

EVALUATION OF VARIATIONAL AND MARKOV CHAIN MONTE CARLO METHODS FOR INFERENCE IN PARTIALLY OBSERVED STOCHASTIC DYNAMIC SYSTEMS.

Y. Shen*, C. Archambeau[†], D. Cornford*, M. Opper[‡], J. Shawe-Taylor[†], and R. Barillec*

*Neural Computing Research Group, Aston University, Birmingham, United Kingdom

[†]Department of Computer Science, University College London, London, United Kingdom

[‡]Artificial Intelligence Group, Technical University Berlin, Berlin, Germany

ABSTRACT

In recent work we have developed a novel variational inference method for partially observed systems governed by stochastic differential equations. In this paper we provide a comparison of the Variational Gaussian Process Smoother with an exact solution computed using a hybrid Monte Carlo approach to path sampling, applied to a stochastic double well potential model. It is demonstrated that the variational smoother provides us a very accurate estimate of mean path while marginal variance is slightly underestimated. We conclude with some remarks as to the advantages and disadvantages of the variational smoother.

1. INTRODUCTION

Stochastic dynamical systems have been used for modelling of real-life systems in various areas ranging from physics to system biology to environmental science. Such systems are often only partially observed, which makes statistical inference in those systems difficult. The inference problem for stochastic dynamical systems usually includes both state and parameter estimation. In this paper, we focus on state estimation and assume that the system equation and its parameters are both known *a priori*. This is known as filtering and/or smoothing in statistical signal processing. It is known that the Kushner-Stratonovich-Pardoux (KSP) equations are the optimal solution to a general filtering/smoothing problem [1, 2, 3]. For linear systems, the filtering part of KSP equations reduces to the well-known Kalman-Bucy filter [4] which is computationally efficient. For non-linear dynamics in general, however, filtering/smoothing is still a challenging problem because a numerical solution to the KSP equations is not feasible for high-dimensional systems. This paper presents a comparison study between a Markov Chain Monte Carlo (MCMC) smoother [5] and the variational smoother which we have recently developed [6].

Mathematically, a stochastic dynamical system is often represented by a stochastic differential equation (SDE) [7]:

$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{x}, t)dt + (\mathbf{2D})^{1/2}(\mathbf{x}, t)d\mathbf{W}(t), \quad (1)$$

where $\mathbf{x}(t) \in \mathcal{R}^d$ is the state vector, $\mathbf{D} \in \mathcal{R}^{d \times d}$ is the so-called diffusion matrix, \mathbf{f} represents a deterministic dynamical process. The driving noise process is represented by a Wiener process $\mathbf{W}(t)$. Eq. 1 is also referred to as a diffusion process. The state is observed via some measurement function $\mathbf{h}(\cdot)$ at discrete times, say $\{t_k\}_{k=1, \dots, M}$. The observations are assumed contaminated by i.i.d Gaussian noise:

$$\mathbf{y}_k = \mathbf{h}(\mathbf{x}(t_k)) + \mathbf{R}^{\frac{1}{2}} \cdot \eta \quad (2)$$

where $\mathbf{y}_k \in \mathcal{R}^{d'}$ is the k -th observation, $\mathbf{R} \in \mathcal{R}^{d' \times d'}$ is the covariance matrix of measurement errors, and ξ represents multivariate white noise. A Bayesian approach to filtering/smoothing is typically adopted in which the posterior distribution

$$p(\mathbf{x}(t) | \{\mathbf{y}_1, \dots, \mathbf{y}_k, t_k < t\})$$

and

$$p(\mathbf{x}([0, T]) | \{\mathbf{y}_1, \dots, \mathbf{y}_M, t_M < T\}),$$

respectively, are to be formulated and estimated. Theoretically, an optimal estimate of $p(\cdot)$ is the solution to the corresponding KSP equations. Computational approaches are based on a variety of approximation schemes or achieved by MCMC sampling methods.

Using Markov Chain Monte Carlo [8], one is able to sample from a posterior distribution exactly. At each step of a MCMC simulation, a new state is proposed and will be accepted or rejected in a probabilistic way. For applications to stochastic dynamical systems, it is also referred to path sampling. A path sampling approach to discrete-time state-space models has been addressed in [5] and references

therein. In those works, a Gibbs-sampler with single-site update was used. To achieve better mixing, several algorithms using multiple-site update are explored in [9]. Recently, a Hybrid Monte Carlo (HMC) algorithm for path sampling is proposed in [10]. We have extended this to include a prior over the initial state. The HMC method updates the entire sample path at each step of path sampling while keeping the acceptance of new paths high. In this work, we first scrutinise the use of HMC for non-linear smoothing and then assess the performance of the variational smoother proposed in [6] by comparing its results with those of HMC.

In contrast to MCMC, all other approaches to non-linear filtering/smoothing, including the one proposed in [6], are based on a particular approximation scheme. The extended Kalman filter is the first attempt to tackle the non-linearity by linearising the dynamics around the currently available state estimate [11]. However, unstable error growth is observed in such linearisation methods [12]. To alleviate this difficulty, the Ensemble Kalman Filter (EnKF) was introduced in [12]. An ensemble of states are integrated forward in time. Therefore, the Kalman gain can be estimated by using the error covariances which are not propagated but calculated from the ensemble of states at each time step. Note that this method drops the linear approximation of non-linear dynamics while keeping the Gaussian assumption of error statistics. The particle filter (PF) proposed in [13] represents a different direction of approximation strategies. Essentially, the posterior density of filtering variables in PF is approximated by a discrete distribution with random support. Each member in the discrete support is called particle and its probability mass is considered as weight.

In essence, the variational smoother in [6] makes a global linear approximation of non-linear dynamics. This implies a Gaussian approximation of the posterior process

$$p(\mathbf{x}([0, T]) | \{\mathbf{y}_1, \dots, \mathbf{y}_M, t_M < T\}).$$

The quality of approximation is measured by Kullback-Leibler (KL) divergence [14] between the true and approximate posterior. The optimal approximate posterior is obtained by minimising the KL divergence. Following this, any statistical inference in the true system is based on the approximate posterior. This method is within the framework of variational approximation for Bayesian inference, which is computationally very efficient and popular in the machine learning community [15].

The structure of this paper is as follows: First, we present a Bayesian treatment of non-linear smoothing. In Sec. 3, the MCMC method is described in detail while we give a summary of the variational smoother in Sec. 4. For detailed

proofs, we refer to [6]. After that, we compare both methods in Sec. 5 by numerical experiments with a double-well potential system. The paper concludes with a discussion.

2. BAYESIAN APPROACH TO NON-LINEAR SMOOTHING

Both for the MCMC method in [10] and for the variational smoother in [6], stochastic differential equations are discretized by using an explicit Euler-Maruyama scheme [7]. The discretized version of Eq. 1 is given by

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{f}(\mathbf{x}_k, t_k)\delta t + (2\mathbf{D})^{1/2}(\mathbf{x}_k, t_k)\sqrt{\delta t} \cdot \xi_k, \quad (3)$$

with $t_k = k \cdot \delta t$, $k = 0, 1, \dots, N$, and a smoothing window from $t = 0$ to $T = N \cdot \delta t$. Note that ξ_k are white noise random variables. An initial state \mathbf{x}_0 needs to be set. There are M observations within the smoothing window, and they are denoted by

$$(t_{k_j}, \mathbf{y}_j)_{j=1, \dots, M} \quad \text{with} \quad \{t_{k_1}, \dots, t_{k_M}\} \subseteq \{t_0, \dots, t_N\}.$$

In the following, we formulate the posterior distribution step by step.

The prior of a diffusion process can be written down as

$$p(\mathbf{x}_0, \dots, \mathbf{x}_N) = p(\mathbf{x}_0) \cdot p(\mathbf{x}_1 | \mathbf{x}_0) \cdot \dots \cdot p(\mathbf{x}_N | \mathbf{x}_{N-1}),$$

where $p(\mathbf{x}_0)$ is the prior of initial states and $p(\mathbf{x}_{k+1} | \mathbf{x}_k)$ with $k = 0, \dots, N - 1$ are transition densities of the diffusion process. For small enough δt , those transition densities can be well approximated by a Gaussian density [16]. Accordingly,

$$p(\mathbf{x}_{k+1} | \mathbf{x}_k) = \mathcal{N}(\mathbf{x}_{k+1} | \mathbf{x}_k + \mathbf{f}(\mathbf{x}_k)\delta t, 2\mathbf{D}\delta t).$$

Therefore, the prior is given by

$$p(\mathbf{x}_0, \dots, \mathbf{x}_N) \propto p(\mathbf{x}_0) \cdot \exp(-\mathcal{H}_{\text{dynamics}}),$$

where

$$\mathcal{H}_{\text{dynamics}} = \sum_{k=0}^{N-1} \frac{\delta t}{4} \left[\frac{\mathbf{x}_{k+1} - \mathbf{x}_k}{\delta t} - \mathbf{f}(\mathbf{x}_k, t_k) \right]^\top \mathbf{D}^{-1} \left[\frac{\mathbf{x}_{k+1} - \mathbf{x}_k}{\delta t} - \mathbf{f}(\mathbf{x}_k, t_k) \right].$$

As we assume that measurement noise is i.i.d. Gaussian random variable, the likelihood is simply given by

$$p(\mathbf{y}_1, \dots, \mathbf{y}_M | \mathbf{x}(t_0), \dots, \mathbf{x}(t_N)) = \exp(-\mathcal{H}_{\text{obs}}),$$

where

$$\mathcal{H}_{\text{obs}} = \frac{1}{2} \sum_{j=1}^M [\mathbf{h}(\mathbf{x}(t_{k_j})) - \mathbf{y}_j]^\top \mathbf{R}^{-1} [\mathbf{h}(\mathbf{x}(t_{k_j})) - \mathbf{y}_j]. \quad (4)$$

In summary, we have the posterior

$$p(\mathbf{x}(t) | \{\mathbf{y}_1, \dots, \mathbf{y}_M\}) \propto p(\mathbf{x}_0) \cdot \exp(-1(\mathcal{H}_{\text{dynamics}} + \mathcal{H}_{\text{obs}})).$$

3. MCMC METHOD

In Hybrid Monte Carlo, the molecular dynamics simulation algorithm is applied to make proposals in a Metropolis-Hastings algorithm, for example,

$$\mathbf{X}^k = \{\mathbf{x}_0^k, \dots, \mathbf{x}_N^k\} \longrightarrow \mathbf{X}^{k+1} = \{\mathbf{x}_0^{k+1}, \dots, \mathbf{x}_N^{k+1}\},$$

at step k . To make a proposal of \mathbf{X}^{k+1} , one simulates a fictitious deterministic system as follows

$$\begin{aligned} \frac{d\mathbf{X}}{d\tau} &= \mathbf{P} \\ \frac{d\mathbf{P}}{d\tau} &= -\nabla_{\mathbf{X}} \hat{\mathcal{H}}(\mathbf{X}, \mathbf{P}) \end{aligned}$$

where $\mathbf{P} = (\mathbf{p}_0, \dots, \mathbf{p}_N)$ represents momentum and $\hat{\mathcal{H}}$ is a fictitious Hamiltonian which is the sum of potential energy \mathcal{H}^{pot} and kinetic energy $\mathcal{H}^{kin} = \frac{1}{2} \sum_{k=1}^N \mathbf{p}_k^2$. For the posterior distribution of non-linear smoothing in Sec. 2, the potential energy is given by

$$\mathcal{H}^{pot} = -\log(p(\mathbf{x}_0)) + \mathcal{H}^{\text{dynamics}} + \mathcal{H}^{\text{obs}}.$$

The above system is initialised by setting $\mathbf{X}(\tau = 0) = \mathbf{X}_k$ and sampling a random number from $\mathcal{N}(0, 1)$ for each component of $\mathbf{P}(\tau = 0)$. After that, one integrates the system equations forward in time with time increment $\delta\tau$ by using leapfrog as follows:

$$\begin{aligned} \mathbf{X}' &= \mathbf{X} + \delta\tau \mathbf{Z} \mathbf{P} + \frac{\delta\tau^2}{2} \mathbf{Z} \mathbf{Z}^\top (-\nabla_{\mathbf{X}} \hat{\mathcal{H}}) \\ \mathbf{P}' &= \mathbf{P} + \frac{\delta\tau}{2} \mathbf{Z}^\top (-\nabla_{\mathbf{X}} \hat{\mathcal{H}} - \nabla_{\mathbf{X}'} \hat{\mathcal{H}}) \end{aligned}$$

where \mathbf{Z} denotes the so-called preconditioning matrix which accelerates the convergence of matrix iterations. The matrix \mathbf{Z} is a circulant matrix which is constructed from the vector $\{1, \exp(-\alpha), \dots, \exp(-\alpha \cdot T)\}$, where α is a tuning parameter. After J iterations, the state $\mathbf{X}(\tau = J\delta\tau)$ is proposed as \mathbf{X}^{k+1} which will be accepted with probability

$$\min\{1, \exp(-\hat{\mathcal{H}}^{k+1} + \hat{\mathcal{H}}^k)\}.$$

4. VARIATIONAL GAUSSIAN PROCESS APPROXIMATION SMOOTHER

The starting point of the Variational Gaussian Process Approximation (VGPA) method is to approximate Eq. 1 by a linear SDE:

$$d\mathbf{x}(t) = \mathbf{f}_L(\mathbf{x}, t)dt + (2\mathbf{D})^{1/2}(\mathbf{x}, t)d\mathbf{W}(t), \quad (5)$$

where

$$\mathbf{f}_L(\mathbf{x}, t) = -\mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t). \quad (6)$$

The matrix $\mathbf{A}(t) \in \mathcal{R}^{d \times d}$ and the vector $\mathbf{b}(t) \in \mathcal{R}^d$ are two variational parameters to be optimised.

The approximation made by Eq. 6 implies that the true posterior process, i.e. $p(\mathbf{x}(t)|\mathbf{y}_1, \dots, \mathbf{y}_M)$, is approximated by a Gaussian Markov process, $q(t)$. If we discretise the linear SDE in the same way as the true SDE, the approximate posterior can be written down as

$$q(\mathbf{x}_0, \dots, \mathbf{x}_N) = q(\mathbf{x}_0) \cdot \prod_{k=0}^{N-1} \mathcal{N}(\mathbf{x}_{k+1} | \mathbf{x}_k + \mathbf{f}_L(\mathbf{x}_k)\delta t, 2\mathbf{D}\delta t).$$

In [6], $q(\mathbf{x}_0)$ is fixed to $\mathcal{N}(\mathbf{x}_0 | \mathbf{m}_0, \mathbf{S}_0)$, and the prior on initial states $p(\mathbf{x}_0)$ is implicitly a uniform distribution.

The optimal $\mathbf{A}(t)$ and $\mathbf{b}(t)$ are obtained by minimising the KL divergence of $q(\cdot)$ and $p(\cdot)$ which is given by

$$\text{KL}[q||p] = \int dq \ln \frac{dq}{dp} = \int_0^T E(t)dt + \text{const.} \quad (7)$$

with $E(t) = E_{sde}(t) + E_{obs}(t)$, where

$$E_{sde}(t) = \frac{1}{4} \langle (\mathbf{f}(\mathbf{x}) - \mathbf{f}_L(\mathbf{x}))^\top \mathbf{D}^{-1} (\mathbf{f}(\mathbf{x}) - \mathbf{f}_L(\mathbf{x})) \rangle_{q_t}$$

and $E_{obs}(t) = \langle \mathcal{H}^{obs} \rangle_{q_t}$ where \mathcal{H}^{obs} is defined as Eq. 4 and $\langle \cdot \rangle_{q_t}$ denotes the expectation w. r. t. the marginal distribution of the approximate posterior process $q(\cdot)$ at time t .

To compute the KL divergence, we introduce two auxiliary variational parameters $\mathbf{m}(t)$ and $\mathbf{S}(t)$ which are the mean and covariance matrix of the marginal distribution q_t . However, the pair $(\mathbf{m}(t), \mathbf{S}(t))$ is not independent of $(\mathbf{A}(t), \mathbf{b}(t))$, leading to the following constraints:

$$\frac{d\mathbf{m}(t)}{dt} = -\mathbf{A}(t)\mathbf{m}(t) + \mathbf{b}(t), \quad (8)$$

and

$$\frac{d\mathbf{S}(t)}{dt} = -\mathbf{A}(t)\mathbf{S}(t) - \mathbf{S}(t)\mathbf{A}^\top(t) + 2\mathbf{D}. \quad (9)$$

Accordingly, we find optimal $(\mathbf{A}(t), \mathbf{b}(t))$, $(\mathbf{m}(t), \mathbf{S}(t))$ by looking for the stationary points of the following Lagrangian

$$\begin{aligned} \mathcal{L} = \int & \{E - \text{tr}\{\Psi(\frac{d\mathbf{S}}{dt} + \mathbf{A}\mathbf{S} + \mathbf{S}\mathbf{A}^\top - 2\mathbf{D})\} \\ & - \lambda^\top (\frac{d\mathbf{m}}{dt} + \mathbf{A}\mathbf{m}) - \mathbf{b}\} dt \end{aligned}$$

where $\Psi(t) \in \mathcal{R}^{d \times d}$ and $\lambda(t) \in \mathcal{R}^d$ are Lagrange multipliers. By definition, $\Psi(T) = 0$ and $\lambda(T) = 0$.

By taking the derivatives of \mathcal{L} with respect to \mathbf{m} , \mathbf{S} , \mathbf{A} and \mathbf{b} , we obtain a system of so-called Euler-Lagrange (EL) equations which the optimal \mathbf{m} , \mathbf{S} , \mathbf{A} , \mathbf{b} , Ψ and λ should fulfil [6]. Hence, the non-linear smoothing problem is reduced to solving a system of first-order differential equations. The equation system is solved iteratively.

We start with an initial guess of \mathbf{m} , \mathbf{S} , \mathbf{A} , \mathbf{b} , Ψ and λ . First, we compute $\mathbf{m}(t)$ and $\mathbf{S}(t)$ by performing standard Gaussian Process regression [17]. Then, we set $\Psi(t) = 0$ and $\lambda(t) = 0$ for all t . Finally, \mathbf{A} and \mathbf{b} are initialised using

$$\mathbf{A}(t) = \left\langle \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right\rangle_{q_t} + \mathbf{D}\Psi(t) \quad (10)$$

$$\mathbf{b}(t) = \langle \mathbf{f}(\mathbf{x}) \rangle_{q_t} + \mathbf{A}(t)\mathbf{m}(t) - 2\mathbf{D}\lambda(t). \quad (11)$$

Note that Eq. (10-11) are derived from the EL-equations.

At iteration i , we first update \mathbf{m} and \mathbf{S} by solving Eq. (8-9) forward in time where \mathbf{A}^i and \mathbf{b}^i are used. Next, Ψ and λ are updated by solving Eq. (12-15) with final condition $\Psi(T) = 0$ and $\lambda(T) = 0$ where \mathbf{m}^{i+1} and \mathbf{S}^{i+1} are used. Note that the data are assimilated at this step. Between two successive observations, we update Ψ and λ by solving

$$\frac{d\Psi(t)}{dt} = 2\Psi(t)\mathbf{A}(t) - \frac{\partial E_{sde}}{\partial \mathbf{S}} \quad (12)$$

$$\frac{d\lambda(t)}{dt} = \mathbf{A}^\top(t)\lambda(t) - \frac{\partial E_{sde}}{\partial \mathbf{m}} \quad (13)$$

When there is an observation at t_{k_j} , $j = 1, \dots, M$, the following jump-conditions apply

$$\Psi(t_{k_j}^+) = \Psi(t_{k_j}^-) - \frac{1}{2}\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H} \quad (14)$$

$$\lambda(t_{k_j}^+) = \lambda(t_{k_j}^-) + \mathbf{H}^\top \mathbf{R}^{-1}(\mathbf{y}_j - \mathbf{H}\mathbf{m}(t_{k_j})). \quad (15)$$

We consider a linear measurement function which is represented by a matrix \mathbf{H} . Eq. (12-15) are derived from the EL-equations [6].

Finally, we compute $\mathbf{A}(t; \mathbf{m}^{i+1}, \mathbf{S}^{i+1}, \Psi^{i+1}, \lambda^{i+1})$ and $\mathbf{b}(t; \mathbf{m}^{i+1}, \mathbf{S}^{i+1}, \Psi^{i+1}, \lambda^{i+1})$ by using Eq. (10-11). To keep the algorithm stable, the update of $\mathbf{A}(t)$ and $\mathbf{b}(t)$ is done by

$$\begin{aligned} \mathbf{A}^{i+1}(t) &= \mathbf{A}^i(t) - \omega\{\mathbf{A}^i(t) - \mathbf{A}(t; \cdot, \cdot, \cdot, \cdot)\} \\ \mathbf{b}^{i+1}(t) &= \mathbf{b}^i(t) - \omega\{\mathbf{b}^i(t) - \mathbf{b}(t; \cdot, \cdot, \cdot, \cdot)\} \end{aligned}$$

where $0 < \omega < 1$. The iteration stops when \mathcal{L} has converged to a constant value.

5. NUMERICAL EXPERIMENTS

The MCMC and variational algorithms are compared on a double-well potential system which is given by

$$\dot{x}(t) = f(x(t)) + \kappa \cdot \xi(t), \quad (16)$$

where

$$f(x) = 4x(1 - x^2)$$

and $\xi(t)$ is white-noise [18]. The parameter κ corresponds to $(2\mathbf{D})^{\frac{1}{2}}$ in Eq. 1 and determines the strength of random fluctuations within the system. This system has two stable states, namely $x = +1$ and $x = -1$. However, random fluctuations could cause a transition of the system from one stable state into another. The average time needed for the occurrence of such events is called the exit time [18]. In this study, we set $\kappa = 0.5$ and the corresponding exit time is about 4000 time units [18]. This provides us some prior knowledge on initial states.

In the numerical experiments, we consider a smoothing window ranging from $t = 0$ to $t = 12.0$. Further, we assume that states x can be observed directly, which makes \mathbf{h} an identity function. Within the smoothing window, we generate two data sets, say A and B , from a sample path which was considered in [18] and [19]. The variance R of measurement errors are 0.04 and 0.36, respectively. Each data set consists of seven data points which are "measured" at times $t_{k_1} = 1.0, \dots, t_{k_M} = 7.0$. Although multiple data sets are generated and analysed for each of those two R -values and for some intermediate R -values, the results of data set A and B are representative and chosen for illustration.

For the MCMC method, Eq. 16 is discretized with time increment $\delta t = 0.1$. The prior on initial states is set to a Gaussian density with mean at $x = +1$ and variance equal to 0.05. This choice is strongly based on our prior knowledge of the system. The tuning parameters of Hybrid Monte Carlo are chosen as follows: $J = 2$, $\delta\tau = 0.005$ and $\alpha = 0.02$. The use of the preconditioning matrix \mathbf{Z} keeps the necessary J small, which makes the simulation computationally more efficient. However, the multiplication of the matrix \mathbf{Z} with various vectors would cost extra computational time. Because of the circulant property of \mathbf{Z} , this part of computational burden is reduced.

For each of 2 data sets, we run a Markov chain of length 5,000,000 and subsample from this chain with sampling interval equal to 1,000. The first 1,000 samples after subsampling are discarded as burn-in period. It turns out that it is insufficient to determine burn-in only by monitoring a summary statistic like energy $\hat{\mathcal{H}}$. On the contrary, one has to monitor the traces of state x at different time points. Particularly, those time points must be chosen from different phases of the smoothing window, for example, transition phases, stationary phases, and the phase before/after the first/last observation.

For the variational GP approximation method, Eq. 16 is discretized with time increment $\delta t = 0.01$. The only tuning

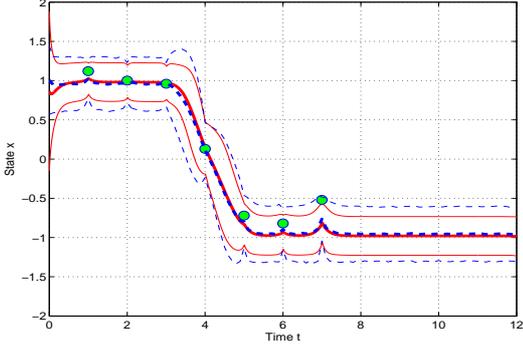


Fig. 1. Comparison of mean-path and marginal-variance estimates between the MCMC- (dashed) and variational (solid) method with a double-well potential system. Filled circles represent 7 observations from data set A , with measurement noise variance equal to 0.04. The mean paths are displayed by thick lines, while each pair of thin lines indicates an envelope of mean path with $2 \times$ standard deviation.

parameter ω is set to 0.15 as in [6]. The number of iterations required for the convergence of VGPA may increase when we extend the smoothing window or add more measurement noise. This is because of a poor estimation of the initial path by standard GP regression.

In Fig. 1 and Fig. 2, the estimates of both mean path and marginal variance are displayed for data set A and B , respectively. In each figure, the results of VGPA (solid lines) are compared with the MCMC results (dashed lines).

For the data sets with relatively small measurement noise, the estimated mean paths of both methods agree with each other very well whereas the estimated conditional variance of VGPA is biased overall but only slightly smaller than that of MCMC. It is also seen that the estimated mean path is slightly biased towards zero during both stationary phases. This can be explained by the fact that although the posterior of x has a distinct mode at $x = +1.0$ before the transition or $x = -1.0$ after the transition, the mode at another stable state is not vanishing. Note that this bias is observed for both VGPA and MCMC. But the effect shown in the MCMC results is slightly more significant than in the VGPA results. Moreover, we see that the mean path stays around the mode at -1.0 long after the last observation. This is also in accordance with the large exit time of the system we consider.

A small dip of mean paths is evident in the results of the VGPA smoother when we look into the initial period of the smoothing window. This is accompanied with large marginal variance $S(0)$. To explain this observation, we run MCMC simulations with increasingly larger prior variance of initial states. As expected, the posterior variance of x_0

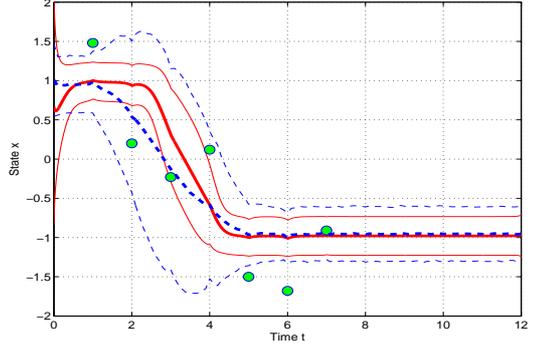


Fig. 2. The same as in Fig. 1 but with data set B (measurement noise variance = 0.36).

increases with its prior variance. Further, it turns out that a similar dip of mean paths appears when the prior variance becomes sufficiently large. It can be understood as follows: Without any data, a double-well systems does show a bimodal probabilistic structure. With a posterior mean of x_0 close to $+1$ and a large value of its posterior variance, the mean path could be further biased towards zero in the initial period where the first observation has little influence. Note that the approximate posterior variance $S(0)$ is not optimised, but held fixed.

Finally, we turn our attention to data set B with very large measurement noise. Note that it is difficult to identify where the transition starts by visual inspection of the data themselves. In contrast, this is possible with data set A . From Fig. 2, we can see that there is significant difference both in mean path and in marginal variance between the MCMC and VGPA smoother, particularly in the period before $t = 5.0$. Due to the ambiguity shown by the data between $t = 2.0$ and $t = 4.0$, the MCMC sampler seems to be exploring the bimodal structure of the posterior distribution. In contrast, the approximate posterior of the VGPA method is fixed to one particular mode at any time. This may explain the difference in mean path between two methods and a significant underestimation of marginal variance for the VGPA smoother.

6. DISCUSSION

By comparing with Markov Chain Monte Carlo, we evaluate a variational method for non-smoothing which was recently proposed in [6]. Both methods are tested on a double-well potential systems. Two data sets with different measurement noise are used to determine the strengths and weaknesses of the novel smoother.

Our investigation is based on the fact that the MCMC

methods provide an exact inference tool for comparison. For data sets with small or moderate measurement noise, it turns out that the VGPA method does produce a very accurate estimate of mean path while the marginal variance is slightly under-estimated. As expected, the variational method is computationally more efficient than MCMC. For the cases we have investigated in this paper, the CPU time needed for the convergence of the VGPA method is roughly equal to that for running 50,000 updating steps in MCMC. Regarding other approximation-based smoothers, it has been reported that Ensemble Kalman smoother fails to reconstruct the transition of a double-well system accurately from a sparse data set [19]. As stated in [19], the failure is due to the fact that in KF and EnKF the propagated states are corrected by a linear interpolation scheme when new data are assimilated.

However, the weakness of the VGPA method is also evident when the ambiguity due to noisy data becomes significant. As many other variational approximation methods, the variational smoother is not good at exploring the multimodal structure of some probability measures. In this paper, the role of prior on initial state is also investigated. It turns out that the estimates of mean path could be biased in the initial phase of the smoothing window, where the first observation has little influence, unless the prior on initial states is incorporated into the variational smoother, which we will shortly do.

Many MCMC algorithms suffer from poor mixing when high-dimensional stochastic complex systems are concerned. Development of efficient MCMC algorithms is always a challenging task. A combination of variational approximation methods and sampling methods would offer a new promising direction to improve the efficiency of MCMC algorithms, which we will shortly investigate.

7. REFERENCES

- [1] H. J. Kushner, “Dynamical equations for optimal filter,” *J. Diff. Eq.*, vol. 3, pp. 179–190, 1967.
- [2] R. L. Stratonovich, “Conditional markov processes,” *Theor. Prob. Appl.*, vol. 5, pp. 156–178, 1960.
- [3] E. Pardoux, “équations du filtrage non linéaire de la prédiction et du lissage,” *Stochastics*, vol. 6, pp. 193–231, 1982.
- [4] R. E. Kalman and R. S. Bucy, “New results in linear filtering and prediction theory,” *J. Basic Eng. D*, vol. 83, pp. 95–108, 1961.
- [5] R. H. Shumway and D. S. Stoffer, *Time Series Analysis and Its Applications*, Springer, New York, 2000.
- [6] C. Archambeau, D. Cornford, M. Opper, and J. Shawe-Taylor, “Gaussian process approximations of stochastic differential equations,” *Journal of Machine Learning Research*, vol. 1, pp. 1–16, 2007.
- [7] P. E. Klöden and Eckhard Platen, *Numerical Solution of Stochastic Differential Equations*, Springer-Verlag, Berlin, 1992.
- [8] C. Andrieu, N. De Freitas, A. Doucet, and M. I. Jordan, “An introduction to mcmc for machine learning,” *Machine Learning*, vol. 50, pp. 5–43, 2003.
- [9] M. Hürzler, *Statistical Methods for General State-Space Models*, Ph.D. thesis, ETH Zürich, 1998.
- [10] F. J. Alexander, G. L. Eyink, and J. M. Restrepo, “Accelerated monte carlo for optimal estimation of time series,” *Journal of Statistical Physics*, vol. 119, pp. 1331–1345, 2005.
- [11] A. Gelb, *Applied Optimal Estimation*, The MIT Press, Cambridge MA, 1974.
- [12] G. Evensen, “Sequential data assimilation with a non-linear quasi-geostrophic model using monte carlo methods to forecast error statistics,” *J. Geophys. Res.*, vol. 99, pp. 10143–10162, 1994.
- [13] G. Kitagawa, “Non-gaussian state space modelling of non-stationary time series,” *J. Am. Statist. Assoc.*, vol. 82, pp. 503–514, 1987.
- [14] S. Kullback and R. A. Leibler, “On information and sufficiency,” *Annal of Mathematical Statistics*, vol. 22, pp. 79–86, 1951.
- [15] D. Saad and M. Opper, Eds., *Advanced Mean Field Methods*, chapter Tutorial on variational approximation methods, pp. 207–219, The MIT Press, Cambridge MA, 2001.
- [16] P. Del Moral D. Crisan and T. J. Lyons, “Interacting particle systems approximations of the kushner-stratonovich equation,” *Advances in Applied Probability*, vol. 31, pp. 819–838, 1999.
- [17] C. E. Rasmussen and C. K. I. Williams, *Gaussian Process for Machine Learning*, The MIT Press, Cambridge MA, 2006.
- [18] R.N. Miller, E. F. Carter, and S. T. Blue, “Data assimilation into nonlinear stochastic models,” *Tellus A*, vol. 51, pp. 167–194, 1999.
- [19] G. L. Eyink, J. M. Restrepo, and F.J. Alexander, “A mean-field approximation in data assimilation for nonlinear dynamics,” *Physica D*, vol. 194, pp. 347–368, 2004.