



Scalable Variational Inference for Stochastic Differential Equations

Master Thesis

Ruifeng Xu

Supervisor: Prof. Dr. Joachim M. Buhmann

Advisors: Stefan Bauer & Nico S. Gorbach

Department of Computer Science, ETH Zurich

Statistical inference of states and parameters of dynamical systems based on noisy, sparse or even incomplete state observations.

Outline

- Dynamical systems
- Motivation & challenges
- Laplace mean-field approximation
- Extension to random dynamical systems
- Experiments
- Conclusion

Dynamical Systems

Ordinary differential equations (ODEs)

A K -dimensional real-valued ODE system is defined as

$$\dot{\mathbf{x}}(t) = \frac{d\mathbf{x}(t)}{dt} = \mathbf{f}(\mathbf{x}(t), \boldsymbol{\theta}) \quad (1)$$

where

$\mathbf{x}(t) = [x_1(t), \dots, x_K(t)]^\top \in \mathbb{R}^K$ are the states at time t ,

$\dot{\mathbf{x}}(t) = [f_1(\mathbf{x}(t), \boldsymbol{\theta}), \dots, f_K(\mathbf{x}(t), \boldsymbol{\theta})]^\top \in \mathbb{R}^K$ are the state derivatives at time t ,

$\mathbf{f} : \mathbb{R}^K \mapsto \mathbb{R}^K$ is the vector fields with parameter $\boldsymbol{\theta} \in \mathbb{R}^M$.

Initial states and parameters determine the future states.

\mathbf{f} may have direct dependency on t , which is suppressed for uncluttered notations.

Stochastic differential equations (SDEs)

Given a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, a K -dimensional SDE system with state-specific, additive Gaussian noises is defined, in the *Itô* form, as

$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{x}(t), \boldsymbol{\theta})dt + \boldsymbol{\Sigma}^{\frac{1}{2}}d\mathbf{W}_t \quad (2)$$

where

$\mathbf{f} : \mathbb{R}^K \mapsto \mathbb{R}^K$ is the deterministic drift function with parameter $\boldsymbol{\theta} \in \mathbb{R}^M$

$\boldsymbol{\Sigma} = \text{diag}(\rho_1^2, \dots, \rho_K^2) \in \mathbb{R}^{K \times K}$ is the diagonal noise covariance matrix,

$\mathbf{W}_t \in \mathbb{R}^K$ is a standard K -dimensional Wiener process.

Each realization is most likely a different *sample path*.

A class of multiplicative noise models can be mapped to this model.

Motivation & Challenges

Motivation

Dynamical systems model various natural phenomena in chemistry, physics, biology, economics, meteorology, etc. For example

- *Protein signalling transduction pathway* models the dynamics among protein species using a set of non-linear differential equations.
- Stochastic *Lorenz 96* model is commonly used in weather forecast.
- Many others . . .

The inference algorithm should be accurate, robust and performant.

Challenges

- Conventional methods requires explicit numerical integrations each time after parameter adaptation, which is slow and not scalable.
- The likelihood surfaces are likely to be multimodal due to nonlinearity within the dynamical systems, which makes parameter search difficult.
- In Bayesian statistics, the marginalization term is intractable and requires approximate inference techniques.

Markov chain Monte Carlo (MCMC) sampling schemes are accurate but computationally expensive and requires onerous convergence analysis.

Laplace Mean-Field Approximation

Noisy observation

Usually, observations $\mathbf{y}(t) = [y_1(t), \dots, y_K(t)]^\top \in \mathbb{R}^K$ are contaminated by noises $\boldsymbol{\varepsilon}(t) = [\varepsilon_1(t), \dots, \varepsilon_K(t)] \in \mathbb{R}^K$ such that

$$\mathbf{y}(t) = \mathbf{x}(t) + \boldsymbol{\varepsilon}(t) \quad (3)$$

For a sequence of observations, we denote

$$\mathbf{Y} = [\mathbf{y}(t_1), \dots, \mathbf{y}(t_N)] \in \mathbb{R}^{K \times N}$$

$$\mathbf{X} = [\mathbf{x}(t_1), \dots, \mathbf{x}(t_N)] \in \mathbb{R}^{K \times N}$$

$$\mathbf{E} = [\boldsymbol{\varepsilon}(t_1), \dots, \boldsymbol{\varepsilon}(t_N)] \in \mathbb{R}^{K \times N}$$

Noisy observation

Assuming *independent and identically distributed (i.i.d.)* state-specific, additive Gaussian noise $\varepsilon(t) \sim \mathcal{N}(\mathbf{0}, \mathbf{D})$ with $\mathbf{D} = \text{diag}(\sigma_1^2, \dots, \sigma_K^2) \in \mathbb{R}^{K \times K}$, then

$$\begin{aligned} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\sigma}) &= \prod_k p(\mathbf{y}_k | \mathbf{x}_k, \sigma_k) \\ &= \prod_k \mathcal{N}(\mathbf{y}_k | \mathbf{x}_k, \sigma_k^2 \mathbf{I}) \end{aligned} \quad (4)$$

where

$\mathbf{y}_k = [y_k(t_1), \dots, y_k(t_N)]^\top \in \mathbb{R}^N$ are the observations for the k -th state over time.

$\mathbf{x}_k = [x_k(t_1), \dots, x_k(t_N)]^\top \in \mathbb{R}^N$ are the values of the k -th state over time.

State prior

Introducing state-specific, independent Gaussian process priors on each \mathbf{x}_k for $k = 1, \dots, K$, then

$$\begin{aligned} p(\mathbf{X}|\boldsymbol{\varphi}) &= \prod_k p(\mathbf{x}_k|\boldsymbol{\varphi}_k) \\ &= \prod_k \mathcal{N}(\mathbf{x}_k|\mathbf{0}, \mathbf{C}_{\boldsymbol{\varphi}_k}) \end{aligned} \tag{5}$$

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

State posterior

Using *Bayes' theorem*, the posterior on X is obtained as

$$\begin{aligned} p(X|Y, \varphi, \sigma) &= \frac{p(X|\varphi)p(Y|X, \sigma)}{\int p(X|\varphi)p(Y|X, \sigma)dX} \\ &= \prod_k p(\mathbf{x}_k|\mathbf{y}_k, \varphi_k, \sigma_k) \\ &= \prod_k \mathcal{N}(\mathbf{x}_k|\boldsymbol{\mu}_k(\mathbf{y}_k), \boldsymbol{\Sigma}_k) \end{aligned} \tag{6}$$

where

$$\begin{aligned} \boldsymbol{\mu}_k(\mathbf{y}_k) &= \mathbf{C}_{\varphi_k}(\mathbf{C}_{\varphi_k} + \sigma_k^2 \mathbf{I})^{-1} \mathbf{y}_k \\ \boldsymbol{\Sigma}_k &= \sigma_k^2 \mathbf{C}_{\varphi_k}(\mathbf{C}_{\varphi_k} + \sigma_k^2 \mathbf{I})^{-1} \end{aligned}$$

Gaussian process response model

Because Gaussian process is closed under differentiation, the joint distribution of \mathbf{x}_k and $\dot{\mathbf{x}}_k$, for $k = 1, \dots, K$, within a finite amount of time points is also Gaussian:

$$\begin{bmatrix} \mathbf{x}_k \\ \dot{\mathbf{x}}_k \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{C}_{\varphi_k} & \mathbf{C}'_{\varphi_k} \\ {}'\mathbf{C}_{\varphi_k} & \mathbf{C}''_{\varphi_k} \end{bmatrix} \right) \quad (7)$$

where

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

Gaussian process response model

The conditional distribution over \dot{X} is given by

$$\begin{aligned} p(\dot{X}|X, \varphi) &= \prod_k p(\dot{x}_k | \mathbf{x}_k, \varphi_k) \\ &= \prod_k \mathcal{N}(\dot{x}_k | \mathbf{m}_k, \mathbf{A}_k) \end{aligned} \quad (8)$$

where

$$\begin{aligned} \mathbf{m}_k &= {}^t C_{\varphi_k} C_{\varphi_k}^{-1} \mathbf{x}_k \\ \mathbf{A}_k &= C''_{\varphi_k} - {}^t C_{\varphi_k} C_{\varphi_k}^{-1} C'_{\varphi_k} \end{aligned}$$

ODE response model

Assuming state-specific, additive Gaussian errors between $\dot{\mathbf{x}}(t)$ and the response from $\mathbf{f}(\mathbf{x}(t), \boldsymbol{\theta})$, we have

$$\begin{aligned} p(\dot{\mathbf{X}}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\gamma}) &= \prod_k p(\dot{\mathbf{x}}_k|\mathbf{X}, \boldsymbol{\theta}, \gamma_k) \\ &= \prod_k \mathcal{N}(\dot{\mathbf{x}}_k|\mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta}), \gamma_k \mathbf{I}) \end{aligned} \quad (9)$$

where

$\dot{\mathbf{x}}_k = [\dot{x}_k(t_1), \dots, \dot{x}_k(t_N)] \in \mathbb{R}^N$ are the derivatives for the k -th state over time.

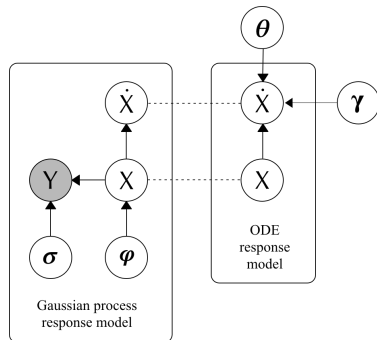
$\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_K]^T \in \mathbb{R}^K$ contains the error variances.

Product of experts

The *product of experts* technique combines (Eq. 8) and (Eq. 9) to obtain

$$p(\dot{X}|X, \varphi, \theta, \gamma) \propto p(\dot{X}|X, \varphi)p(\dot{X}|X, \theta, \gamma) \quad (10)$$

which attains high densities where both $p(\dot{X}|X, \varphi)$ and $p(\dot{X}|X, \theta, \gamma)$ have strong support.



Joint posterior

The joint posterior $p(X, \theta | Y, \varphi, \sigma, \gamma)$ is obtained by

$$\begin{aligned} p(X, \theta | Y, \varphi, \sigma, \gamma) &= \int p(\theta) p(X | Y, \varphi, \sigma) p(\dot{X} | X, \theta, \varphi, \gamma) d\dot{X} \\ &\propto p(\theta) \prod_k [\mathcal{N}(\mathbf{x}_k | \boldsymbol{\mu}_k(\mathbf{y}_k), \boldsymbol{\Sigma}_k) \mathcal{N}(\mathbf{f}_k(X, \theta) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1})] \end{aligned} \quad (11)$$

where

$$\boldsymbol{\Lambda}_k^{-1} = \mathbf{A}_k + \gamma_k \mathbf{I}$$

The “best” parameters θ^* could be estimated using *Maximum a posteriori (MAP)* to yield

$$\begin{aligned} \theta^* &= \arg \max_{\theta} \int p(X, \theta | Y, \varphi, \sigma, \gamma) dX \\ &= \arg \max_{\theta} p(\theta | Y, \varphi, \sigma, \gamma) \end{aligned} \quad (12)$$

which is intractable due to strong non-linear couplings of the states inside the ODEs.

Laplace mean-field approximation

Positing the following factorized proxy distribution:

$$\begin{aligned} Q(X, \theta) &= q(\theta | \eta_\theta, \Xi_\theta) \prod_u q(x_u | \eta_{x_u}, \Xi_{x_u}) \\ &= \mathcal{N}(\theta | \eta_\theta, \Xi_\theta) \prod_u \mathcal{N}(x_u | \eta_{x_u}, \Xi_{x_u}) \end{aligned} \quad (13)$$

Conditional probability $p(\mathbf{x}_u | Y, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})$

Denoting $X_{/\{u\}} = \{\mathbf{x}_o | o = 1, \dots, K \text{ and } o \neq u\}$, for $u = 1, \dots, K$, we have

$$\begin{aligned}
 p(\mathbf{x}_u | Y, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma}) &= \int p(\mathbf{x}_u | Y, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\sigma}) p(\dot{X} | \mathbf{x}_u, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\gamma}) d\dot{X} \\
 &\stackrel{(b)}{=} \int p(\mathbf{x}_u | \mathbf{y}_u, \boldsymbol{\varphi}_k, \sigma_k) p(\dot{X} | X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\gamma}) d\dot{X} \\
 &\propto \mathcal{N}(\mathbf{x}_u | \boldsymbol{\mu}_u(\mathbf{y}_u), \boldsymbol{\Sigma}_u) \prod_k \mathcal{N}(\mathbf{f}_k(X, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1}) \quad (14)
 \end{aligned}$$

where (b) holds because

- $p(\mathbf{x}_u | Y, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\sigma})$ depends only on \mathbf{y}_u due to independent prior assumption,
- $p(\dot{X} | \mathbf{x}_u, X_{/\{u\}}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\gamma})$ is equivalent to $p(\dot{X} | X, \boldsymbol{\varphi}, \boldsymbol{\theta}, \boldsymbol{\gamma})$.

Cost minimization for states

The mean vector and precision matrix of $q(\mathbf{x}_u | \boldsymbol{\eta}_{\mathbf{x}_u}, \boldsymbol{\Xi}_{\mathbf{x}_u})$ for $u = 1, \dots, K$ are given by

$$\begin{aligned}
 \boldsymbol{\eta}_{\mathbf{x}_u} &= \arg \max_{\mathbf{x}_u} \ln [\mathcal{N}(\mathbf{x}_u | \boldsymbol{\mu}_u(\mathbf{y}_u), \boldsymbol{\Sigma}_u) \prod_k \mathcal{N}(\mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1})] \\
 &= \arg \max_{\mathbf{x}_u} [\ln \mathcal{N}(\mathbf{x}_u | \boldsymbol{\mu}_u(\mathbf{y}_u), \boldsymbol{\Sigma}_u) + \sum_k \ln \mathcal{N}(\mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1})] \\
 &= \arg \min_{\mathbf{x}_u} \frac{1}{2} [(\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u))^T \boldsymbol{\Sigma}_u^{-1} (\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u)) + \sum_k (\mathbf{f}_k - \mathbf{m}_k)^T \boldsymbol{\Lambda}_k (\mathbf{f}_k - \mathbf{m}_k)] \\
 &= \arg \min_{\mathbf{x}_u} \text{cost}_{\mathbf{x}_u}(\mathbf{x}_u, \mathbf{X}_{/\{u\}}, \boldsymbol{\theta}, \boldsymbol{\mu}_u(\mathbf{y}_u), \boldsymbol{\Sigma}_u, \mathbf{m}, \boldsymbol{\Lambda}) \tag{15}
 \end{aligned}$$

$$\boldsymbol{\Xi}_{\mathbf{x}_u}^{-1} = \nabla \nabla \text{cost}_{\mathbf{x}_u} |_{\mathbf{x}_u = \boldsymbol{\eta}_{\mathbf{x}_u}} \tag{16}$$

Conditional probability $p(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \gamma)$

For $p(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \gamma)$, we have

$$\begin{aligned}
 p(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \gamma) &\stackrel{(a)}{=} p(\boldsymbol{\theta}|\mathbf{X}, \boldsymbol{\varphi}, \gamma) \\
 &= \int p(\boldsymbol{\theta}) p(\dot{\mathbf{X}}|\mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\theta}, \gamma) d\dot{\mathbf{X}} \\
 &\propto \prod_k \mathcal{N}(\mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1})
 \end{aligned} \tag{17}$$

where (a) holds since $\boldsymbol{\theta}$ depends indirectly on \mathbf{Y} through \mathbf{X} .

Cost minimization for parameters

Denoting $\mathbf{f}_k = \mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta})$, $\mathbf{m} = [\mathbf{m}_1, \dots, \mathbf{m}_K]$, and $\boldsymbol{\Lambda} = [\boldsymbol{\Lambda}_1, \dots, \boldsymbol{\Lambda}_K]$, then mean vector and precision matrix of $p(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{X}, \boldsymbol{\varphi}, \boldsymbol{\gamma})$ are given by

$$\begin{aligned}
 \boldsymbol{\eta}_{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \ln \prod_k \mathcal{N}(\mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1}) \\
 &= \arg \max_{\boldsymbol{\theta}} \sum_k \ln \mathcal{N}(\mathbf{f}_k(\mathbf{X}, \boldsymbol{\theta}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1}) \\
 &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \sum_k (\mathbf{f}_k - \mathbf{m}_k)^T \boldsymbol{\Lambda}_k (\mathbf{f}_k - \mathbf{m}_k) \\
 &= \arg \min_{\boldsymbol{\theta}} \text{cost}_{\boldsymbol{\theta}}(\mathbf{X}, \boldsymbol{\theta}, \mathbf{m}, \boldsymbol{\Lambda})
 \end{aligned} \tag{18}$$

$$\boldsymbol{\Xi}_{\boldsymbol{\theta}}^{-1} = \nabla \nabla \text{cost}_{\boldsymbol{\theta}} |_{\boldsymbol{\theta}=\boldsymbol{\eta}_{\boldsymbol{\theta}}} \tag{19}$$

Inference algorithm

- Initialize using Gaussian process regression
- Repeat until convergence or maximum iteration
 - Update θ while keeping the others fixed
 - For $u = 1, \dots, K$, update x_u while keeping the others fixed
- Calculate precision matrices

Derivation for the gradients and Hessians

Recall that $cost_{x_u}$ for state u is given by

$$cost_{x_u} = \frac{1}{2}[(\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u))^T \boldsymbol{\Sigma}_u^{-1}(\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u)) + \sum_k (\mathbf{f}_k - \mathbf{m}_k)^T \boldsymbol{\Lambda}_k(\mathbf{f}_k - \mathbf{m}_k)]$$

Using matrix derivative and the fact that $\boldsymbol{\Sigma}_u^{-1}$ is symmetric, we have

$$\nabla_{x_u} \frac{1}{2}(\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u))^T \boldsymbol{\Sigma}_u^{-1}(\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u)) = \boldsymbol{\Sigma}_u^{-1} \mathbf{x}_u \quad (20)$$

$$\nabla \nabla_{x_u} \frac{1}{2}(\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u))^T \boldsymbol{\Sigma}_u^{-1}(\mathbf{x}_u - \boldsymbol{\mu}_u(\mathbf{y}_u)) = \boldsymbol{\Sigma}_u^{-1} \quad (21)$$

Derivation for the gradients and Hessians

Using *chain rule* and the fact that Λ_k is symmetric, we have

$$\begin{aligned}
 & \nabla_{x_u} \frac{1}{2} (\mathbf{f}_k - \mathbf{m}_k)^T \Lambda_k (\mathbf{f}_k - \mathbf{m}_k) \\
 &= \begin{bmatrix} \frac{\partial(\mathbf{f}_k)_1}{\partial x_u(t_1)} & \cdots & \frac{\partial(\mathbf{f}_k)_N}{\partial x_u(t_1)} \\ \vdots & \ddots & \vdots \\ \frac{\partial(\mathbf{f}_k)_1}{\partial x_u(t_N)} & \cdots & \frac{\partial(\mathbf{f}_k)_N}{\partial x_u(t_N)} \end{bmatrix} \Lambda_k (\mathbf{f}_k - \mathbf{m}_k) \\
 &\quad - \begin{bmatrix} \frac{\partial(\mathbf{m}_k)_1}{\partial x_u(t_1)} & \cdots & \frac{\partial(\mathbf{m}_k)_N}{\partial x_u(t_1)} \\ \vdots & \ddots & \vdots \\ \frac{\partial(\mathbf{m}_k)_1}{\partial x_u(t_N)} & \cdots & \frac{\partial(\mathbf{m}_k)_N}{\partial x_u(t_N)} \end{bmatrix} \Lambda_k (\mathbf{f}_k - \mathbf{m}_k)
 \end{aligned} \tag{22}$$

Derivation for the gradients and Hessians

The (i, j) -th entry of the Hessian is given by

$$\begin{aligned}
 & \frac{\partial^2 \frac{1}{2} (\mathbf{f}_k - \mathbf{m}_k)^T \mathbf{\Lambda}_k (\mathbf{f}_k - \mathbf{m}_k)}{\partial x_u(t_i) \partial x_u(t_j)} \\
 &= \left[\frac{\partial^2 (\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_i) \partial x_u(t_j)} \quad \dots \quad \frac{\partial^2 (\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_i) \partial x_u(t_j)} \right] \mathbf{\Lambda}_k (\mathbf{f}_k - \mathbf{m}_k) \\
 &+ \left[\frac{\partial (\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_j)} \quad \dots \quad \frac{\partial (\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_j)} \right] \mathbf{\Lambda}_k \begin{bmatrix} \frac{\partial (\mathbf{f}_k - \mathbf{m}_k)_1}{\partial x_u(t_i)} \\ \vdots \\ \frac{\partial (\mathbf{f}_k - \mathbf{m}_k)_N}{\partial x_u(t_i)} \end{bmatrix}
 \end{aligned} \tag{23}$$

Positivity constraint

Let

$$\boldsymbol{\theta} = [\theta_1, \dots, \theta_M]^\top = [e^{\tilde{\theta}_1}, \dots, e^{\tilde{\theta}_M}]^\top \quad (24)$$

Since the exponential function is monotonic, we can first find

$$\begin{aligned} \tilde{\boldsymbol{\theta}}^* &= \arg \min_{\tilde{\boldsymbol{\theta}}} \text{cost}_{\boldsymbol{\theta}}(\mathbf{X}, e^{\tilde{\boldsymbol{\theta}}}, \mathbf{m}, \boldsymbol{\Lambda}) \\ &= \arg \min_{\tilde{\boldsymbol{\theta}}} \sum_k \ln \mathcal{N}(\mathbf{f}_k(\mathbf{X}, e^{\tilde{\boldsymbol{\theta}}}) | \mathbf{m}_k, \boldsymbol{\Lambda}_k^{-1}) \end{aligned} \quad (25)$$

and then obtain $\boldsymbol{\theta}^*$ as

$$\boldsymbol{\theta}^* = [e^{\tilde{\theta}_1^*}, \dots, e^{\tilde{\theta}_M^*}]^\top \quad (26)$$

Since $e^r > 0$ for any $r \in \mathbb{R}$, we essentially restrict $\boldsymbol{\theta}^*$ to positive values.

Positivity constraint

The positivity constraint on states can be achieved by transforming $cost_{x_u}$ to $cost_{\tilde{x}_u}$, where we define for $u = 1, \dots, K$

$$\mathbf{x}_u = [x_u(t_1), \dots, x_u(t_N)]^\top = [e^{\tilde{x}_u(t_1)}, \dots, e^{\tilde{x}_u(t_N)}]^\top \quad (27)$$

Using chain rule, we have for $u = 1, \dots, K$ and $n = 1, \dots, N$ the following:

$$\begin{aligned} \frac{d\tilde{x}_u(t_n)}{dt} &= \frac{d \ln x_u(t_n)}{dt} \\ &= \frac{1}{x_u(t_n)} \frac{dx_u(t_n)}{dt} \\ &= \frac{f_u(e^{\tilde{x}(t_n)}, \boldsymbol{\theta})}{e^{\tilde{x}(t_n)}} \end{aligned} \quad (28)$$

Positivity constraint

Caveats

- The covariance matrix for the new variables cannot be transformed back to the original variables easily.
- Probabilistic interpretability is lost since proper change of random variables requires the evaluation of the Jacobian determinant, which is computationally expensive.

Extension to Random Dynamical Systems

Random ordinary differential equations (RODEs)

Given a complete probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let

$(\zeta_t)_{t \in [0, T]} \in \mathbb{R}^W$ be a stochastic process with continuous sample paths,
 $f : \mathbb{R}^K \times \mathbb{R}^W \mapsto \mathbb{R}^K$ be a continuous function.

For all $\omega \in \Omega$, a K -dimensional RODE defined as

$$\frac{dx(t)}{dt} = f(x(t), \zeta_t(\omega)) \quad (29)$$

is a *non-autonomous* ODE system

$$\dot{x}(t) = \frac{dx(t)}{dt} = F_\omega(x, t) = f(x(t), \omega_t) \quad (30)$$

RODEs

An example of a scalar RODE with additive noise is given by

$$\frac{dx(t)}{dt} = -x + \cos(W_t(\omega)) \quad (31)$$

where W_t is a one-dimensional Wiener process.

To ensure the existence of a unique solution for the initial value problem on the finite time interval $[0, T]$, we assume that f is infinitely differentiable in x , and hence, it is locally *Lipschitz* in x .

RODEs

Since ζ_t is usually *Hölder continuous* in t , the vector fields of F_ω are continuous but not differentiable in t for every $\omega \in \Omega$.

- Numerical schemes, e.g. the *Runge-Kutta* method, fail to achieve high order of convergence.
- The solution paths of the RODEs are once differentiable, which implies that the gradient matching model can be ideally applied.

Solve many RODEs sample paths deterministically in parallel to obtain an *ensemble* solution.

- The Laplace mean-field approximation is computationally very efficient.

Doss-Sussmann/Imkeller-Schmalfuss correspondence

Any finite dimensional SDE system with additive noise can be transformed to an equivalent RODE system and vice versa as follows:

$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{x}(t), \boldsymbol{\theta})dt + d\mathbf{W}_t \Leftrightarrow \frac{dz(t)}{dt} = \mathbf{f}(z(t) + \mathbf{O}_t, \boldsymbol{\theta}) + \mathbf{O}_t \quad (32)$$

where $z(t) = \mathbf{x}(t) - \mathbf{O}_t$ and \mathbf{O}_t is a stationary stochastic *Ornstein-Uhlenbeck* process defined as

$$d\mathbf{O}_t = -\mathbf{O}_t dt + d\mathbf{W}_t \quad (33)$$

Using the scalar RODE in (Eq. 31) as an example,

$$d \begin{pmatrix} x_t \\ y_t \end{pmatrix} = \begin{pmatrix} -x_t + \cos(y_t) \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} dW_t \Leftrightarrow \frac{dx(t)}{dt} = -x + \cos(W_t(\omega)) \quad (34)$$

Inference algorithm

-
- 1: Transform the SDEs into RODEs.
 - 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 using Laplace mean-field approximation
 - 5: **end for**
 - 6: Average the estimation results from all the sample paths.
-

Experiments

Lotka-Volterra model

Two non-linear ODEs to model the interaction between prey and predator.

$$\begin{aligned}\dot{x}(t) &= \alpha x(t) - \beta x(t)y(t) \\ \dot{y}(t) &= \delta x(t)y(t) - \gamma y(t)\end{aligned}\tag{35}$$

where

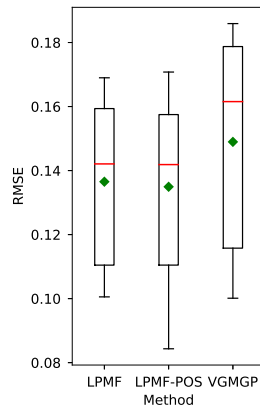
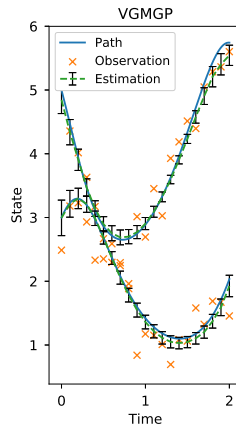
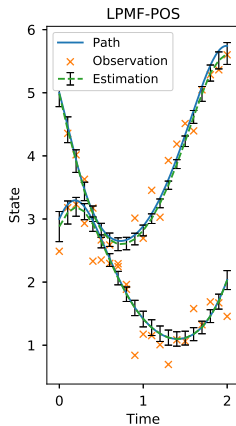
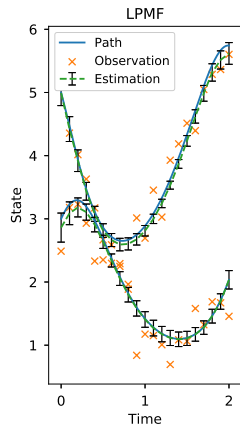
$x(t), y(t) \in \mathbb{R}_{\geq 0}$ are the populations of prey and predator at time t .

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

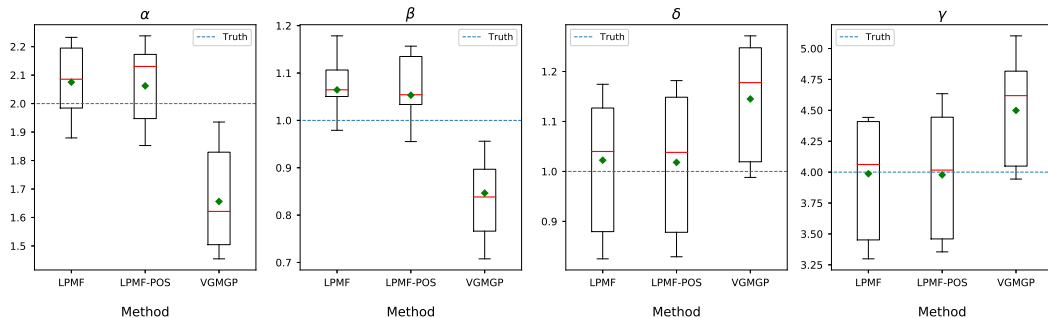
Experimental setup

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

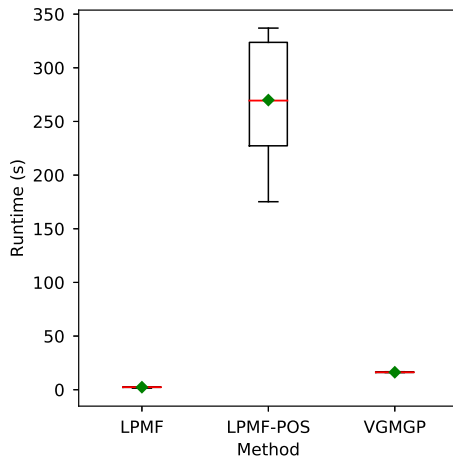
State estimation



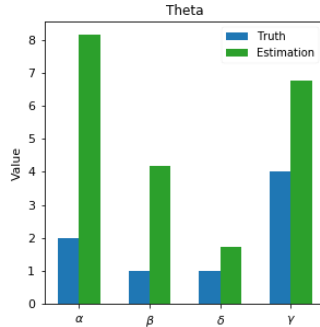
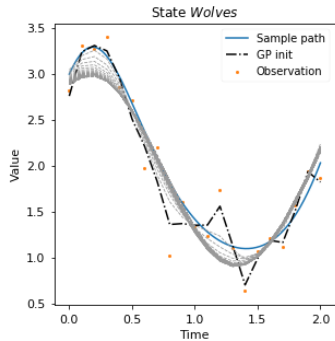
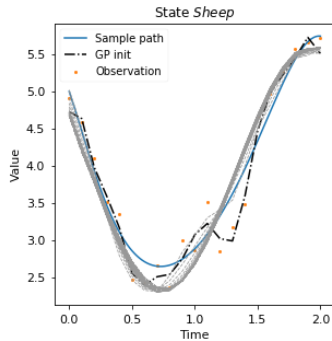
Parameter estimation



Runtime performance



Positivity constraint on states



Protein signalling transduction pathway

A signal transduction cascade model to describe the dynamics among protein species.

$$\begin{aligned}
 \dot{S} &= -k_1 \times S - k_2 \times S \times R + k_3 \times RS \\
 d\dot{S} &= k_1 \times S \\
 \dot{R} &= -k_2 \times S \times R + k_3 \times RS + V \times \frac{Rpp}{K_m + Rpp} \\
 \dot{RS} &= k_2 \times S \times R - k_3 \times RS - k_4 \times RS \\
 \dot{Rpp} &= k_4 \times RS - V \times \frac{Rpp}{K_m + Rpp}
 \end{aligned} \tag{36}$$

where

$S, dS, R, RS, Rpp \in \mathbb{R}_{\geq 0}$ are the concentrations of proteins.

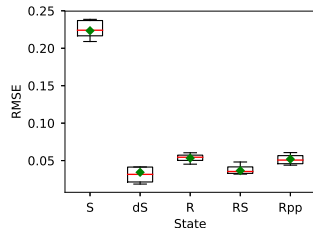
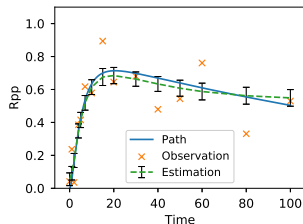
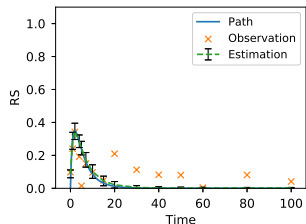
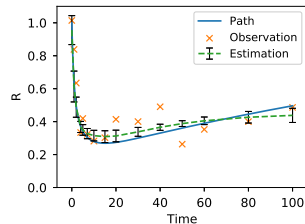
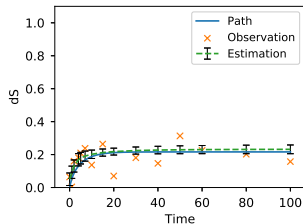
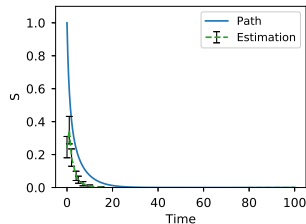
$k_1, k_2, k_3, k_4, K_m, V \in \mathbb{R}^+$ are kentic parameters

Experimental setup

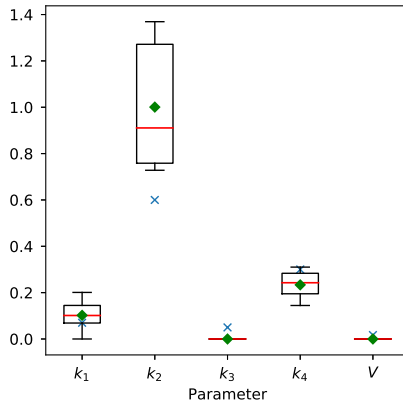
K	K_{obs}	t_0	t_T	δt	$k_1, k_2, k_3, k_4, V, Km$	σ_k^2
5	4	0	100	0.05	0.07, 0.6, 0.05, 0.3, 0.017, 3	0.01

The observations, in total 15, are collected at time points 0, 1, 2, 4, 5, 7, 10, 15, 20, 30, 40, 50, 60, 80 and 100.

State estimation



Parameter estimation



Lorenz 96 model

A minimalistic weather forecast model. A K -dimensional deterministic Lorenz 96 dynamical system is defined for $k = 1, \dots, K$, state-wise as follows:

$$\dot{x}_k(t) = (x_{k+1}(t) - x_{k-2}(t))x_{k-1}(t) - x_k(t) + F \quad (37)$$

where

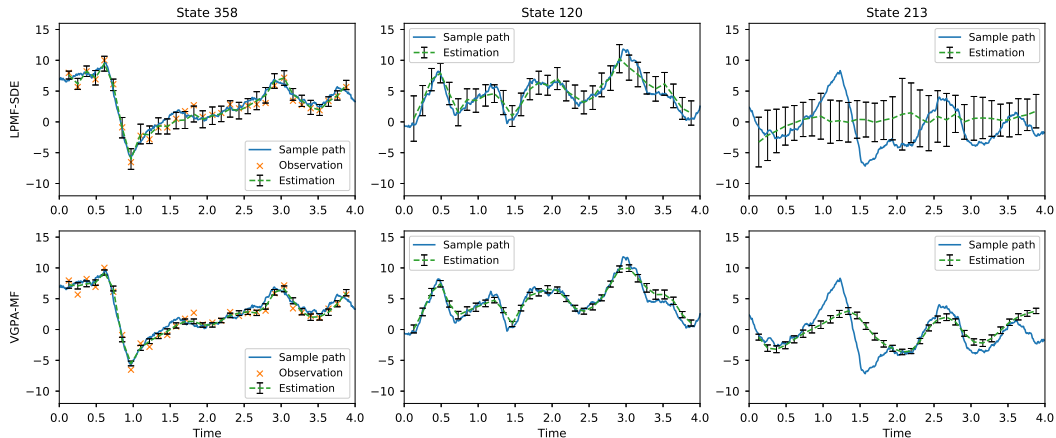
$$\begin{aligned} x_{-1}(t) &= x_{K-1}(t) \\ x_0(t) &= x_K(t) \\ x_{K+1}(t) &= x_1(t) \end{aligned}$$

$$F \in \mathbb{R} \text{ controls the behavior of the system} \quad (38)$$

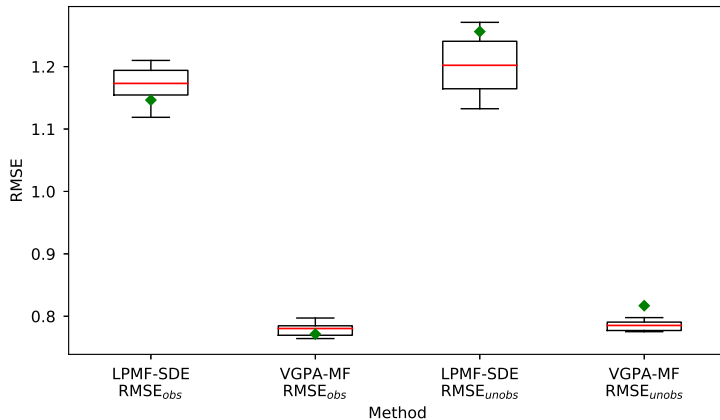
Experimental setup

K	K_{obs}	t_0	t_T	δt	F	σ_k^2	ρ_k^2	$freq_{obs}$
500	325 (65%)	0	4	0.01	8	1	4	8

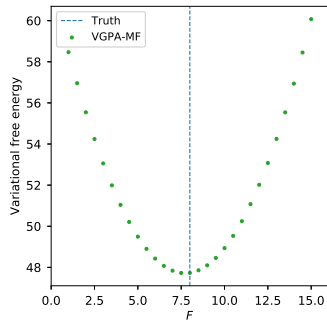
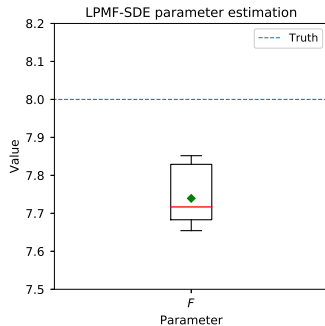
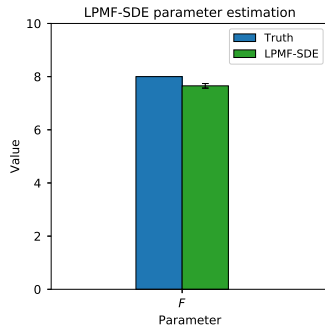
State estimation



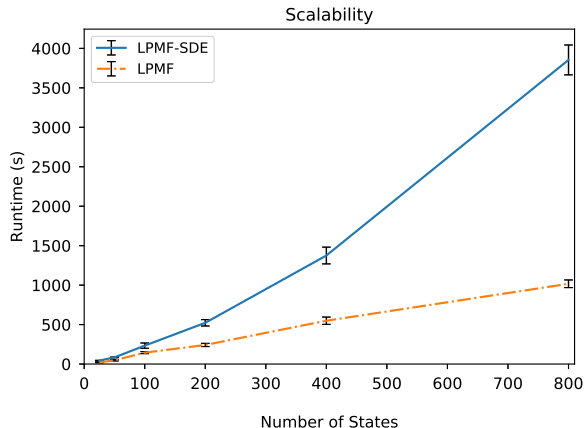
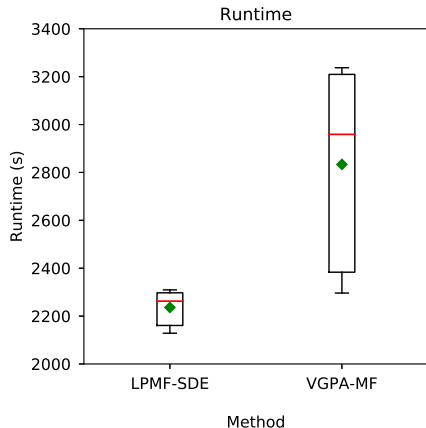
State estimation



Parameter estimation



Runtime performance



Lorenz 63 model

Low dimensional mathematical model for thermal convection in the atmosphere. The vector field of the deterministic Lorenz 63 is defined as follows:

$$\mathbf{f}(\mathbf{x}(t), \boldsymbol{\theta}) = \begin{bmatrix} \dot{x}(t) \\ \dot{y}(t) \\ \dot{z}(t) \end{bmatrix} = \begin{bmatrix} \sigma(y(t) - x(t)) \\ x(t)(\rho - z(t)) - y(t) \\ x(t)y(t) - \beta z(t) \end{bmatrix} \quad (39)$$

where

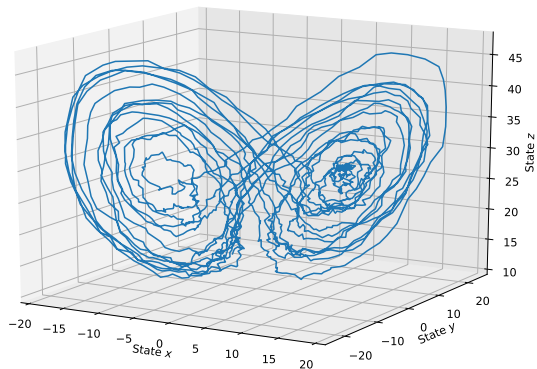
$\mathbf{x}(t) = [x(t), y(t), z(t)]^\top \in \mathbb{R}^3$ is the state vector.

$\boldsymbol{\theta} = [\sigma, \rho, \beta]^\top \in \mathbb{R}^3$ is the parameter vector

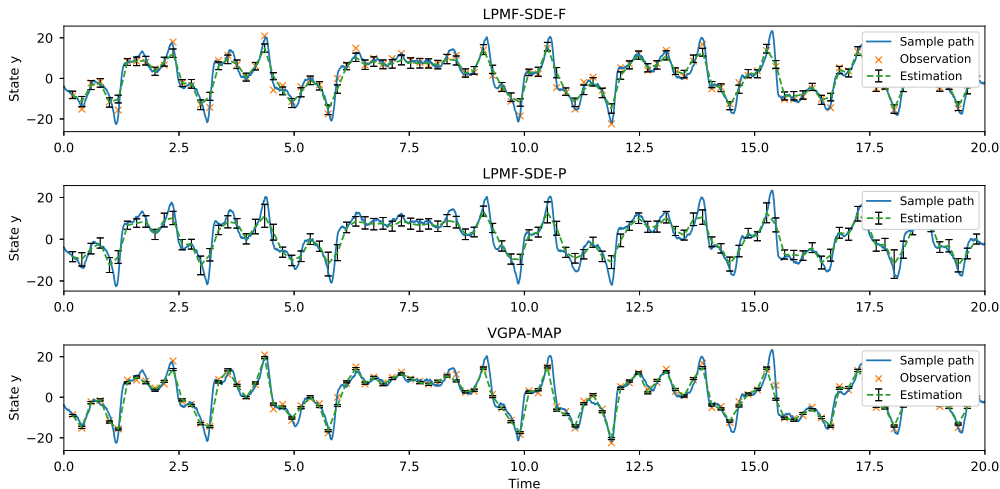
Experimental setup

K	K_{obs}	t_0	t_T	δt	σ, ρ, β	σ_k^2	ρ_k^2	$freq_{obs}$
3	3 or 2	0	20	0.01	10, 28, $\frac{8}{3}$	2	10	5

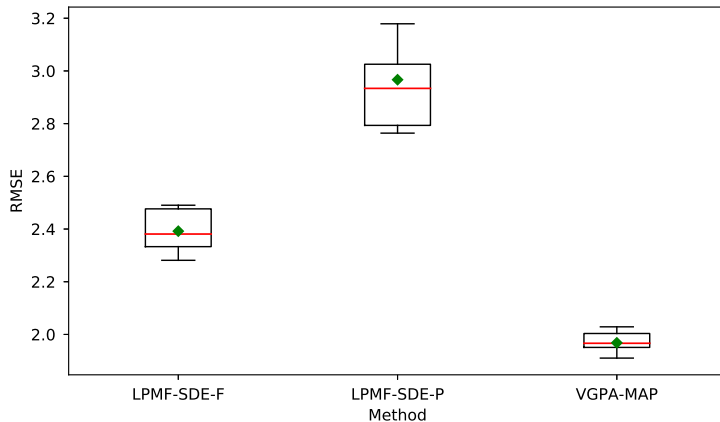
Sample path



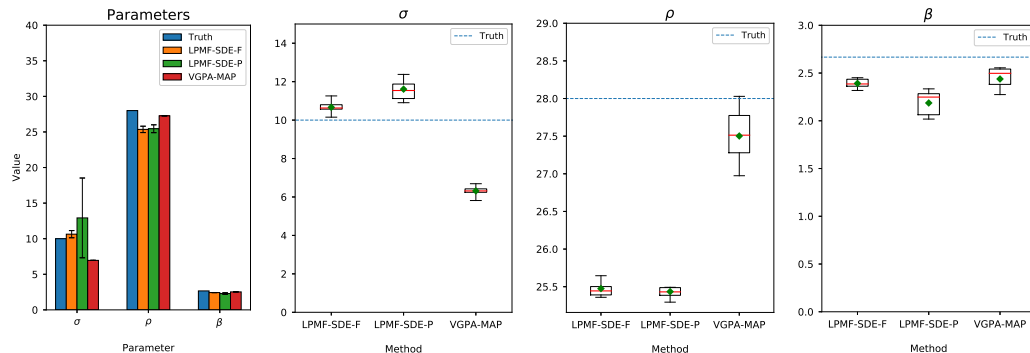
State estimation



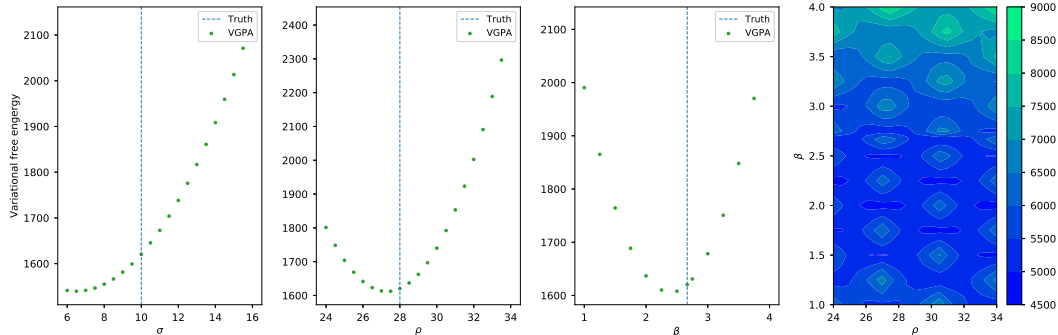
State estimation



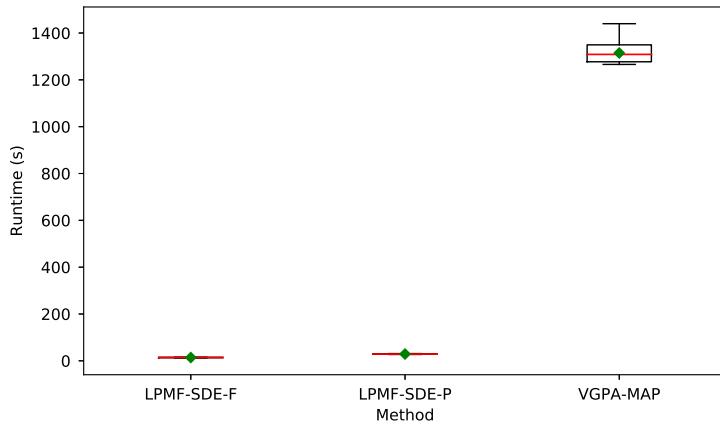
Parameter estimation



Parameter estimation



Runtime performance



Conclusion

Contributions

- Implemented a brand new Python inference library.
- Proposed the Laplace mean-field approximation for ODEs.
- Derived analytical solution for the gradients and Hessians.
- Improved the flexibility of the ODE inference algorithm by
 - relaxing the structural assumption on the dynamical systems.
 - adding support for positivity constraints on states and parameters.
- Developed an efficient parallel inference pipeline for SDEs.

Future work

- Add support to infer the diffusion noise.
- Add support to select hyperparameters.
- Improve the numerical optimization strategy.
- Add GPU implementation.

Q & A