

Gaussian filtering and variational approximations for Bayesian smoothing in continuous-discrete stochastic dynamic systems

Juha Ala-Luhtala^a, Simo Särkkä^b, Robert Piché^a

^a*Tampere University of Technology, Finland*

^b*Aalto University, Finland*

Abstract

The Bayesian smoothing equations are generally intractable for systems described by nonlinear stochastic differential equations and discrete-time measurements. Gaussian approximations are a computationally efficient way to approximate the true smoothing distribution. In this work, we present a comparison between two Gaussian approximation methods. The Gaussian filtering based Gaussian smoother uses a Gaussian approximation for the filtering distribution to form an approximation for the smoothing distribution. The variational Gaussian smoother is based on minimizing the Kullback–Leibler divergence of the approximate smoothing distribution with respect to the true distribution. The results suggest that for highly nonlinear systems, the variational Gaussian smoother can be used to iteratively improve the Gaussian filtering based smoothing solution. We also present linearization and sigma-point methods to approximate the intractable Gaussian expectations in the Variational Gaussian smoothing equations. In addition, we extend the variational Gaussian smoother for certain class of systems with singular diffusion matrix.

Keywords: Bayesian smoothing; Gaussian approximation; Variational inference

1. Introduction

The continuous-discrete system refers to a system whose process dynamics are governed by a continuous-time stochastic differential equation (SDE) and whose measurements are taken at discrete time instants. Bayesian filtering and smoothing equations give the solution to the problem of estimating the state of the system from the noisy measurements. Computing the filtering and smoothing distributions involves solving the related partial differential equations [1, 2], and is only tractable for linear-Gaussian systems (and some other special cases [3]). In this paper, we consider two different approaches for computing a Gaussian approximation for the smoothing distribution.

The Gaussian approximation for the filtering and smoothing equations is well known in the literature. The first Gaussian approximations were based on linearization using the Taylor series based methods [4, 1]. The Taylor series based methods can be seen as a special case of the more general Gaussian filtering and smoothing framework, where different filters and smoothers arise based on the numerical method for computing the Gaussian expectations [5, 6]. A different approach for Gaussian smoothing was considered in [7, 8, 9], where a Gaussian approximation is sought by approximating the stochastic process giving the smoothed distribution with a linear process. The method is based on the fixed-form variational Bayes approximation [10] and minimizes the Kullback–Leibler divergence of the approximate distribution with respect to the true distribution.

The variational Gaussian approximation is considered further in [11, 12, 13, 14]. Shen *et al.* [11] compared the variational approximation to a Monte Carlo Markov Chain (MCMC) solution for the one dimensional double well system and found that the variational method performed comparatively to the MCMC solution when the uncertainties in the measurements were not so large as to cause the true posterior to be multimodal. In [12] the variational Gaussian smoothing solution is used as a proposal distribution in an MCMC method to improve the efficiency of the algorithm. The variational MCMC method was found to outperform the hybrid Monte Carlo method for sparsely observed diffusion processes. Vrettas *et al.* [13, 14] considered approximating the variational Gaussian smoothing equations using a radial basis function representation for the variational parameters. By imposing a certain structure for the variational parameter functions, the overall number of parameters to be optimized can be reduced.

The variational Gaussian smoothing equations derived by Archambeau *et al.* [7, 8, 9] require a non-singular effective diffusion matrix. Our first goal is to extend the variational Gaussian approximation to a certain class of singular models by considering an alternative derivation based on Girsanov’s theorem.

Another problem in the variational Gaussian algorithm is the need to compute Gaussian expectations over nonlinear functions. Previous works on the variational approximation [7, 8, 11, 13, 14] have not presented details for computing the Gaussian expectations for general nonlinear systems. Our second goal is to extend the variational Gaussian smoothing method for general nonlinear systems by considering numerical approximations for the Gaussian expectations. The treatment is similar to the one in [6], where Taylor series based linearization, cubature and unscented transform based sigma-point methods and Gauss–Hermite quadrature were used to compute Gaussian expectations in the Gaussian smoothers based on the classical Gaussian filtering framework.

In addition, we provide in this paper a comparison between the variational Gaussian approximation and the Gaussian smoothers presented in [6]. Using a suitable change of variables the Gaussian smoothing equations can be converted to a variational form similar to the variational Gaussian smoothing equations. The computation of the Gaussian filtering based smoothing solution is numerically stable and provides good initial conditions for the variables in the variational Gaussian smoothing algorithm. This might help to overcome the problems

reported by Vrettas *et al.* [13] in the initialization of the variational Gaussian algorithm for high dimensional systems. Also, we study if the variational Gaussian smoother can be used to iteratively improve the results from the Gaussian filtering based smoother.

The organization of this paper is as follows. First we present the variational Gaussian smoother and the Gaussian smoother based on the Gaussian filtering framework for the continuous-discrete system. Here we extend the variational Gaussian smoothing equations for a certain class of singular models. Next we compare theoretically the Gaussian smoothers by presenting a conversion of the Gaussian filtering based smoothing equations to the variational form by a suitable change of variables. This also provides the initial values for the variational Gaussian smoothing algorithm. The problem of numerically computing the Gaussian expectations in the variational Gaussian smoothing equations is treated in the next section. We present Taylor series linearization and sigma-point methods to approximate the Gaussian expectations in the variational Gaussian smoothing equations. The paper concludes with two synthetic-data examples that are used to compare the Gaussian smoothers, and also provides comparison of the different numerical methods for the computation of the Gaussian expectations.

1.1. Problem statement

The continuous-discrete system considered in this paper is given by

$$dx = f(x, t) dt + L(t)d\beta(t), \quad (1)$$

$$y_k = h_k(x(t_k)) + v_k, \quad (2)$$

where $x(t)$ is the state, $f(x(t), t)$ is the drift term, and $\beta(t)$ is a Brownian motion stochastic process with diffusion matrix $Q(t)$. The effective diffusion matrix for the process is given by

$$\Sigma(t) = L(t)Q(t)L^T(t). \quad (3)$$

The initial conditions are assumed to be normally distributed $x(t_0) \sim N(m_0, P_0)$. The measurement noise $\{v_k\}$ is a zero mean Gaussian white noise sequence with covariance matrix R_k . The measurement noise $\{v_k\}$, process noise $\beta(t)$ and initial conditions x_0 are assumed to be mutually independent.

Let y_1, \dots, y_K be the measurements taken at discrete time instants t_1, \dots, t_K . The solution to the Bayesian smoothing problem is the posterior distribution

$$p(x(t) | y_1, \dots, y_K), \quad t \in [t_0, t_k]. \quad (4)$$

In this paper, we concentrate on Gaussian approximations for the smoothing distribution. That is, the smoothing distribution is approximated as

$$p(x(t) | y_1, \dots, y_K) \approx N(x(t) | m(t), P(t)), \quad (5)$$

where $m(t)$ is the mean function and $P(t)$ is autocovariance $P(t, t')$ at $t = t'$ for the approximating distribution. The smoothing problem now reduces to finding the expressions for the mean and covariance functions.

2. Gaussian smoothing for continuous-discrete systems

2.1. Variational Gaussian approximation

The variational Gaussian approximation for the continuous-discrete smoothing problem was derived by Archambeau *et al.* [7, 8, 9]. The method is based on approximating the smoothing process with a linear process

$$dx = [-A(t) + b(t)]dt + \sqrt{\Sigma(t)}d\beta(t), \quad (6)$$

where $A(t)$ and $b(t)$ are parameters of the approximation and $\beta(t)$ is a Brownian stochastic process with identity diffusion matrix. The solution to the linear SDE (6) is a Gaussian process. The marginal density at each time is given by $q(x(t)) = N(x(t) | m(t), P(t))$, where the mean and covariance are computed from the ordinary differential equations

$$\frac{d}{dt}m(t) = -A(t)m(t) + b(t) \quad (7)$$

$$\frac{d}{dt}P(t) = -A(t)P(t) - P(t)A^T(t) + \Sigma(t). \quad (8)$$

The parameters $A(t)$ and $b(t)$ are computed by minimizing the Kullback–Leibler (KL) divergence of the probability law \mathbb{Q}_X of the approximating process with respect to the probability law $\mathbb{P}_{X|Y}$ of the true smoothing process. The KL-divergence is given by [7, 8, 9]

$$\text{KL}(\mathbb{Q}_X \parallel \mathbb{P}_{X|Y}) = \int_{t_0}^{t_K} \mathbb{E}_q \left[e(x(t), t) + \sum_{k=1}^K u_k(x(t))\delta(t - t_k) \right] dt, \quad (9)$$

where

$$e(x(t), t) = \frac{1}{2}[f(x(t), t) + A(t)x(t) - b(t)]^T \Sigma^{-1}(t)[f(x(t), t) + A(t)x(t) - b(t)], \quad (10)$$

$$u_k(x(t)) = \frac{1}{2}[y_k - h_k(x(t))]^T R_k^{-1}[y_k - h_k(x(t))]. \quad (11)$$

The expectations $\mathbb{E}_q[\cdot]$ are computed with respect to the marginal distribution $q(x(t))$ of the approximating process.

The KL-divergence is minimised using the mean and covariance differential equations (7) and (8) as constraints. Introducing Lagrange multiplier functions $\lambda(t)$ and $\Psi(t)$ for the constraints, the objective function is given by

$$\begin{aligned} \mathcal{F}(A, b, m, P) = \int_{t_0}^{t_K} \left\{ \mathbb{E}_q \left[e(x, t) + \sum_{k=1}^K u_k(x)\delta(t - t_k) \right] \right. \\ \left. - \lambda(t)^T \left(\frac{dm}{dt} + A(t)m(t) - b(t) \right) \right. \\ \left. - \text{tr} \left[\Psi(t) \left(\frac{dP}{dt} + A(t)P(t) + P(t)A^T(t) - \Sigma(t) \right) \right] \right\} dt. \quad (12) \end{aligned}$$

For notational convenience, we use $x = x(t)$. The Euler–Lagrange equations for the optimal $A(t)$ and $b(t)$ can then be written as [7,8]

$$\frac{d}{dt}\lambda(t) = A^T(t)\lambda(t) - \nabla_m \mathbb{E}_q[e(x, t)] \quad (13)$$

$$\frac{d}{dt}\Psi(t) = \Psi(t)A(t) + A^T(t)\Psi(t) - \nabla_P \mathbb{E}_q[e(x, t)] \quad (14)$$

$$A(t) = -\mathbb{E}_q[F_x(x, t)] + 2\Sigma(t)\Psi(t) \quad (15)$$

$$b(t) = \mathbb{E}_q[f(x, t)] + A(t)m(t) - \Sigma(t)\lambda(t). \quad (16)$$

At observation times, the Lagrange multiplier functions satisfy

$$\lambda(t_k^+) = \lambda(t_k^-) + \nabla_m \mathbb{E}_q[u_k(x)], \quad (17)$$

$$\Psi(t_k^+) = \Psi(t_k^-) + \nabla_P \mathbb{E}_q[u_k(x)]. \quad (18)$$

To find a solution satisfying the Euler–Lagrange equations, we propose here a slight modification of the iterative algorithm given by Archambeau *et al.* [7, 8]. Given the previous estimates $A^{(k)}(t)$ and $b^{(k)}(t)$, the mean and covariance differential equations (7) and (8) are solved forward in time from t_0 which gives $m^{(k+1)}(t)$ and $P^{(k+1)}(t)$. Using these, the Lagrange differential equations (13) and (14) are then solved backward in time from t_K which gives $\lambda^{(k+1)}(t)$ and $\Psi^{(k+1)}(t)$. New estimates $A^{(k+1)}(t)$ and $b^{(k+1)}(t)$ are then computed by using a damped fixed-point update

$$A^{(k+1)}(t) = A^{(k)}(t) + \gamma_k \left(A(t) - A^{(k)}(t) \right) \quad (19)$$

$$b^{(k+1)}(t) = b^{(k)}(t) + \gamma_k \left(b(t) - b^{(k)}(t) \right), \quad (20)$$

where $A(t)$ and $b(t)$ are computed using Equations (15) and (16) respectively, and $\gamma_k \in (0, 1)$ is a damping parameter that is used to prevent numerical instabilities caused by too large updates. Instead of using constant damping parameter as in [7], we propose to select the parameter γ_k using backtracking line search so that the objective function and therefore also the KL-divergence is reduced at each time step. In a computer implementation of the variational Gaussian smoother, the values of the functions $A(t)$ and $b(t)$ are computed and stored at discrete time points. The differential equations can be solved using any standard numerical solver such as Euler or Runge–Kutta methods.

2.2. Variational Gaussian smoothing for singular models

Note that to compute the term $e(x, t)$, a nonsingular effective diffusion matrix $\Sigma(t)$ is required. We extend here the variational Gaussian smoother to singular models of the form

$$\frac{dx_1}{dt} = F_1(t)x \quad (21)$$

$$dx_2 = f_2(x, t) dt + \sqrt{\Sigma_2(t)} d\beta, \quad (22)$$

where $\Sigma_2(t)$ is a nonsingular diffusion matrix.

Denote by n_1 and n_2 the dimensions of the state vectors x_1 and x_2 respectively. We now seek an approximate smoothing process of the form

$$\frac{dx_1}{dt} = F_1(t)x \quad (23)$$

$$dx_2 = (-A(t)x + b(t))dt + \sqrt{\Sigma_2(t)}d\beta, \quad (24)$$

where $A(t)$ and $b(t)$ are the variational parameters. This gives a Gaussian process, where mean $m(t)$ and covariance $P(t)$ follow the differential equations

$$\frac{d}{dt}m(t) = -\tilde{A}(t)m(t) + \tilde{b}(t) \quad (25)$$

$$\frac{d}{dt}P(t) = -\tilde{A}(t)P(t) - P(t)\tilde{A}^T(t) + \Sigma(t), \quad (26)$$

with

$$\tilde{A}(t) = \begin{bmatrix} -F_1(t) \\ A(t) \end{bmatrix}, \quad \tilde{b}(t) = \begin{bmatrix} 0_{n_1 \times 1} \\ b(t) \end{bmatrix}, \quad \Sigma(t) = \begin{bmatrix} 0_{n_1 \times n_1} & 0_{n_1 \times n_2} \\ 0_{n_2 \times n_1} & \Sigma_2(t) \end{bmatrix}. \quad (27)$$

The KL-divergence term for the singular system can be computed using Girsanov's theorem and is given by (see [15] and Appendix A for details)

$$\text{KL}(\mathbb{Q}_X \parallel \mathbb{P}_{X|Y}) = \int_{t_0}^{t_K} \mathbb{E}_q \left[e(x, t) + \sum_{k=1}^K u_k(x) \delta(t - t_k) \right] dt, \quad (28)$$

where $u_k(x)$ is given by eq. (11) and

$$e(x, t) = \frac{1}{2} [f_2(x, t) + A(t)x(t) - b(t)]^T \Sigma_2^{-1}(t) [f_2(x, t) + A(t)x(t) - b(t)]. \quad (29)$$

The objective function is given by

$$\begin{aligned} \mathcal{F}(A, b, m, P) = & \int_{t_0}^{t_K} \left\{ \mathbb{E}_q \left[e(x, t) + \sum_{k=1}^K u_k(x) \delta(t - t_k) \right] \right. \\ & - \lambda(t)^T \left(\frac{dm}{dt} + \tilde{A}(t)m(t) - \tilde{b}(t) \right) \\ & \left. - \text{tr} \left[\Psi(t) \left(\frac{dP}{dt} + \tilde{A}(t)P(t) + P(t)\tilde{A}^T(t) - \Sigma(t) \right) \right] \right\} dt \quad (30) \end{aligned}$$

and the solution satisfies

$$\frac{d}{dt}\lambda(t) = \tilde{A}^T(t)\lambda(t) - \nabla_m \mathbb{E}_q[e(x, t)] \quad (31)$$

$$\frac{d}{dt}\Psi(t) = \Psi(t)\tilde{A}(t) + \tilde{A}^T(t)\Psi(t) - \nabla_P \mathbb{E}_q[e(x, t)] \quad (32)$$

$$A(t) = -\mathbb{E}_q[F_{2,x}(x, t)] + 2\Sigma_2(t)M\Psi(t) \quad (33)$$

$$b(t) = \mathbb{E}_q[f_2(x, t)] + A(t)m(t) - \Sigma_2(t)M\lambda(t), \quad (34)$$

where $F_{2,x}(x, t)$ is the Jacobian of $f_2(x, t)$ and the matrix M , which selects the relevant part of the Lagrange parameter functions, is given by

$$M = \begin{bmatrix} 0_{n_1 \times n_2} & I_{n_2 \times n_2} \end{bmatrix}. \quad (35)$$

At observation times, the Lagrange multipliers satisfy the boundary conditions given by eqs. (17)-(18).

2.3. Gaussian filtering based Gaussian smoother

The Gaussian filtering based Gaussian smoother uses the Gaussian approximation for the filtering distribution to form the Gaussian approximation for the smoothing distribution [6].

The Gaussian filtering (or Gaussian assumed density filtering) approach is well known in the literature (see e.g. [16, 1, 17]) and uses the approximation

$$p(x(t) | y_1, \dots, y_k) \approx N(x(t) | m_f(t), P_f(t)), \quad (36)$$

where $p(x(t) | y_1, \dots, y_k)$ is the filtering distribution. The mean and covariance function are recursively computed using the following prediction and update steps. In the prediction step the mean and covariance functions are propagated from time t_{k-1} to time t_k using

$$\frac{dm_f}{dt} = \mathbb{E}_f[f(x, t)] \quad (37)$$

$$\frac{dP_f}{dt} = \mathbb{E}_f[(x - m_f)f(x, t)^T] + \mathbb{E}_f[f(x, t)(x - m_f)^T] + \Sigma(t). \quad (38)$$

In the update step, the information from the latest measurement y_k is used to update the predicted estimates $m(t_k^-)$ and $P(t_k^-)$ using equations

$$S_k = \mathbb{E}_f [(h_k(x) - \mathbb{E}_f[h_k(x)])(h_k(x) - \mathbb{E}_f[h_k(x)])^T] + R_k \quad (39)$$

$$K_k = \mathbb{E}_f [(x - m_f(t_k^-))(h_k(x) - \mathbb{E}_f[h_k(x)])^T] S_k^{-1} \quad (40)$$

$$m_f(t_k) = m_f(t_k^-) + K_k (y_k - \mathbb{E}_f[h_k(x)]) \quad (41)$$

$$P_f(t_k) = P_f(t_k^-) - K_k S_k K_k^T. \quad (42)$$

Särkkä and Sarmavuori [6] extend the general Gaussian filtering ideas also to smoothing problems and derive three types of Gaussian smoothers labeled as type I, type II and type III. The type I smoother is derived by using the Gaussian approximation for the filtering distribution to approximate the exact partial differential equations for the smoothed mean and covariance. The type II smoother is derived by discretizing the dynamic model, applying the discrete-time smoothing equations and then taking the limit as the discretization time approaches zero. Formulating the type II smoothing equations into a computationally efficient form gives the type III smoother.

The type I smoothing equations are numerically quite sensitive, which causes the type I smoother to diverge quite often compared to the type II and type

III smoothers. Also, the type I smoother is computationally more demanding and was not found to be clearly better than the type II and type III smoothers in the synthetic-data example considered in [6]. For this reason we decided to concentrate on the type II and type III smoothers in this paper.

The type II smoothing equations are given by

$$\frac{dm_s}{dt} = \mathbb{E}_f[f(x, t)] + [\mathbb{E}_f[f(x, t)(x - m_f)^T] + \Sigma(t)] P_f^{-1}(m_s - m_f) \quad (43)$$

$$\begin{aligned} \frac{dP_s}{dt} &= (\mathbb{E}_f[f(x, t)(x - m_f)^T] + \Sigma) P_f^{-1} P_s \\ &\quad + P_s P_f^{-1} (\mathbb{E}_f[f(x, t)(x - m)^T]^T + \Sigma(t)) - \Sigma(t). \end{aligned} \quad (44)$$

If the Jacobian $F_x(x, t)$ of $f(x, t)$ is available, we can alternatively use

$$\mathbb{E}_f[f(x, t)(x - m_f)^T] = \mathbb{E}_f[F_x(x, t)] P_f \quad (45)$$

in the filtering and smoothing equations. The expectations in the smoothing equations are with respect to the filtering density, which means that they can be computed already during the filtering stage. This is used in the type III smoothing equations, which reformulate the filtering and type II smoothing equations so that no expectations need to be computed during the smoothing stage (see [6] for details).

2.4. Differences between the Gaussian smoothers

In this section, we study the differences between the variational and Gaussian filtering based Gaussian smoothers. First we introduce a change of variables that converts the type II Gaussian smoothing equations to a form similar to the variational Gaussian smoothing equations. This conversion is similar to the results from Rauch, Tung and Striebel [18], where it was shown that the Rauch–Tung–Striebel smoothing equations are formally equivalent to the smoothing equations presented by Bryson and Frazier [4].

The conversion is achieved by using the change of variables

$$\lambda(t) = -P_f^{-1}(t)(m_s(t) - m_f(t)), \quad \Psi(t) = -\frac{1}{2}(P_f^{-1}(t) - P_s^{-1}(t)). \quad (46)$$

Inserting $\lambda(t)$ and $\Psi(t)$ to the Equations (43)-(44) and computing the time derivatives of $\lambda(t)$ and $\Psi(t)$ gives (see Appendix B)

$$\frac{d}{dt} m_s(t) = -A(t)m_s(t) + b(t) \quad (47)$$

$$\frac{d}{dt} P_s(t) = -A(t)P(t) - P(t)A^T(t) + \Sigma(t) \quad (48)$$

$$\frac{d}{dt} \lambda(t) = A^T(t)\lambda(t) - 2\Psi(t)\Sigma(t)\lambda(t) \quad (49)$$

$$\frac{d}{dt} \Psi(t) = \Psi(t)A(t) + A^T(t)\Psi(t) - 2\Psi(t)\Sigma(t)\Psi(t), \quad (50)$$

where

$$A(t) = -\mathbb{E}_f[F_x(x, t)] + 2\Sigma(t)\Psi(t) \quad (51)$$

$$b(t) = \mathbb{E}_f[f(x, t)] + \mathbb{E}_f[F_x(x, t)](m_s(t) - m_f(t)) + A(t)m_s(t) - \Sigma(t)\Psi(t). \quad (52)$$

For linear measurement function $h_k(x) = H_k x$, the measurement update for $\lambda(t)$ and $\Psi(t)$ in the variational form of the type II smoother is the same as in the variational Gaussian smoother and is given by

$$\lambda(t_k^+) = \lambda(t_k^-) + H_k^T R_k^{-1}(m_s(t_k) - y_k) \quad (53)$$

$$\Psi(t_k^+) = \Psi(t_k^-) + \frac{1}{2} H_k^T R_k^{-1} H_k. \quad (54)$$

For nonlinear measurement function, the measurement updates for $\lambda(t)$ and $\Psi(t)$ are in general not equal to the variational Gaussian smoother update Equations (17) and (18).

Similarities to the variational Gaussian smoothing equations are evident from Equations (47)-(52). Note that Equations (51) and (52) for the parameters $A(t)$ and $b(t)$ are otherwise similar to the variational Gaussian smoothing Equations (15) and (16), but the function $f(x, t)$ is replaced with statistical linearization with respect to the filtering distribution:

$$f(x, t) \approx \mathbb{E}_f[f(x, t)] + \mathbb{E}_f[F_x(x, t)](x - m_f(t)). \quad (55)$$

Furthermore, using the statistical linearization (55) to approximate the gradients in Equations (13) and (14) gives

$$\nabla_m \mathbb{E}[e(x, t)] \approx 2\Psi(t)\Sigma(t)\lambda(t) \quad (56)$$

$$\nabla_P \mathbb{E}[e(x, t)] \approx 2\Psi(t)\Sigma(t)\Psi(t). \quad (57)$$

That is, the differential equations (49) and (50) can be seen as an approximation to the exact differential equations (13) and (14) in the variational Gaussian smoother. This suggests that for linear measurements, the type II Gaussian smoother can be seen to approximate the variational Gaussian smoother by using statistical linearization with respect to the filtering distribution.

We can use the type II Gaussian smoothing solution as an initial iterand for the variational parameters by computing $A^{(0)}(t)$ and $b^{(0)}(t)$ using equations (51) and (52), where $\lambda(t)$ and $\Psi(t)$ are given by (46). In [13] it was noted that the iterative solution of the variational Gaussian smoothing equations is sensitive to the initial values of the variational parameters $A(t)$ and $b(t)$. This way, the variational Gaussian smoother can be seen as an iterative way to improve the type II Gaussian smoothing solution. The benefit of using the variational Gaussian smoother to improve the type II smoother results is studied further in the synthetic-data examples.

3. Computation of Gaussian expectations in the variational Gaussian smoother

The Gaussian smoothers considered in this paper require computations of Gaussian expectations over arbitrary nonlinear functions. For some simple models these can be computed analytically, but for many models no analytical expression exists. The computation of Gaussian expectations for the Gaussian filtering based smoothers is presented in [6]. In this work we concentrate on computing the expectations in the variational Gaussian smoothing equations and present the extended, cubature, unscented and Gauss–Hermite forms of the variational Gaussian smoother.

Gaussian expectations need to be computed in the differential equations (13) and (14) for the functions $\lambda(t)$ and $\Psi(t)$, in Equations (15) and (16) for the variational parameter functions $A(t)$ and $b(t)$ and in the measurement update equations (17) and (18). In order to avoid computing derivatives of the drift function $f(x, t)$ and of the measurement function $h_k(x)$, the gradients with respect to $m(t)$ and $P(t)$ can be written in the form (see Appendix C)

$$\nabla_m \mathbb{E}[e(x, t)] = P^{-1} \mathbb{E}[e(x, t)(x - m)] \quad (58)$$

$$\nabla_P \mathbb{E}[e(x, t)] = \frac{1}{2} P^{-1} \mathbb{E}[e(x, t)(x - m)(x - m)^T] P^{-1} - \frac{1}{2} \mathbb{E}[e(x, t)] P^{-1}. \quad (59)$$

For the measurement updates, the gradients are computed similarly with $e(x, t)$ replaced by $u_k(x)$.

If the Jacobians $F_x(x, t)$ and $H_{k,x}(x)$ of $f(x, t)$ and $h_k(x)$ are available, the gradients can be alternatively written in the form

$$\nabla_m \mathbb{E}[e(x, t)] = \mathbb{E}[e_x(x, t)] \quad (60)$$

$$\nabla_P \mathbb{E}[e(x, t)] = \frac{1}{2} P^{-1} \mathbb{E}[e_x(x, t)(x - m)^T]^T. \quad (61)$$

where

$$e_x(x, t) = \nabla_x e(x, t) = [F_x(x, t) + A(t)]^T \Sigma^{-1}(t) [f(x, t) + A(t)x(t) - b(t)]. \quad (62)$$

To compute the gradients in the measurement update, $e_x(x, t)$ is replaced with

$$u_{k,x}(x) = \nabla_x u_k(x) = H_{k,x}^T(x) R_k^{-1} [h_k(x) - y_k]. \quad (63)$$

3.1. Taylor series based linearization

In the extended Kalman filter and smoother, the Gaussian expectations are computed by using a first order Taylor series linearization. Proceeding similarly, we use the following approximations for the extended variational Gaussian smoother:

$$f(x, t) \approx f(m, t) + F_x(m, t)(x(t) - m(t)) \quad (64)$$

$$h_k(x) \approx h_k(m) + H_{k,x}(m)(x(t) - m(t)), \quad (65)$$

where $F_x(m, t)$ and $H_{k,x}(m)$ are the Jacobians of $f(x, t)$ and $h_k(x)$ respectively evaluated at the current mean estimate $m(t)$. Using this approximation, the expectations in Equations (15) and (16) are given by

$$\mathbb{E}[f(x, t)] \approx f(m, t), \quad (66)$$

$$\mathbb{E}[F_x(x, t)] \approx F_x(m, t). \quad (67)$$

The expectations needed in the gradients are given by

$$\nabla_m \mathbb{E}[e(x, t)] \approx (F_x(m, t) + A(t))^T \Sigma^{-1}(t) (f(m, t) + A(t)m(t) - b(t)) \quad (68)$$

$$\nabla_P \mathbb{E}[e(x, t)] \approx \frac{1}{2} (F_x(m, t) + A(t))^T \Sigma^{-1}(t) (F_x(m, t) + A(t)) \quad (69)$$

and for the measurement update

$$\nabla_m \mathbb{E}[u_k(x)] \approx H_{k,x}^T(m) R_k^{-1}(t) (h_k(m) - y_k) \quad (70)$$

$$\nabla_P \mathbb{E}[u_k(x)] \approx \frac{1}{2} H_{k,x}^T(m) R_k^{-1}(t) H_{k,x}(m). \quad (71)$$

3.2. Sigma-point methods

The general sigma-point rule computes the Gaussian expectations using the approximation

$$\mathbb{E}[g(x, t)] \approx \sum_i W^{(i)} g(m + \sqrt{P} \xi_i, t), \quad (72)$$

where the weights $W^{(i)}$ and vectors ξ_i are chosen depending on the used sigma-point method. In this paper we consider the cubature, unscented and Gauss-Hermite sigma-point methods. The cubature method uses $2n$ sigma-points with vectors

$$\xi_i = \begin{cases} \sqrt{n} e_i, & i = 1, \dots, n \\ -\sqrt{n} e_{i-n}, & i = n+1, \dots, 2n \end{cases} \quad (73)$$

and weights $W^{(i)} = 1/(2n)$ for all $i = 1, \dots, 2n$. The unscented transform uses $2n+1$ sigma-points with vectors

$$\xi_0 = 0, \quad \xi_i = \begin{cases} \sqrt{\lambda + n} e_i, & i = 1, \dots, n \\ -\sqrt{\lambda + n} e_{i-n}, & i = n+1, \dots, 2n. \end{cases}, \quad (74)$$

where $\lambda = \alpha^2(n + \kappa) - n$ and α , β and κ are parameters of the method. The weights are defined to be

$$W_m^{(0)} = \frac{\lambda}{n + \lambda}, \quad W_m^{(i)} = \frac{1}{2(n + \lambda)}, \quad i = 1, \dots, 2n \quad (75)$$

$$W_c^{(0)} = W_m^{(0)} + 1 - \alpha^2 + \beta, \quad W_c^{(i)} = W_m^{(i)}, \quad i = 1, \dots, 2n. \quad (76)$$

The weights W_m are used in approximating the transformed mean and W_c in the covariance approximation. The cubature method is a special case of the unscented transform with parameters $\alpha = 1$, $\beta = \kappa = 0$. The Gauss-Hermite

integration uses s^n sigma-points, where s is a parameter that gives the order of the used Hermite polynomial. Details for computing the vectors ξ_i and weights $W^{(i)}$ are given in [17, 5].

The sigma-point approximation for expectations in Equations (15) and (16) is given by

$$\mathbb{E}[f(x, t)] \approx \sum_i W^{(i)} f(m + \sqrt{P} \xi_i, t) \quad (77)$$

$$\mathbb{E}[F_x(x, t)] = \mathbb{E}[f(x, t)(x - m)^T] P^{-1} \approx \sum_i W^{(i)} f(m + \sqrt{P} \xi_i, t) \xi_i^T \sqrt{P}^{-1}. \quad (78)$$

The general sigma-point approximations for the gradients are given by

$$\nabla_m \mathbb{E}[e(x, t)] \approx \sum W^{(i)} e(m + \sqrt{P} \xi_i, t) \sqrt{P}^{-1} \xi_i \quad (79)$$

$$\nabla_P \mathbb{E}[e(x, t)] \approx \frac{1}{2} \sum W^{(i)} e(m + \sqrt{P} \xi_i, t) \sqrt{P}^{-T} (\xi_i \xi_i^T - I) \sqrt{P}^{-1}. \quad (80)$$

The expectations needed in the observation updates (17) and (18) are computed with $e(x, t)$ replaced by $u_k(x)$.

The sigma-point approximation for the alternative forms (60) and (61) of the gradients are given by

$$\nabla_m \mathbb{E}[e(x, t)] \approx \sum W^{(i)} e_x(m + \sqrt{P} \xi_i), \quad (81)$$

$$\nabla_P \mathbb{E}[e(x, t)] \approx \frac{1}{2} \sum W^{(i)} \sqrt{P}^{-T} \xi_i e_x^T(m + \sqrt{P} \xi_i, t). \quad (82)$$

The measurement updates are computed similarly, with $e_x(x, t)$ replaced by $u_{k,x}(x)$. The cubature, unscented and Gauss–Hermite forms of the variational Gaussian smoother are then obtained by using the corresponding choice for the weights $W^{(i)}$ and vectors ξ_i .

For a linear drift function $f(x, t)$, the term $\mathbb{E}[e(x, t)(x - m)(x - m)^T]$ is a fourth order polynomial. For this reason, the sigma-point rules that are only accurate up to a third order monomial (cubature and unscented rule) give generally a poor approximation of this expectation. Therefore, for cubature and unscented sigma-point methods, the use of the alternative form given by Equations (81) and (82) is recommended.

4. Numerical experiments

The Gaussian smoothers are compared using two different synthetic-data experiments. The tests are done by running first the Gaussian filtering based Gaussian smoother (GFGS) and then using the result as initial conditions for the variational Gaussian smoother (VGS). The VGS iteration is terminated when the absolute change in the KL-divergence between successive iterations is less than 10^{-3} .

In the first experiment, a one dimensional double well system is used. The same system was also used to demonstrate the VGS in [7, 11, 13]. A 5-dimensional reentry problem is used for the second experiment. This system was used to test the continuous-discrete unscented Kalman filter in [19] and demonstrates the use of VGS for singular systems. Also, for this system the computation of the needed Gaussian expectations is not possible analytically.

The metrics used to compare the estimates given by the Gaussian smoothers are the root mean square error (RMSE), negative log-likelihood (NLL) and 95%-consistency. The RMSE and NLL are given by equations

$$\text{RMSE} = \sqrt{\frac{1}{t_K - t_0} \int_{t_0}^{t_K} \|x(t) - m(t)\|^2 dt} \quad (83)$$

$$\text{NLL} = \frac{1}{t_K - t_0} \int_{t_0}^{t_K} \ln N(x(t) | m(t), P(t)) dt, \quad (84)$$

where $x(t)$ is the true state, $m(t)$ is the estimated mean and $P(t)$ is the estimated covariance. The 95%-consistency is defined as the fraction of times the true state is inside the 95% ellipsoid of $N(m(t), P(t))$. The values of the continuous time metrics are computed using the values of $m(t)$ and $P(t)$ computed at discrete time points.

The different approximation methods for computing the Gaussian expectations in the smoothing equations are also compared. The tested methods are labeled as

- EXT: The method using the Taylor series based linearization.
- CT: Cubature rule based sigma-point method.
- UT: Unscented rule based sigma-point method with parameter values $\alpha = 1$, $\beta = 2$ and $\kappa = 0$.
- G-H: Gauss-Hermite series based sigma-point method with order 3.

The methods labeled CT2, UT2 and G-H2 use the respective sigma-point method with the alternative formulation given in Equations (81) and (82).

4.1. Double well

The double well system is given by

$$dx = 4x(1 - x^2)dt + \sqrt{\sigma}d\beta, \quad y_k = x(t_k) + v_k, \quad (85)$$

where the measurement noise v_k is zero-mean Gaussian with variance R . The prior distribution for the double well system is non-Gaussian and multimodal, but the smoothing distribution can be reasonably well approximated with a Gaussian provided that the measurement variance is not too large. The modes are located at $x = 1$ and $x = -1$ and for sufficiently large value of the process noise parameter σ , there is frequent transition from one mode to the other.

Table 1: The mean 95%-consistencies for the GFGS and VGS for different values of measurement variance. The Gaussian expectations in the smoothing equations are computed analytically.

	Measurement variance R			
	0.02	0.1	0.5	2.5
GFGS	0.93	0.89	0.81	0.91
VGS	0.91	0.88	0.80	0.74

The Gaussian smoothers are compared for 4 different values for the measurement variance R . For each value of the measurement variance R , a data set of 100 Monte Carlo simulations is generated using Euler–Maruyama discretization with time step $\Delta t = 0.01$ from $t_0 = 0$ to $t = 10$. The process noise parameter is chosen to be $\sigma = 1$, which is sufficiently large to cause frequent transition between the modes. The initial state is chosen as $x(0) \sim \mathcal{N}(0, 1)$.

For this system, the Gaussian expectations needed in the Gaussian smoothers can be computed analytically (see e.g. [7]). For both Gaussian smoothers, the differential equations are solved using 4th order Runge–Kutta method with time step 0.01. Average number of 35 iterations was observed for the VGS when using the GFGS results as initial conditions. Also, we observed a much faster convergence using this initialization than with the naive initialization using just the initial conditions.

The boxplots of the RMSE and NLL results for different values of the measurement variance R are shown in Figure 1. For comparison, a reference smoothing solution is also computed using a finite difference approximation of the exact Bayesian smoothing equations [20]. For the reference solution, the NLL is computed using a finite difference approximation for the smoothing density.

For small values of the measurement variance, the VGS results are very close to the reference solution and clearly outperform the GFGS approach in terms of RMSE and NLL. The relatively poor RMSE values for the GFGS are due to the poor estimation of the transitions between the two modes. A typical time series for measurement variance $R = 0.1$ is shown in Figure 2.

For measurement variance $R = 2.5$, the true posterior is bimodal and the VGS tends to have the estimated mean close to one of the modes with relatively small variance. This results in very large NLL values for the VGS, when the true path is not close to the mode. In comparison, the GFGS tends to have the mean close to zero with the 95%-confidence region covering both modes. This shows as very good NLL values and smaller spread of the RMSE values compared to the VGS.

The mean 95%-consistency results over the 100 Monte Carlo simulations are shown in Table 1. From the consistency results, we see that in general the VGS tends to underestimate the variance compared to the GFGS. This is especially clear in the $R = 2.5$ case. The underestimation of the variance is a general property of the variational type of approximations [21, p. 431].

The different approximation methods for the Gaussian expectations were

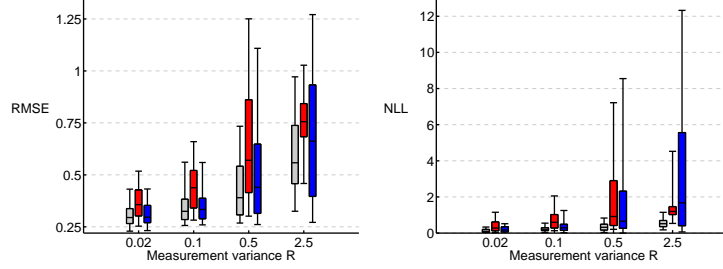


Figure 1: The RMSE (left) and NLL (right) for the reference (grey), GFGS (red) and VGS (blue) for different values of the measurement variance R . The Gaussian expectations in the smoothing equations are computed analytically. The boxplots show the 5%, 25%, 50%, 75% and 95% quantiles.

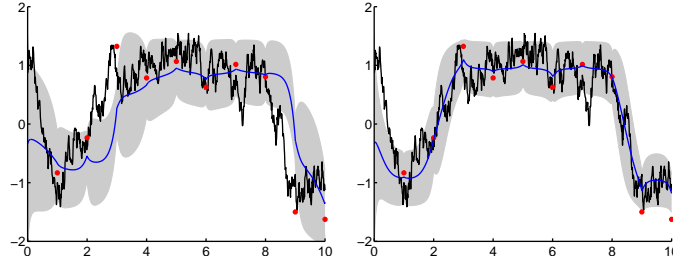


Figure 2: The estimated mean (blue) and 95%-confidence region (shaded) for GFGS (left) and VGS (right). True state (black) and measurements (red dot) are also shown. The Gaussian expectations in the smoothing equations are computed analytically.

compared for measurement variance $R = 0.02$ and using the same data set of 100 Monte Carlo simulations as in the first experiment. The VGS failed to converge to a solution on 5 cases using the EXT method, but no failures were observed using the other methods.

The boxplots of the RMSE and NLL values for the GFGS and VGS when using the different methods to approximate the Gaussian expectations are shown in Figure 3. The VGS using the EXT and G-H methods clearly improve the results of the corresponding GFGS in terms of RMSE and NLL. Numerical problems were observed for the VGS using CT and UT rules, which resulted in poor RMSE and NLL values compared to the GFGS results. Using the alternative formulation of CT2 and UT2 works clearly better and slightly improve the results of the GFGS using the CT and UT methods. The VGS using G-H2 method gives results nearly identical to the VGS using the exact Gaussian expectations.

The mean 95%-consistency values over the 100 Monte Carlo simulations are shown in Table 2. The VGS using the EXT method shows a slight improvement in the consistency compared to the corresponding GFGS result. For the sigma-point methods, the differences are smaller with GFGS giving in general slightly better consistency results.

Table 2: The mean 95%-consistencies for the GFGS and VGS using different methods to compute the Gaussian expectations.

	Integration method						
	EXT	CT	UT	G-H	CT2	UT2	G-H2
GFGS	0.83	0.94	0.94	0.93	0.93	0.93	0.93
VGS	0.86	0.93	0.94	0.93	0.94	0.94	0.91

For comparison, we also included CT2, UT2 and G-H2 versions for the GFGS that use Equation (45) to compute the Gaussian expectations. There seems to be no significant improvement in using the alternative formulation for the GFGS. Also, this increases the computational load, since also the Jacobian needs to be evaluated for each sigma-point.

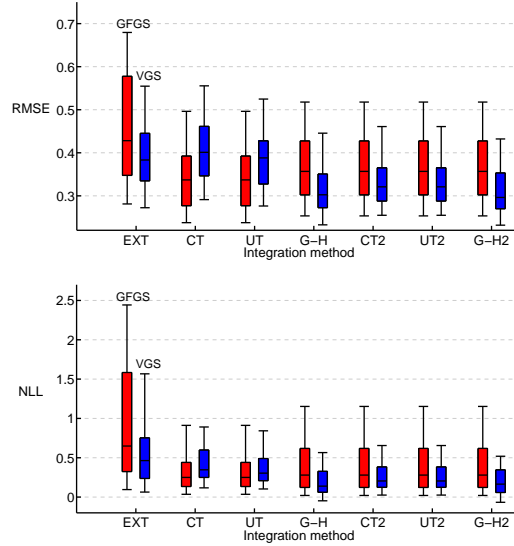


Figure 3: The RMSE (top) and NLL (bottom) for the GFGS (red) and VGS (blue) using different methods to compute the Gaussian expectations. The boxplots show the 5%, 25%, 50%, 75% and 95% quantiles.

4.2. Reentry

The state $x = [r, v, \alpha]^T$ of the reentry problem consists of the vehicle's position r and velocity v in a 2-dimensional coordinate system and a parameter α of its aerodynamic properties. The dynamics are given by

$$d \begin{bmatrix} r(t) \\ v(t) \\ \alpha(t) \end{bmatrix} = \begin{bmatrix} 0_{2 \times 2} & I_{2 \times 2} & 0_{1 \times 1} \\ G(x, t)I_{2 \times 2} & D(x, t)I_{2 \times 2} & 0_{1 \times 1} \\ 0_{1 \times 2} & 0_{1 \times 2} & 0_{1 \times 1} \end{bmatrix} \begin{bmatrix} r(t) \\ v(t) \\ \alpha(t) \end{bmatrix} dt + \begin{bmatrix} 0_{2 \times 3} \\ I_{3 \times 3} \end{bmatrix} d\beta(t), \quad (86)$$

where the gravity related force term $G(x, t)$ and drag related force term $D(x, t)$ are given by

$$G(x, t) = -\frac{Gm_0}{\|r(t)\|^3}, \quad (87)$$

$$D(x, t) = -\beta_0 e^\alpha \exp\left\{\frac{R_0 - \|r(t)\|}{H_0}\right\} \|v(t)\|. \quad (88)$$

The diffusion matrix for the Brownian motion $\beta(t)$ is

$$Q(t) = \begin{bmatrix} 2.4064 \cdot 10^{-5} & 0 & 0 \\ 0 & 2.4064 \cdot 10^{-5} & 0 \\ 0 & 0 & 1 \cdot 10^{-5} \end{bmatrix}. \quad (89)$$

The effective diffusion matrix for this model is singular and the dynamic model can be written in the form of Equations (21) and (22). The values $\beta_0 = -0.59783$, $H_0 = 13.406$, $Gm_0 = 3.9860 \cdot 10^5$ and $R_0 = 6374$ are used as typical values for the parameters [19] (see [22]).

A radar located at $s = [s_x, s_y]^T$ periodically measures the range and bearing of the vehicle with 1 Hz frequency. The measurement model is given by

$$y_k = \begin{bmatrix} \|r(t_k) - s\| \\ \tan^{-1}\left(\frac{r_2(t_k) - s_y}{r_1(t_k) - s_x}\right) \end{bmatrix} + v_k, \quad (90)$$

where the measurement noise v_k is zero-mean Gaussian with covariance matrix

$$R_k = \begin{bmatrix} 1 \cdot 10^{-3} & 0 \\ 0 & 1.7 \cdot 10^{-3} \end{bmatrix}. \quad (91)$$

The state trajectory and noisy measurements are simulated from $t_0 = 0$ to $t_K = 200$ using Euler–Maruyama discretization with time-step $\Delta t = 0.01$. The initial state is drawn from a Gaussian prior with mean and covariance given by

$$m(t_0) = [6500.4 \quad 349.14 \quad -1.8093 \quad -6.7967 \quad 0.6932]^T, \quad (92)$$

$$P(t_0) = \begin{bmatrix} 10^{-6} \cdot I_{4 \times 4} & 0_{4 \times 1} \\ 0_{1 \times 4} & 0 \end{bmatrix}. \quad (93)$$

For this model the computation of the Gaussian expectations needed in the GFGS and VGS is not possible analytically. The Gaussian expectations were computed using the EXT, CT, UT and G-H integration rules. For the VGS it was necessary to use the alternative formulation of CT2 and UT2 rules, since the CT and UT rules caused numerical problems and failure of the algorithm to converge. For G-H this was not a problem. The differential equations were solved using the standard 4-stage Runge–Kutta method with integration step of 0.1. The initial mean and covariance for the smoothers are given by Equation (93), where we used $m_5(t_0) = 0$ and $P_{5,5}(t_0) = 1$ for the unknown aerodynamic parameter.

The boxplots of RMSE and NLL results for 100 Monte Carlo simulations are shown in Figure 4. No clear difference can be seen between GFGS and VGS methods, or between the different methods for computing the Gaussian expectations. The NLL values are slightly better for the GFGS, especially for the position. This is the result of the slightly underestimated variance for the position when using the VGS method. The more compact variance estimate of the VGS can also be seen from the mean 95%-consistency results in Table 3. On average, only 5 iterations were needed in the VGS before convergence. Also, only a small decrease of KL-divergence was observed, which explains the small difference between the methods.

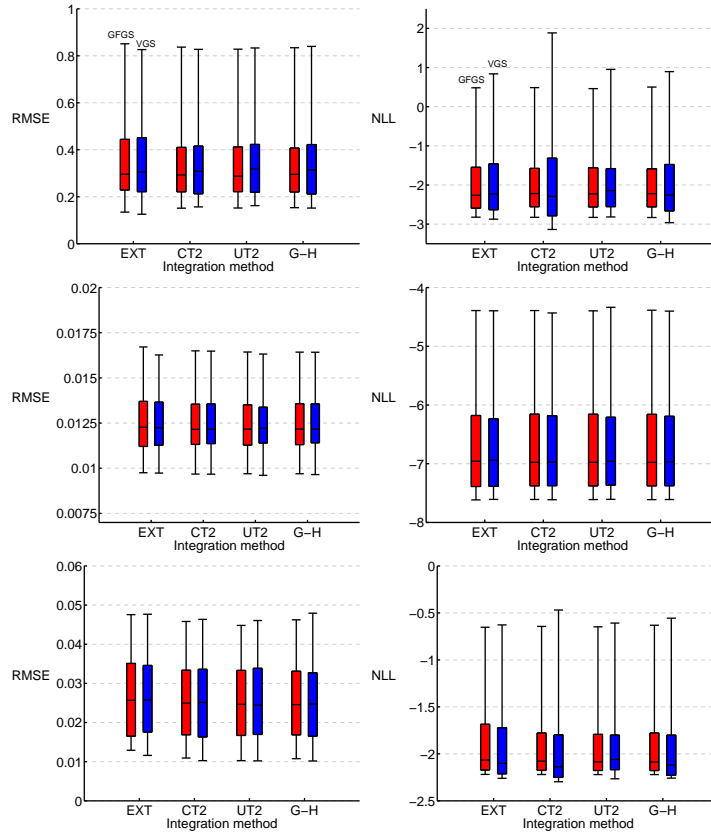


Figure 4: The RMSE and NLL values for the Gaussian filtering based (red) and variational (blue) Gaussian smoothers using different methods to compute the Gaussian expectations. The top row shows the results for position, the middle row for velocity and the bottom row for the aerodynamic parameter. The boxplots show the 5%, 25%, 50%, 75% and 95% quantiles.

Table 3: The mean 95%-consistencies for the Gaussian filtering based and variational Gaussian smoothers.

	Position	Velocity	Parameter
EXT-GFGS	0.94	0.95	0.95
EXT-VGS	0.91	0.95	0.95
CT-GFGS	0.95	0.95	0.96
CT2-VGS	0.88	0.94	0.94
UT-GFGS	0.95	0.95	0.96
UT2-VGS	0.94	0.95	0.96
GH-GFGS	0.95	0.95	0.96
GH-VGS	0.92	0.95	0.95

5. Discussion and conclusions

Compared to the Gaussian filtering based Gaussian smoother the variational Gaussian smoother is more complex to implement and is computationally heavier, since each iteration requires approximately the same amount of computations as one run of the filtering and smoothing equations. The Gaussian filtering based Gaussian smoother provides good initial conditions for the variational Gaussian smoother and could solve the problems with initialization mentioned in [13]. The variational Gaussian smoother provides Gaussian approximation that is optimal in the sense that it minimizes the Kullback–Leibler divergence of the approximating distribution with respect to the true distribution. However, the examples considered here show that this will not always improve the estimate with respect to other commonly used metrics.

Using statistical linearization with respect to the filtering distribution in the variational Gaussian smoothing equations gives formally the Gaussian filtering based smoother as a special case. This suggests that for highly nonlinear systems, the variational Gaussian smoother could better capture the nonlinearities. This is demonstrated in the numerical experiment for the double well system, where the variational Gaussian smoother clearly improves the Gaussian filtering based smoother estimate for small measurement variances. However, no clear improvement was observed in the second numerical experiment for the reentry problem. The reason could be that the nonlinearities in the reentry system are not high enough to gain benefit from using the variational Gaussian smoother equations. A drawback of the variational Gaussian smoother is that it tends to underestimate the variance compared to the Gaussian filtering based smoother.

For general nonlinear systems the Gaussian expectations can be computed using Taylor series based linearization or standard sigma-point methods. The Taylor series based variational Gaussian smoother linearizes the drift and measurement function with respect to the current mean estimate and is therefore similar in idea to the iterated Extended Kalman smoother [23]. Using unscented transform and cubature rule based sigma-point methods in the variational Gaussian smoother resulted in some numerical problems. This is caused since the

unscented and cubature methods are only accurate up to third order monomials, but even for linear systems some of the expectations are over fourth order polynomials. The numerical problems can be reduced by computing the Jacobians of the drift and measurement functions and using an alternative form for the expectations. Also, higher order unscented transform could be used for these expectations (see [19, 24]). Third order Gauss–Hermite integration rule worked well for both examples considered in this paper, but the computational cost is high especially for high dimensional systems. For some high dimensional systems the computational cost of the Gauss–Hermite method could be reduced by using Rao–Blackwellisation [25].

Acknowledgments

The first author received financial support from the Tampere University of Technology Doctoral Programme in Engineering and Natural Sciences. The second author would like to thank the Academy of Finland.

References

- [1] A. H. Jazwinski, Stochastic processes and filtering theory, Dover, 1970.
- [2] C. T. Leondes, J. B. Peller, E. B. Stear, Nonlinear smoothing theory, Systems Science and Cybernetics, IEEE Transactions on 6 (1) (1970) 63–71.
- [3] F. E. Daum, Exact finite-dimensional nonlinear filters, Automatic Control, IEEE Transactions on 31 (7) (1986) 616–622.
- [4] A. Bryson, M. Frazier, Smoothing for linear and nonlinear dynamic systems, in: Proceedings of the Optimum System Synthesis Conference, 1963, pp. 353–364.
- [5] Y. Wu, D. Hu, M. Wu, X. Hu, A numerical-integration perspective on Gaussian filters, Signal Processing, IEEE Transactions on 54 (8) (2006) 2910–2921.
- [6] S. Särkkä, J. Sarmavuori, Gaussian filtering and smoothing for continuous-discrete dynamic systems, Signal Process. 93 (2) (2013) 500–510.
- [7] C. Archambeau, D. Cornford, M. Opper, J. Shawe-Taylor, Gaussian process approximations of stochastic differential equation, Proceedings of the, Journal of Machine Learning Research: Workshop and Conference 11 (2007) 1–16.
- [8] C. Archambeau, M. Opper, Y. Shen, D. Cornford, J. S. Shawe-taylor, Variational inference for diffusion processes, in: J. Platt, D. Koller, Y. Singer, S. Roweis (Eds.), Advances in Neural Information Processing Systems 20, Curran Associates, Inc., 2008, pp. 17–24.

- [9] C. Archambeau, M. Opper, Approximate inference for continuous-time Markov processes, in: *Bayesian Time Series Models*, Cambridge University Press, 2011, pp. 125–140.
- [10] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [11] Y. Shen, C. Archambeau, D. Cornford, M. Opper, J. Shawe-Taylor, R. Barillec, A comparison of variational and Markov chain Monte Carlo methods for inference in partially observed stochastic dynamic systems, *Journal of Signal Processing Systems* 61 (1) (2010) 51–59.
- [12] Y. Shen, D. Cornford, M. Opper, C. Archambeau, Variational Markov chain Monte Carlo for Bayesian smoothing of non-linear diffusions, *Computational Statistics* 27 (1) (2012) 149–176.
- [13] M. D. Vrettas, D. Cornford, M. Opper, Y. Shen, A new variational radial basis function approximation for inference in multivariate diffusions, *Neurocomputing* 73 (7) (2010) 1186–1198.
- [14] M. D. Vrettas, D. Cornford, M. Opper, Estimating parameters in stochastic systems: A variational Bayesian approach, *Physica D: Nonlinear Phenomena* 240 (23) (2011) 1877–1900.
- [15] S. Särkkä, T. Sottinen, Application of Girsanov theorem to particle filtering of discretely observed continuous-time non-linear systems, *Bayesian Analysis* 3 (3) (2008) 555–584.
- [16] H. Kushner, Dynamical equations for optimal nonlinear filtering, *Journal of Differential Equations* 3 (2) (1967) 179 – 190.
- [17] K. Ito, K. Xiong, Gaussian filters for nonlinear filtering problems, *Automatic Control, IEEE Transactions on* 45 (5) (2000) 910 – 927.
- [18] H. E. Rauch, F. Tung, C. T. Striebel, Maximum likelihood estimates of linear dynamic systems, *American Institute of Aeronautics and Astronautics (AIAA) Journal* 3 (8) (1965) 1445–1450.
- [19] S. J. Julier, J. K. Uhlmann, Unscented filtering and nonlinear estimation, *Proceedings of the IEEE* 92 (3) (2004) 401–422.
- [20] S. Särkkä, *Recursive Bayesian inference on stochastic differential equations*, Ph.D. thesis, Helsinki University of Technology (2006).
- [21] D. J. C. MacKay, *Information Theory, Inference and Learning Algorithms*, Cambridge University Press, 2003.
- [22] J. Austin, C. Leondes, Statistically linearized estimation of reentry trajectories, *Aerospace and Electronic Systems, IEEE Transactions on* AES-17 (1) (1981) 54–61.

- [23] B. M. Bell, The iterated Kalman smoother as a Gauss-Newton method, SIAM Journal on Optimization 4 (3) (1994) 626–636.
- [24] U. N. Lerner, Hybrid Bayesian networks for reasoning about complex systems, Ph.D. thesis, Stanford University (2002).
- [25] P. Closas, C. Fernández-Prades, The marginalized square-root Quadrature Kalman filter, in: Signal Processing Advances in Wireless Communications (SPAWC), 2010 IEEE Eleventh International Workshop on, 2010, pp. 1–5.

Appendix A. Computing the KL-divergence using Girsanov’s theorem

This section presents the derivation of the KL-divergence term in Equation (28) for the system with singular effective diffusion matrix. The KL-divergence term can be partitioned to [9]

$$\text{KL}(\mathbb{Q}_X \parallel \mathbb{P}_X | \mathcal{Y}) = \text{KL}(\mathbb{Q}_X \parallel \mathbb{P}_X) - \sum_{k=1}^K \mathbb{E}_q [\ln p(y_k | x_k)], \quad (\text{A.1})$$

where \mathbb{P}_X corresponds to the joint probability law of the stochastic processes $x_1(t)$ and $x_2(t)$ and \mathbb{Q}_X to the joint probability law of the stochastic processes $s_1(t)$ and $s_2(t)$ defined by

$$\begin{aligned} \frac{dx_1}{dt} &= F_1(t)x, \\ dx_2 &= f_2(x(t), t) + d\beta(t) \\ \frac{ds_1}{dt} &= F_1(t)s, \\ ds_2 &= g_2(x(t), t) + d\beta(t), \end{aligned}$$

where $\beta(t)$ is Brownian motion with diffusion matrix $Q(t)$ with respect to measure \mathbb{P}_X .

The processes $s_1(t)$ and $s_2(t)$ are weak solutions to the original system under the measure \mathbb{Q}_X that is defined through the Radon-Nikodym derivative [15]

$$\mathbb{E} \left[\frac{d\mathbb{Q}_X}{d\mathbb{P}_X} \mid \mathcal{F}_t \right] = Z(t),$$

where \mathcal{F}_t is the natural filtration of the Brownian motion $\beta(t)$ and

$$\begin{aligned} Z(t) &= \exp \left[\int_0^t \{f_2(s_1(t), s_2(t), t) - g_2(s_1(t), s_2(t), t)\}^T d\beta(t) \right. \\ &\quad \left. - \frac{1}{2} \int_0^t \{f_2(s_1(t), s_2(t), t) - g_2(s_1(t), s_2(t), t)\}^T Q^{-1}(t) \right. \\ &\quad \left. \{f_2(s_1(t), s_2(t), t) - g_2(s_1(t), s_2(t), t)\}^T \right] . \end{aligned}$$

The KL-divergence term in the right side of Equation (A.1) is then given by

$$\begin{aligned} \text{KL}(\mathbb{Q}_X || \mathbb{P}_X) &= \mathbb{E}_{\mathbb{Q}_X} [-\ln Z(t)] \\ &= \frac{1}{2} \int_0^t \mathbb{E}_q \left[\{f_2(s_1(t), s_2(t), t) - g_2(s_1(t), s_2(t), t)\}^T Q^{-1}(t) \right. \\ &\quad \left. \{f_2(s_1(t), s_2(t), t) - g_2(s_1(t), s_2(t), t)\} \right]. \end{aligned}$$

Inserting $g_2(s_1(t), s_2(t), t) = -A(t)s(t) + b(t)$ gives then the desired KL-divergence.

Appendix B. Converting the Gaussian filtering based Gaussian smoother to the variational form

Here we present the derivation of the variational form of the Gaussian filtering based Gaussian smoother. This is achieved by using the change of variables:

$$\lambda = -P_f^{-1}(m_s - m_f), \quad \Psi = -\frac{1}{2} \left(P_f^{-1} - P_s^{-1} \right).$$

Computing the time derivatives of the new variables and inserting the filtering and smoothing differential equations (37)-(38) and (43)-(44) gives

$$\begin{aligned} \frac{d}{dt} \lambda &= P_f^{-1} \left(\frac{d}{dt} P_f \right) P_f^{-1} (m_s - m_f) - P_f^{-1} \left(\frac{d}{dt} m_s - \frac{d}{dt} m_f \right) \\ &= P_f^{-1} (\mathbb{E}_f[F_x(x)] P_f + P_f \mathbb{E}[F_x^T(x)] + Q) P_f^{-1} (m_s - m_f) \\ &\quad - P_f^{-1} (\mathbb{E}_f[f(x)] + \mathbb{E}_f[F_x(x)] (m_s - m_f) + Q P_f^{-1} (m_s - m_f) - \mathbb{E}_f[f(x)]) \\ &= -\mathbb{E}_f[F_x(x)]^T \lambda \end{aligned}$$

and

$$\begin{aligned} \frac{d}{dt} \Psi &= \frac{1}{2} P_f^{-1} \left(\frac{d}{dt} P_f \right) P_f^{-1} - \frac{1}{2} P_s^{-1} \left(\frac{d}{dt} P_s \right) P_s^{-1} \\ &= \frac{1}{2} P_f^{-1} (\mathbb{E}_f[F_x(x)] P_f + P_f \mathbb{E}[F_x(x)]^T + Q) P_f^{-1} \\ &\quad - \frac{1}{2} P_s^{-1} (\mathbb{E}_f[F_x(x)] P_s + Q P_f^{-1} P_s + P_s \mathbb{E}[F_x^T(x)] + P_s P_f^{-1} Q - Q) P_s^{-1} \\ &= -\Psi \mathbb{E}_f[F_x(x)] - \mathbb{E}_f[F_x(x)]^T \Psi + 2\Psi Q \Psi. \end{aligned}$$

Inserting $A(t) = -\mathbb{E}_f[F_x(x)] + 2\Sigma(t)\Psi(t)$ to the above equations gives

$$\begin{aligned} \frac{d}{dt} \lambda(t) &= A^T(t) \lambda(t) - 2\Psi(t) \Sigma(t) \lambda(t) \\ \frac{d}{dt} \Psi(t) &= \Psi(t) A(t) + A^T(t) \Psi(t) - 2\Psi(t) \Sigma(t) \Psi(t). \end{aligned}$$

Appendix C. Gradients with respect to mean and covariance

In this section we derive the expressions for the gradients with respect to the mean vector m and covariance matrix P of a Gaussian expectation over a scalar function $e(x)$. These are used to form the sigma-point approximations for the Gaussian expectations in the variational Gaussian smoothing equations. Computing the gradients gives

$$\begin{aligned}\nabla_m \mathbb{E}[e(x)] &= \nabla_m \left[\int e(x) \mathcal{N}(x | m, P) dx \right] = \int e(x) \nabla_m \mathcal{N}(x, | m, P) dx \\ &= \int e(x) \mathcal{N}(x, | m, P) P^{-1} (x - m) dx \\ &= P^{-1} \mathbb{E}[e(x)(x - m)] = \mathbb{E}[\nabla_x e(x)]\end{aligned}$$

and

$$\begin{aligned}\nabla_P \mathbb{E}[e(x)] &= \nabla_P \left[\int e(x) \mathcal{N}(x | m, P) dx \right] = \int e(x) \nabla_P \mathcal{N}(x | m, P) dx \\ &= \int e(x) \mathcal{N}(x | m, P) \frac{1}{2} [P^{-1}(x - m)(x - m)^T P^{-1} - P^{-1}] dx \\ &= \frac{1}{2} P^{-1} \mathbb{E}[e(x)(x - m)(x - m)^T] P^{-1} - \frac{1}{2} \mathbb{E}[e(x)] P^{-1} \\ &= \frac{1}{2} P^{-1} \mathbb{E}[\nabla_x e(x)(x - m)^T].\end{aligned}$$