

COMPUTATIONAL EFFICIENCY OF INFERENCE METHODS FOR STOCHASTIC DIFFERENTIAL EQUATION DYNAMICS

ALAN AMIN

1. Introduction. Stochastic differential equations with respect to a Wiener process (SDEs) are useful implicit descriptions of certain Markov processes. They have been applied to modelling financial markets [5, 1], projecting the paths of spacecraft [1], and describing biological phenomenon [8]. They also can be used to build stochasticity into state space models [9].

Inference of SDEs has remained a challenge - in particular, the task of computationally efficient calculation of gradients of a loss function. Early methods included Monte Carlo forward differentiation of parameters of the drift and noise functions, however, these scale poorly with the number of parameters, which can be very high when these functions are high dimensional or if a deeper or wider neural network is necessary for example. A more computationally efficient method in this case is Monte Carlo backpropogation through the SDE solver given a realization of the Wiener process [?]. However, this method requires one to save the forward path of the SDE which can be memory intensive for long time series. Adjoint methods instead define a backwards SDE that one may solve to calculate gradients. Early adjoint methods considered each step of an Euler-Maruyama SDE solver a constrained optimization problem and were able to define the backwards SDEs by taking derivatives: one effort defined the backwards SDE for use in calculating monte carlo estimates of the gradient [6], while another used quadrature to approximately integrate over all paths [2]. Recently, the adjoint method has been extended for the use of any Stratonovich SDE solver [9].

Here I explore these ideas by implementing a high performance version of both this later adjoint method, and backpropogation. To achieve low level control of the SDE solver, I implement my own solver using the Euler-Heun method. I compare the accuracy of the gradients of this adjoint method, and backpropogation as well as benchmark their computational cost and ability to fit data. I will demonstrate the convergence of my solver and gradient estimates to the correct values. Then I will show that my implementation of gradient estimation using backpropogation is more computationally efficient and results in higher quality inference than my implementation of the adjoint method for gradient calculation. This investigation represents my edification in the solving and inference of SDEs, as well as an investigation into the factors that effect the utility of these two gradient estimation methods.

2. Results.

2.1. Euler-Heun method solver.

2.1.1. Setup. I consider here the SDE, from time 0 to T ,

$$dX_t = f(X_t, t)dt + g(X_t, t) \circ dW_t$$

where X_t has d dimensions and the Wiener process W has m dimensions. To investigate how the architecture of an SDE solver may affect the performance of gradient estimation techniques, and to easily implement backpropogation, I implement my own SDE solver.

The naive implementation of a solver for a Stratonovich integral is known as the Euler-Heun. In probability, as the mesh size of the grating $0 = t_0 < \dots < t_n = T$

goes to 0,

$$\sum_{i=0}^{n-1} \frac{g(X_{t_{i+1}}, t_{i+1}) + g(X_{t_i}, t_i)}{2} (W_{t_{i+1}} - W_{t_i}) \rightarrow \int_0^T g(X_t, t) \circ dW_t.$$

42 The Euler-Heun (EH) method SDE solver approximates the integral for a time step
 43 δt via a second order Runge-Kutta (RK) method to approximate the average of the
 44 drift and diffusion evaluated at the start and stop time of each step:

$$\begin{aligned} 45 \quad H_{i+1} &= X_i + f(X_i, t_i)\delta t + g(X_i, t_i)\Delta W_i \\ 46 \quad X_{i+1} &= X_i + \frac{\delta t}{2} (f(X_i, t_i) + f(H_{i+1}, t_{i+1})) \\ 47 \quad &+ \frac{1}{2} (g(X_i, t_i) + g(H_{i+1}, t_{i+1})) \Delta W_i \\ 48 \quad \Delta W_i &\sim \sqrt{\delta t} N(0, 1). \end{aligned}$$

50 The EH method is strong order 0.5 - if $X_T^{\delta t}$ is approximated using the EH method
 51 with step size δt , then

$$52 \quad (2.1) \quad \|X_T^{\delta t} - X_T\|_{L^2} \lesssim \delta t^{0.5}.$$

53 **2.1.2. Results.** Throughout, I will use one layer neural networks with tanh ac-
 54 tivation as the drift and diffusion functions f and g . As well, I will consider a two
 55 dimensional SDE with 2 dimensions of diffusion, i.e. $d = m = 2$. g will not be assumed
 56 diagonal or commutative. To allow for efficient evaluation of the neural networks, I
 57 used the @einsum macro, which avoids boundschecking by checking compatibility of
 58 the dimensions of tensors before the multiplication. I used function barriers when
 59 unpacking parameters to allow the function to specialize on the input type (functions
 60 f_NN! and g_NN! in the code). Finally, to avoid repeated memory allocation for eval-
 61 uations of intermediate values in the evaluation of f or g , I preallocate the memory to
 62 store these values; I then pass the variables associated to this memory as a parameter
 63 to the functions f_NN! and g_NN! and mutate it to these intermediate values (variables
 64 t1 and t2 in the code).

65 I implement a high performance EH solver for SDEs that takes as input, a starting
 66 point, a parametrized drift function, a parametrized diffusion function, and parame-
 67 ters (function my_SDEsolve! in code) (Fig 1A). To minimize memory allocation per
 68 step, I preallocated all the memory necessary to store intermediate values for each
 69 step and mutate these variables (function EH! in code).

70 One step of the solver, given the step size of the Wiener process, was benchmarked
 71 as having 0 memory allocations and taking 1.790 μs on my machine (section "Time
 72 Euler-Heun" in the code). Much of this computational complexity came from the
 73 need to apply the costly tanh non-linearity: the time was almost halved to 1.080
 74 μs when using a ReLU activation. Below however, I only consider the use of the
 75 tanh non-linearity unless stated otherwise. To evaluate the convergence of my solver,
 76 I used the Euler-Heun solver implemented in DifferentialEquations.jl with step size
 77 2^{-16} to approximate X_T (in this case I used $T = 1$). Since the solver had a fixed step
 78 size, it was easy to save its noise and use it in my solver to calculate $X_T^{\delta t}$ with step
 79 sizes δt ranging from 2^{-7} to 2^{-12} (section "EH is order 0.5" in the code). Empirical
 80 estimation of the L^2 distance reveal that my method converges to the correct solution
 81 and a linear regression reveals that it is of approximate order 0.522, similar to the
 82 theoretical 0.5 described in equation 2.1.

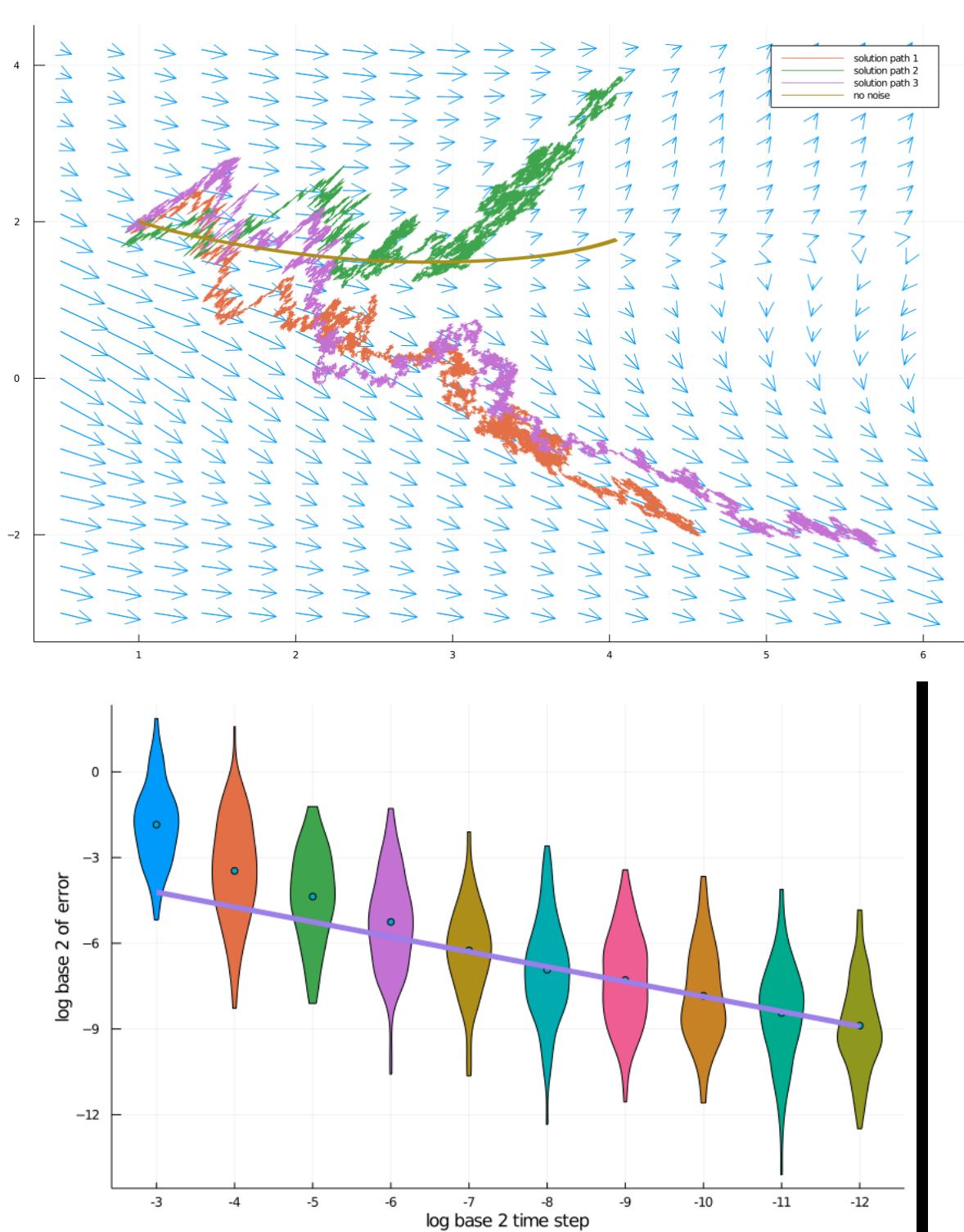


FIG. 1. A: Example paths starting at the point $(1, 2)$ given by my Euler-Heun solver given a instantiation of the drift (blue arrows) and diffusion functions. The solution to the ordinary differential equation without the drift is also shown. In this example, the stochasticity causes some paths to cross a separatrix and depart significantly from the ODE solution. B: Converge of the L^2 error defined in equation 2.1. A linear regression of the last 6 points gives a slope of -0.522 and an R^2 of 0.995 , close to the theoretical -0.5 . The violin plots contain 100 point each.

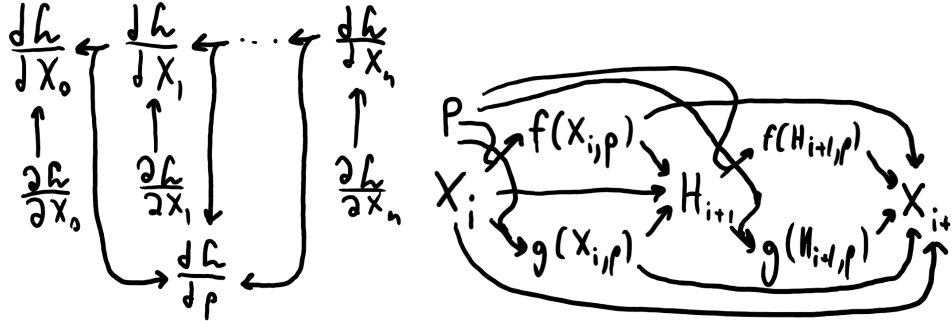


FIG. 2. A: A schematic of the calculation of the gradient of the loss with respect to the parameters and starting point by backpropagation through the solver. B: The computation graph of a single Euler-Heun step. I reverse it to calculate the pullback.

83 2.2. Backpropagation through the solver.

84 **2.2.1. Setup.** Here I consider the problem of approximating gradients of some
 85 loss function \mathcal{L} by backpropagation through the solver. Say $(X_i)_{i=0}^n$ are the steps
 86 of the SDE iterator, i.e. $X_{i+1} = EH(X_i, p)$ where EH is one step of the solver.
 87 Call $\mathcal{B}_{EH(X,p)}^p$ the p component of the backpropagation of the Euler-Heun step and
 88 $\mathcal{B}_{EH(X,p)}^X$ the X component. Then

$$89 \quad (2.2) \quad \frac{d\mathcal{L}}{dX_i} = \frac{\partial \mathcal{L}}{\partial X_i} + \mathcal{B}_{EH(X_i,p)}^X \left(\frac{d\mathcal{L}}{dX_{i+1}} \right)$$

$$90 \quad (2.3) \quad \frac{d\mathcal{L}}{dp} = \sum_{i=1}^n \mathcal{B}_{EH(X_{i-1},p)}^p \left(\frac{d\mathcal{L}}{dX_i} \right)$$

92 (Fig 2A). To calculate $\mathcal{B}_{EH(X,p)}^X$ and $\mathcal{B}_{EH(X,p)}^p$ I consider the computation graph in
 93 fig 2B.

94 **2.2.2. Results.** To generate data to define a loss function, I initialize the pa-
 95 rameters of f and g to random values and sample values from a path of the solution
 96 (section "Gen data" in the code). Below I consider the ℓ^2 loss with respect to this
 97 data.

98 When writing the pullback of the neural networks, I again used the @einsum
 99 macro (functions pullback_f! and pullback_g! in the code). When these functions
 100 were called, I again preallocate memory to store intermediate calculations and pass
 101 these variables as parameters. I also again use function barriers when unpacking
 102 parameters.

103 I implement a high performance pullback of the SDE solver that takes as input,
 104 the forward path, a parametrized drift function, a parametrized diffusion function,
 105 pullbacks of the drift and diffusion that add to the derivatives with respect to the
 106 parameters instead of overwrite them, and parameters (function my_SDEbackprop!
 107 in code).

108 One iteration of the Euler-Heun pullback, given the step size of the Wiener
 109 process, was benchmarked as having 0 memory allocations and taking 18.799 μs on my
 110 machine. The derivatives calculated by the backpropagation algorithm are validated
 111 by those calculated by finite differencing (section "Backprop" in code) (Fig 3). Small
 112 differences between the gradient calculated by backpropagation and that calculated

113 by finite differencing are likely due to the finite differencing algorithm only returning
 114 values up to a certain precision or precision loss when calculating gradients through
 115 a tanh function; in particular, the step size doesn't affect the agreement between the
 116 two methods.

117 **2.3. Adjoint method.**

2.3.1. Setup. I now implement the adjoint method of gradient estimation de-
 scribed in [9]. Call $\mathcal{B}_{f(X,p)}$ the pullback of the drift function and $\mathcal{B}_{g_i(X,p)}$ the pullback
 of the i -th column diffusion function. Consider the backwards SDE defined in [9]

$$d(A_t, B_t) = \mathcal{B}_{f(X,p)}(A_t) + \sum_{i=1}^m \mathcal{B}_{g_i(X,p)}(A_t) \circ d\tilde{W}_{i,t}$$

118 With $A_T = \frac{\partial \mathcal{L}}{\partial X_T}$, $B_T = \vec{0}$ and A jumping at time t' by a magnitude of $\frac{\partial \mathcal{L}}{\partial X_{t'}}$. Then,
 119 by the results in [9], $A_0 = \frac{d\mathcal{L}}{dX_0}$ and $B_0 = \frac{d\mathcal{L}}{dp}$. It is also shown that the forward path
 120 X satisfies the backwards SDE

121 (2.4)
$$dX_t = -f(X_t, p)dt - g(X_t, p) \circ d\tilde{W}_t.$$

122 Crucially, the stochastic integrals in these results are of the Stranovich interpretation.
 123 Thus, the EH method may be used to evaluate them.

Solving such a backwards SDE is identical to solving the forwards SDE

$$dY_t = -f(Y_t, p)dt - g(Y_t, p) \circ dV_t$$

124 where $Y_t = X_{T-t}$ and $V_t = -W_{T-t}$. The major difficulty in implementing this model
 125 is that the forwards and backwards noise must be the same. One solution for fixed
 126 step-size solvers is to save the path of the noise and reuse it when going backwards.
 127 [9] avoids this by implementing using what they call a Brownian tree. Here, a single
 128 random seed determines the value of the Wiener process path at all points $(i2^{-m})_{i=0}^{2^m}$
 129 for a given m . The computation cost of querying the Brownian tree scales linearly
 130 with m .

131 **2.3.2. Results.** I use the implementation of the Brownian tree at <http://github.com/SciML/DiffEqNoiseProcess.jl/pull/65>.

132 To implement a high performance version of the adjoint method, I simply plugged
 133 in the pullback function for the neural networks f and g into my SDE solver (section
 134 "Backwards SDE" in the code). I show that the forward and reverse solutions to the
 135 SDE are the same as described by equation 2.4 (Fig 4). Furthermore, I show that
 136 the gradients of the adjoint method agree with those calculated by finite differencing,
 137 with the agreement becoming better as the step size used in the solver decreases, and
 138 thus the solution becomes more accurate (Fig 3).

139 At a step size of 2^{-12} , the Brownian tree algorithm takes $5.750 \mu s$ and makes
 140 100 allocations of 10 kB. This is almost an order of magnitude more resources than a
 141 step of the solver. Since the Euler-Heun method is a fixed step size solver and saving
 142 the path does not incur a significant memory cost, hereafter, I opt instead to simply
 143 save the noise from the forward pass for use in the backwards pass. As well, instead
 144 of using the forward solution as calculated in the backwards pass, I simply save the
 145 forwards pass, incurring minimal memory cost, as pass this path to the backwards
 146 path; I expect this to increase the stability of the solver.

148 **2.4. Comparing the adjoint method and backpropogation.**

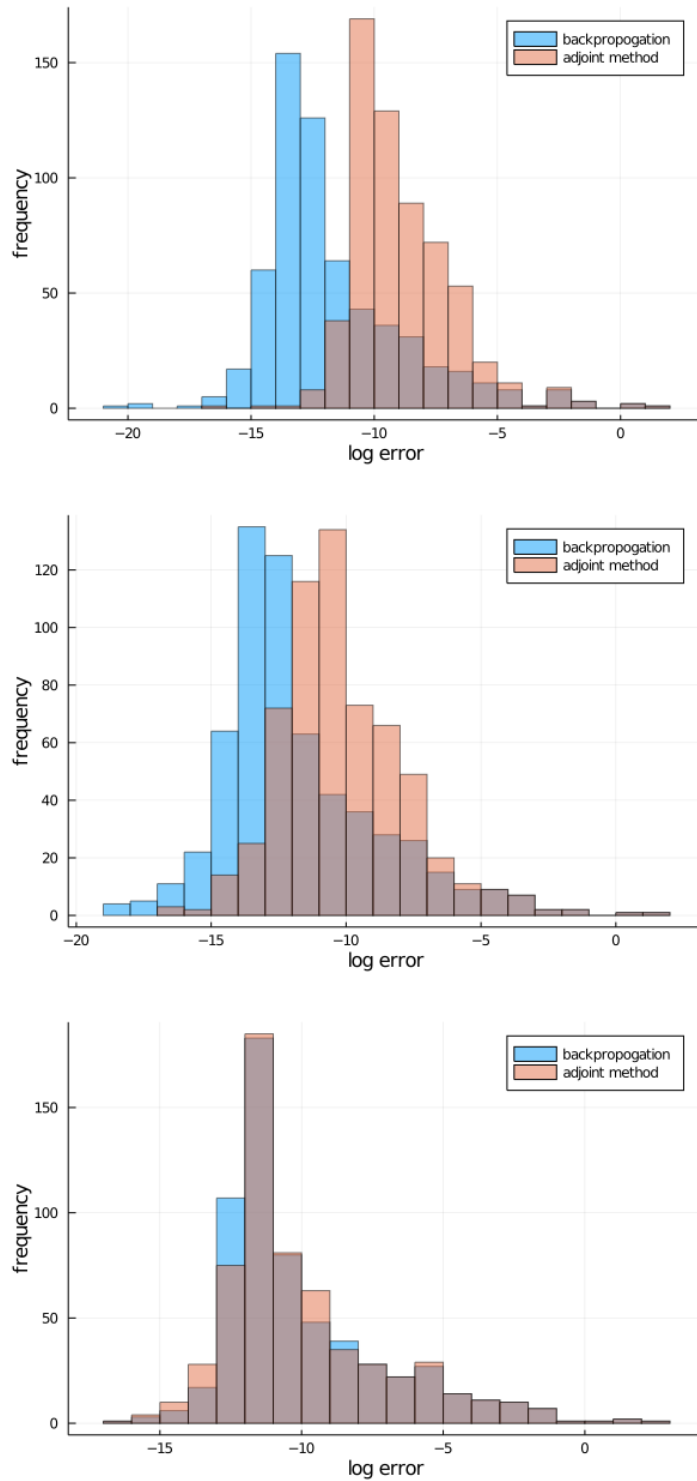


FIG. 3. Fractional error of parameter gradients calculated using back propagation or the adjoint method in comparison to those calculated by finite differencing for step sizes A: 2^{-10} , B: 2^{-12} , C: 2^{-14} . These plots are each for a single realization of the Wiener process and plot all parameters of f , g and the starting point. Notice the x axis scaling is not the same between plots (sorry!).

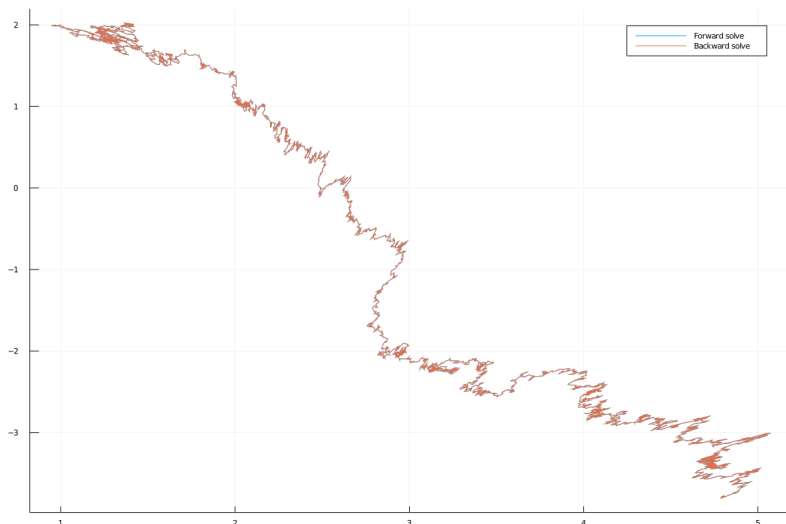


FIG. 4. Shown here is the forwards and backwards solution to an SDE for a given realization of the Wiener process. The two paths are superimposed and thus visually identical.

149 **2.4.1. Performance.** Here I will compare the computational complexity and
 150 accuracy of the gradient calculation by the adjoint method and by backpropagation
 151 through the solver.

152 For a step size of 2^{-12} I was able to optimize the backpropagation method to
 153 take 86.573 ms while the adjoint method took 554.085 ms with an order of magnitude
 154 more allocations. The most time consuming step in the later case is the calculation
 155 of the pullback of the diffusion function. In particular, my implementation updates
 156 the diffusion matrix column by column requiring one to take slices of the diffusion
 157 matrix. These views are allocated to the heap. This represents an opportunity for
 158 optimization by predefining the views and reusing them step to step.

159 As was shown above, both methods agree with parameter gradient estimated
 160 by finite differencing with the agreement increasing for the adjoint method as the
 161 step size decreased. The result is similar when comparing the gradients of the two
 162 methods directly (Fig 5A). In principle, the gradients calculated by finite differencing
 163 and backpropagation are not exactly the gradients of the loss function as the forward
 164 solve is inexact. However, there is no reason *a priori* to believe the adjoint method
 165 should have more accurate gradient estimates; in fact, one would expect that the
 166 error from the backwards solve add to that of the forwards solve to give an even less
 167 accurate estimate. Thus, those values in Fig 5A may be interpreted as lower bound for
 168 the error of the gradients. Clearly, due to the increased computational efficiency and
 169 more accurate gradient estimates, my implementation of backpropagation through
 170 the solver is superior to my implementation of the adjoint method.

171 As a curiosity, it is worth mentioning that path of $(A_t)_t$ as calculated by the
 172 adjoint method is very similar to $(\frac{d\mathcal{L}}{dX_{T-t}})_t$ (Fig 5B), although it is expected that
 173 these two paths have similar start and end points $(\frac{d\mathcal{L}}{dX_T})_t$ and $(\frac{d\mathcal{L}}{dX_0})_t$ respectively.

174 **2.4.2. Inference.** Here I will investigate how these two method perform in an
 175 inference task. I will attempt to minimize the ℓ^2 error to the data generated above
 176 via (stochastic) gradient descent (section "Training loop" in the code). I had better

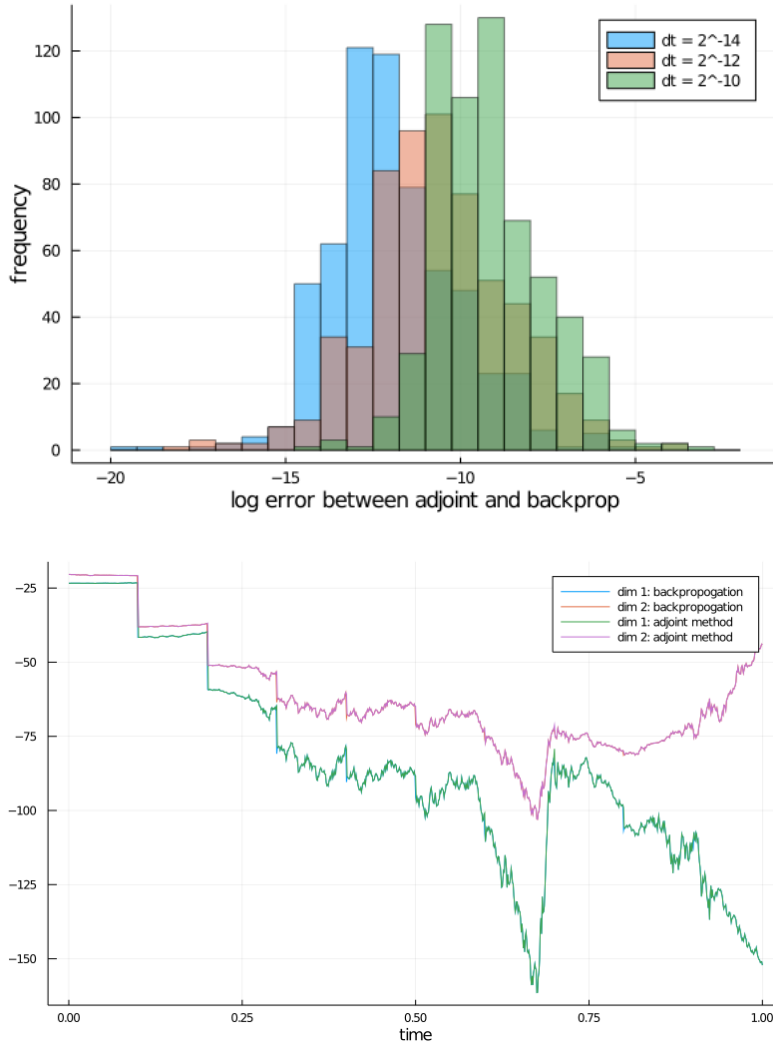


FIG. 5. A: Fractional error of parameter gradients calculated using back propagation in comparison to those calculated by adjoint method for step sizes 2^{-10} , 2^{-12} , 2^{-14} . B: Two dimensions of the paths of the adjoints A_t and $(\frac{dL}{dX_T-t})_t$ for a particular realization of the Wiener process. The two paths are superimposed and visually identical.

177 results using a ReLU activation in the neural networks than tanh so below I will
 178 consider fitting f and g with ReLU activations. As well, inference of the starting point
 179 made the problem much more difficult, so in this setting, I considered the starting
 180 point known.

181 Because of the stochasticity of the gradients, inference was very difficult. Parame-
 182 ters would explode or the magnitude of g would stabilize at a large value if the learning
 183 rate was too high. To avoid divergence, it was necessary to average the gradients of
 184 the parameters over large batches of realizations of the Wiener process. Here I used
 185 a batch size of 100 and a learning rate of 0.0005. I also observed better results when
 186 initializing with f and g near 0 but with the values of the first hidden layer positive,

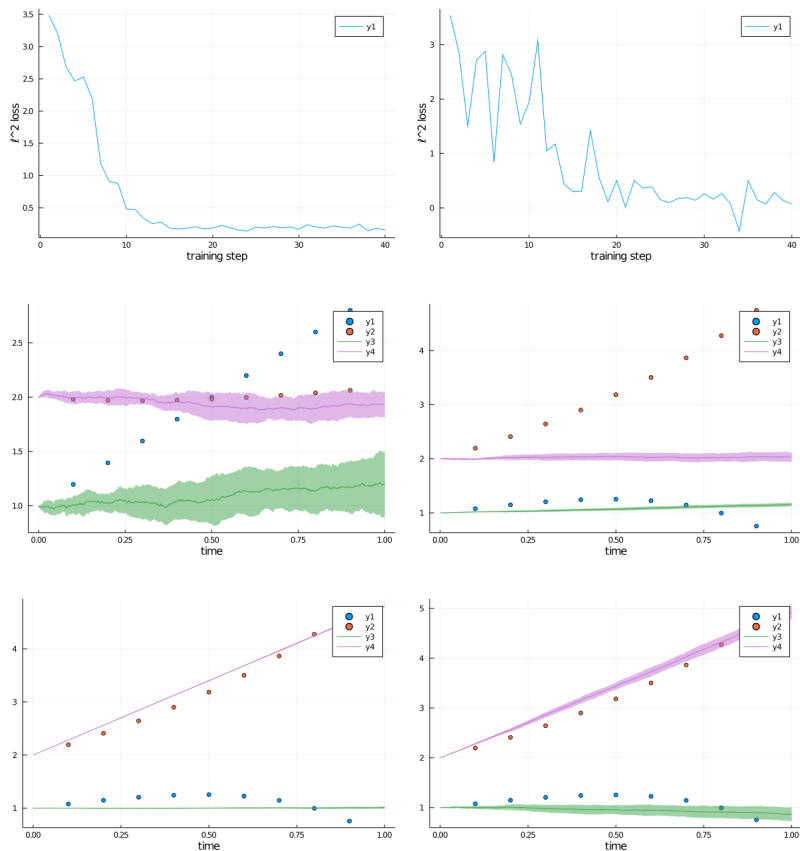


FIG. 6. The first column shows inference using backpropogation through the solver to estimate gradients while the second shows the use of the adjoint method to estimate gradients. For each set of data, both models were trained for 40 steps, the first column at a learning rate of 0.0005 and the second at 0.0003. The first row shows the loss decreasing stochastically with increased training. The middle and last row show the data (points) and the average and standard deviation of 10 paths of the SDE; the middle row shows paths from the SDE with untrained f and g , and the last row uses the trained f and g .

187 so as not to be set to 0 by the activation function.

188 I thorough comparison of the two methods would require tuning of the inference
 189 procedure which is outside the scope of this report. However, I noticed that when
 190 using gradients calculated by backpropogation, inference with the above parameters
 191 would usually converge while when using the adjoint method the sometimes diverged.
 192 Decreasing the learning rate slightly to 0.0003 usually resulted in convergence for
 193 the adjoint method. However, both methods were eventually able to converge to
 194 solutions (Fig 6). The solution arrived at when using the adjoint method had a
 195 higher magnitude g when stochasticity is not necessary to explain the data and had a
 196 much more stochastic loss curve; this is possibly a by-product of less accurate gradient
 197 estimates. Neither solution fit the data satisfactorily despite the loss plateauing, likely
 198 because of many of the nodes of the neural network being turned off by the ReLu
 199 activation. Further tuning of the architectures of f and g would likely lead to higher
 200 quality fits.

201 **3. Discussion.** Here I have implemented the EH SDE solver, backpropagation
 202 through the solver, and implemented the adjoint method of calculating integrals de-
 203 scribed in [9]; I was able to use these techniques to fit an SDE to data. My preliminary
 204 results show that the use of the Brownian tree represents a significant computational
 205 cost in the evaluation of gradients using the adjoint method, and that the more ac-
 206 curate gradient estimates of the backpropogation method resulted in higher quality
 207 inference.

208 While backpropogation through an EH step should in principle cost more function
 209 evaluations (one extra pullback of f) the particulars of my implementation meant that
 210 the backpropogation method was significantly more computationally efficient than the
 211 adjoint method for gradient calculation. To evaluate accurately the computational
 212 complexities of these models, a more careful implementation of the adjoint method is
 213 necessary.

214 Multiple ideas important to the comparison of the adjoint method and backpro-
 215 pagation were not explored in this report.

216 1) Most important is the selection of the SDE solver.

217 a) The EH method has fixed timesteps. This allowed me to save the forward path
 218 and noise at the time-points visited in the forwards pass for use in the backwards
 219 path, significantly decreasing the computational complexity of the adjoint method
 220 while likely increasing its stability. Many solvers use adaptive step sizes however to
 221 increase solution stability. This negates the benefit in stability and speed I saw in my
 222 analysis.

223 b) When efficiency is measured as number of calculations required to reach a
 224 certain error rate, the stochastic Runge-Kutta (SRK) method described in [12] leads
 225 to very efficient solvers; in particular, methods that are more efficient than the EH
 226 method. By running these SRK methods with two tableau, one may calculate esti-
 227 mated error to create adaptive SDE solvers as well [11]. These family of methods are
 228 more efficient given more stringent requirements on the form of the diffusion term g :
 229 the SRA methods are the most efficient, assuming additive noise (i.e. $g(x, t) = g(t)$);
 230 the SRID methods are less efficient, only assuming the diffusion function is diagonal,
 231 i.e. the noise in each dimension is independent; the SRIC methods are less efficient as-
 232 suming a commutativity property; and finally, the general SRI methods, which make
 233 no assumption about the form of the noise, are least efficient as they must calculate
 234 iterated stochastic integrals [13]. The form of the noise for the backwards SDE may
 235 differ from that of the forwards, requiring a less efficient solver. In particular, when
 236 the forward SDE has diagonal noise, the backwards SDE is only guaranteed to have
 237 commutative noise. However, when the forwards SDE has additive noise, the pull
 238 back of the noise does as well. Thus it is possible that the adjoint method may only
 239 be optimal for certain SDEs.

240 c) Here I only considered simple neural SDEs without stiffness. Stiff SDEs require
 241 the application of solvers with large stability regions. To achieve this, implicit methods
 242 have often been used. Backpropogation through implicit methods is extremely costly.
 243 This fact has previously been described as an advantage of the adjoint method for
 244 gradient estimation [3]. However, recent work has also developed stability optimised
 245 SRK methods SOSRI and SOSRI2 [10]; these methods are stable enough to solve stiff
 246 SDEs, are significantly more computationally efficient than implicit methods, and are
 247 much easier to backpropogate through. As well, these methods are restricted to SDEs
 248 with diagonal noise however and thus cannot be used for the backwards pass of the
 249 adjoint method. This potentially puts the adjoint method at a disadvantage when
 250 handling stiff SDEs.

d) Here I considered an SDE that did not need to be run with small time steps for the solver to get accurate estimates of the solution. SDEs where smaller step sizes are necessary would increase the memory burden, making the adjoint method more appealing. However, adaptive methods that decrease the step size of the solver commensurately with the estimated error of the solver may decrease the number of steps in the path and thus incur a smaller burden on memory. As well, more computationally efficient methods, with fewer calls to the Wiener process and function evaluations would incur a smaller memory cost to the solver. Finally, not every intermediate calculation need be saved to implement backpropagation: some may be recalculated, such as H_{i+1} in my implementation of backpropagation.

2) One potential solution to the instability of inference of SDEs is regularization by penalizing the magnitude of higher order derivatives of the solutions or the derivatives of the dynamics as suggested in [4, 7]. While the former solution is not easily translatable to the SDE case where solutions are a.s. nowhere differentiable, the latter may lead to more stable inference. This could make the instability of the adjoint method less of a detriment during inference.

3) In [9] it was recommended that latent SDEs be trained using the evidence lower bound (ELBO) as an objective. This objective avoids the collapse of the noise during training as seen in Fig 6. To make sense of the ELBO, one must be able to compute the KL divergence to a prior. Using Girsanov's theorem, one may find the KL divergence between a prior SDE with the same diffusion, but possibly different drift, and the learned SDE (when the diffusions are different, it is possible that the supports of the laws of the solutions are non-overlapping and thus have divergence infinity - for example solutions for SDEs with different additive diffusions a.s. have paths with different quadratic variation). The Bayesian setup would normally require that one fix the prior, this would result in fixing the diffusion. Optimizing the ELBO with respect to the diffusion is, however, still mathematically possible ([9] section 9.7). How the treatment of the diffusion as a learnable function or hyper-parameter affects the quality of the latent representation is unexplored; it is thus possible that the adjoint method may be superior in the situation where one must only learn the drift.

Both methods of gradient calculation explored in this report have their theoretical advantages. However, an enumeration of use case and class of SDE is complex and results are subject to change as new SDE solvers are implemented. However, such challenges must be overcome to progress towards scalable inference of SDEs.

REFERENCES

- [1] R. F. BASS, *Stochastic Processes*, Cambridge University Press, 2011.
- [2] H. S. BHAT AND R. W. MADUSHANI, *Nonparametric adjoint-based inference for stochastic differential equations*, Proceedings - 3rd IEEE International Conference on Data Science and Advanced Analytics, DSAA 2016, (2016), pp. 798–807, <https://doi.org/10.1109/DSAA.2016.69>.
- [3] R. T. Q. CHEN, Y. RUBANOVA, J. BETTENCOURT, AND D. DUVENAUD, *Neural Ordinary Differential Equations*, in NeurIPS, 2018, https://doi.org/10.1007/978-3-662-55774-7_3, <https://arxiv.org/abs/arXiv:1806.07366v5>.
- [4] C. FINLAY, J.-H. JACOBSEN, L. NURBEKYAN, AND A. M. OBERMAN, *How to train your neural ODE: the world of Jacobian and kinetic regularization*, (2020), <http://arxiv.org/abs/2002.02798>, <https://arxiv.org/abs/2002.02798>.
- [5] M. GILES AND P. GLASSERMAN, *Smoking adjoints: fast monte carlo greeks*, Risk, (2006), pp. 88–92, [#0](http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Smoking+adjoints:+fast+Monte+Carlo+Greeks).

- 301 [6] B. P. GROSS, *Applications of the Adjoint Method in Stochastic Financial Modelling*, (2015).
302 [7] J. KELLY, J. BETTENCOURT, M. J. JOHNSON, AND D. DUVENAUD, *Learning differential equations*
303 *that are easy to solve*, arXiv, 3 (2020), <https://arxiv.org/abs/2007.04504>.
304 [8] M. KIMURA, *On the probability of fixation of mutant genes in a population.*, Genetics, 47
305 (1962), pp. 713–719.
306 [9] X. LI, T.-K. L. WONG, R. T. Q. CHEN, AND D. DUVENAUD, *Scalable Gradients for Stochastic*
307 *Differential Equations*, in Proceedings of the 23rd International Conference on Artificial
308 Intelligence and Statistics (AISTATS), vol. 108, 2020, <http://arxiv.org/abs/2001.01328>,
309 <https://arxiv.org/abs/2001.01328>.
310 [10] C. RACKAUCKAS, *Stability-Optimized High Order Methods and Stiffness Detection for Pathwise*
311 *Stiff Stochastic Differential Equations* *ods for Diagonal Noise SDEs*.
312 [11] C. RACKAUCKAS AND Q. NIE, *Adaptive methods for stochastic differential equations via natu-*
313 *ral embeddings and rejection sampling with memory*, Discrete and Continuous Dynamical
314 Systems - Series B, 22 (2017), pp. 2731–2761, <https://doi.org/10.3934/dcdsb.2017133>.
315 [12] A. RÖSSLER, *Runge-Kutta methods for Stratonovich stochastic differential equation systems*
316 *with commutative noise*, Journal of Computational and Applied Mathematics, 164-165
317 (2004), pp. 613–627, <https://doi.org/10.1016/j.cam.2003.09.009>.
318 [13] M. WIKTORSSON, *Joint characteristic function and simultaneous simulation of iterated Itô*
319 *integrals for multiple independent Brownian motions*, Annals of Applied Probability, 11
320 (2001), pp. 470–487, <https://doi.org/10.1214/aoap/1015345301>.