
Differentiable Bayesian inference for SDEs using a pathwise series expansion of Brownian motion

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 By invoking a pathwise series expansion of Brownian motion, we propose to
2 approximate a stochastic differential equation (SDE) with an ordinary differential
3 equation (ODE). This allows us to reformulate Bayesian inference for a SDE as the
4 parameter estimation task for an ODE. Unlike a nonlinear SDE, the likelihood for
5 an ODE model is tractable and its gradient can be obtained using adjoint sensitivity
6 analysis. This reformulation allows us to use an efficient sampler, such as NUTS,
7 that rely on the gradient of the log posterior. Applying the reparameterisation trick,
8 variational inference can also be used for the same estimation task. We illustrate
9 the proposed method on two biological SDE models. We obtain similar parameter
10 estimates when compared to data augmentation techniques. Finally, we introduce a
11 Bayesian neural SDE model to test out approach for a system identification task.
12 We obtain competitive performance on a benchmark motion capture dataset when
13 compared to latent ODE models.

14 1 Introduction

15 The solution to a stochastic differential equation – a diffusion process – that evolves randomly over
16 time is a flexible and useful tool in modelling stochastic systems. Naturally, SDEs are widely used in
17 modelling continuously evolving systems of real valued quantities subjected to noise or stochastic
18 fluctuations.

19 The task of estimating the parameters of a SDE observed at discrete times is highly challenging. One
20 has to deal with the estimation of a high dimensional latent diffusion in addition to the governing
21 parameters. Moreover, the exact transition densities given by the forward Kolmogorov equation,
22 required to calculate the likelihood, are intractable for a nonlinear SDE. In a Bayesian paradigm
23 the estimation task requires repeated sampling of the parameters as well as the latent diffusion,
24 wherein the proposals for the latter rely on the Euler-Maruyama discretisation of the SDE. The
25 primary challenge here is the non-trivial sampling of the latent diffusion path. Additionally any
26 sampling scheme has to also deal with the strong correlation between the diffusion path and the
27 model parameters. Slow convergence of MCMC incurs a high computational cost due to the iterative
28 nature of sampling a discretised diffusion path.

29 With the widespread availability of automatic differentiation (AD) techniques, Bayesian inference
30 is increasingly carried out using Markov chain Monte Carlo (MCMC) algorithms that traverse the
31 parameter space based on gradient of the target density. Such algorithms have proposal generating
32 mechanism that can rapidly explore the parameter space when compared to traditional random-
33 walk MCMC. Moreover, optimisation based alternatives to MCMC such as black-box variational
34 inference has the potential to further expedite the inference process. AD has made it possible to apply
35 such algorithms to many complex models with a differentiable target density. Many probabilistic

programming platforms, that rely on AD, includes such algorithms and have vastly automated the entire process of Bayesian inference.

Application of the aforementioned efficient algorithms to a SDE is highly non-trivial, as this requires the extension of AD to a SDE. Firstly, existing approaches (Gobet and Munos, 2005; Yang and Kushner, 1991; Giles and Glasserman, 2006) to differentiate a diffusion process generally scales poorly with the parameter dimension, or with memory. Additionally, the intractable likelihood of a SDE is often approximated using non-differentiable simulation methods, which in turn renders any attempt to use the gradient of the target density futile.

In case of an ODE, in contrast to a SDE, the likelihood is tractable and there exists an efficient numerical sensitivity analysis technique that lends itself to AD (Chen et al., 2018). Thus, we propose to approximate a SDE by an equivalent ODE using a pathwise truncated series expansion of Brownian motion (Lyons et al., 2012, 2014). The resulting ODE contains the same SDE parameters in addition to extra parameters that are the expansion coefficients. We estimate all these parameters jointly where the marginal density of the model parameters is the desired quantity of interest. We carry out this estimation task using the No-U-Turn sampler as well as using black-box variational inference, both inherently capable of handling the inflated parameter dimension efficiently.

2 Related work

Series expansion of Brownian motion for inference, in conjunction with a Gibbs sampling scheme, was first introduced in Lyons et al. (2012) for systems with additive noise. Lyons et al. (2014) applied series expansion for filtering in nonlinear state-space models. By marrying series expansion with adjoint sensitivity analysis, our approach can be applied to both additive and multiplicative noise models, the latter is prevalent in biology. Moreover, our approach is targeted towards differentiable inference, which accommodates faster and more efficient inference algorithms.

AD for SDEs was introduced in Li et al. (2020) by deriving an adjoint sensitivity analysis technique for Stratonovich SDEs, thus enabling more efficient AD for a SDE, when compared to approaches in Gobet and Munos (2005); Yang and Kushner (1991); Giles and Glasserman (2006). We instead focused on approximating a SDE with an ODE, the latter being fully integrated with AD (Chen et al., 2018; Ghosh et al., 2021). Additionally, working with an ODE we can apply familiar tools for inference and model building. Unlike Chen et al. (2018), our approach can handle both classical inference as well as system identification (say using a Neural SDE).

Variational inference for a SDE, exploiting AD, was recently proposed in Ryder et al. (2018). Li et al. (2020) used the path integral approach for this purpose. In contrast to these approaches our method do not require the usage of a neural network for inference. Thus, our approach is more conducive for mechanistic modelling, especially the calibration task.

ODE approximation for an SDE can also be derived using a Gaussian approximation of the transition density (Golightly and Gillespie, 2013; Fearnhead et al., 2014). Our approach rely on a pathwise approximation of the diffusion. So, in our approach we can accommodate non-Gaussian transition densities.

3 Background

We begin by first introducing the Bayesian inference framework for a diffusion process described by a SDE. Consider a K -dimensional diffusion process that satisfies the following SDE:

$$d\mathbf{X}_t = \mathbf{a}(\mathbf{X}_t, \boldsymbol{\theta})dt + \sqrt{\mathbf{B}(\mathbf{X}_t, \boldsymbol{\theta})}d\mathbf{W}_t, \quad \mathbf{X}_0 = \mathbf{x}_0, \quad (1)$$

where \mathbf{X}_t denotes the value of the process at time t , \mathbf{a} is a K -dimensional *drift* vector, \mathbf{B} a $K \times K$ diffusion matrix and the driving noise \mathbf{W}_t is a K -dimensional Brownian motion which is treated in the Itô sense. Both the drift and diffusion matrix depends on an unknown parameter vector $\boldsymbol{\theta} \in \mathbb{R}^D$. If the initial condition \mathbf{x}_0 is unknown then this becomes an element of $\boldsymbol{\theta}$.

We assume that $\mathbf{a}(\cdot)$ and $\mathbf{B}(\cdot)$ are sufficiently regular functions such that Equation (1) has a weak non-explosive solution (Oksendal, 2013).

Consider a set of noisy experimental observations $\mathbf{y} \in \mathbb{R}^{M \times K}$ observed at M experimental time points, $\{t_i\}_{i=1}^M$, for the K states. Within the Bayesian inferential paradigm we want to place a

85 prior distribution on the unknown parameters $p(\theta)$ and intend to obtain the corresponding posterior
 86 distribution

$$p(\mathbf{X}, \theta | \mathbf{y}) \propto p(\mathbf{y} | \mathbf{X}, \theta) p(\mathbf{X}) p(\theta) \quad (2)$$

87 of the latent path \mathbf{X} and the unknown parameters θ , where $p(\mathbf{y} | \mathbf{X}, \theta)$ is the likelihood and $p(\mathbf{X})$ is
 88 the distribution of the diffusion \mathbf{X} defined from the SDE given by Equation (1). For a nonlinear SDE,
 89 $p(\mathbf{X})$ is intractable. To work around this intractability an Euler-Maruyama discretisation is generally
 90 used, yielding a Gaussian approximation to the transition density. Inference proceeds by generating
 91 samples of \mathbf{X} and θ .

92 Within an MCMC scheme, samples of \mathbf{X} and θ can be drawn by using a Gibbs sampling (Golightly
 93 and Wilkinson, 2008) scheme. As an alternative to such Gibbs proposal mechanism both \mathbf{X} and θ
 94 can be jointly updated using a particle MCMC scheme wherein a particle filter is used to draw the
 95 diffusion path and evaluate an unbiased estimate of the likelihood. Particle MCMC methods, despite
 96 their higher computational cost, produce a faster mixing Markov chain and thus have become the state-
 97 of-the-art for inference in SDEs. Nonetheless, both these methods are inherently non-differentiable
 98 and rely crucially on a random-walk proposal for updating θ .

99 4 Pathwise series approximation

100 Within an interval $[0, T]$ a standard Brownian motion has the following series expansion (Lyons et al.,
 101 2012, 2014):

$$W_t = \sum_{i=1}^{\infty} \left(\int_0^T \phi_i(s) dW(s) \right) \int_0^t \phi_i(u) dW(u), \quad (3)$$

102 where $\{\phi_i\}_{i \geq 1}$ is an orthonormal basis of $L^2[0, T]$. For example this can be the Karhunen–Loève
 103 (KL) expansion of a Gaussian process given by

$$\phi_i(t) = (2/T)^{1/2} \cos\{(2i-1)\pi t/2T\}. \quad (4)$$

104 We will use the shorthand $Z_i = \int_0^T \phi_i(s) dW(s)$. Since the basis functions $\{\phi_i\}$ are deterministic
 105 and orthonormal, it follows from standard results of Itô calculus that $Z_i \sim \mathcal{N}(0, 1)$.

106 The infinite series in Equation 3 can be truncated after N terms to obtain an approximation of \hat{W}_t
 107 Brownian motion. Taking derivative with respect to time we obtain the following approximation to
 108 the differential of Brownian motion given by

$$\frac{d\hat{W}}{dt} = \sum_{i=1}^N Z_i \phi_i(t). \quad (5)$$

109 Consider now a scalar SDE driven by the approximate Brownian motion:

$$dX_t = a(X_t, \theta) dt + b(X_t, \theta) d\hat{W}_t, \quad (6)$$

110 where a, b are some scalar drift and diffusion term, that ensure a weak solution of this SDE exists. As
 111 $N \rightarrow \infty$ the solution of the above ODE X_t will converge to the solution had it been driven by the
 112 exact Brownian motion. The seminal work of Wong and Zakai (1965) shows that as $N \rightarrow \infty$ the
 113 solution actually converges to the solution of a Stratonovich SDE given by

$$dX_t = a(X_t, \theta) dt + b(X_t, \theta) \circ dW_t. \quad (7)$$

114 Thus, we can approximate the above Stratonovich SDE as

$$dX_t = a(X_t, \theta) dt + b(X_t, \theta) \circ d\hat{W}_t, \quad (8)$$

115 which is just the ODE:

$$\frac{d\hat{X}}{dt} = a(\hat{X}, \theta) + b(\hat{X}, \theta) \sum_{i=1}^N Z_i \phi_i(t). \quad (9)$$

116 Similarly, we can approximate a K -dimensional Brownian motion by applying the truncated series
 117 approximation to each of the K dimensions. Substituting such a K -dimensional approximation in
 118 Equation (1) we obtain the following ODE:

$$\frac{d\hat{\mathbf{X}}}{dt} = \mathbf{a}(\hat{\mathbf{X}}, \theta) + \mathbf{B}(\hat{\mathbf{X}}, \theta) \sum_{i=1}^N \mathbf{Z}_i \Phi_i(t), \quad (10)$$

where $\mathbf{Z}_i, \Phi_i \in \mathbb{R}^K$. In the multivariate case the convergence to the Stratonovich SDE is not guaranteed in general (Lyons et al., 2014). However, if one chooses to expand the Brownian motion using Haar wavelets then convergence is indeed guaranteed (Lyons et al., 2014; McShane, 2020). The latter is the Levy-Ciecelski construction of Brownian motion. In our experiments we did not observe a failure to converge to the Stratonovich limit while using the Karhunen-Loeve basis functions. Similar findings were reported in Lyons et al. (2012).

The first few coefficients in general captures the large scale oscillations whereas the remaining coefficients determine the small-scale (high frequency) oscillations. Thus, even without using a really large number of coefficients a good approximation to the true time t marginal distribution $p(\hat{\mathbf{X}}(t))$ can be achieved. In Figure 1 we compare the marginal density of the states of a Lotka-Volterra model (see section 8.1) obtained through the Euler-Maruyama and the series expansion introduced here. Clearly, with even 10 coefficients both the approximations match well.

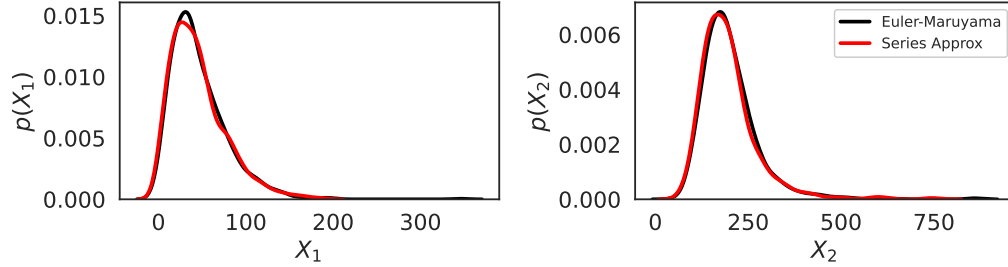


Figure 1: We compare the marginal density of the Lotka-Volterra model states at time $T=30$. We used 1000 samples to evaluate these densities.

130

131 5 Inference using MCMC

In order to use the ODE approximation of the SDE for estimating the SDE parameters θ , we first assume that \mathbf{y} is the noisy observation, at the M time-points, of the solution $\hat{\mathbf{X}}(\theta) := \hat{\mathbf{X}}(\theta, \mathbf{Z}, \mathbf{x}_0) \in \mathbb{R}^K$ of the ODE given by Equation (5), which approximates the actual diffusion path $\hat{\mathbf{X}} \approx \mathbf{X}$, where $\mathbf{Z} := \{\mathbf{Z}_i\}_{i=1}^N$ and \mathbf{x}_0 is the initial value. Note that for a state-dependent diffusion term the Itô and Stratonovich formulation differ in their drift functions. Since the series approximation converges to the Stratonovich SDE, thus an Itô SDE needs to be converted to its equivalent Stratonovich form before applying the approximation. This can be done by modifying the drift term (see appendix A for details). The task then is to infer the posterior distribution of all the unknown parameters of the ODE: $\theta, \mathbf{Z}, \mathbf{x}_0$. Since the random variables \mathbf{Z} are distributed as standard Normal, we use the same as the prior on these variables. Thus, we have the prior: $p(\mathbf{Z}_i) = \mathcal{N}(\mathbf{0}, \mathbb{I}_{K \times K})$.

We can now write the joint posterior distribution of $(\theta, \mathbf{Z}, \mathbf{x}_0)$ as

$$p(\theta, \mathbf{Z}, \mathbf{x}_0 | \mathbf{y}) \propto p(\mathbf{y} | \hat{\mathbf{X}}(\theta)) p(\mathbf{Z}) p(\theta, \mathbf{x}_0), \quad (11)$$

where the likelihood $p(\mathbf{y} | \hat{\mathbf{X}}(\theta))$ is now both tractable and differentiable, with respect to $\theta, \mathbf{Z}, \mathbf{x}_0$.

We can obtain samples from this posterior distribution using a MCMC, sampler that targets gradient of the log posterior: $\nabla_{\theta, \mathbf{Z}, \mathbf{x}_0} \{\log p(\theta, \mathbf{Z}, \mathbf{x}_0 | \mathbf{y})\}$, such as the No-U-Turn (Hoffman and Gelman, 2014) sampler. The marginal density $p(\theta | \mathbf{y})$ can be easily obtained by collecting the corresponding samples. We can also obtain the approximation to the latent diffusion's sample paths, conditioned on the data, $\hat{\mathbf{X}}(\theta) |_{\theta, \mathbf{Z}, \mathbf{x}_0 \sim p(\theta, \mathbf{Z}, \mathbf{x}_0 | \mathbf{y})}$ by solving the ODE using the parameter samples.

149 6 Variational inference

Let us first introduce the shorthand $\Theta := (\theta, \mathbf{Z}, \mathbf{x}_0)$ to denote the vector containing all the unknown quantities that we want to estimate. Using variational inference we can approximate $p(\Theta | \mathbf{y})$ with a

tractable distribution $q(\Theta|\lambda)$ from a family of distributions $q(\cdot|\lambda)$, indexed by λ , by maximising the evidence lower bound (ELBO) (Jordan et al., 1999) given by

$$\mathcal{L}(\lambda) = \mathbb{E}[\log p(y|\hat{X}(\Theta))p(\Theta)] - \mathbb{E}[\log q(\Theta|\lambda)] \quad (12)$$

where the above expectations are with respect to $q(\Theta|\lambda)$. If gradient of the ELBO, w.r.t the variational parameter λ , is available then variational inference can be formulated as a simple gradient descent problem as follows:

$$\lambda \leftarrow \lambda + \gamma \nabla_{\lambda} \mathcal{L}(\lambda), \quad (13)$$

where γ is a learning rate. However, for an ODE the above expectations are intractable. Following Ghosh et al. (2021), we apply the reparameterisation trick (Kingma et al., 2014; Rezende et al., 2014; Titsias and Lázaro-Gredilla, 2014) to obtain a Monte Carlo (MC) estimate, using L samples, of the gradient of the ELBO:

$$\nabla_{\lambda} \mathcal{L}_{MC}(\lambda) = \frac{1}{L} \sum_{l=1}^L \nabla_{\Theta} \left[\left\{ \log p(y|\hat{X}(\Theta^l))p(\Theta^l) - \log q(\Theta^l|\lambda) \right\} \right] \nabla_{\lambda} g(\lambda, \epsilon^{(l)}), \quad (14)$$

where Θ^l is the output of an invertible, differentiable function $g(\lambda, \epsilon^{(l)})$ and $\epsilon^{(l)} \sim p(\epsilon)$ —a parameter free distribution. Substituting this MC estimate in Equation (13), we can find the optimal λ^* using the following stochastic optimisation update:

$$\lambda \leftarrow \lambda + \gamma \nabla_{\lambda} \mathcal{L}_{MC}(\lambda). \quad (15)$$

6.1 Choice of the approximation

To place prior distributions with support on positive reals we need to transform the support of Θ to the unconstrained real line \mathbb{R}^D : $T: \mathbb{R}_{>0}^D \rightarrow \mathbb{R}^D$, and subsequently obtain a transformed parameter vector $\xi = T(\Theta)$. The posterior density $p(\Theta|y)$, with the above transformation, is given by

$$p(\Theta|y) \propto p(y|X(T^{-1}(\xi)))p(T^{-1}(\xi)) \left| \det J_{T^{-1}}(\xi) \right|, \quad (16)$$

where $J_{T^{-1}}(\xi)$ is the Jacobian of the inverse of T . This transformation lets us choose an approximating distribution $q(\xi|\lambda)$ with unconstrained support, such as a Gaussian.

Gaussian approximation: Parameters of a nonlinear ODE such as Equation (10) are often strongly correlated. Moreover due the nature of our approximation we also expect correlations to exist between the expansion coefficients \mathbf{Z} and Θ . For nonlinear ODEs, a full-rank Gaussian density as the variational approximation was shown in Ghosh et al. (2021) to be able to capture the correlation structure of the posterior amply. Thus, following Ghosh et al. (2021) we also chose a full-rank Gaussian approximation: $q(\xi|\lambda) = \mathcal{N}(\xi|\mu, \Sigma)$, where $\lambda = (\mu, \Sigma)$ are the variational parameters. To ensure that the covariance matrix Σ remains positive semidefinite, we parameterise the covariance using Cholesky factorisation, $\Sigma = LL'$. To ensure uniqueness we take the logarithm of the diagonal elements of L . The required reparameterisation, $\xi = g(\lambda, \epsilon)$, then simply follows as the affine transform $\xi = \mu + L\epsilon$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$.

Normalizing flows: The full-rank Gaussian, although is able to capture strong correlations, lacks richness to adequately approximate a complex and possibly multi-modal posterior density if we model the drift and diffusion functions using a neural network, as required for a system identification task (see section 9). A series of invertible nonlinear transformations, a normalizing flow, can be applied to a simple base density (say a Gaussian) to generate a richer approximating density. In particular we use a neural autoregressive flow which has a compact parameterisation: block neural autoregressive flow (BNAF) (De Cao et al., 2020), and is proven to be an universal density approximator.

7 Gradient evaluations

To carry out inference we need to obtain the gradient of a scalar function: i) the log posterior in the MCMC case or, ii) the ELBO for variational inference with respect to the model parameters Θ or the variational parameters λ . This in turn requires the propagation of gradients through the ODE. This can be efficiently carried out using the *adjoint sensitivity* analysis (Chen et al., 2018; Rackauckas et al., 2018), which is the continuous formulation of reverse-mode automatic differentiation. We

193 describe this for obtaining the gradient of the ELBO w.r.t Θ . Note that the downstream gradients,
 194 $\frac{d\Theta}{d\lambda}$, can be trivially obtained using AD. Gradient of the prior and likelihood densities can be obtained
 195 analogously.

196 Consider a cost function, such as the ELBO, that depends on the ODE solution at the measurement
 197 times:

$$C(\hat{\mathbf{X}}) = \sum_i c(\hat{\mathbf{X}}_{t_i}), \quad (17)$$

198 where the sum appears due the factorisation of the likelihood over the time axis. In *adjoint sensitivity*
 199 analysis (Rackauckas et al., 2018) the gradient of the above scalar-valued cost function $C(\cdot)$, whose
 200 input is the ODE solution, can be computed directly. The first step is to solve a backwards ODE, the
 201 adjoint problem:

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t)^\top \frac{\partial \mathbf{f}}{\partial \mathbf{X}}, \quad (18)$$

202 where we use the shorthand \mathbf{f} to denote the velocity field of the ODE in Equation (10).
 203 Furthermore, at each experimental time point t_i this backward ODE is perturbed by $\frac{\partial c(\hat{\mathbf{X}}_{t_i})}{\partial \mathbf{X}}$. The
 204 gradient of the cost function with respect to the ODE parameters can be evaluated by another
 205 quadrature as follows:

$$\frac{dC}{d\Theta} = \mathbf{a}(t_1)^\top \frac{\partial \mathbf{f}(\hat{\mathbf{X}}_{t_1}, \Theta)}{\partial \Theta} + \sum_i \int_{t_i}^{t_{i+1}} \mathbf{a}(t)^\top \frac{\partial \mathbf{f}}{\partial \Theta} dt. \quad (19)$$

206 Note that a continuous solution of the system and the adjoint states is required for the integration
 207 above. Alternatively, the system, the adjoint and the cost function ODEs can be solved simultaneously
 208 backward in time (Chen et al., 2018). We have used this reverse-time simultaneous-solution technique
 209 introduced in Chen et al. (2018).

210 8 Benchmarking: parameter estimation in mechanistic models

211 We begin by testing the efficacy of the series approximation for the task of parameter estimation. In
 212 this regard we use two biological SDEs: (i) the Lotka-Volterra model, and (ii) the stochastic SIR
 213 epidemic model. For the Lotka-Volterra model we have used simulated data but for the SIR model
 214 we used real data. For both these models we first apply a particle marginal Metropolis-Hastings
 215 (PMMH) algorithm (Andrieu et al., 2010; Golightly and Wilkinson, 2011) to obtain a *gold standard*
 216 estimate of the posterior of θ , \mathbf{X} that we treat as a proxy for the true unknown posterior. We used
 217 a Bootstrap particle filter (Gordon et al., 1995; Golightly and Wilkinson, 2011) which uses the
 218 Euler-Maruyama discretisation as the proposal (Golightly and Wilkinson, 2011) for \mathbf{X} and we update
 219 θ using an adaptive random-walk Metropolis-Hastings algorithm. Having run the PMMH, we then
 220 compare with this gold standard the estimates of the posterior obtained by applying MCMC using the
 221 No-U-Turn (NUTS) algorithm (Hoffman and Gelman, 2014) and variational inference, both using the
 222 proposed series approximation: which we denote as **SA-ODE**. For, variational inference we used
 223 the RMSprop (Tieleman and Hinton, 2016) optimisation algorithm with a step size of 10^{-3} . Our
 224 choice of RMSprop is motivated by the findings in Ghosh et al. (2021). By VI we denote variational
 225 inference.

226 For both the models we have used $N = 10$ basis functions, each being the KL expansion function
 227 (see Equation (4)). In appendix B we have compared the posterior estimates for the Lotka-Volterra,
 228 as obtained by NUTS, between the KL and the Haar wavelet basis. Although we obtain similar
 229 estimates using both type of basis functions, the wavelet bases require stricter error tolerances for
 230 the ODE solver. Furthermore, we have also carried out a sensitivity analysis to the choice of N , the
 231 results can be found in appendix B. Finally, we have set the value of T in Equation (4) to be the
 232 end-point of the chosen time interval for each of the models.

233 We used the `numpyro` probabilistic programming library to apply NUTS and variational inference for
 234 the ease of comparison. Furthermore, we used `google_jax` implementation of the Dormand-Prince
 235 adaptive ODE solver for integrating the **SA-ODE**. This solver provides an implementation of adjoint
 236 sensitivity that is needed for applying NUTS and variational inference. For the PMMH algorithm
 237 we implemented a vectorised particle filter in `google_jax` and implemented the adaptive MCMC in
 238 Python. The code is available at [retracted](#)

8.1 Stochastic Lotka-Volterra model

The stochastic Lotka–Volterra model (Wilkinson, 2018) has been widely used for benchmarking (see Fearnhead et al. (2014); Giagos (2010)). This model describes a population comprising of two competing species: *predators* which die with rate c_2 and reproduce with rate c_1 by consuming prey, which in turn reproduce with rate c_3 . This system can be defined using the following drift and diffusion terms (see (Wilkinson, 2018) for derivation):

$$a(\mathbf{X}_t, \boldsymbol{\theta}) = \begin{bmatrix} c_1 X_1 - c_2 X_1 X_2 \\ c_2 X_1 X_2 - c_3 X_2 \end{bmatrix}, \quad B(\mathbf{X}_t, \boldsymbol{\theta}) = \begin{bmatrix} c_1 X_1 & c_2 X_1 X_2 \\ -c_2 X_1 X_2 & c_3 X_2 + c_2 X_1 X_2 \end{bmatrix}, \quad (20)$$

where we denote by X_1, X_2 the prey and predator species respectively. The parameter vector is the rate constants: $\boldsymbol{\theta} = (c_1, c_2, c_3)$, and the task is to estimate these given a noise corrupted sample path from the above system. We generate such a sample path using the Euler-Maruyama discretisation, with initial values $\mathbf{x}_0 = (100., 100.)$, between the time interval $[0 : 0.1 : 50]$. A set of 10 evenly spaced values from this path corrupted with Gaussian noise with $\sigma = 10$ constitute the observations \mathbf{y} . Following Golightly and Wilkinson (2011) we consider \mathbf{x}_0, σ to be known and thus we are left with estimating the rate constants and the expansion coefficients.

The likelihood, for the series approximation, is a Gaussian, $p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{Z}, \sigma) = \prod_i \mathcal{N}(X(t_i; \boldsymbol{\theta}, \mathbf{Z}), \sigma^2 \mathbb{I})$, where \mathbb{I} is a 2×2 identity matrix. We place a $\text{Gamma}(1, 2)$ prior on all the parameters. We placed the following priors on the rates: $c_1 \sim \text{Beta}(2, 1)$ $c_2 \times 100 \sim \text{Half } \mathcal{N}(0, 1)$ and $c_3 \sim \text{Beta}(1, 2)$. The priors on the expansion coefficients are just standard Gaussians.

We ran two chains of PMMH, each for 100,000 iterations from slightly differing initial states. We discarded the first 50,000 iterations as burnin for each chain and thinned accordingly to have 1000 samples representing the gold standard posterior estimate. For the series approximation two chains of NUTS is run for 1000 iterations after an initial 500 warmup iterations. The NUTS samples, from the two chains, are then thinned to obtain 1000 samples. VI with the series approximation was run for 30,000 iterations with $L = 1$. For plotting and summarising the posterior distributions, and comparing to the gold standard we used 1000 samples from the variational approximation for this and the subsequent (SIR model) example.

We present summaries of the posterior marginals in Table 1. Plots of the posterior diffusion paths are given in appendix C. We notice close match between the posteriors estimated with the series approximation to that of the PMMH. Interestingly, the variational inference estimate did not show underestimation of posterior variance (similar observations were made in Ghosh et al. (2021)). NUTS produced a relative ESS (averaged over the three parameters) of 0.89 and PMMH produced 0.05. We repeated the inference process using two more realisations of the artificial noise. Additional results are summarised in appendix C.

8.2 The SIR compartmental model

The SIR model (Anderson et al., 1992) of infectious disease models the number of susceptible (S), infected (I), and recovered (R) people in a population subjected to an epidemic. The stochastic version of the SIR model, for a population of N people, can be defined using an SDE with the following drift and diffusion terms (see Fuchs (2013) for derivation):

$$a(\mathbf{X}_t, \boldsymbol{\theta}) = \begin{bmatrix} -\beta SI \\ \beta SI - \gamma I \end{bmatrix}, \quad B(\mathbf{X}_t, \boldsymbol{\theta}) = \begin{bmatrix} \beta SI & -\beta SI \\ -\beta SI & \beta SI + \gamma I \end{bmatrix}, \quad (21)$$

where the infection β and recovery γ rates are unknown parameters. Also, we have $N = S(t) + I(t) + R(t)$. We used this SDE to model an outbreak of influenza at a boys boarding school in 1978 (Jackson et al., 2013). This particular dataset was previously used for benchmarking in Ryder et al. (2018). This data consists of the number of infections I_{obs} for a period of 14 days. The population size is $N = 763$. In addition to β, γ we also estimated the fractional initial susceptibility, $s_0 = S(t=0)/N$, assuming the initial recovered fraction $r_0 = 0$ and thus $i_0 = 1 - s_0$. As this is count data we have used a Poisson likelihood while using the series approximation, $p(y(t)|\beta, \gamma, s_0, \mathbf{Z}) = \text{Poisson}(I(t))$, and placed the following priors: $\beta, \gamma \sim \text{Gamma}(2, 2)$ and $s_0 \sim \text{Beta}(2, 1)$.

For this problem we ran two chains of PMMH, each for 200,000 iterations, again started with slightly different initial states. In this case we discarded the first 100,000 iterations as burnin for each chain

Table 1: We summarise the **mean \pm standard deviation** of the posterior distribution of each parameter for the two models. **VI** and **NUTS** are using the **SA-ODE** model. These were run on a 3.6 GHz machine with 16 GB memory

| THE SIR MODEL | | | | |
|-------------------------------------|------------|---------------------|---------------------|---------------------|
| θ | TRUE VALUE | PMMH | VI | NUTS |
| β | – | 1.8427 ± 0.0719 | 1.8069 ± 0.1319 | 1.8479 ± 0.1413 |
| γ | – | 0.4875 ± 0.0190 | 0.4849 ± 0.0278 | 0.4851 ± 0.0258 |
| s_0 | – | 0.9964 ± 0.0010 | 0.9957 ± 0.0010 | 0.9959 ± 0.0014 |
| THE STOCHASTIC LOTKA-VOLTERRA MODEL | | | | |
| c | TRUE VALUE | PMMH | VI | NUTS |
| c_1 | 0.5 | 0.5081 ± 0.0261 | 0.4846 ± 0.0205 | 0.4961 ± 0.0219 |
| $100 \times c_2$ | 0.25 | 0.2492 ± 0.0113 | 0.2417 ± 0.0092 | 0.2454 ± 0.0096 |
| c_3 | 0.3 | 0.2965 ± 0.0159 | 0.2872 ± 0.0115 | 0.2924 ± 0.0126 |

and thinned accordingly to obtain 1000 samples that represent the gold standard. We used the same setup, as was used in the previous example, for applying NUTS and VI with the series approximation. We present summaries of the posterior marginals in Table 1. The model fit plot is shown in appendix C. In this example NUTS produced a relative ESS (again averaged over the three parameters) of 0.95 and PMMH produced 0.05.

9 System identification task

So far we have concentrated on estimating the parameters of an SDE whose drift and diffusion terms are known. This is the context of classical inference in all mechanistic modelling scenarios. In the case of system identification we may have no knowledge about the drift or the diffusion functions. That is we are given a time-series dataset and the goal is then fit a black-box function that can predict the response at future time points. For this type of task Chen et al. (2018) proposed an ODE model whose velocity field is given by a neural network. Several extensions to this modelling framework have been proposed recently. We introduce here a (Stratonovich) SDE version of such a Neural ODE model, where the drift and the diffusion functions are modelled with neural networks. Such an SDE is given by

$$d\mathbf{X}_t = \mathbf{f}(\mathbf{X}_t, \boldsymbol{\eta}_m)dt + \mathbf{g}(\mathbf{X}_t, \boldsymbol{\eta}_s) \circ d\mathbf{W}_t, \quad \mathbf{X}_0 = \mathbf{x}_0, \quad (22)$$

where \mathbf{f}, \mathbf{g} are neural networks with parameters $\boldsymbol{\eta}_m$ and $\boldsymbol{\eta}_s$ respectively. By applying the series approximation we get the following **SA-ODE**:

$$\frac{d\hat{\mathbf{X}}}{dt} = \mathbf{f}(\hat{\mathbf{X}}_t, \boldsymbol{\eta}_m) + \mathbf{g}(\hat{\mathbf{X}}_t, \boldsymbol{\eta}_s) \sum_{i=1}^N \mathbf{Z}_i \Phi_i(t), \quad (23)$$

where $\{\mathbf{Z}_i\}_{i=1}^N$ are expansion coefficients. Notice that if we make the diffusion term zero then we recover the usual Neural ODE (Chen et al., 2018). The diffusion term can be treated as an additional regularisation from a Neural ODE perspective. Interestingly, the series approximation provides a straightforward approach to regularise Neural ODE models.

Unlike a latent ODE (a time-series VAE) formulation as proposed in Chen et al. (2018)], we intend to use the variational inference formulation that we have setup in this paper. That is we place priors on $\boldsymbol{\eta}_m, \boldsymbol{\eta}_s$ and try to infer their posterior distribution. Once, a variational approximation to the joint posterior $p(\boldsymbol{\eta}_m, \boldsymbol{\eta}_s, \mathbf{Z}|\mathbf{y})$ is learnt, we can simply generate future predictions using the posterior predictive distribution. In this case we have used a BNAF as the variational approximation (see appendix D for details).

Here our goal is to investigate that whether the diffusion term and its series approximation can regularise the Neural ODE enough to stop overfitting, without any other additional regularisation.

Table 2: Average MSE on future time points.

| MODEL | TEST ERROR | REFERENCE |
|-------------------------|-------------------------------------|------------------------|
| GPDM | 126.46 ± 34 | WANG ET AL. (2007) |
| VGPLVM | 142.18 ± 1.92 | DAMIANOU ET AL. (2011) |
| DTsBN-S | 80.21 ± 0.04 | GAN ET AL. (2015) |
| NPODE | 45.74 | HEINONEN ET AL. (2018) |
| NEURAL ODE | 87.23 ± 0.02 | CHEN ET AL. (2018) |
| ODE ² VAE | 93.07 ± 0.72 | YILDIZ ET AL. (2019) |
| ODE ² VAE-KL | 15.99 ± 4.16 | YILDIZ ET AL. (2019) |
| BnSA-ODE | 41.21 ± 36.07 | THIS WORK |

9.1 Motion capture task

To demonstrate the efficacy of our Bayesian neural SA-ODE (which we denote as **BnSA-ODE**) model we have used a dataset extracted from the CMU motion capture library. This dataset consists of 43 walking sequences of several subjects, each of which is fitted separately. Heinonen et al. (2018) used this dataset to demonstrate an ODE model, npODE, whose velocity field was modelled using a Gaussian process. In Yildiz et al. (2019) this dataset was used to demonstrate a second order extension, ODE²VAE, to the latent ODE model of Chen et al. (2018). In Yildiz et al. (2019) the same preprocessing as was done in Wang et al. (2007) was used which left a 50-dimensional time-series. This preprocessed dataset was made available in Yildiz et al. (2019) and this is what we have used. Additionally, we have scaled the data to be within a $[-1, 1]$ range. We followed the convention in Yildiz et al. (2019) of using the first two third of each sequence for training/validation and the rest for testing purpose i.e evaluating a MSE score.

We considered a 3-dimensional state-space and used a decoder network that maps the latent state-space to the 50-dimensional output. Moreover, we consider a Gaussian measurement model and we estimate the noise standard deviation on which we place a Half $\mathcal{N}(0, 1)$ prior. We place a standard Gaussian prior on all the neural network parameters. We used a shallow architecture with 10 hidden units for the drift, diffusion and decoder networks (see appendix D for details). We retained the same time interval for numerical solution of the **BnSA-ODE** as was used in Yildiz et al. (2019). For the series expansion we have used $N = 15$ basis functions and due to the need of extrapolation we set the value of T for the KL basis functions (see Equation (4)) to be twice the end-point of the training interval. We ran VI for 1200 iterations with the RMSProp where the step-size was set to 10^{-2} . We set $L = 40$. Running further iterations results in over-fitting. We evaluated the MSE between the data and the mean of the posterior predictive distribution at the test time-points for each sequence. We evaluated the posterior predictive distribution pointwise using 50 samples from the variational approximation.

BnSA-ODE performs better than the Neural ODE on average and only ODE²VAE with a KL penalty produces smaller average test error. However, **BnSA-ODE** performs poorly for some sequences and that is why it has a much higher variance of the test error. One way to enhance the **BnSA-ODE** model is to use a deeper architecture. Alternatively, **BnSA-ODE** can be integrated within a latent ODE framework.

10 Conclusion

We presented a method to carry out inference for a nonlinear SDE by approximating the SDE with an ODE. We do this using a truncated series expansion of Brownian motion. For the parameter estimation task, of mechanistic models, this method produce close approximation to the estimate obtained using a particle MCMC algorithm. Which is considered the state-of-the-art for this type of inference problems. Furthermore, our approach lets us easily convert a Neural ODE to a Neural SDE and vice-versa. We exploited this facet of our approach by approximating a Neural SDE using the series expansion for a system identification task. In future we want extend the proposed approach to develop a latent SDE model using the VAE formulation.

References

- Roy M Anderson, B Anderson, and Robert M May. *Infectious diseases of humans: dynamics and control*. Oxford university press, 1992.
- Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3): 269–342, 2010.
- Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in neural information processing systems*, pages 6571–6583, 2018.
- Andreas C. Damianou, Michalis K. Titsias, and Neil D. Lawrence. Variational gaussian process dynamical systems. In *Proceedings of the 24th International Conference on Neural Information Processing Systems*, NIPS’11, page 2510–2518, 2011.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. Block neural autoregressive flow. In *Uncertainty in Artificial Intelligence*, pages 1263–1273. PMLR, 2020.
- Paul Fearnhead, Vasileios Giagos, and Chris Sherlock. Inference for reaction networks using the linear noise approximation. *Biometrics*, 70(2):457–466, 2014.
- Christiane Fuchs. *Inference for diffusion processes: with applications in life sciences*. Springer Science & Business Media, 2013.
- Zhe Gan, Chunyuan Li, Ricardo Henao, David Carlson, and Lawrence Carin. Deep temporal sigmoid belief networks for sequence modeling. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’15, page 2467–2475, 2015.
- Sanmitra Ghosh, Paul Birrell, and Daniela De Angelis. Variational inference for nonlinear ordinary differential equations. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 2719–2727. PMLR, 13–15 Apr 2021.
- Vasileios Giagos. *Inference for auto-regulatory genetic networks using diffusion process approximations*. PhD thesis, Lancaster University, 2010.
- Mike Giles and Paul Glasserman. Smoking adjoints: Fast monte carlo greeks. *Risk*, 19(1):88–92, 2006.
- Emmanuel Gobet and Rémi Munos. Sensitivity analysis using itô–malliavin calculus and martingales, and application to stochastic optimal control. *SIAM Journal on control and optimization*, 43(5): 1676–1713, 2005.
- Andrew Golightly and Colin S Gillespie. Simulation of stochastic kinetic models. In *In Silico Systems Biology*, pages 169–187. Springer, 2013.
- Andrew Golightly and Darren J Wilkinson. Bayesian inference for nonlinear multivariate diffusion models observed with error. *Computational Statistics & Data Analysis*, 52(3):1674–1693, 2008.
- Andrew Golightly and Darren J Wilkinson. Bayesian parameter inference for stochastic biochemical network models using particle markov chain monte carlo. *Interface focus*, 1(6):807–820, 2011.
- Neil Gordon, David Salmond, and Craig Ewing. Bayesian state estimation for tracking and guidance using the bootstrap filter. *Journal of Guidance, Control, and Dynamics*, 18(6):1434–1443, 1995.
- Markus Heinonen, Cagatay Yildiz, Henrik Mannerström, Jukka Intosalmi, and Harri Lähdesmäki. Learning unknown ode models with gaussian processes. In *International Conference on Machine Learning*, pages 1959–1968. PMLR, 2018.
- Matthew D Hoffman and Andrew Gelman. The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *Journal of Machine Learning Research*, 15(1):1593–1623, 2014.

399 Charlotte Jackson, Emilia Vynnycky, Jeremy Hawker, Babatunde Olowokure, and Punam Mangtani.
400 School closures and influenza: systematic review of epidemiological studies. *BMJ open*, 3(2),
401 2013.

402 Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to
403 variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.

404 D P Kingma, Max Welling, et al. Auto-encoding variational bayes. In *Proceedings of the International
405 Conference on Learning Representations (ICLR)*, volume 1, 2014.

406 Xuechen Li, Ting-Kam Leonard Wong, Ricky T. Q. Chen, and David Duvenaud. Scalable gradients
407 for stochastic differential equations. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings
408 of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108
409 of *Proceedings of Machine Learning Research*, pages 3870–3882. PMLR, 26–28 Aug 2020.

410 Simon M. J. Lyons, Amos J. Storkey, and Simo Särkkä. The coloured noise expansion and parameter
411 estimation of diffusion processes. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher
412 J. C. Burges, Léon Bottou, and Kilian Q. Weinberger, editors, *Advances in Neural Information
413 Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012.
414 Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States*, pages
415 1961–1969, 2012.

416 Simon MJ Lyons, Simo Särkkä, and Amos J Storkey. Series expansion approximations of brownian
417 motion for non-linear kalman filtering of diffusion processes. *IEEE Transactions on Signal
418 Processing*, 62(6):1514–1524, 2014.

419 EJ McShane. Stochastic differential equations and models of random processes. In *Contributions to
420 Probability Theory*, pages 263–294. University of California Press, 2020.

421 Bernt Oksendal. *Stochastic differential equations: an introduction with applications*. Springer
422 Science & Business Media, 2013.

423 Christopher Rackauckas, Yingbo Ma, Vaibhav Dixit, Xingjian Guo, Mike Innes, Jarrett Revels, Joakim
424 Nyberg, and Vijay Ivaturi. A comparison of automatic differentiation and continuous sensitivity
425 analysis for derivatives of differential equation solutions. *arXiv preprint arXiv:1812.01892*, 2018.

426 Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and
427 approximate inference in deep generative models. In *Proceedings of the 31st International
428 Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages
429 1278–1286. PMLR, 2014.

430 Tom Ryder, Andrew Golightly, A. Stephen McGough, and Dennis Prangle. Black-box variational
431 inference for stochastic differential equations. In *Proceedings of the 35th International Conference
432 on Machine Learning*, pages 4423–4432, 2018.

433 T. Tieleman and G Hinton. Lecture 6.5-rmsprop: divide the gradient by a running average of its
434 recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 4, 2016.

435 Michalis Titsias and Miguel Lázaro-Gredilla. Doubly stochastic variational bayes for non-conjugate
436 inference. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32
437 of *Proceedings of Machine Learning Research*, pages 1971–1979. PMLR, 2014.

438 Jack M Wang, David J Fleet, and Aaron Hertzmann. Gaussian process dynamical models for human
439 motion. *IEEE transactions on pattern analysis and machine intelligence*, 30(2):283–298, 2007.

440 Darren J Wilkinson. *Stochastic modelling for systems biology*. CRC press, 2018.

441 Eugene Wong and Moshe Zakai. On the Convergence of Ordinary Integrals to Stochastic Integrals.
442 *The Annals of Mathematical Statistics*, 36(5):1560 – 1564, 1965. doi: 10.1214/aoms/1177699916.
443 URL <https://doi.org/10.1214/aoms/1177699916>.

444 Jichuan Yang and Harold J Kushner. A monte carlo method for sensitivity analysis and parametric
445 optimization of nonlinear stochastic systems. *SIAM journal on control and optimization*, 29(5):
446 1216–1249, 1991.

447 Cagatay Yildiz, Markus Heinonen, and Harri Lähdesmäki. ODE2VAE: deep generative second order
 448 odes with bayesian neural networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer,
 449 Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information*
 450 *Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019,*
 451 *NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13412–13421, 2019.

452 Checklist

453 The checklist follows the references. Please read the checklist guidelines carefully for information on
 454 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
 455 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
 456 the appropriate section of your paper or providing a brief inline description. For example:

- 457 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 458 • Did you include the license to the code and datasets? **[No]** The code and the data are
 459 proprietary.
- 460 • Did you include the license to the code and datasets? **[N/A]**

461 Please do not modify the questions and only use the provided macros for your answers. Note that the
 462 Checklist section does not count towards the page limit. In your paper, please delete this instructions
 463 block and only keep the Checklist section heading above along with the questions/answers below.

464 1. For all authors...

- 465 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 466 contributions and scope? **[Yes]**
- 467 (b) Did you describe the limitations of your work? **[Yes]**
- 468 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 469 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 470 them? **[Yes]**

471 2. If you are including theoretical results...

- 472 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 473 (b) Did you include complete proofs of all theoretical results? **[N/A]**

474 3. If you ran experiments...

- 475 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 476 mental results (either in the supplemental material or as a URL)? **[Yes]**
- 477 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 478 were chosen)? **[Yes]**
- 479 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 480 ments multiple times)? **[Yes]**
- 481 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 482 of GPUs, internal cluster, or cloud provider)? **[Yes]**

483 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 484 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
- 485 (b) Did you mention the license of the assets? **[N/A]**
- 486 (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
- 487 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 488 using/curating? **[N/A]**
- 489 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 490 information or offensive content? **[N/A]**

491 5. If you used crowdsourcing or conducted research with human subjects...

- 492 (a) Did you include the full text of instructions given to participants and screenshots, if
 493 applicable? **[N/A]**

- 494 (b) Did you describe any potential participant risks, with links to Institutional Review
 495 Board (IRB) approvals, if applicable? [N/A]
 496 (c) Did you include the estimated hourly wage paid to participants and the total amount
 497 spent on participant compensation? [N/A]

498 **A Converting from Itô to Stratonovich formulation**

499 The Stratonovich SDE given by

$$dX_t^i = \tilde{a}(X_t^i, \boldsymbol{\theta})dt + \sum_{j=1}^K \sqrt{b^{i,j}(X_t, \boldsymbol{\theta})} \circ dW_t^j, \quad i, j = 1, \dots, K, \quad (24)$$

500 with the same solutions as the K -dimensional Ito SDE driven by a K -dimensional Wiener process
 501 given by

$$dX_t^i = a(X_t^i, \boldsymbol{\theta})dt + \sum_{j=1}^K \sqrt{b^{i,j}(X_t, \boldsymbol{\theta})} dW_t^j, \quad i, j = 1, \dots, K, \quad (25)$$

502 has a drift coefficient that is defined component-wise as

$$\tilde{a}(X_t^i, \boldsymbol{\theta})dt = a(X_t^i, \boldsymbol{\theta})dt + \sum_{k=1}^K \sum_{j=1}^K \sqrt{b^{k,j}(X_t, \boldsymbol{\theta})} \frac{\partial \{\sqrt{b^{i,j}(X_t, \boldsymbol{\theta})}\}}{\partial X_t^k}. \quad (26)$$

503 The above partial derivatives can be evaluated either through hand-calculation or using AD itself.
 504 The later is what we did in our experiments.

505 **B Wavelet basis function**

506 For the multivariate case uniform convergence of the SA-ODE to the corresponding Stratonovich SDE
 507 is guaranteed if one chooses Haar wavelets as an orthonormal basis in which to expand the driving
 508 Brownian motion (Lyons et al., 2014). Thus, we wanted to compare the approximation achieved by
 509 using the KL expansion with that of Haar wavelets. But, before we delve into the comparison, let us
 510 briefly introduce the Haar wavelets: which is a complete orthonormal basis of $L^2[0, T]$.

511 The Haar wavelets are parameterised by two natural numbers: the scale $n \geq 0$ and the shift
 512 $0 \leq k < 2^n$. The first wavelet is defined as

$$\psi_{0,0}(t) = \begin{cases} 1 & 0 \leq t < \frac{T}{2} \\ -1 & \frac{T}{2} \leq t \leq T \\ 0 & \text{otherwise} \end{cases}. \quad (27)$$

513 Further wavelets are defined by rescaling $\psi_{0,0}$ so that it is non-zero only on some sub-interval of
 514 $[0, T]$, while ensuring that the wavelet still has unit norm. In general,

$$\psi_{n,k}(t) = \frac{2^{n/2}}{\sqrt{T}} \psi_{0,0}(2^n - kt). \quad (28)$$

515 Thus, $\psi_{1,0}$ is a copy of $\psi_{0,0}$ restricted to $[0, T/2]$, and $\psi_{1,1}$ is a copy restricted to $[T/2, T]$. Further-
 516 more, we add the constant function $\psi_* = \frac{1}{\sqrt{T}}$ to form a complete basis. To be consistent with the
 517 notation introduced in section 4 we set $\phi_1 = \psi_*$, $\phi_2 = \psi_{0,0}$, $\phi_3 = \psi_{1,0}$ and so on.

518 To compare the wavelet and KL basis we ran NUTS and VI on the simulated dataset used in section
 519 8.1, with the same algorithmic settings retained. We set $N = 10$ and $T = 50$ as was done in section
 520 8.1. Marginal densities of the parameters are plotted in Figure 2. Although we get similar estimates,
 521 as in (Lyons et al., 2014, 2012), following which we used KL expansion throughout, the wavelets
 522 required stricter error tolerances for the ODE solver. We thus recommend the usage of a stiff solver
 523 when using the wavelet basis. Extension of AD for a stiff solver can be done using the custom op
 524 creation method of Ghosh et al. (2021). However, a JIT compiled solver, provided with Jax, is faster.
 525

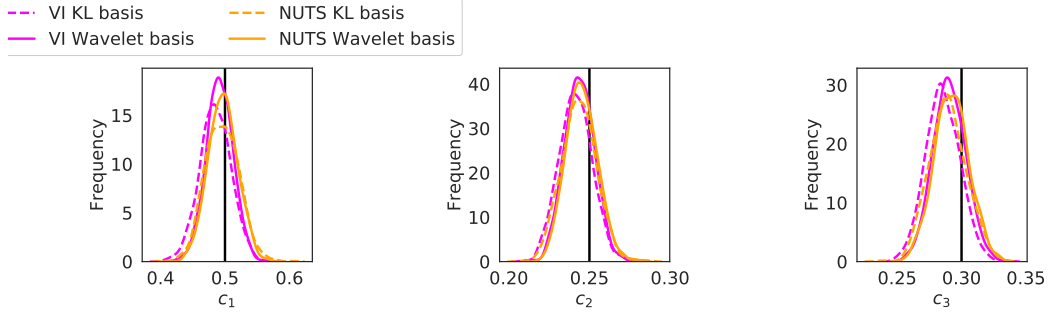


Figure 2: We compare the marginal density of the Lotka-Volterra parameters as obtained by NUTS and VI, using the wavelets (solid lines) and KL (dashed lines) basis. The black vertical lines indicate the true parameter values.

526 B.1 Sensitivity to truncation of the series: the choice of N

527 To carry out this sensitivity analysis we use three realisations of the simulated dataset. The first
 528 one is used in section 8.1 and the estimates on the rest are summarised in section C.1. We ran the
 529 NUTS algorithm with $N = 3, 5, 8, 10$ respectively and measured the maximum mean discrepancy
 530 (MMD) to the corresponding estimates obtained by PMMH for each dataset. We used the KL basis
 531 throughout and retained all algorithmic settings for NUTS as in section 8.1. Figure 3 compares the
 MMD for specific choices of N . It is apparent that with $N = 8$ the MMD plateaus.

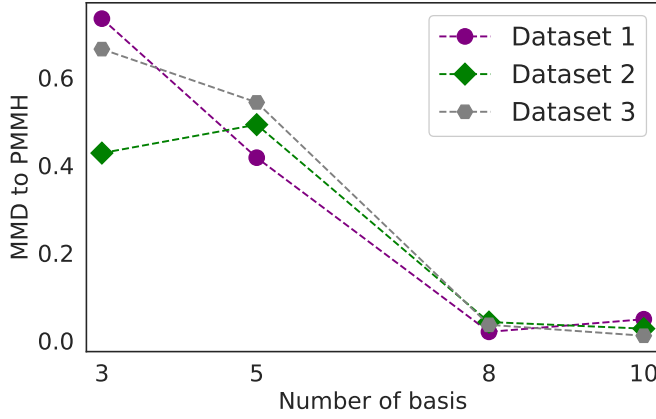


Figure 3: We compare the effect of increasing the number of expansion terms N on the MMD between the PMMH estimate and the estimates obtained from running NUTS, with $N = 3, 5, 8, 10$, for the Lotka-Volterra model. All estimates used 1000 samples from the posterior.

532

533 C Plots of the latent diffusion for Lotka-Volterra

534 Here we evaluate the posterior predictive distribution of the latent diffusion $\hat{X}^*(\theta)|_{\theta, \mathbf{Z}, \mathbf{x}_0 \sim p(\theta, \mathbf{Z}, \mathbf{x}_0 | \mathbf{y})}$,
 535 using the SA-ODE approximation, on a finer grid t^* than the observations. The posterior predictive
 536 distribution is evaluated pointwise using samples obtained from running NUTS or drawn from the VI
 537 approximation. We used the posterior distributions estimated using the simulated dataset that was
 538 used in section 8.1. The mean and the 95% credible intervals of \hat{X}^* are shown in Figure 4.

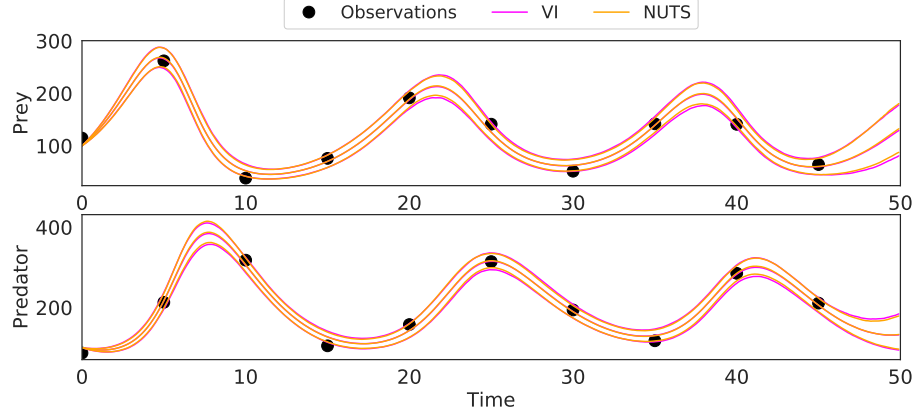


Figure 4: Mean and 95% credible intervals of the posterior predictive distribution of the latent diffusion \hat{X}^* , for the Lotka-Volterra model.

Table 3: We summarise the **mean \pm standard deviation** of the posterior distribution of each parameter for the The Stochastic Lotka-Volterra model. **VI** and **NUTS** are using the **SA-ODE** model.

| THE STOCHASTIC LOTKA-VOLTERRA MODEL: DATASET 2 | | | | |
|--|------------|---------------------|---------------------|---------------------|
| c | TRUE VALUE | PMMH | VI | NUTS |
| c_1 | 0.5 | 0.5194 ± 0.0273 | 0.5026 ± 0.0205 | 0.5093 ± 0.0224 |
| $100 \times c_2$ | 0.25 | 0.2557 ± 0.0122 | 0.2477 ± 0.0092 | 0.2494 ± 0.0095 |
| c_3 | 0.3 | 0.2937 ± 0.0163 | 0.2883 ± 0.0125 | 0.2906 ± 0.0132 |
| THE STOCHASTIC LOTKA-VOLTERRA MODEL: DATASET 3 | | | | |
| c | TRUE VALUE | PMMH | VI | NUTS |
| c_1 | 0.5 | 0.5172 ± 0.0213 | 0.4970 ± 0.0159 | 0.5210 ± 0.0204 |
| $100 \times c_2$ | 0.25 | 0.2540 ± 0.0080 | 0.2432 ± 0.0067 | 0.2514 ± 0.0076 |
| c_3 | 0.3 | 0.3140 ± 0.0117 | 0.3048 ± 0.0101 | 0.3130 ± 0.0111 |

539 C.1 Additional results for Lotka-Volterra

540 To benchmark the methods with this model we used simulated data for which the corresponding
541 estimates were summarised in Table 1. In addition to this dataset we generated two more datasets
542 using two new realisation of the artificial noise corruption. The results for these additional datasets
543 are summarised in Table 3. All the algorithmic and model specific settings were kept the same as was
544 used in section 8.1.

545 C.2 Model fit plot for the SIR model

546 We evaluated the posterior predictive distribution $p(y^*|y)$ on a finer time grid t^* , using the **SA-ODE**.
547 Samples of the posterior predictive distribution were evaluated pointwise using the posterior estimates
548 obtained from running NUTS and VI. Figure 5 summarises the mean and 95% credible intervals of
549 $p(y^*|y)$.

550 D BNAF as approximating distribution

551 In BNAF each dimension of the transformed parameter ξ (see section 6.1), as a bijection from
552 a standard Normal variable, is modelled using a neural network that (autoregressively) takes the
553 preceding dimensions as inputs. Such bijections can be chained together (applied repeatedly) to form
554 a flow. The dense layers of these networks have weight matrices constructed as a block-diagonal

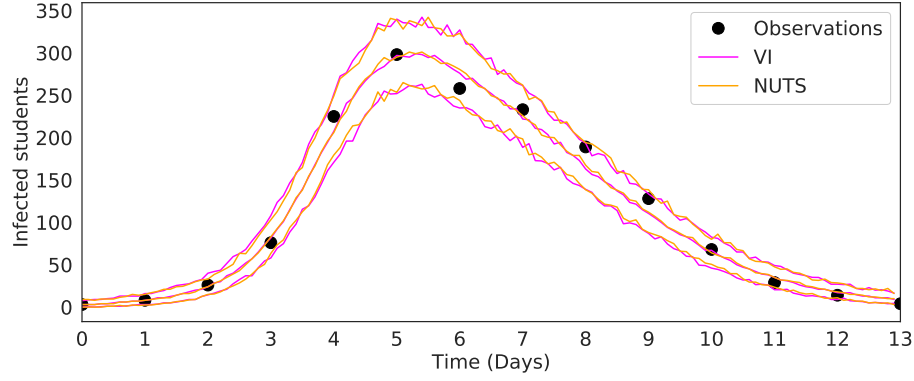


Figure 5: Mean and 95% credible intervals of the posterior predictive distribution for the SIR model.

matrix with the diagonal elements being $a \times b$ matrices themselves, for some chosen $a, b > 1$. We set $a, b = 8$. We used the validation data to choose between $a, b = 5, 8, 10$. $a, b = 8$ produced the best prediction errors. Furthermore, we set the number of flows to 1.

D.1 Architecture of the BnSA-ODE

We used a neural network with a single hidden layer to model the drift, diffusion and the decoder function. The architecture is shown in Figure 6. We used the validation data to experiment with 5, 10, 15, 25, 50 hidden units. We obtained best prediction errors on the future frames using 10 units.

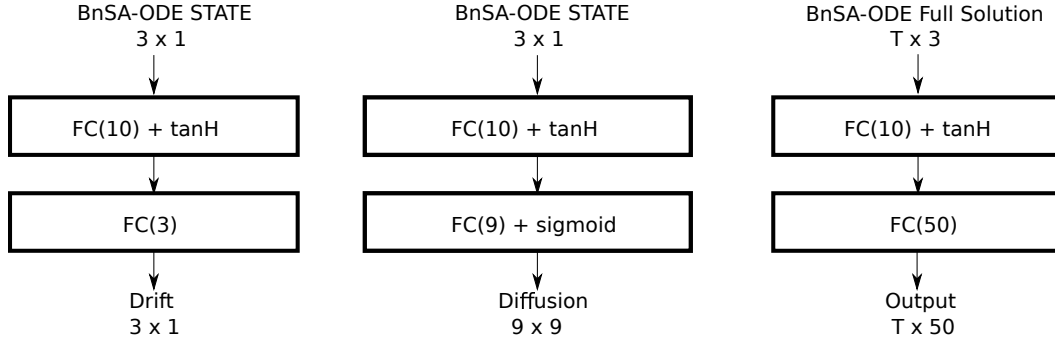


Figure 6: The **BnSA-ODE** architecture for the drift, diffusion and decoder networks that produced the lowest average validation error. We set the number of latent states $K = 3$. Note that by T here we denote the number of time points for training/validation/test data.