

### VGMGP algorithm

To close this section, the pseudocode of their algorithm is summarized in Algorithm 1. The algorithm first uses Gaussian process regression to obtain a smoothed estimate for the states and the distribution on the state derivatives. As noted by Gorbach et al. (2017), the Gaussian process priors come in naturally in cases where there exist unobserved states.

---

**Algorithm 1** Pseudocode for the VGMGP algorithm.

---

```

1: for  $u = 1, \dots, K$  do
2:   Initialize  $\mu_u(\mathbf{y}_u)$ ,  $\Sigma_u$ ,  $\mathbf{m}_u$  and  $\Lambda_u$  using Gaussian process regression
3: end for
4:
5: while not converged or maximum iteration not reached do
6:   for  $u = 1, \dots, K$  do
7:     Update the mean and variance of  $\hat{q}_{\psi_u}$  using  $\hat{\theta}^{(i)}$ 
8:   end for
9:    $\hat{\theta}^{(i+1)} = \arg \max_{\theta} \mathbb{E}_Q[\sum_k \ln \mathcal{N}(f_k(\mathbf{X}, \theta) | \mathbf{m}_k, \Lambda_k^{-1})]$ 
10: end while

```

---

In the second step, the algorithm uses the variational gradient matching framework described before to maximize the lower bounds iteratively. In every iteration, the states and parameters are optimized in an *expectation-maximization (EM)* (Bishop, 2006, Section 9.4) fashion until convergence is reached or the maximum allowed number of iterations is reached.

---

# Laplace Mean-Field Approximation

---

One limitation of the previous VGMGP methodology is that analytical variational lower bounds are obtainable only if the structural assumption on the ODEs is satisfied. Although many dynamical systems such the Lotka-Volterra model (Section 6.2), the Lorenz 63 model (Section 6.5) and the Lorenz 96 model (Section 6.4) fulfill this requirement, it would be valuable to devise a more general solution without constraining the structure of the ODEs. Moreover, for models of biochemical or physical interactions such as the protein signaling transduction pathway model (Section 6.3), sometimes the states or the parameters need to be positive in order to give meaningful results. This positivity constraint is also not supported by VGMGP. Lastly, as can be seen from (Eq. 3.41) and (Eq. 3.42), the analytical solution requires lots of book-keeping, due to the evaluation of the expectations, which renders the implementation cumbersome and error prone.

As an extension to the previous variational approach, this chapter derives another approximation scheme based on Laplace approximation. Section 4.1 give a brief review about the general Laplace approximation technique. In Section 4.2, it is applied to the VGMGP model to derive a new solution called *Laplace mean-field (LPMF)*, which relaxes the assumption about the ODE structure. Similar to the variational approach, this solution also transforms the original inference problem into an optimization problem. However, due to the relaxation on the ODE structure, the optimization objectives can in general no longer be optimized in closed-form, and hence, numerical solutions have to be used. As will be shown in Section 4.3, the gradients and even the Hessians of the state objectives can be computed efficiently, which allows us to rely on second-order optimization techniques and enables the algorithm to scale to large-scale dynamical systems. Viewing the LPMF optimization objectives as risk functions, we can further use the reparameterization trick to enforce positivity constraints on the states and the ODE parameters, which will be discussed in Section 4.4. The strengths and

weaknesses of the LPMF solution are examined empirically and discussed in Chapter 6.

## 4.1 Laplace approximation

This section reviews the *Laplace approximation* technique in the context of approximating an unknown multivariate probability distribution based on (Bishop, 2006, Section 4.4) and (MacKay, 2003, Section 27).

Suppose we are interested in the following probability distribution

$$p(\mathbf{x}) = \frac{\tilde{p}(\mathbf{x})}{Z} \quad (4.1)$$

where  $\mathbf{x} \in \mathbb{R}^D$ ,  $\tilde{p}(\mathbf{x})$  is known and  $Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$  is the normalization constant assumed to be intractable. The Laplace approximation of  $p$  is a Gaussian distribution  $q$  such that its mean is centered on a mode of the  $p$  (MacKay, 2003).

Assuming that  $p$  has a peak at the point  $\mathbf{x}_0$ , then the gradient  $\nabla p(\mathbf{x}_0) = 0$ , or equivalently,  $\nabla \tilde{p}(\mathbf{x}_0) = 0$ . The second-order *Taylor expansion* of  $\ln \tilde{p}(\mathbf{x})$  around  $\mathbf{x}_0$  is given by

$$\begin{aligned} \ln \tilde{p}(\mathbf{x}) &\approx \ln \tilde{p}(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla \tilde{p}(\mathbf{x}_0) - \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \tilde{\mathbf{H}}(\mathbf{x} - \mathbf{x}_0) \\ &= \ln \tilde{p}(\mathbf{x}_0) - \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \tilde{\mathbf{H}}(\mathbf{x} - \mathbf{x}_0) \end{aligned} \quad (4.2)$$

where  $\tilde{\mathbf{H}}$  is the negative of the Hessian matrix  $\mathbf{H}$  at  $\mathbf{x}_0$  such that

$$\tilde{H}_{ij} = -H_{ij} = -\frac{\partial^2}{\partial x_i \partial x_j} \ln \tilde{p}|_{\mathbf{x}=\mathbf{x}_0} \quad (4.3)$$

or

$$\tilde{\mathbf{H}} = -\mathbf{H} = -\nabla \nabla \ln \tilde{p}(\mathbf{x}_0) \quad (4.4)$$

Exponentiating both sides of (Eq. 4.2), we have

$$\tilde{p}(\mathbf{x}) \approx \tilde{p}(\mathbf{x}_0) \exp \left[ -\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \tilde{\mathbf{H}}(\mathbf{x} - \mathbf{x}_0) \right] \quad (4.5)$$

which is of quadratic form and can be normalized by inspection to obtain the following Gaussian distribution:

$$\begin{aligned} q(\mathbf{x}) &= \frac{|\tilde{\mathbf{H}}|^{\frac{1}{2}}}{(2\pi)^{\frac{D}{2}}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \tilde{\mathbf{H}}(\mathbf{x} - \mathbf{x}_0) \right] \\ &= \mathcal{N}(\mathbf{x}|\mathbf{x}_0, \tilde{\mathbf{H}}^{-1}) \end{aligned} \quad (4.6)$$



## 4.2 Laplace mean-field approximation

As discussed in Section 3.3, the conditional distributions  $p(\theta|Y, X, \boldsymbol{\varphi}, \gamma)$  (Eq. 3.30) and  $p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \sigma, \gamma)$  (Eq. 3.29) are both Gaussians provided that the ODEs fulfill the structural assumption described by (Eq. 3.28). If we relax the constraint on the structure of the ODEs, the distributions  $p(\theta|Y, X, \boldsymbol{\varphi}, \gamma)$  and  $p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \sigma, \gamma)$  are no longer Gaussians and cannot be normalized in closed-form anymore.

From the gradient matching model, we have

$$\begin{aligned} p(\theta|Y, X, \boldsymbol{\varphi}, \gamma) &\stackrel{(a)}{=} p(\theta|X, \boldsymbol{\varphi}, \gamma) \\ &\propto p(\theta) \int p(\dot{X}|X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma) d\dot{X} \\ &\propto \prod_k \mathcal{N}(f_k(X, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1}) \end{aligned} \quad (4.7)$$

where (a) holds since  $\boldsymbol{\theta}$  depends indirectly on  $Y$  through  $X$  (Gorbach et al., 2017), and  $p(\dot{X}|X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma)$  is the product of experts result (Eq. 3.19).

Similarly, for each state  $x_u$  with  $u = 1, \dots, K$ , we have

$$\begin{aligned} p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \sigma, \gamma) &= p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \sigma) \int p(\dot{X}|x_u, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma) d\dot{X} \\ &\stackrel{(b)}{\propto} \mathcal{N}(x_u | \boldsymbol{\mu}_u(\mathbf{y}_u), \boldsymbol{\Sigma}_u) \prod_k \mathcal{N}(f_k(X, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1}) \end{aligned} \quad (4.8)$$

where (b) holds because  $p(x_u|Y, X_{/\{u\}}, \boldsymbol{\varphi}, \sigma)$  depends only on  $\mathbf{y}_u$  as the consequence of the independent Gaussian process prior assumption on states (Eq. 3.14), and  $p(\dot{X}|x_u, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma)$  is equivalent to  $p(\dot{X}|X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma)$ .

### Cost functions

If we follow the same factorization assumption over the states and the parameters as Gorbach et al. (2017), and require each component of the proxy distribution  $Q(X, \boldsymbol{\theta})$  to be Gaussian, then  $Q(X, \boldsymbol{\theta})$  is of the following form

$$\begin{aligned} Q(X, \boldsymbol{\theta}) &= q(\boldsymbol{\theta} | \boldsymbol{\eta}_\theta, \boldsymbol{\Xi}_\theta) \prod_u q(x_u | \boldsymbol{\eta}_{x_u}, \boldsymbol{\Xi}_{x_u}) \\ &= \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\eta}_\theta, \boldsymbol{\Xi}_\theta) \prod_u \mathcal{N}(x_u | \boldsymbol{\eta}_{x_u}, \boldsymbol{\Xi}_{x_u}) \end{aligned} \quad (4.9)$$

Based on (Eq. 4.7) and the Laplace approximation technique reviewed be-

fore, the mean of the Gaussian proxy  $q(\theta|\eta_\theta, \Xi_\theta)$  can be found by

$$\begin{aligned}\eta_\theta &= \arg \max_{\theta} \sum_k \ln \mathcal{N}(f_k(X, \theta) | m_k, \Lambda_k^{-1}) \\ &= \arg \min_{\theta} \frac{1}{2} \sum_k (f_k - m_k)^T \Lambda_k (f_k - m_k) \\ &= \arg \min_{\theta} \text{cost}_\theta(X, \theta, m, \Lambda)\end{aligned}\tag{4.10}$$

where  $f_k = f_k(X, \theta)$ ,  $m = [m_1, \dots, m_K]$  and  $\Lambda = [\Lambda_1, \dots, \Lambda_K]$ . The precision matrix  $\Xi_\theta^{-1}$  is then the Hessian of the objective function (Eq. 4.10) evaluated at the optimal point  $\eta_\theta$

$$\Xi_\theta^{-1} = \nabla \nabla \text{cost}_\theta |_{\theta=\eta_\theta}\tag{4.11}$$

Similarly, the mean of the Gaussian proxy  $q(x_u | \eta_{x_u}, \Xi_{x_u})$  for  $u = 1, \dots, K$  can be found using (Eq. 4.8) as

$$\begin{aligned}\eta_{x_u} &= \arg \max_{x_u} [\ln \mathcal{N}(x_u | \mu_u(y_u), \Sigma_u) + \sum_k \ln \mathcal{N}(f_k(X, \theta) | m_k, \Lambda_k^{-1})] \\ &= \arg \min_{x_u} \frac{1}{2} [(x_u - \mu_u(y_u))^T \Sigma_u^{-1} (x_u - \mu_u(y_u)) + \sum_k (f_k - m_k)^T \Lambda_k (f_k - m_k)] \\ &= \arg \min_{x_u} \text{cost}_{x_u}(x_u, X_{/\{u\}}, \theta, \mu_u(y_u), \Sigma_u, m, \Lambda)\end{aligned}\tag{4.12}$$

The corresponding precision matrix  $\Xi_{x_u}^{-1}$  is given by

$$\Xi_{x_u}^{-1} = \nabla \nabla \text{cost}_{x_u} |_{x_u=\eta_{x_u}}\tag{4.13}$$

### LPMF algorithm

To conclude, given the initialization from the Gaussian processes, the optimization step (Eq. 4.10) first updates the estimation of the parameters, and then for each state, the optimization procedure (Eq. 4.19) improves state estimation. These two steps are executed in turn iteratively until convergence of the estimation or the maximum allowed number of iterations is reached. Lastly, the covariance matrices of the estimation results are calculated using (Eq. 4.11) and (Eq. 4.13). The pseudocode for the LPMF algorithm is given in Algorithm 2.

Note that when the ODEs satisfy the structural assumption about the parameters  $\theta$ , closed-form solutions for  $\eta_\theta$  and  $\Xi_\theta$  can be obtained using (Eq. 3.31) and (Eq. 3.32) respectively, which requires rewriting the ODEs as linear combination of the parameters plus a term that is independent of them. Similarly, analytical solutions for  $\eta_{x_u}$  and  $\Xi_{x_u}$  can be derived from (Eq. 3.33) and (Eq. 3.34) for  $u = 1, \dots, K$  when the structural assumption about the

---

**Algorithm 2** Pseudocode for the LPMF algorithm.
 

---

```

1: for  $u = 1, \dots, K$  do
2:   Initialize  $\mu_u(y_u)$ ,  $\Sigma_u$ ,  $m_u$  and  $\Lambda_u$  using Gaussian process regression
3: end for
4:
5: while not converged or maximum iteration not reached do
6:    $\eta_\theta = \arg \min_\theta \text{cost}_\theta(\eta_X, \theta, m, \Lambda)$ 
7:   for  $u = 1, \dots, K$  do
8:      $\eta_{x_u} = \arg \min_{x_u} \text{cost}_{x_u}(x_u, \eta_{X/\{u\}}, \eta_\theta, \mu_u(y_u), \Sigma_u, m, \Lambda)$ 
9:   end for
10: end while
11:
12:  $\Xi_\theta^{-1} = \nabla \nabla \text{cost}_\theta(\eta_X, \eta_\theta, m, \Lambda)$ 
13: for  $u = 1, \dots, K$  do
14:    $\Xi_{x_u}^{-1} = \nabla \nabla \text{cost}_{x_u}(\eta_{x_u}, \eta_{X/\{u\}}, \eta_\theta, \mu_u(y_u), \Sigma_u, m, \Lambda)$ 
15: end for
    
```

---

states is fulfilled. This also requires rewriting of the ODEs as a linear combination of the states plus an independent term, which is more cumbersome than the previous case for dynamical systems with large number of states.

### 4.3 Derivation for the gradients and Hessians

This section derives the gradient and Hessian for the cost function  $\text{cost}_{x_u}$  (Eq. 4.19) and shows that they can be evaluated efficiently.

Recall that the cost function  $\text{cost}_{x_u}$  for state  $u$  is given by

$$\text{cost}_{x_u} = \frac{1}{2} [(x_u - \mu_u(y_u))^T \Sigma_u^{-1} (x_u - \mu_u(y_u)) + \sum_k (f_k - m_k)^T \Lambda_k (f_k - m_k)]$$

We are interested in the gradient and Hessian of  $\text{cost}_{x_u}$  w.r.t state  $x_u$ , i.e.  $\nabla_{x_u} \text{cost}_{x_u}$  and  $\nabla \nabla_{x_u} \text{cost}_{x_u}$ . Since differentiation is a linear operation, we can look at the gradient and Hessian of each component of the cost function separately.

Using matrix derivative and the fact that  $\Sigma_u^{-1}$  is symmetric, the gradient and Hessian of the first term in  $\text{cost}_{x_u}$  are

$$\nabla_{x_u} \frac{1}{2} (x_u - \mu_u(y_u))^T \Sigma_u^{-1} (x_u - \mu_u(y_u)) = \Sigma_u^{-1} x_u \quad (4.14)$$

and

$$\nabla \nabla_{x_u} \frac{1}{2} (x_u - \mu_u(y_u))^T \Sigma_u^{-1} (x_u - \mu_u(y_u)) = \Sigma_u^{-1} \quad (4.15)$$

#### 4. LAPLACE MEAN-FIELD APPROXIMATION

Using the *chain rule* and the fact that  $\Lambda_k$  is symmetric, for each component from the second term in  $\text{cost}_{x_u}$ , the gradient is given by

$$\begin{aligned}
& \nabla_{x_u} (\mathbf{f}_k - \mathbf{m}_k)^T \Lambda_k (\mathbf{f}_k - \mathbf{m}_k) \\
&= \begin{bmatrix} \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_1)} & \dots & \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_1)} \\ \vdots & \ddots & \vdots \\ \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_N)} & \dots & \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_N)} \end{bmatrix} \Lambda_k (\mathbf{f}_k - \mathbf{m}_k) \\
&= \begin{bmatrix} \frac{\partial(\mathbf{f}_k)_1}{\partial x_u(t_1)} & \dots & \frac{\partial(\mathbf{f}_k)_N}{\partial x_u(t_1)} \\ \vdots & \ddots & \vdots \\ \frac{\partial(\mathbf{f}_k)_1}{\partial x_u(t_N)} & \dots & \frac{\partial(\mathbf{f}_k)_N}{\partial x_u(t_N)} \end{bmatrix} \Lambda_k (\mathbf{f}_k - \mathbf{m}_k) \\
&\quad - \begin{bmatrix} \frac{\partial(\mathbf{m}_k)_1}{\partial x_u(t_1)} & \dots & \frac{\partial(\mathbf{m}_k)_N}{\partial x_u(t_1)} \\ \vdots & \ddots & \vdots \\ \frac{\partial(\mathbf{m}_k)_1}{\partial x_u(t_N)} & \dots & \frac{\partial(\mathbf{m}_k)_N}{\partial x_u(t_N)} \end{bmatrix} \Lambda_k (\mathbf{f}_k - \mathbf{m}_k) \tag{4.16}
\end{aligned}$$

Since the states do not appear inside the vector field across time points, the first matrix above is necessarily a diagonal matrix. If the  $\mathbf{f}_k$  is independent of  $x_u$ , then this matrix is a zero matrix. Since  $\mathbf{m}_k = \mathbf{C}_{\varphi_k}'' \mathbf{C}_{\varphi_k}^{-1} \mathbf{x}_k$  is a linear combination of  $\mathbf{x}_k$ , the second matrix is the transpose of  $\mathbf{C}_{\varphi_k}'' \mathbf{C}_{\varphi_k}^{-1}$  if  $k = u$  and is a zero matrix otherwise.

The Hessian of  $(\mathbf{f}_k - \mathbf{m}_k)^T \Lambda_k (\mathbf{f}_k - \mathbf{m}_k)$  w.r.t.  $x_u$  also exhibits such sparse structure. In brief, the  $(i, j)$ -th entry of the Hessian is given by

$$\begin{aligned}
& \frac{\partial^2 (\mathbf{f}_k - \mathbf{m}_k)^T \Lambda_k (\mathbf{f}_k - \mathbf{m}_k)}{\partial x_u(t_i) \partial x_u(t_j)} \\
&= \begin{bmatrix} \frac{\partial^2(\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_i) \partial x_u(t_j)} & \dots & \frac{\partial^2(\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_i) \partial x_u(t_j)} \end{bmatrix} \Lambda_k (\mathbf{f}_k - \mathbf{m}_k) \\
&\quad + \begin{bmatrix} \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_j)} & \dots & \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_j)} \end{bmatrix} \Lambda_k \begin{bmatrix} \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_i)} \\ \vdots \\ \frac{\partial(\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_i)} \end{bmatrix} \tag{4.17}
\end{aligned}$$

The overall gradients and Hessian are just the sums of the respective components discussed above. With proper handling of multiplication involving diagonal matrices and caching of constant terms, the gradient and Hessian of the cost function  $\text{cost}_{x_u}$  can be very efficiently calculated, as will be shown empirically in Chapter 6.

Although the gradient and Hessian of  $\text{cost}_{\theta}$  (Eq. 4.10) w.r.t.  $\theta$  do not have such structure, they can still be evaluated relatively efficiently since the number of parameters is usually small in comparison to the number of states. In

this work, both symbolic tools and machine learning libraries supporting auto-differentiation are employed to calculate the gradient and Hessian of  $cost_\theta$ .

## 4.4 Positivity constraints

An advantage of the LPMF solution is that we can treat the optimization objectives  $cost_\theta$  and  $cost_{x_u}$  for  $u = 1, \dots, K$  as risk functions, and hence the name “cost” is used in the notations. This allows us to apply the following reparameterization trick to enforce a positivity constraint on the parameters and the states.

From Section 4.2, the optimal estimation for the parameters and states is given by

$$\theta^* = \arg \min_{\theta} cost_\theta(X, \theta, m, \Lambda) \quad (4.18)$$

$$x_u^* = \arg \min_{x_u} cost_{x_u}(x_u, X_{/\{u\}}, \theta, \mu_u(y_u), \Sigma_u, m, \Lambda) \quad (4.19)$$

for  $u = 1, \dots, K$ . Suppose we desire the parameters to be positive values. Instead of inferring  $\theta$  directly, the reparameterization trick transforms the cost function  $cost_\theta$  into a new cost function  $cost_{\tilde{\theta}}$ , where we define  $\theta = [\theta_1, \dots, \theta_M]^\top = [e^{\tilde{\theta}_1}, \dots, e^{\tilde{\theta}_M}]^\top$ . Because the exponential function is monotonic, we can therefore first find  $\tilde{\theta}^*$  that minimizes the new cost

$$\begin{aligned} \tilde{\theta}^* &= \arg \min_{\tilde{\theta}} cost_\theta(X, e^{\tilde{\theta}}, m, \Lambda) \\ &= \arg \min_{\tilde{\theta}} \sum_k \ln \mathcal{N}(f_k(X, e^{\tilde{\theta}}) | m_k, \Lambda_k^{-1}) \end{aligned} \quad (4.20)$$

and then exponentiate it element-wise to obtain an estimation for  $\theta^* = [e^{\tilde{\theta}_1^*}, \dots, e^{\tilde{\theta}_M^*}]^\top$ , which would correspond to the configuration that minimizes the original  $cost_\theta$ . Since  $e^r > 0$  for any  $r \in \mathbb{R}$ , we essentially restrict  $\theta^*$  to only contain positive values.

Analogously, the positivity constraint on states can be achieved by transforming the cost function  $cost_{x_u}$  into its equivalent  $cost_{\tilde{x}_u}$ , where we define  $x_u = [e^{\tilde{x}_u(t_1)}, \dots, e^{\tilde{x}_u(t_N)}]^\top$  for  $u = 1, \dots, K$ . Using the chain rule, the vector field after the transformation is given by

$$\frac{d\tilde{x}(t)}{dt} = \frac{f(e^{\tilde{x}(t)}, \theta)}{e^{\tilde{x}(t)}} \quad (4.21)$$

Note that the application of the positivity constraint on states is flexible in the sense that not all the states are required to be constrained at the same

time. Also, the constraint can theoretically be applied to both the parameters and the states at the same time.

To distinguish from the unconstrained Laplace mean-field approximation, we use the term LPMP-POS to refer to this extension. The pseudocode for this extension is essentially the same as Algorithm 2 except for the replacement of the optimization functions and variables with the extra step to transform the inference result back to the original variables at the end.

One caveat for this approach is that the covariance matrix calculated for the transformed variables cannot be transformed back to the original variables easily. Also, interpretability from a probabilistic point of view is lost since exponentiation is a non-linear operation and rigorous treatment of the change of random variables would require the evaluation of the relevant Jacobian determinant, which is computationally very expensive.

---

## Extension to Random Dynamical Systems

---

As another important family of differential equations, *random ordinary differential equations (RODEs)* are closely related to both ODEs and SDEs. A system of RODEs is simply a set of ODEs with a stochastic process in its vector field functions (Kloeden and Jentzen, 2007), while an SDE system can be analyzed using its RODE counterpart (Sussmann, 1978; Imkeller and Schmalfuss, 2001). Since RODEs are pathwise ODEs, the Laplace mean-field approximation described in Chapter 4 can then be used to infer the states and parameters of the RODEs, or equivalently, the corresponding SDEs.

This chapter is organized as follows. Section 5.1 gives a brief introduction to RODEs. In Section 5.3, the Laplace mean-field approximation technique is applied to RODEs to devise an *ensemble-like* (Murphy, 2012, Section 16.6) solution to infer the states and parameters for diffusion processes. We demonstrate the performance and accuracy of the solution empirically by comparing with other state-of-art techniques in Chapter 6.

### 5.1 Random ordinary differential equations

Adopting the definition from Kloeden and Jentzen (2007), RODEs are simply ODEs with a stochastic process in their vector fields. Let  $f : \mathbb{R}^K \times \mathbb{R}^W \mapsto \mathbb{R}^K$  be a continuous function, and  $(\zeta_t)_{t \in [0, T]}$  be an  $\mathbb{R}^W$  valued stochastic process with continuous sample paths defined on the complete probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . For all  $\omega \in \Omega$ , a  $K$ -dimensional RODE defined as

$$\frac{dx(t)}{dt} = f(x(t), \zeta_t(\omega)) \quad (5.1)$$

is a *non-autonomous* ODE system

$$\dot{x}(t) = \frac{dx(t)}{dt} = F_\omega(x, t) = f(x(t), \omega_t) \quad (5.2)$$

An example (Grüne and Kloeden, 2001) of a scalar RODE with additive noise is given by

$$\frac{dx(t)}{dt} = -x + \cos(W_t(\omega)) \quad (5.3)$$

where  $W_t$  is a one-dimensional Wiener process. A RODE example with multiplicative noise can be defined similarly but is not considered in this work.

Following Kloeden and Jentzen (2007), to ensure the existence of a unique solution for the initial value problem defined in Section 2.1 on the finite time interval  $[0, T]$ , we typically assume that  $f$  is arbitrarily smooth, i.e. it is infinitely differentiable in its variables, and thus is locally *Lipschitz* in  $x$ . Since the stochastic process is usually only *Hölder continuous* in time, the vector field of the non-autonomous ODEs  $F_\omega(x, t)$  is therefore continuous but not differentiable in time for every fixed realization of  $\omega \in \Omega$ .

Because for every fixed realization of  $\omega \in \Omega$ , the RODEs (Eq. 5.1) turn into a deterministic ODE system (Eq. 5.2), one approach to solving the RODEs is to use sampling methods to first obtain many sample paths, and then solving each sample path deterministically. In order to cover the statistics of the solution, a massive number of ODEs must be solved efficiently. A high performance related study was conducted by Riesinger et al. (2016), where the sample paths are solved in parallel on a modern GPU cluster.

## 5.2 Doss-Sussmann/Imkeller-Schmalzfuss correspondence

Since a system of RODEs can be analyzed pathwise using deterministic calculus, it offers an opportunity to study its related SDEs as discussed below.

First of all, it has been shown by Jentzen and Kloeden (2011) that any RODE system with a Wiener process can be expressed as its equivalent SDE system. Using the scalar RODE in (Eq. 5.3) as an example, its SDE formulation is described as

$$d \begin{pmatrix} x_t \\ y_t \end{pmatrix} = \begin{pmatrix} -x_t + \cos(y_t) \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} dW_t \quad (5.4)$$

Similarly, any finite dimensional SDE system can be transformed into its equivalent ODEs by utilizing the *Doss-Sussmann/Imkeller-Schmalzfuss* (Sussmann, 1978; Imkeller and Schmalzfuss, 2001) correspondence. Specifically, for SDE models with additive noise, the statement from (Jentzen and Kloeden, 2011, Chapter 2) is stated as the following proposition.



**Proposition 5.1** *Any finite dimensional SDE system can be transformed into an equivalent RODE system and vice versa as follows:*

$$dx(t) = f(x(t), \theta)dt + dW_t \iff \frac{dz(t)}{dt} = f(z(t) + O_t, \theta) + O_t \quad (5.5)$$

where  $z(t) = x(t) - O_t$  and  $O_t$  is a stationary stochastic Ornstein-Uhlenbeck process defined as

$$dO_t = -O_t dt + dW_t \quad (5.6)$$

### 5.3 Laplace mean-field for random dynamical systems

As discussed previously, although RODE sample paths can be analyzed using deterministic calculus, the existence of the stochastic process causes traditional numerical schemes such as the *Euler* and the *Runge-Kutta* methods (Butcher, 2016) to fail to achieve their usual order of convergence when applied to RODEs (Grüne and Kloeden, 2001). In the past, improved numerical solutions such as the integral versions of the implicit Taylor-like expansions (Kloeden and Jentzen, 2007) have been proposed to achieve better result.

On the other hand, since the solution paths of the RODEs are once differentiable, the gradient matching model can be ideally applied to the RODEs. Moreover, the computational efficiency of the LPMF method allows a large number of RODE sample paths to be solved simultaneously. Lastly, by using the Doss-Sussmann/Imkeller-Schalfuss correspondence described before, we derive the following ensemble gradient matching algorithm, denoted as LPMF-SDE, to infer the states and the parameters of the SDEs without requiring any stochastic calculus.

---

**Algorithm 3** Pseudocode for the LPMF-SDE algorithm.

---

- 1: Transform the SDEs into RODEs using (Eq. 5.5).
  - 2: Generate  $N_{paths}$  RODEs sample paths using each time an independently generated Ornstein-Uhlenbeck process sample path.
  - 3: **for**  $i = 1, \dots, N_{paths}$  **do**
  - 4:     Estimate the states and parameters using Algorithm 2.
  - 5: **end for**
  - 6: Average the estimation results from all the sample paths.
- 

Since the experiments are conducted with the Lorenz 63 and the Lorenz 96 models, the estimation of the parameters are carried out by extending the

closed-forms solutions from (Eq. 3.31) and (Eq. 3.32) as follows:

$$\boldsymbol{\eta}_\theta = \boldsymbol{\Omega}_\theta \sum_k \mathbf{B}_{\theta k}^T \boldsymbol{\Lambda}_k (\mathbf{m}_k - \mathbf{b}_{\theta k} - \mathbf{O}_k) \quad (5.7)$$

$$\boldsymbol{\Xi}_\theta^{-1} = \sum_k \mathbf{B}_{\theta k}^T \boldsymbol{\Lambda}_k \mathbf{B}_{\theta k} \quad (5.8)$$

where

$$\mathbf{m}_k = \mathbf{C}_{\varphi_k}'' \mathbf{C}_{\varphi_k}^{-1} \mathbf{z}_k \quad (5.9)$$

$$\boldsymbol{\Lambda}_k^{-1} = \mathbf{C}_{\varphi_k}'' - {}^t \mathbf{C}_{\varphi_k} \mathbf{C}_{\varphi_k}^{-1} \mathbf{C}_{\varphi_k}' + \gamma_k \mathbf{I} \quad (5.10)$$

Similar to the notations introduced in Section 2.1,  $\mathbf{z}_u \in \mathbb{R}^N$  and  $\mathbf{Z} \in \mathbb{R}^{K \times N}$  refer to the states of the RODE sample path, while  $\mathbf{O}$  refer to the states of the Ornstein-Uhlenbeck process. Note that the rewriting of the vector field as a linear combination of the parameters  $\boldsymbol{\theta}$  and an extra term should satisfy

$$\mathbf{B}_{\theta k}(\mathbf{Z} + \mathbf{O})\boldsymbol{\theta} + \mathbf{b}_{\theta k}(\mathbf{Z} + \mathbf{O}) + \mathbf{O} = \mathbf{f}_k(\mathbf{Z} + \mathbf{O}, \boldsymbol{\theta}) + \mathbf{O} \quad (5.11)$$

for  $k = 1, \dots, K$ .

Due to the complexity of expressing as a linear combination of the states and an extra term, the states of the RODE sample path are estimated using gradient-descent by minimizing the adapted cost function (Eq. 4.19) as follows:

$$\text{cost}_{\mathbf{z}_u} = \ln [\mathcal{N}(\mathbf{z}_u | \boldsymbol{\mu}(\mathbf{y}_u), \boldsymbol{\Sigma}_u) \prod_k \mathcal{N}(\mathbf{f}_k(\mathbf{Z} + \mathbf{O}, \boldsymbol{\theta}) + \mathbf{O} | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1})] \quad (5.12)$$

for  $u = 1, \dots, K$ , where

$$\boldsymbol{\mu}_u(\mathbf{y}_u) = \mathbf{C}_{\varphi_u} (\mathbf{C}_{\varphi_u} + \sigma_u^2 \mathbf{I})^{-1} \mathbf{y}_u \quad (5.13)$$

$$\boldsymbol{\Sigma}_u = \sigma_u^2 \mathbf{C}_{\varphi_u} (\mathbf{C}_{\varphi_u} + \sigma_u^2 \mathbf{I})^{-1} \quad (5.14)$$

---

# Experiments

---

This chapter examines the estimation accuracy, runtime performance and scalability of the general LPMF method and its extensions introduced in Chapter 4 and Chapter 5 by comparing with the other state-of-art inference techniques empirically. In Section 6.1, the implementation of the algorithms is discussed. Then four dynamical systems with different dimensionality and complexity are used for the experiments. Section 6.2 and Section 6.3 consider the state and parameter estimation for deterministic dynamical systems, while Section 6.4 and Section 6.5 consider two random dynamical systems.

## 6.1 Implementation

The source code of this work is implemented from the ground up using the general purpose programming language Python <sup>3</sup><sup>1</sup>. The MATLAB<sup>2</sup> code for the VGMGP algorithm from Gorbach et al. (2017) is used as the blueprint during implementation. Several important open-source packages have made this Python solution possible:

- The SciPy<sup>3</sup> package provides a powerful optimization module and other I/O utilities to read and write files in MATLAB format.
- The NumPy<sup>4</sup> package is the backbone for linear algebra operations and random variable generation.
- The SymPy<sup>5</sup> package is used for symbolic mathematics when implementing the general interface for dynamical systems and kernel func-

---

<sup>1</sup><https://www.python.org/>

<sup>2</sup><https://www.mathworks.com/products/matlab.html>

<sup>3</sup><https://www.scipy.org/>

<sup>4</sup><http://www.numpy.org/>

<sup>5</sup><http://www.sympy.org/en/index.html>

**Table 6.1:** Experimental setup for the Lotka-Volterra model. The system dimension is denoted by  $K$  and the number of observable dimensions is  $K_{obs}$ . Based on the parameter values  $\alpha$ ,  $\beta$ ,  $\delta$ , and  $\gamma$ , the ODEs are integrated from time  $t_0$  to  $t_T$  with a step size of  $\delta t$ . The observation noise variance  $\sigma_k^2$  is assumed to be the identical for each state. For each time unit,  $freq_{obs}$  denotes the number of observations to be collected, which are equally distributed over the time line.

$K$	$K_{obs}$	$t_0$	$t_T$	$\delta t$	$\alpha, \beta, \delta, \gamma$	$\sigma_k^2$	$freq_{obs}$
2	2	0	2	0.01	2, 1, 1, 4	0.1	10

tions. It is also used to calculate the gradients and Hessians of the cost functions for part of the solution.

- The matplotlib<sup>6</sup> package helps to produce publication quality plotting.
- The Jupyter Notebook<sup>7</sup> is used as the GUI interface to combine code, visualizations and documentation in a sharable format.
- The sdeint<sup>8</sup> package provides the numerical methods to solve SDEs.
- The TensorFlow<sup>9</sup> package is a machine learning library popular among the deep learning community, which provides auto-differentiation support for part of the solution.

## 6.2 Lotka-Volterra model

The first deterministic dynamical system examined in this chapter is the *Lotka-Volterra* model (Lotka, 1932), which is frequently used in ecology to describe the interaction between the prey species and the predator species over time. The model consists of two first-order, nonlinear differential equations where the states  $x(t), y(t) \in \mathbb{R}_{\geq 0}$  are the populations of the prey and the predator respectively at time point  $t$ . The ODEs of the model are given by

$$\begin{aligned}\dot{x}(t) &= \alpha x(t) - \beta x(t)y(t) \\ \dot{y}(t) &= \delta x(t)y(t) - \gamma y(t)\end{aligned}\tag{6.1}$$

where  $\alpha, \beta, \delta, \gamma \in \mathbb{R}^+$  are the parameters controlling the dynamics.

### Experimental setup

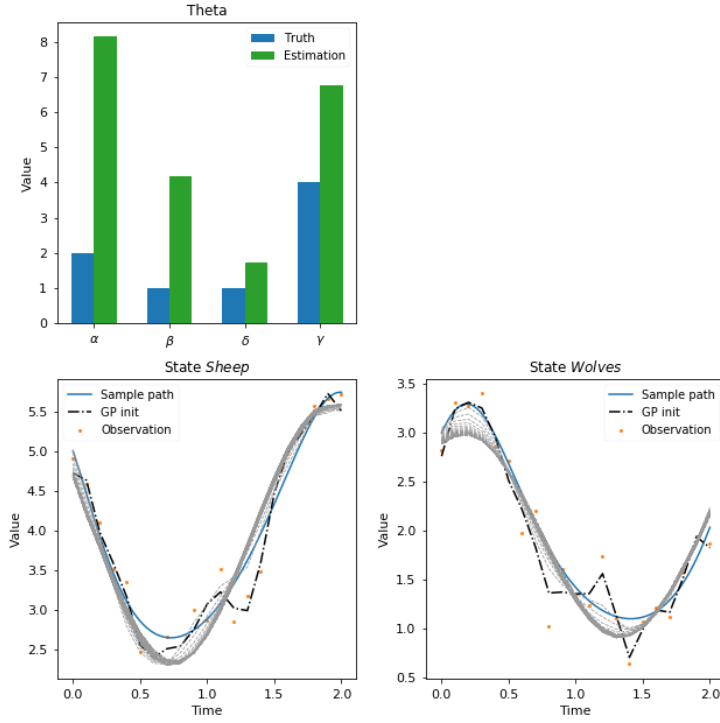
Table 6.1 shows the setup for the experiment. The experiment is repeated 10 times with each time an independently collected observation set. Since

<sup>6</sup><http://matplotlib.org/>

<sup>7</sup><http://jupyter.org/>

<sup>8</sup><https://github.com/mattja/sdeint>

<sup>9</sup><https://www.tensorflow.org/>



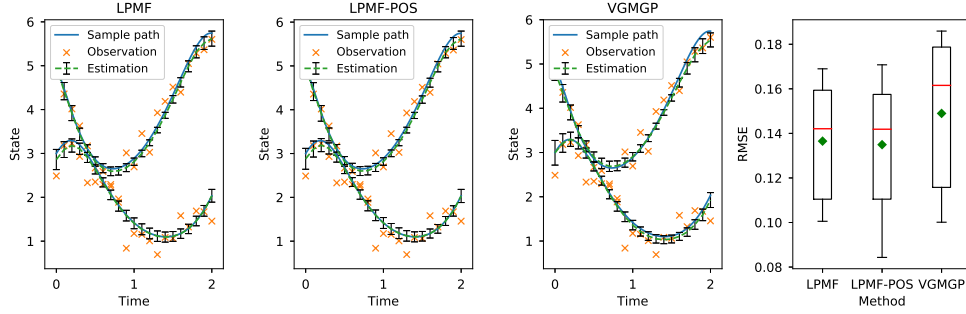
**Figure 6.1:** Inference for the Lotka-Volterra model fails when both the states and parameters are constrained to be positive.

the LPMF method is derived from the VGMGP method, in the following, we compare the results using both methods. The LPMF method is run first without any positivity constraint and then it is run again with positivity constraint on the parameters, which are referred to as LPMF and LPMP-POS respectively. It would be interesting to constrain both the states and the parameters to be positive, but the inference fails as shown in Figure 6.1. The reason for that is unclear yet due to time constraints and requires further investigation. In order to provide a fair comparison in terms of runtime, both methods are deployed on a desktop with an Intel i5 quad-core CPU and 16 GB of memory.

### State estimation

Figure 6.2 shows the results after the 10 independent runs. For illustration purposes, the observations from one run is also plotted to indicate the noise level. The dotted green lines are obtained by averaging the means of the state estimation for the 10 runs, while the error bars indicate one standard deviation of the means. The figure shows that the state estimation results are very close to each other and is almost identical to the ground truth.

## 6. EXPERIMENTS



**Figure 6.2:** State estimation results for the Lotka-Volterra model. The left plot shows the result using LPMF; the left middle plot shows the result using LPMF-POS; the right middle plot shows the result using VGMGP; the right plot summarizes the RMSE after 10 independent runs. The error bars in the first three plots indicate one standard deviation.

To quantify the accuracy, the root mean square error (RMSE) is used and is defined as follows:

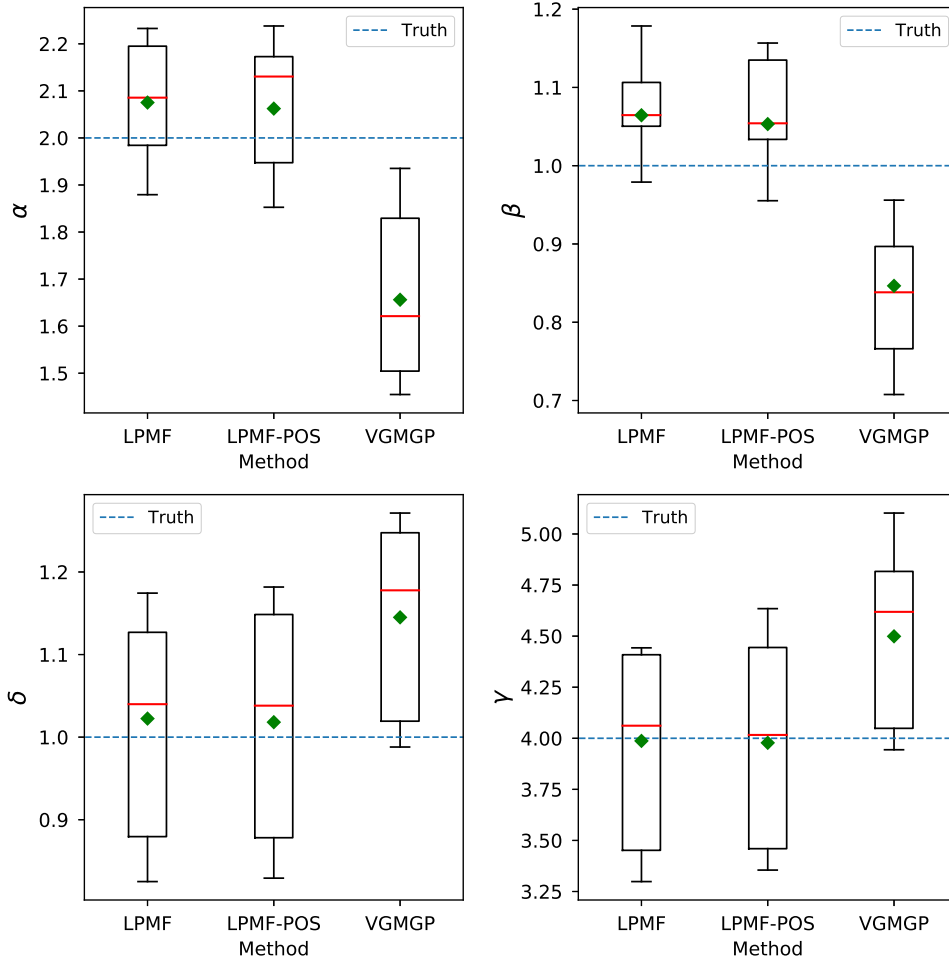
$$RMSE = \frac{1}{K} \sum_k \sqrt{\frac{1}{N} \sum_{n=1}^N (\hat{x}_k(t_n) - x_k(t_n))^2} \quad (6.2)$$

where  $K$  indicates the number of states,  $N$  is the total number of observations for each dimension, and  $\hat{x}_k(t_n)$  and  $x_k(t_n)$  are the predicted and true values for the  $k$ -th state at time point  $t$  respectively. The RMSEs are very close to each other with the LPMF method having slightly lower error.

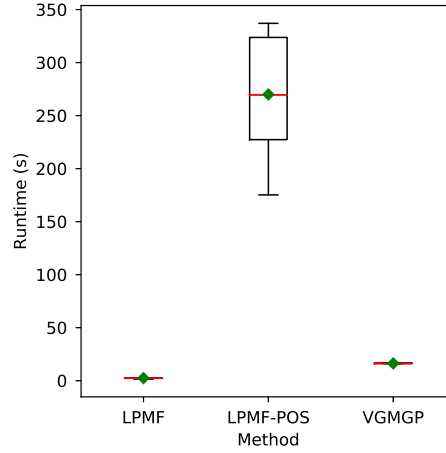
### Parameter estimation

The estimation results for the parameters of the Lotka-Volterra model are shown in Figure 6.3. The LPMF method again achieves better results than the VGMGP algorithm, even with the positivity constraint on the parameters. The mean values for the prediction from LPMF are also very close to the true parameter values.

To explain this, first note that both methods assume the decoupling of states from the other states and the decoupling of the states and the parameters when constructing the proxy distribution  $Q$ . However, the VGMGP method further assumes the decoupling of the same states across time points, which is not the case for LPMF. Since the ODEs of the Lotka-Volterra model satisfies the structural assumption, if the conditional distributions in (Eq. 3.29) and (Eq. 3.30) are indeed Gaussian, then Laplace approximation is expected to correctly find the mode the distribution.



**Figure 6.3:** Parameter estimation results for the Lotka-Volterra model. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles. The true parameter value is shown as the dotted blue line.



**Figure 6.4:** Runtime performance for the Lotka-Volterra model. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles.

### Runtime performance

In terms of runtime, both the LPMF method when no positivity constraint is imposed and the VGMGP method are extremely fast. On average, the LPMF method finishes in 2.4 seconds and the VGMGP method completes in 16.3 seconds. Since LPMF is implemented in Python while VGMGP is implemented in MATLAB, exact comparison is infeasible. In general, the LPMF is likely to be more efficient since the evaluation of the expectations (Eq. 3.41) and (Eq. 3.42) is not required. Moreover, the states are inferred by minimizing the cost function for the LPMF method in this experiment. If closed-form solutions are used, it would be expected to be even faster.

Lastly, it is unclear what causes the slow down of the LPMF algorithm after the introduction of the positivity constraint. Given the time constraints and in order to achieve fast prototyping, the gradients and Hessians of the cost functions with positivity constraints are obtained from symbolic libraries. This is in contrast to the highly vectorized implementation when no positivity constraint is enforced. Part of the reason is probably due to the inefficient implementation, but it requires further investigation.

## 6.3 Protein signalling transduction pathway

As already mentioned in Section 1.1, the biochemical *protein signalling transduction pathway* (Vyshemirsky and Girolami, 2007) is a signal transduction cascade model describing the dynamics among protein species. The model



**Table 6.2:** Experimental setup for the protein signaling transduction pathway model. The system dimension is denoted by  $K$  and the number of observable dimensions is  $K_{obs}$ . Based on the parameters  $k_1, k_2, k_3, k_4, V, Km$ , the ODEs are solved from time  $t_0$  to  $t_T$  with a step size of  $\delta t$ . The observation noise variance  $\sigma_k^2$  is assumed to be the identical for each state.

$K$	$K_{obs}$	$t_0$	$t_T$	$\delta t$	$k_1, k_2, k_3, k_4, V, Km$	$\sigma_k^2$
5	5	0	100	0.05	0.07, 0.6, 0.05, 0.3, 0.017, 3	0.01

can be represented by the following 5-dimensional ODEs:

$$\begin{aligned}
 \dot{S} &= -k_1 \times S - k_2 \times S \times R + k_3 \times RS \\
 \dot{dS} &= k_1 \times S \\
 \dot{R} &= -k_2 \times S \times R + k_3 \times RS + V \times \frac{Rpp}{K_m + Rpp} \\
 \dot{RS} &= k_2 \times S \times R - k_3 \times RS - k_4 \times RS \\
 \dot{Rpp} &= k_4 \times RS - V \times \frac{Rpp}{K_m + Rpp}
 \end{aligned} \tag{6.3}$$

where the input signal is the concentration level of the protein  $S$ , which can either bind to the protein  $R$  to form the complex  $RS$ , or activate it into its phosphorylated form  $Rpp$ , or degrade into  $dS$ . The protein  $Rpp$  can be deactivated. The conversion between  $Rpp$  and  $R$  is governed by the *Michaelis-Menten kinetic law* with parameters  $V$  and  $K_m$ , while the rest of the interactions are defined by the *Mass Action kinetic law* with their respective parameters  $k_1, k_2, k_3$  and  $k_4$ .

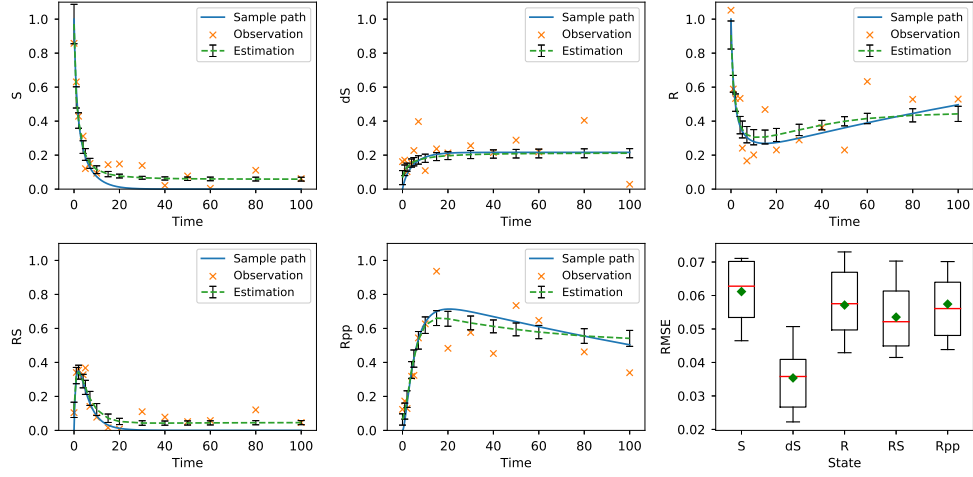
Different from other dynamical models in this chapter, the state  $Rpp$  and the parameter  $K_m$  both violate the structural assumption on the ODEs. It is also a difficult benchmark system due to the large number of parameters and its sensitivity to noise. In this section, we use the LPMF method with positivity constraint on the parameters to run the experiment.

Table 6.2 shows the experimental setup. Since the states flattens out towards the end, the observation time points are set to 0, 1, 2, 4, 5, 7, 10, 15, 20, 30, 40, 50, 60, 80 and 100. Because the observations in this experiment are required to be positive, when the noise is too big to generate negative observations, it is resampled until the positivity requirement is satisfied.

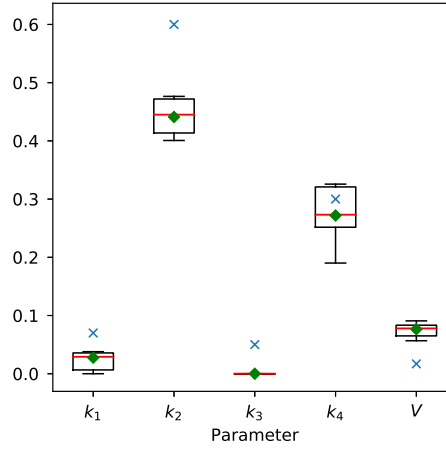
#### Direct inference on $Rpp$

In this experiment, the parameter  $K_m$  is not inferred but all the states are directly inferred. This setup includes the  $Rpp$  state that violates the structural assumption, which is not possible using the VGMGP method. The results after 10 independent repetitions are shown in Figure 6.5 and Figure 6.6. For

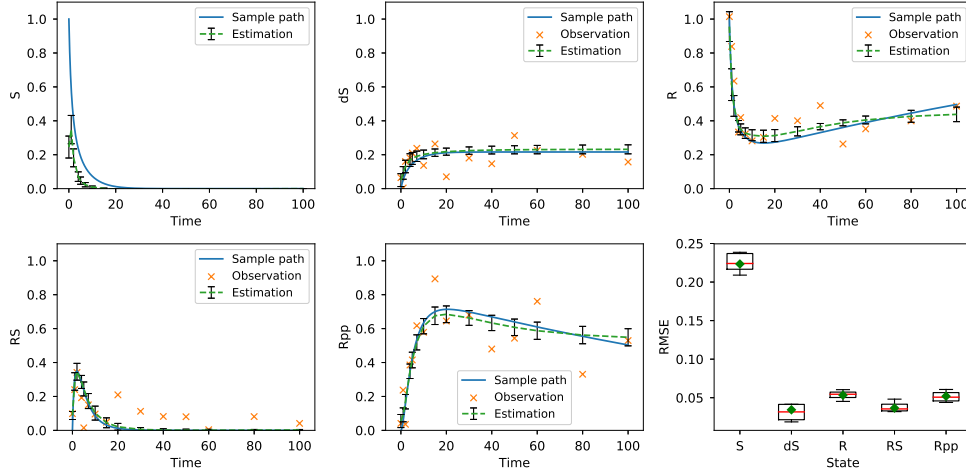
## 6. EXPERIMENTS



**Figure 6.5:** State estimation results for the protein signaling transduction pathway model with the parameter  $K_m$  set to constant. The ground truth is shown as the blue line. The observations, when available, are shown as the orange crosses. The estimation is shown as the dotted green line. The error bars in the first five plots indicate one standard deviation. For the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles.



**Figure 6.6:** Parameter estimation results for the protein signaling transduction pathway model with the parameter  $K_m$  set to constant. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles. The true parameter values are indicated by the blue crosses.



**Figure 6.7:** State estimation results for the protein signaling transduction pathway model with the parameter  $K_m$  set to constant and the state  $S$  unobservable. The ground true is shown as the blue line. The observations, when available, are shown as the orange crosses. The estimation is shown as the dotted green line. The error bars in the first five plots indicate one standard deviation. For the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles.

state estimation, the overall trend of the states is captured by the learning algorithm. Since there are much fewer observations toward the end of the time point, the prediction becomes worse. The parameter estimation result is not ideal in comparison to the states.

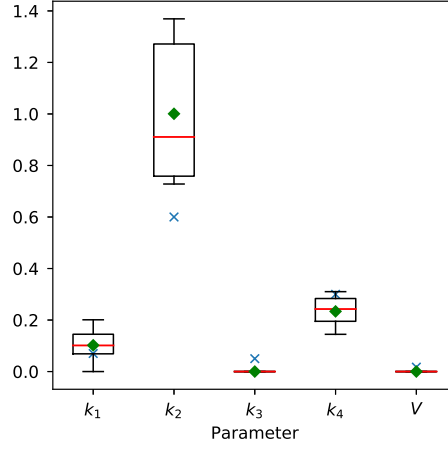
### Partial observation

To increase the difficulty of the problem, the observations on the state  $S$  is masked out. The results after 10 independent repetitions is shown in Figure 6.7 and Figure 6.8. Given even less observations, the accuracy of both state and parameter estimation become worse. From Figure 6.7 it can be seen that gradient model indeed does provide some guidance on the estimation of the state  $S$ .

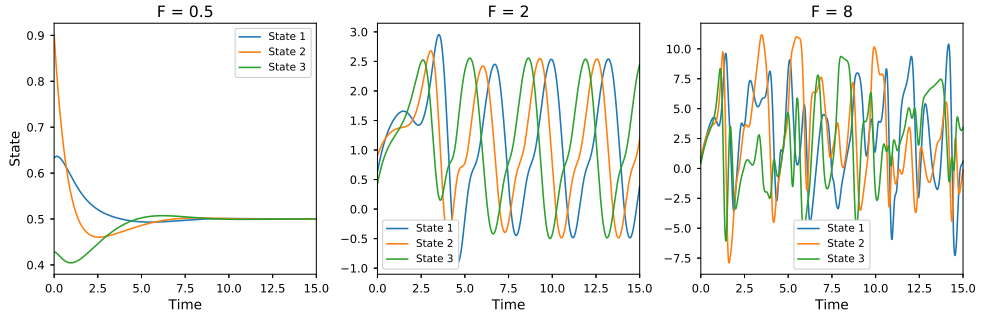
## 6.4 Lorenz 96 model

As a minimalistic weather forecast model, the *Lorenz 96* model (Lorenz, 1996) is another widely used benchmark system due to its chaotic behavior under certain configurations and its flexibility to scale to large numbers of states. A  $K$ -dimensional deterministic Lorenz 96 dynamical system is defined for

## 6. EXPERIMENTS



**Figure 6.8:** Parameter estimation results for the protein signaling transduction pathway model with the parameter  $K_m$  set to constant and the state  $S$  unobservable. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles. The true parameter values are indicated by the blue crosses.



**Figure 6.9:** Trajectories of the 1st, 2nd and 3rd dimensions of a 10-dimensional deterministic Lorenz 96 model with different parameter values. The numerical integration is carried out from time 0 to 15 with a step size of 0.01.

$k = 1, \dots, K$ , state-wise as follows:

$$\dot{x}_k(t) = (x_{k+1}(t) - x_{k-2}(t))x_{k-1}(t) - x_k(t) + F \quad (6.4)$$

where  $x_{-1}(t) = x_{K-1}(t)$ ,  $x_0(t) = x_K(t)$ ,  $x_{K+1}(t) = x_1(t)$ , and  $F \in \mathbb{R}$  is the parameter controlling the behavior of the system. When  $F < 0.895$ , the states decay into a steady value equal to  $F$ ; when  $F$  is between 0.895 and 4.0, the states are periodic; when  $F \geq 4.0$ , the system exhibits chaotic behavior (Vrettas et al., 2015). An example of the trajectories of a 10-dimensional deterministic Lorenz 96 model is shown in Figure 6.9.

The experiments in this section are mainly concerned about the stochastic

**Table 6.3:** Experimental setup to generate sample paths and to collect observations for the stochastic Lorenz 96 model. The system dimension is denoted by  $K$  and the number of observable dimensions is  $K_{obs}$ . Based on the drift parameter  $F$  and the diffusion noise, the sample paths are generated from time  $t_0$  to  $t_T$  with a step size of  $\delta t$ . The observation noise variance  $\sigma_k^2$  and the diffusion noise variance  $\rho_k^2$  are assumed to be the identical for each state. For each time unit,  $freq_{obs}$  denotes the number observations to be collected, which are equally distributed over the time line.

$K$	$K_{obs}$	$t_0$	$t_T$	$\delta t$	$F$	$\sigma_k^2$	$\rho_k^2$	$freq_{obs}$
500	325 (65%)	0	4	0.01	8	1	4	8

Lorenz 96 model, which can be obtained by simply including a noise process as described in Section 2.2. Experiments about its deterministic counterpart can be found in Gorbach et al. (2017). In the following, we compare the results from the LPMF-SDE and the VGPA-MF methods.

### Experimental setup

Following Vrettas et al. (2015), Table 6.3 shows the experimental setup to generate sample paths and to collect observations. The sample paths are generated by the VGPA-MF MATLAB code, which uses the first order Euler-Maruyama scheme. To avoid exploiting the special structure of the drift function (Eq. 6.4), the observed states are selected randomly. After observations are collected, the data files are converted into the format compatible with the LPMF-SDE source code.

Note that the system dimension is scaled down from 1000 to 500 due to time constraints. Accordingly, the number of observed states is scaled down proportionally. Nevertheless, it is still very interesting to quantify the scalability of the LPMF-SDE solution, which will be examined in detail in a separate experiment later.

For each SDE sample path, the LPMF-SDE method averages individual inference results from 100 independently generated RODE sample paths to obtain an ensemble solution. The RODE sample paths are estimated in parallel on the ETH Euler cluster<sup>10</sup>. Only 1 CPU core and 2 GB of memory are allocated to each inference. For this experiment, the RODE states are estimated using the gradient-based method while the parameters are calculated in closed-form.

On the other hand, the VGPA-MF code is deployed on a desktop with an Intel i5 quad-core CPU and 16 GB of memory. We also note that thread pooling is used to utilize all the CPU resources for the VGPA-MF implementation, which is not the case for LPMF-SDE.

<sup>10</sup><https://scicomp.ethz.ch/wiki/Euler>

### State estimation

Figure 6.10 shows an example of the state inference result for the stochastic Lorenz 96 model using both methods. For observed states, the mean values of the predictions from both methods are very close to the ground truth, while the VGPA-MF method tends to have lower variance. For unobserved states, the accuracy of both methods depends on the amount of available information that is directly related to them. As indicated by the drift function (Eq. 6.4), the states are only coupled with their neighbors for the Lorenz 96 model. Therefore, estimation on an unobserved state can be reasonably good when its neighbors are observed, e.g. state 120 has observed neighbors 121, 122, 123, etc. Estimation can fail when the neighbors of that state are also unobserved so that they are unable to provide sufficient information to the algorithm, e.g. the true trajectory of state 213 cannot be recovered since states 211 to 215 are all unobserved. As expected, the uncertainty about the estimation for unobserved states is higher than that for the observed states, which is indicated by the wider error bars.

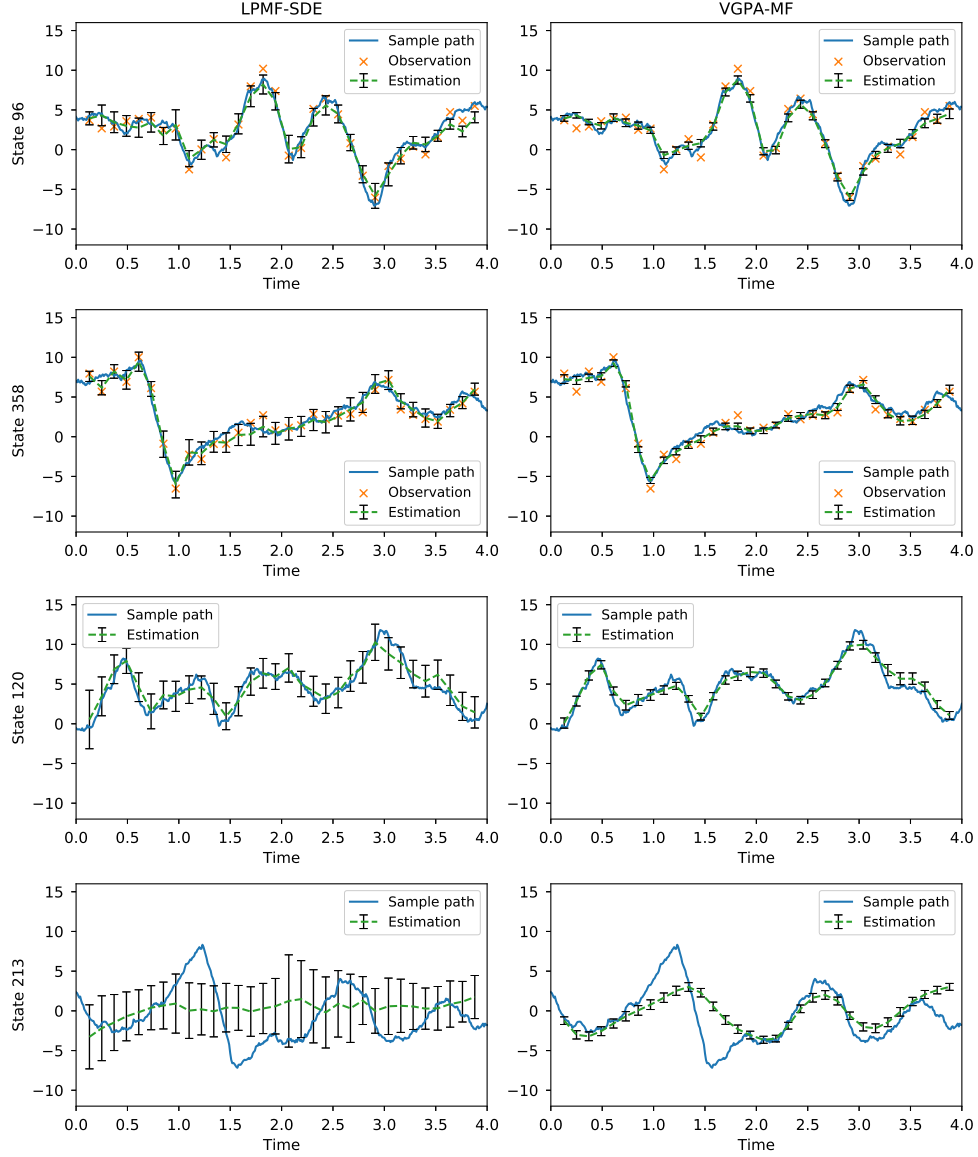
In order to better quantify the difference between both methods, 10 independent Lorenz 96 SDE sample paths are generated. For each SDE sample path, a set of observations is collected, and then the above experiment is repeated. The root mean square error (RMSE) is used as the accuracy measure, and the RMSEs for observed and unobserved states are considered separately. The formulas for the RMSEs are

$$RMSE_{obs} = \frac{1}{K_{obs}} \sum_{k \in \mathcal{S}_{obs}} \sqrt{\frac{1}{N} \sum_{n=1}^N (\hat{x}_k(t_n) - x_k(t_n))^2} \quad (6.5)$$

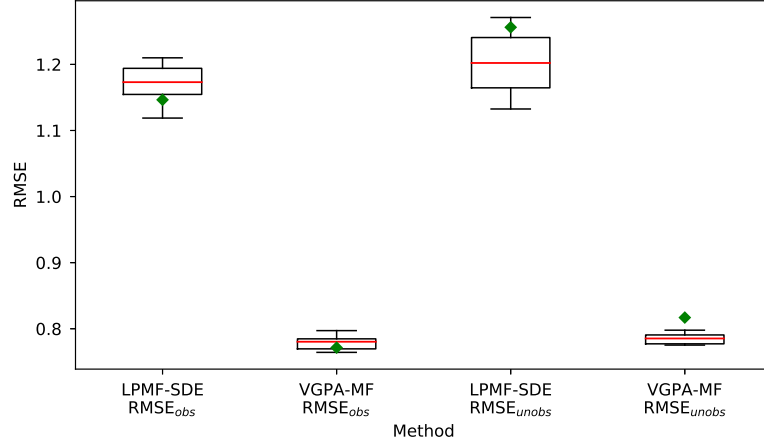
$$RMSE_{unobs} = \frac{1}{K_{unobs}} \sum_{k \in \mathcal{S}_{unobs}} \sqrt{\frac{1}{N} \sum_{n=1}^N (\hat{x}_k(t_n) - x_k(t_n))^2} \quad (6.6)$$

where  $K_{obs}$  and  $K_{unobs}$  indicate the number of observed and unobserved states respectively with  $\mathcal{S}_{obs}$  and  $\mathcal{S}_{unobs}$  as the corresponding sets containing the state indices,  $N$  is the total number of observations for each dimension, and  $\hat{x}_k(t_n)$  and  $x_k(t_n)$  are the predicted and true values for the  $k$ -th state at time point  $t$  respectively.

The results of the 10 independent repetitions are summarized in Figure 6.11. At first glance, it is surprising that the VGPA-MF method has such low RMSE values. This is partially because in the MATLAB code we have obtained, only state inference is carried out and the true drift parameter  $F$  is given to the inference procedure. However, the LPMF-SDE method is estimating both the states and parameters simultaneously. Note that the VGPA-MF solution is fully capable of jointly estimating the states and parameters. But as at the time of this writing, we have not obtained the code



**Figure 6.10:** State estimation results for 4 selected states from the 500-dimensional stochastic Lorenz 96 model. Among the 500 states, 325 are observed. The left column shows the result using the LPMF-SDE method, while the right column contains the result based on the VGPA-MF algorithm. In the four rows, state 96 and its neighboring states are directly observed; state 358 and only one side of its neighbors are directly observed; state 120 is not observed but its neighbors are; state 213 and its neighbors are all unobserved. The ground truth is shown as a solid blue line. The mean of estimation is drawn as a dotted green line together with the error bars indicating one standard deviation. Observations, when available, are indicated by the orange crosses.



**Figure 6.11:** Box plot for the RMSEs of state estimation over 10 independent repetitions using sample paths from the 500-dimensional stochastic Lorenz 96 model with 325 observable states. The RMSEs of observed and unobserved states are calculated separately. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles.

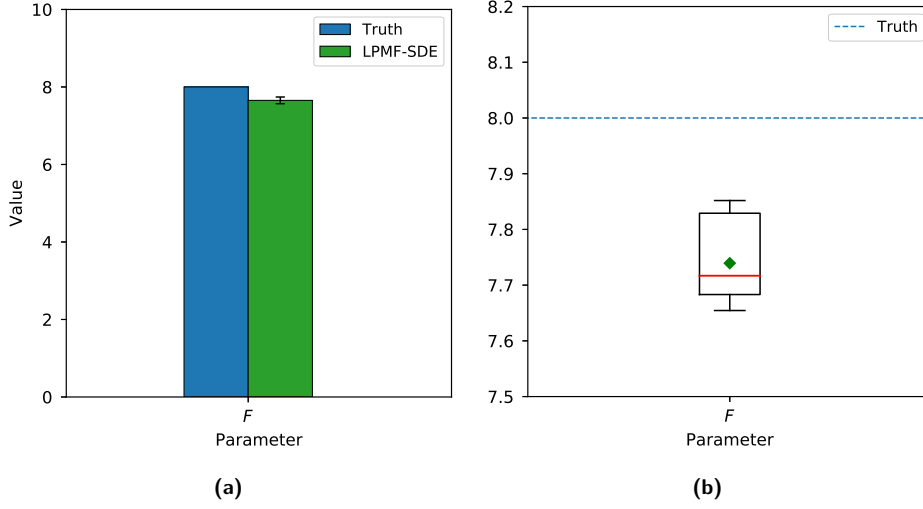
that computes the gradient of their variational free energy function w.r.t. the parameters in order to complete the required experiments. Moreover, as claimed by Vrettas et al. (2015), their solution can also estimate the diffusion noise, which is an advantage over the LPMF-SDE method and can be considered as a possible future extension to this work. Considering the sparsity of the observations and that the observation noise variances are set to 1, we see some potential for the LPMF-SDE solution.

One shortcoming of the LPMF-SDE method in this experiment is that it only considers a single optimal value of the cost function found using gradient-based method when making predictions. Depending on the initialization conditions and the complexity of the cost function, which is subject to the complexity of drift function, sometimes the optimization procedure be trapped in a local optimum. As a consequence, even though the mean can be recovered reasonably well, the variance of the prediction is less under control. Since currently only out-of-package second-order optimization techniques such as the *truncated Newton method*, the *dog-leg trust-region algorithm*, etc. (Nocedal and Wright, 2006), are used, further improvement on the optimization strategy can also be considered as an extension to the current work.

### Parameter estimation

This section only discusses parameter inference results using the LPMF-SDE method due to the lack of relevant files from the VGPA-MF MATLAB





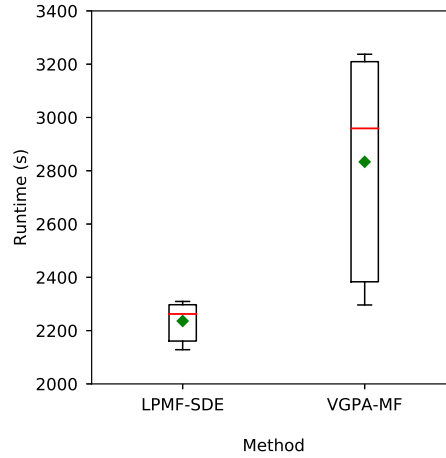
**Figure 6.12:** Parameter estimation result using the LPMF-SDE algorithm for the 500-dimensional stochastic Lorenz 96 model with 325 observable states. (a) Estimation result for the SDE sample path shown in Figure 6.10. The blue bar indicates the true parameter value. The green bar shows the mean of estimation by averaging individual results from 100 RODE sample paths. The error bar indicates one standard deviation. (b) Box plots for parameter estimation over 10 independent repetitions. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles. The dotted blue line indicates the true parameter value.

code as mentioned previously. The result for a single SDE sample path is illustrated in Figure 6.12a, while the result from the 10 dependent repetitions is summarized in Figure 6.12b.

Figure 6.12a shows that the estimation of the drift parameter  $F$  is very close to the true value 8, and the variance is much lower in comparison to state estimation. This is due to the fact that the Lorenz 96 model has only one drift parameter, which appears inside the drift function for each state. Hence, the gradient matching model is able to obtain much more information about the drift parameter, despite the existence of unobserved states. Figure 6.12b further shows that the mean values of parameter estimation over the 10 independent repetitions are narrowly distributed. But in general, the LPMF-SDE method tends to underestimate the parameter value.

### Runtime performance

For the SDE sample path shown in Figure 6.10, the average runtime for the LPMF-SDE algorithm to solve one RODE sample path is 2309 seconds with a standard deviation of around 128 seconds, while the runtime for the VGPA-MF method takes 3232 seconds. As the current VGPA-MF MATLAB



**Figure 6.13:** Box plot for the runtimes over 10 repetitions using independent sample paths from the stochastic Lorenz 96 model. The median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles.

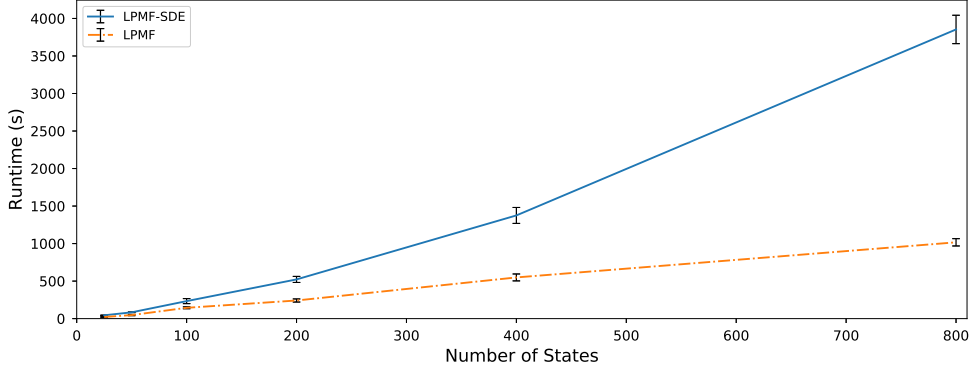
code only runs the state smoothing algorithm, it is expected to be much more computationally intensive when parameter estimation is carried out simultaneously, since the forward-backward loop has to be entered many times until convergence.

Figure 6.13 summarizes the average runtimes for LPMF to solve one RODE sample path and the runtimes for VGPA-MF to complete one iteration of state smoothing over the 10 independent repetitions. It shows that the LPMF-SDE solution performs consistently faster and more stably than its counterpart, as indicated by lower median runtime value and the narrower distribution.

To conclude, the LPMF-SDE algorithm requires much less resources and also runs faster for inference on one RODE sample path. However, the ensemble nature of the LPMF-SDE strategy relies on averaging results from a large number RODE estimations, which could increase the computational requirements significantly. Nevertheless, given the availability of computer farms and the ease of setting up distributed inference pipeline nowadays, the parallelism requirement for the LPMF-SDE solution can be easily satisfied.

### Scalability

As mentioned by Gorbach et al. (2017), inference on the deterministic Lorenz 96 model with 1000 states using their VGMGP solution completes within 400 seconds on average. However, the LPMF-SDE method, which extends



**Figure 6.14:** Scalability of the LPMF-SDE algorithm versus its deterministic counterpart, i.e. LPMF, as the dimensionality of the system increases. For the LPMF-SDE method, the blue line indicates the average interpolated runtime to inference one RODE sample path over 10 independent RODE sample paths by connecting the measurements at 25, 50, 100, 200, 400 and 800. For the LPMF method, the orange line is obtained by connecting the averages of 10 independent runs on the same ODE trajectory with different observations at the same measurement points. The error bars in both setups indicate one standard deviation.

their methodology, takes on average more than 2200 seconds for one RODE sample path with only 500 states, as shown in Figure 6.13. It is therefore interesting to re-examine the scalability of the LPMF-SDE solution.

First, we noticed during implementation that the performance of the LPMF method for ODEs is comparable to the VGMGP method. The major difference between LPMF-SDE and LPMF is the introduction of a stochastic Ornstein-Uhlenbeck process into the vector field of the ODEs. To reveal the influence of the stochastic process on the inference procedure, the following experiment compares side-by-side the performance of the LPMF-SDE method running on the stochastic Lorenz 96 model with the performance of the LPMF method running on its deterministic counterpart as the number of states increases.

Specifically, measurements are taken with a system dimensionality of 25, 50, 100, 200, 400, and 800 with the number of observed states kept at 65 percent of the total number of states. Using the same setup in Table 6.3, except that the diffusion noise is not used for the Lorenz 96 ODEs, the LPMF-SDE algorithm is tested with 10 independent RODE sample paths while the LPMF method is tested with 10 independent observation sets from the same ODE trajectory for each system dimensionality. The experiments are all conducted on the Euler cluster with the same hardware configuration.

As can be seen from Figure 6.14, benefiting from the efficient implementation of the gradient and Hessian evaluation subroutines, the gradient-based

LPMF method incurs only a slightly higher performance penalty than the closed-form solution from VGMGP. Also, the runtime seems to increase linearly as the number of states increases. On the other hand, the performance of the LPMF-SDE method degrades much faster than LPMF, and the gap between them become larger as the number of states increases. The cause of this phenomenon requires further investigation but a plausible explanation is the complication of the optimization objective after the introduction of the stochastic process into the vector field.

### 6.5 Lorenz 63 model

The last dynamical system examined in this chapter is the stochastic version of the *Lorenz 63* model (Lorenz, 1963), which is a low dimensional mathematical model for thermal convection in the atmosphere. The vector field of the deterministic Lorenz 63 is defined as follows:

$$f(x(t), \theta) = \begin{bmatrix} \dot{x}(t) \\ \dot{y}(t) \\ \dot{z}(t) \end{bmatrix} = \begin{bmatrix} \sigma(y(t) - x(t)) \\ x(t)(\rho - z(t)) - y(t) \\ x(t)y(t) - \beta z(t) \end{bmatrix} \quad (6.7)$$

where the state vector is given by  $x(t) = [x(t), y(t), z(t)]^\top \in \mathbb{R}^3$ , and  $\theta = [\sigma, \rho, \beta]^\top \in \mathbb{R}^3$  is the parameter vector controlling the system behavior. Although consisting of only 3 states, this model is highly nonlinear and exhibits chaotic behavior under certain parameter configurations.

Correspondingly, the stochastic Lorenz 63 model with a state-specific, additive noise process is given by

$$dx(t) = f(x(t), \theta)dt + \Sigma^{\frac{1}{2}}dW_t \quad (6.8)$$

where  $f(x(t), \theta)$  is defined in (Eq. 6.7),  $\Sigma \in \mathbb{R}^{3 \times 3}$  is a diagonal matrix containing the diffusion noise variance, and  $W_t$  is a 3-dimensional standard Wiener process. Using the parameter set  $[10, 28, \frac{8}{3}]^\top$ , which is well-known for its resulting chaotic behavior, a sample path is shown in Figure 6.15.

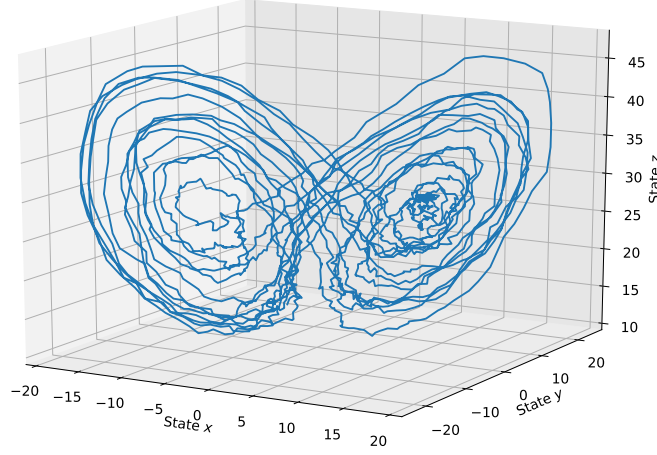
#### Comparison algorithm

In order to make comparison, a minimalistic MAP estimation of the drift parameters based on (Vrettas et al., 2011, Table 3) is self-implemented<sup>11</sup> by extending the Python source code<sup>12</sup> for the VGPA algorithm. In the following, the abbreviation VGPA-MAP is used to refer to this extension.

The extension consists of an inner loop and an outer loop. The inner loop enhances state estimation by running the VGPA smoothing algorithm to

<sup>11</sup><https://github.com/ruiixu23/VGPA>

<sup>12</sup><https://github.com/vrettasm/VGPA>



**Figure 6.15:** Sample path for the stochastic Lorenz 63 model generated based on the parameter values  $[\sigma, \rho, \beta]^\top = [10, 28, \frac{8}{3}]^\top$  and the diffusion noise variances  $\rho_1^2 = \rho_2^2 = \rho_3^2 = 10$ . The integration is performed from time 0 to 20 with a step size of 0.01.

compute the optimal approximate state posterior, based on the current estimation of the parameters. The outer loop then takes a gradient step to update the parameters. This procedure is repeated until either the gradient over the parameters vanishes or the state estimation from the inner loop does not improve further significantly. Details of the algorithm can be found in (Vrettas et al., 2011, Section 5.2).

Note that Vrettas et al. (2011) claim that the aforementioned algorithm can be similarly applied to the estimation of the diffusion noise covariance matrix  $\Sigma$ , which is not implemented here as the current LPMF-SDE method only supports inference on the drift parameters. Also, the extension adopts a simple gradient update strategy with a fixed learning rate, which may lead to suboptimal result but can nonetheless be used as a baseline.

### Experimental setup

Adopting the experimental setup from Vrettas et al. (2015), the configuration used to generate sample paths and to collect observations is shown in Table 6.4. Specifically, the sample paths are generated by the VGPA source code using the first order Euler-Maruyama method with a small step size to achieve higher accuracy. After observations are collected, the data files are transformed into the format compatible with the LPMF-SDE method.

In order to provide a fair comparison, both methods are deployed to the Eu-

**Table 6.4:** Experimental setup to generate sample paths and to collect observations for the stochastic Lorenz 63 model. The system dimension is denoted by  $K$  and the number of observable states is denoted by  $K_{obs}$ . Given the drift parameters  $\sigma, \rho, \beta$  and the diffusion noise, the sample paths are generated from time  $t_0$  to  $t_T$  with a step size of  $\delta t$ . The variances of the observation noise  $\sigma_k^2$  and the diffusion noise  $\rho_k^2$  are assumed to be identical for each state. For each time unit,  $freq_{obs}$  denotes the number of observations to be collected, which are equally distributed over the time line.

$K$	$K_{obs}$	$t_0$	$t_T$	$\delta t$	$\sigma, \rho, \beta$	$\sigma_k^2$	$\rho_k^2$	$freq_{obs}$
3	3 or 2	0	20	0.01	10, 28, $\frac{8}{3}$	2	10	5

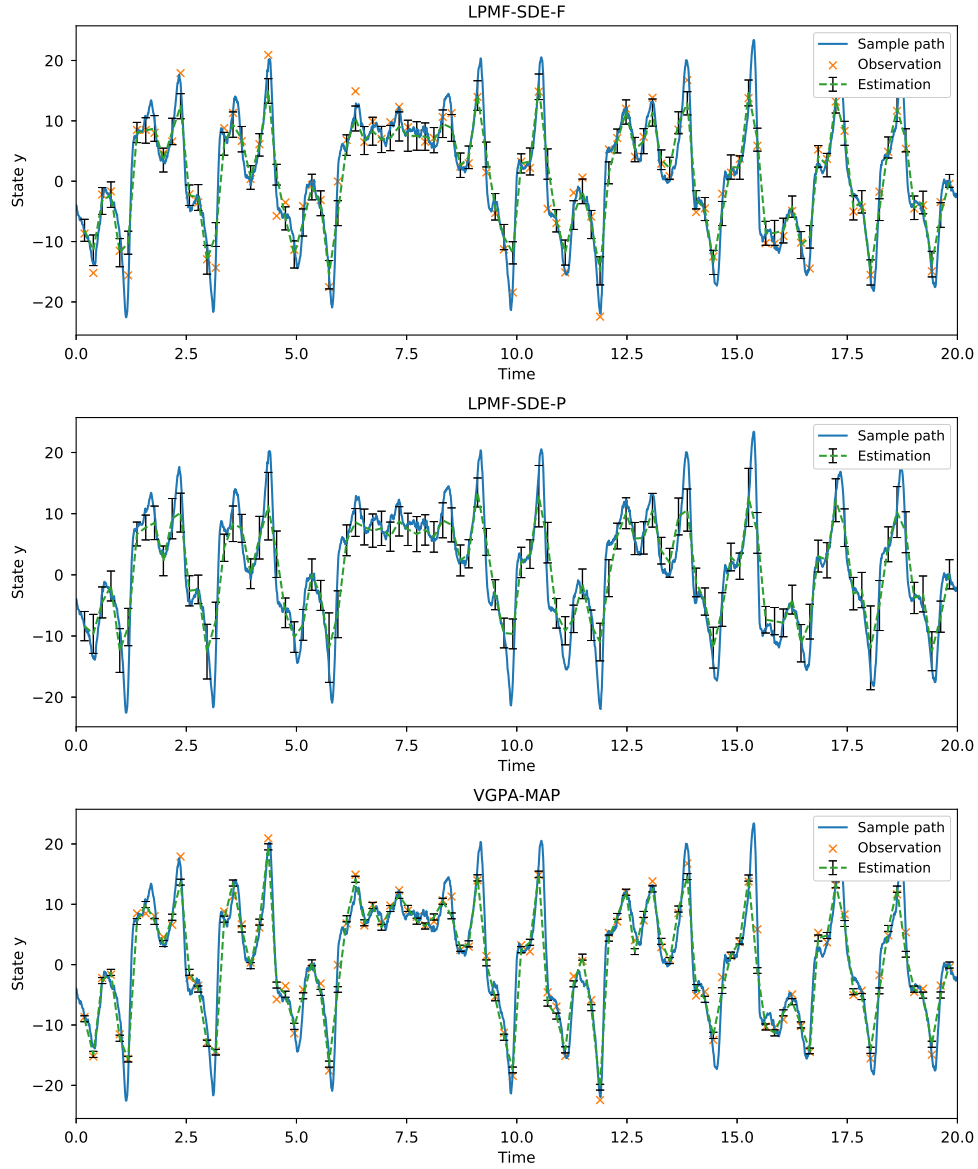
ler cluster using the same hardware configuration described in Section 6.4. For LPMF-SDE, 100 independent RODE sample paths are solved to obtain an ensemble result for one SDE sample path. The states are optimized using a gradient-based method while the parameters are calculated analytically with mirroring of negative parameters<sup>13</sup>. For VGPA-MAP, the learning rate is fixed to 0.001, the maximum number of iterations is 250, and the stopping criteria are the decrease in the total variational free energy and the  $L2$  norm of the change to the drift parameters with a threshold of  $10^{-6}$ . We noticed that all the experiments from VGPA-MAP reached the maximum number of 250 iterations. After manual examination, the result at the 80th iteration is used for the following comparisons.

It would be interesting to compare these two methods with one unobserved state. Unfortunately, the smoothing result from the VGPA algorithm with unobserved state is not ideal from preliminary examination. Therefore, a comparison with full state observability is provided. But in order to demonstrate that the LPMF-SDE solution is capable of handling partial observations, an additional experiment is conducted by masking out the observations for state  $y$ . To distinguish the two scenarios, we use LPMF-SDE-F and LPMF-SDE-P to refer to the fully and partially observable cases respectively.

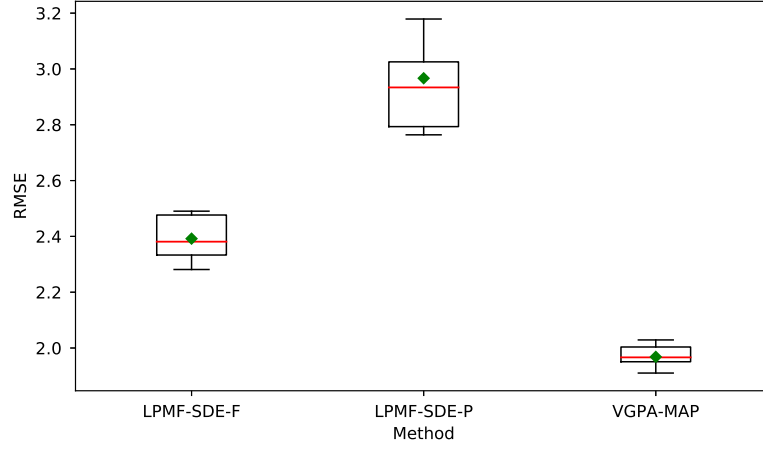
### State estimation

The estimation for state  $y$  from one SDE sample path is shown in Figure 6.16. Similar to the result for the stochastic Lorenz 96 model in Section 6.4, the general trend of the state is successfully captured by both methods with minor differences on the details. The VGPA-MAP method seem to be better at points where the state changes dramatically to form high and low peaks. The variance of the estimation is also lower than the LPMF-SDE approach. One the other hand, it is worth noting that the LPMF-SDE method even functions quite well when state  $y$  is unobserved but with naturally higher uncertainty, especially around the peaks. Since state  $y$  is unobserved, the

<sup>13</sup>The mirroring of negative parameters is a heuristic solution. A better approach would be the enforcement of positivity constraint.



**Figure 6.16:** Estimation for state  $y$  of the stochastic Lorenz 63 model. The top and middle plots are the results for the fully and partially observable cases using the LPMF-SDE method respectively. The bottom plot contains the result for the fully observable case using the VGPA-MAP extension. The ground truth is shown as a solid blue line. The mean of the estimation is drawn as a dotted green line together with the error bars indicating one standard deviation. Observations, when available, are indicated by the orange crosses.



**Figure 6.17:** Box plot for the RMSEs of state estimation over 10 independent SDE sample paths from the stochastic Lorenz 63 model. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles.

only source of information to estimate it is from the drift function (Eq. 6.7), which demonstrates the potential of the gradient matching framework.

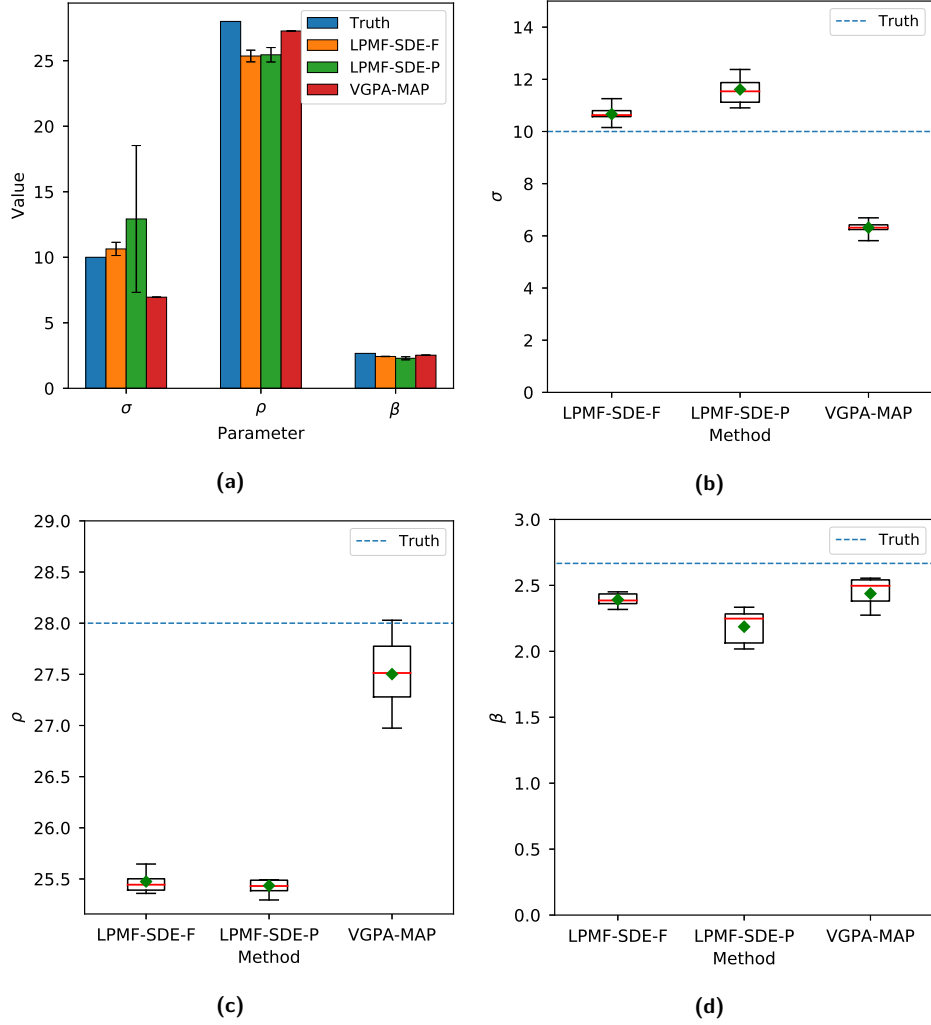
After repeating the above experiment 10 times using each time an independently generated SDE sample path and observation set, the RMSEs for state estimation is computed and summarized in Figure 6.17. Since there are only three states, the RMSEs over all states are considered together. The figure shows that the RMSE of state estimation using VGPA-MAP is slightly below the observation noise variance, which outperforms LPMF-SDE with full state observability by around 0.4 on average. On the other hand, the RMSE using LPMF-SDE when state  $y$  is unobserved is still reasonably low, considering that the values of the state range from -20 to 20, and the inference is carried out over a much longer time period than all previous experiments.

### Parameter estimation

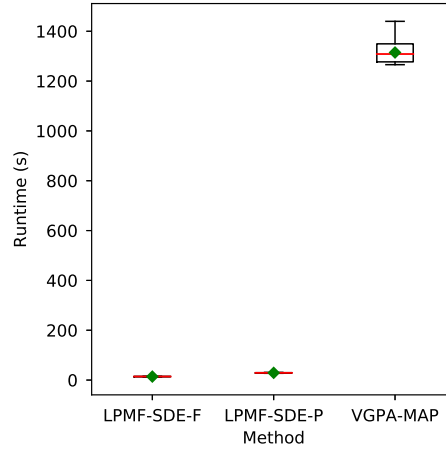
For parameter estimation, the results are presented in Figure 6.18. Overall, the estimation based on both methods seem to be on par with each other with relatively low variance when all the states are observed. Given that the VGPA-MAP algorithm achieves better accuracy when estimating the states, it would be expected that the parameter estimation to also be better. The non-optimal performance is probably due to the simple gradient update strategy as mentioned before.

If we look at Figure 6.18a in detail, the large variance around the estimation for parameter  $\sigma$  when state  $y$  is unobserved is noticeable. To explain this,





**Figure 6.18:** Parameter estimation result for the stochastic Lorenz 63 model. (a) Estimation result for the SDE sample path shown in Figure 6.16 with the error bar indicating one standard deviation. (b), (c), and (d) are the box plots for parameters  $\sigma$ ,  $\rho$  and  $\beta$  after 10 independent repetitions respectively. In the box plot, the median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles. The dotted blue line indicates the true parameter value.



**Figure 6.19:** Box plot for the runtimes over 10 repetitions using independent sample paths from the stochastic Lorenz 63 model. The median is indicated by the red line, while the mean is shown as the green diamond. The box shows the lower and upper quartiles, while the whiskers are the 5th and 95th percentiles.

first note that  $\sigma$  appears only inside the first equation of (Eq. 6.7) together with state  $x$  and  $y$ . Since  $y$  is not observable, the only source of information to estimate  $\sigma$  is  $x$ . This is in contrast to the other two parameters  $\rho$  and  $\beta$ , where both states  $x$  and  $z$  are used to estimate them. Lastly, the means of the predicted parameters over the 10 independent runs are generally concentrated except for parameter  $\rho$  inferred by the VGPA-MAP algorithm, as shown in Figure 6.18b to Figure 6.18c.

### Runtime performance

Figure 6.19 shows the distribution of the average time to solve one RODE sample path and the runtime of the VGPA-MAP algorithm to infer one SDE sample path after 10 independent repetitions. For each RODE sample, LPMF-SDE takes a fraction of the time required by VGPA-MAP for one SDE sample path. Even if the RODEs are not solved in parallel, the total runtime required by LPMF-SDE is still comparable to that of VGPA-MAP. This further demonstrates the runtime efficiency of the LPMF-SDE scheme.

---

## Conclusion

---

This work examines the problem of state and parameter estimation in deterministic and random dynamical systems given noisy, sparse or even incomplete observations. A mean-field Laplace approximation solution is proposed to address the intractability of the posterior distribution and has been shown empirically to be robust and scalable. Further, it relaxes the structural assumption imposed on the dynamical systems from previous work and introduces positivity constraints on the states and parameters. Based on the correspondence between RODEs and SDEs, a highly efficient parallel inference technique is devised to address problems involving diffusion processes.



---

## Bibliography

---

- Archambeau, C., Cornford, D., Opper, M., and Shawe-Taylor, J. Gaussian process approximations of stochastic differential equations. In *Gaussian Processes in Practice*, pages 1–16, 2007.
- Archambeau, C., Opper, M., Shen, Y., Cornford, D., and Shawe-taylor, J. S. Variational inference for diffusion processes. In *Advances in Neural Information Processing Systems*, pages 17–24, 2008.
- Babtie, A. C., Kirk, P., and Stumpf, M. P. H. Topological sensitivity analysis for systems biology. *Proceedings of the National Academy of Sciences*, 111(52): 18507–18512, 2014.
- Bishop, C. M. *Pattern recognition and machine learning*. Springer, 2006.
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017. doi: 10.1080/01621459.2017.1285773.
- Butcher, J. C. *Numerical methods for ordinary differential equations*. John Wiley & Sons, 2016.
- Calderhead, B., Girolami, M., and Lawrence, N. D. Accelerating bayesian inference over nonlinear differential equations with gaussian processes. In *Advances in neural information processing systems*, pages 217–224, 2009.
- Dondelinger, F., Husmeier, D., Rogers, S., and Filippone, M. Ode parameter inference using adaptive gradient matching with gaussian processes. In *Artificial Intelligence and Statistics*, pages 216–228, 2013.
- Ellner, S. P. and Guckenheimer, J. *Dynamic models in biology*. Princeton University Press, 2011.
- Gardiner, C. W. *Stochastic methods: a handbook for the natural and social sciences*. Springer, 2009.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep learning*. MIT press, 2016.

- Gorbach, N. S., Bauer, S., and Buhmann, J. M. Mean-field variational inference for gradient matching with gaussian processes. *arXiv preprint arXiv:1610.06949*, 2016.
- Gorbach, N. S., Bauer, S., and Buhmann, J. M. Scalable variational inference for dynamical systems. *arXiv preprint arXiv:1705.07079*, 2017.
- Grüne, L. and Kloeden, P. Pathwise approximation of random ordinary differential equations. *BIT Numerical Mathematics*, 41(4):711–721, 2001.
- Higham, D. J. An algorithmic introduction to numerical simulation of stochastic differential equations. *SIAM review*, 43(3):525–546, 2001.
- Hinton, G. E. Products of experts. In *ICANN’99*, volume 1, pages 1–6, 1999.
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- Imkeller, P. and Schmalfuss, B. The conjugacy of stochastic and random differential equations and the existence of global attractors. *Journal of Dynamics and Differential Equations*, 13(2):215–249, 2001.
- Jensen, J. L. W. V. Sur les fonctions convexes et les inégalités entre les valeurs moyennes. *Acta mathematica*, 30(1):175–193, 1906.
- Jentzen, A. and Kloeden, P. E. *Taylor approximations for stochastic partial differential equations*. SIAM, 2011.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- Kalnay, E. *Atmospheric modeling, data assimilation and predictability*. Cambridge university press, 2003.
- Kloeden, P. E. and Jentzen, A. Pathwise convergent higher order numerical schemes for random ordinary differential equations. In *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, volume 463, pages 2929–2944. The Royal Society, 2007.
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. M. Automatic differentiation variational inference. *Journal of Machine Learning Research*, 18(14):1–45, 2017.
- Kullback, S. and Leibler, R. A. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- Lorenz, E. The butterfly effect. *World Scientific Series on Nonlinear Science Series A*, 39:91–94, 2000.
- Lorenz, E. N. Deterministic nonperiodic flow. *Journal of the atmospheric sciences*, 20(2):130–141, 1963.

- Lorenz, E. N. Predictability: A problem partly solved. In *Proc. Seminar on predictability*, volume 1, 1996.
- Lotka, A. J. The growth of mixed populations: two species competing for a common food supply. *Journal of the Washington Academy of Sciences*, 22 (16/17):461–469, 1932.
- Macdonald, B. and Husmeier, D. Gradient matching methods for computational inference in mechanistic models for systems biology: a review and comparative analysis. *Frontiers in bioengineering and biotechnology*, 3, 2015.
- MacKay, D. J. *Information theory, inference and learning algorithms*. Cambridge university press, 2003.
- Murphy, K. P. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- Nocedal, J. and Wright, S. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer New York, 2006. ISBN 9780387227429.
- Øksendal, B. *Stochastic Differential Equations: An Introduction with Applications*. Springer, 2013.
- Petersen, K. B. and Pedersen, M. S. *The matrix cookbook*. 2012.
- Ramsay, J. O., Hooker, G., Campbell, D., and Cao, J. Parameter estimation for differential equations: a generalized smoothing approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(5):741–796, 2007.
- Rasmussen, C. E. and Williams, C. K. I. *Gaussian processes for machine learning*, volume 1. MIT press, 2006.
- Riesinger, C., Neckel, T., and Rupp, F. Solving random ordinary differential equations on gpu clusters using multiple levels of parallelism. *SIAM Journal on Scientific Computing*, 38(4):C372–C402, 2016.
- Sussmann, H. J. On the gap between deterministic and stochastic ordinary differential equations. *The Annals of Probability*, pages 19–41, 1978.
- Vrettas, M. D., Cornford, D., and Opper, M. Estimating parameters in stochastic systems: A variational bayesian approach. *Physica D: Nonlinear Phenomena*, 240(23):1877–1900, 2011.
- Vrettas, M. D., Opper, M., and Cornford, D. Variational mean-field algorithm for efficient inference in large systems of stochastic differential equations. *Physical Review E*, 91(1):012148, 2015.
- Vysheirsky, V. and Girolami, M. A. Bayesian ranking of biochemical system models. *Bioinformatics*, 24(6):833–839, 2007.

## BIBLIOGRAPHY

---

Wainwright, M. J., Jordan, M. I., et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.





Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

## Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

**Title of work** (in block letters):

SCALABLE VARIATIONAL INFERENCE FOR  
STOCHASTIC DIFFERENTIAL EQUATIONS

**Authored by** (in block letters):

*For papers written by groups the names of all authors are required.*

**Name(s):**

XU

**First name(s):**

RUIFENG

With my signature I confirm that

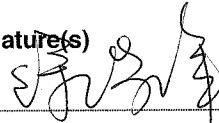
- I have committed none of the forms of plagiarism described in the 'Citation etiquette' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

**Place, date**

Zürich, 01.09.2017

**Signature(s)**



*For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.*