
Learning Stochastic Dynamical Systems via Bridge Sampling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We develop algorithms to automate discovery of stochastic dynamical system
2 models from noisy, vector-valued time series. By discovery, we mean learning
3 both a nonlinear drift vector field and a diagonal diffusion matrix for an Itô stochastic
4 differential equation in \mathbb{R}^d . We parameterize the vector field using tensor
5 products of Hermite polynomials, enabling the model to capture highly nonlinear
6 and/or coupled dynamics. We solve the resulting estimation problem using
7 expectation maximization (EM). This involves two steps. We augment the data
8 via diffusion bridge sampling, with the goal of producing time series observed at
9 a higher frequency than the original data. With this augmented data, the resulting
10 expected log likelihood maximization problem reduces to a least squares problem.
11 Through experiments on systems with dimensions one through three, we show
12 that this EM approach enables accurate estimation for multiple time series with
13 possibly irregular observation times. We study how the EM method performs as a
14 function of the noise level in the data, the volume of data, and the amount of data
15 augmentation performed.

16 1 Introduction

17 Traditional mathematical modeling in the sciences and engineering often has as its goal the devel-
18 opment of equations of motion that describe observed phenomena. Classically, these equations of
19 motion usually took the form of deterministic systems of ordinary or partial differential equations
20 (ODE or PDE, respectively). Especially in systems of contemporary interest in biology and finance
21 where intrinsic noise must be modeled, we find stochastic differential equations (SDE) used instead
22 of deterministic ones. Still, these models are often built from first principles, after which the model's
23 predictions (obtained, for instance, by numerical simulation) are compared against observed data.

24 Recent years have seen a surge of interest in using data to automate discovery of ODE, PDE, and
25 SDE models. These machine learning approaches complement traditional modeling efforts, using
26 available data to constrain the space of plausible models, and shortening the feedback loop linking
27 model development to prediction and comparison to real observations. We posit two additional
28 reasons to develop algorithms to learn SDE models. First, SDE models—including the models
29 considered here—have the capacity to model highly nonlinear, coupled stochastic systems, including
30 systems whose equilibria are non-Gaussian and/or multimodal. Second, SDE models often allow for
31 interpretability. Especially if the terms on the right-hand side of the SDE are expressed in terms of
32 commonly used functions (such as polynomials), we can obtain a qualitative understanding of how
33 the system's variables influence, regulate, and/or mediate one other.

34 In this paper, we develop an algorithm to learn SDE models from high-dimensional time series. To
35 our knowledge, this is the most general expectation maximization (EM) approach to learning an
36 SDE with multidimensional drift vector field and diagonal diffusion matrix. Prior EM approaches

were restricted to one-dimensional SDE [8], or used a Gaussian process approximation, linear drift approximation, and approximate maximization [21]. To develop our method, we use diffusion bridge sampling as in [28, 29], which focused on Bayesian nonparametric methods for SDE in \mathbb{R}^1 . After augmenting the data using bridge sampling, we are left with a least-squares problem, generalizing the work of [6] from the ODE to the SDE context.

In the literature, variational Bayesian methods are the only other SDE learning methods that have been tested on high-dimensional problems [31]. These methods use approximations consisting of linear SDE with time-varying coefficients [1], kernel density estimates [2], or Gaussian processes [3]. In contrast, we parameterize the drift vector field using tensor products of Hermite polynomials; as mentioned above, the resulting SDE has much higher capacity than linear and/or Gaussian process models.

Many other techniques explored in the statistical literature focus on scalar SDE [4, 12, 13, 30].

As mentioned, differential equation discovery problems have attracted considerable recent interest. A variety of methods have been developed to learn ODE [6, 7, 16, 22, 24, 25, 27] as well as PDE [17, 18, 20, 23]. Unlike many of these works, we do not focus on model selection and/or regularization; if needed, our methods can be combined with model selection procedures developed in the ODE context [10, 11].

2 Problem Setup

Let W_t denote Brownian motion in \mathbb{R}^d —informally, an increment dW_t of this process has a multivariate normal distribution with zero mean vector and covariance matrix Idt . Let X_t denote an \mathbb{R}^d -valued stochastic process that evolves according to the Itô SDE

$$dX_t = f(X_t)dt + \Gamma dW_t. \quad (1)$$

For rigorous definitions of Brownian motion and SDE, see [5, 32]. The nonlinear vector field $f : \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the *drift* function, and the $d \times d$ matrix Γ is the *diffusion* matrix. To reduce the number of model parameters, we assume $\Gamma = \text{diag } \gamma$.

Our goal is to develop an algorithm that accurately estimates the functional form of f and the vector γ from time series data.

Parameterization. We parameterize f using Hermite polynomials. The n -th Hermite polynomial takes the form

$$H_n(x) = (\sqrt{2\pi}n!)^{-1/2}(-1)^n e^{x^2/2} \frac{d^n}{dx^n} e^{-x^2/2} \quad (2)$$

Let $\langle f, g \rangle_w = \int_{\mathbb{R}} f(x)g(x) \exp(-x^2/2) dx$ denote a weighted L^2 inner product. Then, $\langle H_i, H_j \rangle_w = \delta_{ij}$, i.e., the Hermite polynomials are orthonormal with respect to the weighted inner product. In fact, with respect to this inner product, the Hermite polynomials form an orthonormal basis of $L_w^2(\mathbb{R}) = \{f : \langle f, f \rangle_w < \infty\}$.

Now let $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{Z}_+^d$ denote a multi-index. We use the notation $|\alpha| = \sum_j \alpha_j$ and $x^\alpha = \prod_j (x_j)^{\alpha_j}$ for $x = (x_1, \dots, x_d) \in \mathbb{R}^d$. For $x \in \mathbb{R}^d$ and a multi-index α , we also define

$$H_\alpha(x) = \prod_{j=1}^d H_{\alpha_j}(x_j). \quad (3)$$

We write $f(x) = (f_1(x), \dots, f_d(x))$ and then parameterize each component

$$f_j(x) = \sum_{m=0}^M \sum_{|\alpha|=m} \beta_\alpha^j H_\alpha(x). \quad (4)$$

We see that the maximum degree of $H_\alpha(x)$ is $|\alpha|$. Hence we think of the double sum in (4) as first summing over degrees and then summing over all terms with a fixed maximum degree. We say maximum degree because, for instance, $H_2(z) = (z^2 - 1)/(\sqrt{2\pi}2)^{1/2}$ contains both degree 2 and degree 0 terms.

76 There are $\binom{m+d-1}{d-1}$ possibilities for a d -dimensional multi-index α such that $|\alpha| = m$. Summing
 77 this from $m = 0$ to M , there are $\widetilde{M} = \binom{M+d}{d}$ total multi-indices in the double sum in (4). Let (i)
 78 denote the i -th multi-index according to some ordering. Then we can write

$$f_j(x) = \sum_{i=1}^{\widetilde{M}} \beta_{(i)}^j H_{(i)}(x). \quad (5)$$

79 Essentially, we parameterize f using tensor products of Hermite polynomials.

80 **Data.** We consider our data $\mathbf{x} = \{x_j\}_{j=0}^L$ to be direct observations of X_t at discrete points in time
 81 $\mathbf{t} = \{t_j\}_{j=0}^L$. Note that these time points do not need to be equispaced. In the derivation that follows,
 82 we will consider the data (\mathbf{t}, \mathbf{x}) to be one time series. Later, we indicate how our methods generalize
 83 naturally to multiple time series, i.e., repeated observations of the same system.

84 To achieve our estimation goal, we apply expectation maximization (EM). We regard \mathbf{x} as the incom-
 85 plete data. Let $\Delta t = \max_j(t_j - t_{j-1})$ be the maximum interobservation spacing. We think of the
 86 missing data \mathbf{z} as data collected at a time scale $h \ll \Delta t$ fine enough such that the transition density
 87 of (1) is approximately Gaussian. To see how this works, let $\mathcal{N}(\mu, \Sigma)$ denote a multivariate normal
 88 with mean vector μ and covariance matrix Σ . Now discretize (1) in time via the Euler-Maruyama
 89 method with time step $h > 0$; the result is

$$\widetilde{X}_{n+1} = \widetilde{X}_n + f(\widetilde{X}_n)h + h^{1/2}\Gamma Z_{n+1}, \quad (6)$$

90 where $Z_{n+1} \sim \mathcal{N}(0, I)$ is a standard multivariate normal, independent of X_n . This implies that

$$(\widetilde{X}_{n+1} | \widetilde{X}_n = v) \sim \mathcal{N}(v + f(v)h, h\Gamma^2). \quad (7)$$

91 As h decreases, $\widetilde{X}_{n+1} | \widetilde{X}_n = v$ —a Gaussian approximation—will converge to the true transition
 92 density $X_{(n+1)h} | X_{nh} = v$, where X_t refers to the solution of (1).

93 **Diffusion Bridge.** To augment or complete the data, we employ diffusion bridge sampling, using a
 94 Markov chain Monte Carlo (MCMC) method that goes back to [15, 19]. Let us describe our version
 95 here. We suppose our current estimate of $\theta = (\beta, \gamma)$ is given. Define the diffusion bridge process to
 96 be (1) conditioned on both the initial value x_i at time t_i , and the final value x_{i+1} at time t_{i+1} . The
 97 goal is to generate sample paths of this diffusion bridge. By a sample path, we mean $F - 1$ new
 98 samples $\{z_{i,j}\}_{j=1}^{F-1}$ at times $t_i + jh$ with $h = (t_{i+1} - t_i)/F$.

99 To generate such a path, we start by drawing a sample from a Brownian bridge with the same
 100 diffusion as (1). That is, we sample from the SDE

$$d\widehat{X}_t = \Gamma dW_t \quad (8)$$

101 conditioned on $\widehat{X}_{t_i} = x_i$ and $\widehat{X}_{t_{i+1}} = x_{i+1}$. This Brownian bridge can be described explicitly

$$\widehat{X}_t = \Gamma(W_t - W_{t_i}) + x_i - \frac{t - t_i}{t_{i+1} - t_i}(\Gamma(W_{t_{i+1}} - W_{t_i}) + x_i - x_{i+1}) \quad (9)$$

102 Here $W_0 = 0$ (almost surely), and $W_t - W_s \sim \mathcal{N}(0, (t - s)I)$ for $t > s \geq 0$.

103 Let \mathbb{P} denote the law of the diffusion bridge process, and let \mathbb{Q} denote the law of the Brownian bridge
 104 (9). Using Girsanov's theorem [14], we can show that

$$\frac{d\mathbb{P}}{d\mathbb{Q}} = C \exp \left(\int_{t_i}^{t_{i+1}} f(\widehat{X}_s)^T \Gamma^{-2} d\widehat{X}_s - \frac{1}{2} \int_{t_i}^{t_{i+1}} f(\widehat{X}_s)^T \Gamma^{-2} f(\widehat{X}_s) ds \right), \quad (10)$$

105 where the constant C depends only on x_i and x_{i+1} . The left-hand side is a Radon-Nikodym deriva-
 106 tive, equivalent to a density or likelihood; the ratio of two such likelihoods is the accept/reject ratio
 107 in the Metropolis algorithm [26].

108 Putting the above pieces together yields the following Metropolis algorithm to generate diffusion
 109 bridge sample paths. Fix $F \geq 2$ and $i \in \{0, \dots, L - 1\}$. Assume we have stored the previous
 110 Metropolis step, i.e., a path $\mathbf{z}^{(\ell)} = \{z_{i,j}^{(\ell)}\}_{j=1}^{F-1}$.

- 111 1. Use (9) to generate samples of \widehat{X}_t at times $t_i + jh$, for $j = 1, 2, \dots, F - 1$ and $h =$
 112 $(t_{i+1} - t_i)/F$. This is the proposal $\mathbf{z}^* = \{z_{i,j}^*\}_{j=1}^{F-1}$.
 113 2. Numerically approximate the integrals in (10) to compute the likelihood of the proposal.
 114 Specifically, we compute

$$p(\mathbf{z}^*)/C = \sum_{j=0}^{F-1} f(z_{i,j}^*)^T \Gamma^{-2} (z_{i,j+1}^* - z_{i,j}^*) \\ - \frac{h}{4} \sum_{j=0}^{F-1} [f(z_{i,j}^*)^T \Gamma^{-2} f(z_{i,j}^*) + f(z_{i,j+1}^*)^T \Gamma^{-2} f(z_{i,j+1}^*)]$$

115 We have discretized the stochastic $d\widehat{X}_s$ integral using Itô's definition, and we have dis-
 116 cretized the ordinary ds integral using the trapezoidal rule.

- 117 3. Accept the proposal with probability $p(\mathbf{z}^*)/p(\mathbf{z}^{(\ell)})$ —note the factors of C cancel. If the
 118 proposal is accepted, then set $\mathbf{z}^{(\ell+1)} = \mathbf{z}^*$. Else set $\mathbf{z}^{(\ell+1)} = \mathbf{z}^{(\ell)}$.

119 We initialize this MCMC algorithm with a Brownian bridge path and use post-burn-in steps as the
 120 diffusion bridge samples we seek.

121 **Expectation Maximization (EM).** Let us now give details to justify the intuition expressed above,
 122 that employing the diffusion bridge to augment the data on a fine scale will enable estimation. Let
 123 $\mathbf{z}^{(r)} = \{z_{i,j}^{(r)}\}_{j=1}^{F-1}$ be the r -th diffusion bridge sample path. We interleave this sampled data together
 124 with the observed data \mathbf{x} to create the completed time series

$$\mathbf{y}^{(r)} = \{y_j^{(r)}\}_{j=1}^N,$$

125 where $N = LF + 1$. By interleaving, we mean that $y_{1+iF}^{(r)} = x_i$ for $i = 0, 1, \dots, L$, and that
 126 $y_{1+j+iF}^{(r)} = z_{i,j}^{(r)}$ for $j = 1, 2, \dots, F - 1$ and $i = 0, 1, \dots, L - 1$. With this notation, we can more
 127 easily express the EM algorithm. Let us assume that we currently have access to $\boldsymbol{\theta}^{(k)}$, our estimate
 128 of the parameters after k iterations. If $k = 0$, we set $\boldsymbol{\theta}^{(0)}$ equal to an initial guess. Then we follow
 129 two steps:

- 130 1. For the expectation step, we first generate an ensemble of R diffusion bridge sample paths.
 131 Interleaving as above, this yields R completed time series $\mathbf{y}^{(r)}$ for $r = 1, \dots, R$. In what
 132 follows, we will use an average over this ensemble to approximate the expected value. Let
 133 h_j denote the elapsed time between observations y_j and y_{j+1} . Using the completed data,
 134 the temporal discretization (6) of the SDE, the Markov property, and property (7), we have

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(k)}) = \mathbb{E}_{\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^{(k)}} [\log p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta})] \quad (11)$$

$$\approx \frac{1}{R} \sum_{r=1}^R \log p(\mathbf{y}^{(r)} \mid \boldsymbol{\theta}) \\ = \frac{1}{R} \sum_{r=1}^R \sum_{n=1}^{N-1} \log p(y_{n+1}^{(r)} \mid y_n^{(r)}, \boldsymbol{\theta}) \\ = -\frac{1}{R} \sum_{r=1}^R \sum_{n=1}^{N-1} \left[\sum_{j=1}^d \frac{1}{2} \log(2\pi h_n \gamma_j^2) \right. \\ \left. + \frac{1}{2h_n} \left\| \Gamma^{-1} \left(y_{n+1}^{(r)} - y_n^{(r)} - h_n \sum_{\ell=1}^{\widetilde{M}} \beta_{(\ell)} H_{(\ell)}(y_n^{(r)}) \right) \right\|_2^2 \right]. \quad (12)$$

- 135 2. For the M step, we maximize in stages

$$\beta^{(k+1)} = \arg \max_{\beta} Q((\beta, \gamma^{(k)}), \boldsymbol{\theta}^{(k)}) \\ \gamma^{(k+1)} = \arg \max_{\gamma} Q((\beta^{(k+1)}, \gamma), \boldsymbol{\theta}^{(k)})$$

136 The maximization over β is a least squares problem. Reinterpreting H as a matrix, the
 137 solution is given by forming the matrix

$$\mathcal{M}_{k,\ell} = \frac{1}{R} \sum_{r=1}^R \sum_{n=1}^{N-1} h_n H_{(k)}^T(y_n^{(r)}) \Gamma^{-2} H_{(\ell)}(y_n^{(r)}) \quad (13)$$

138 and the vector

$$\rho_k = \frac{1}{R} \sum_{r=1}^R \sum_{n=1}^{N-1} H_{(k)}^T(y_n^{(r)}) \Gamma^{-2} (y_{n+1}^{(r)} - y_n^{(r)}). \quad (14)$$

139 We then solve the system $\mathcal{M}\beta = \rho$ for β . Now that we have β , we maximize over γ . The
 140 solution can be obtained in closed form

$$\gamma_i^2 = \frac{1}{RN} \sum_{r=1}^R \sum_{n=1}^{N-1} h_n^{-1} ((y_{n+1}^{(r)} - y_n^{(r)} - h_n \sum_{\ell=1}^{\widetilde{M}} \beta_{(\ell)} H_{(\ell)}(y_n^{(r)})) \cdot e_i)^2 \quad (15)$$

141 where e_i is the i^{th} canonical basis vector in \mathbb{R}^d .

142 We iterate the above two steps until $\|\theta^{(k+1)} - \theta^{(k)}\| / \|\theta^{(k)}\| < \delta$ for some tolerance $\delta > 0$.

143 When the data consists of multiple time series $\{\mathbf{t}^{(i)}, \mathbf{x}^{(i)}\}_{i=1}^S$, everything scales accordingly. For
 144 instance, we create an ensemble of R diffusion bridge samples for each of the S time series. If we
 145 index the resulting completed time series appropriately, we simply replace R by RS in (13), (14),
 146 and (15) and keep everything else the same.

147 There are three sources of error in the above algorithm. The first relates to replacing the expectation
 148 by a sample average; the induced error should, by the law of large numbers, decrease as $R^{-1/2}$. The
 149 second stems from the approximate nature of the computed diffusion bridge samples—as indicated
 150 above, we use numerical integration to approximate the Girsanov likelihood. The third source of
 151 error is in using the Gaussian transition density to approximate the true transition density of the
 152 SDE. Both the second and third sources of error vanish in the $F \rightarrow \infty$ limit [9].

153 3 Experiments

154 We present a series of increasingly higher-dimensional experiments with synthetic data. To generate
 155 this data, we start with a known stochastic dynamical system of the form (1). Using Euler-Maruyama
 156 time stepping starting from a randomly chosen initial condition, we march forward in time from
 157 $t = 0$ to a final time $t = 10$.

158 In all examples, we step forward internally at a time step of $h = 0.0001$, but for the purposes of
 159 estimation, we only use data sampled every 0.1 units of time, discarding 99.9% of the simulated
 160 trajectory. We use a fine internal time step to reduce, to the extent possible, numerical error in the
 161 simulated data. We save the data on a coarse time scale to test the proposed EM algorithm.

162 To study how the EM method performs as a function of noise strength, data volume, and data aug-
 163 mentation, we perform four sets of experiments. When we run EM, we randomly generate the initial
 164 guess $\beta^{(0)} \sim \mathcal{N}(\mu = 0, \sigma^2 = 0.5)$. We set the EM tolerance parameter $\delta = 0.01$. The only reg-
 165 ularization we include is to threshold β —values less than 0.01 in absolute value are reset to zero.
 166 Finally, in the MCMC diffusion bridge sampler, we use 10 burn-in steps and then create an ensemble
 167 of size $R = 100$.

168 To quantify the error between the estimated $\tilde{\beta}$ and the true β , we apply the Frobenius norm:

$$\varepsilon = \sqrt{\sum_i \|\beta_{(i)} - \tilde{\beta}_{(i)}\|^2} \quad (16)$$

169 The $\tilde{\beta}$ coefficients are the Hermite coefficients of the estimated drift vector field f . For each example
 170 system, we compute the true Hermite coefficients β by multiplying the true ordinary polynomial
 171 coefficients by a change-of-basis matrix that is easily computed.

172 We test the method using stochastic systems in dimensions $d = 1, 2, 3$. In 1D, we use

$$dX_t = (1 + X_t - X_t^2)dt + \gamma dW_t.$$

173 In 2D, we use a stochastic Duffing oscillator with no damping or driving:

$$dX_{0,t} = X_{1,t}dt + \gamma_0 dW_{0,t}, \quad dX_{1,t} = (-X_{0,t} - X_{0,t}^3)dt + \gamma_1 dW_{1,t}$$

174 For the 3D case, we consider the stochastic, damped, driven Duffing oscillator:

$$\begin{aligned} dX_{0,t} &= X_{1,t}dt + \gamma_0 dW_{0,t} \\ dX_{1,t} &= (X_{0,t} - X_{0,t}^3 - 0.3X_{1,t} + 0.5 \cos(X_{2,t}))dt + \gamma_1 dW_{1,t} \\ dX_{2,t} &= 1.2dt + \gamma_2 dW_{2,t} \end{aligned}$$

175 In what follows, we refer to these systems as the 1D, 2D, and 3D systems.

176 **Experiment 1: Varying Number of Time Series.** Here we vary data volume by stepping the
 177 number S of time series from $S = 1$ to $S = 10$. Each time series has length $L + 1 = 101$. The
 178 results, as plotted in Figures 1 and 2, show that increasing S leads to much better estimates of β . As
 179 a rule of thumb, the results indicate that at least $S \geq 4$ time series are needed for accurate estimation.

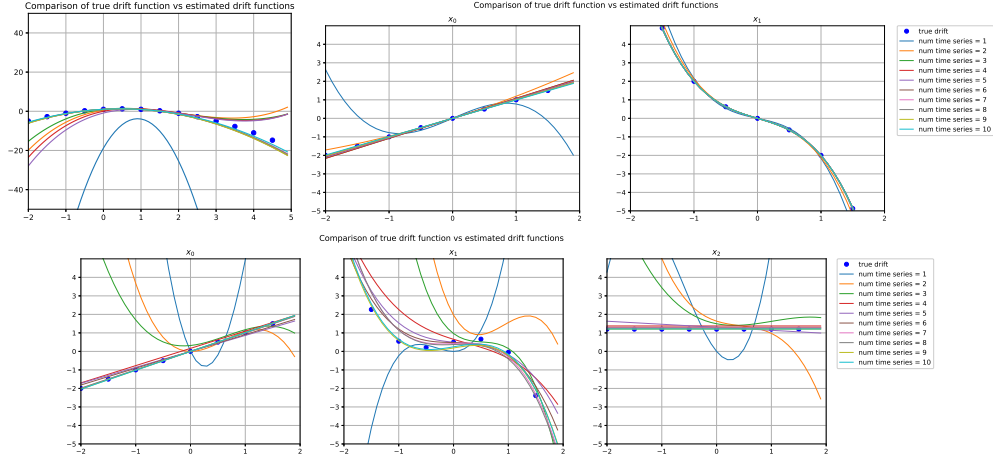


Figure 1: As we increase the number S of time series used to learn the drift, the estimated drift more closely approximates the ground truth. From top to bottom, left to right, we have plotted estimated and true drifts for the 1D, 2D, and 3D systems. For the 1D and 2D systems, the true drifts depend on only one variable. For the $dX_{1,t}$ component of the 3D system, we have plotted the dependence of the drifts on X_0 only, keeping X_1 and X_2 fixed at 0.

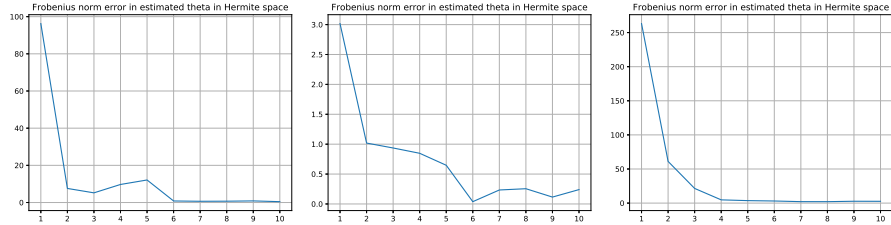


Figure 2: As we increase the number S of time series used to learn the drift, the Frobenius norm error between estimated and true drifts—see (16)—decreases significantly. From left to right, we have plotted results for the 1D, 2D, and 3D systems.

180 **Experiment 2: Varying Length of Time Series.** Here we vary data volume by stepping the length
 181 $L + 1$ of the time series from $L + 1 = 11$ to $L + 1 = 101$, keeping the number of time series fixed
 182 at $S = 10$. Also note that in this experiment, observation times strictly between the initial and
 183 final times are chosen randomly. In Figure 3, we have plotted the estimated and true drifts for only
 184 the 3D system; in Figure 4, we have plotted the error (16) for all three systems. Comparing with
 185 Experiment 1, we see that randomization of the observation times improves estimation. That is, even
 186 with $L + 1 = 11$ data points per time series, we obtain accurate estimates.

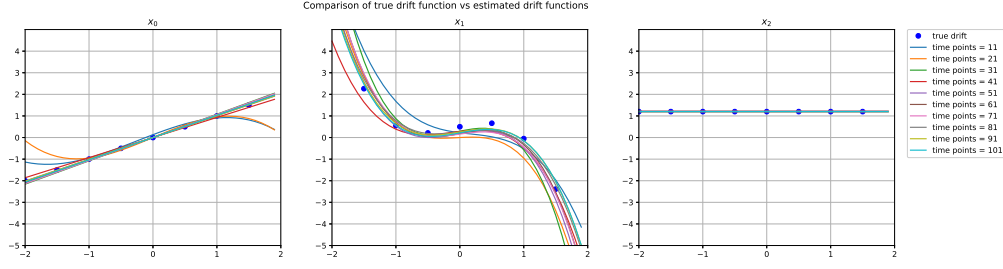


Figure 3: We plot true and estimated drifts for the 3D system as a function of increasing time series length L . The three components of the vector field are plotted as in the third row of Figure 1. The results show that randomization of observation times compensates for a small value of L , enabling accurate estimation.

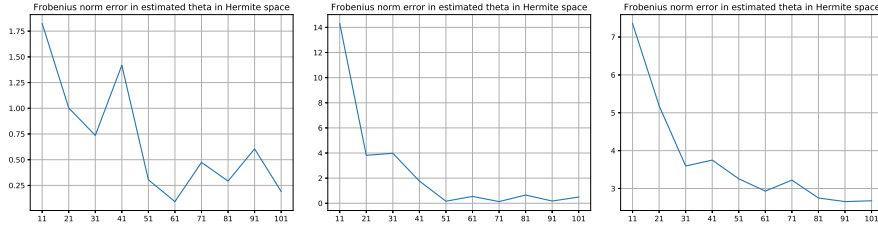


Figure 4: As we increase the length L of each time series used for learning, the Frobenius norm error between estimated and true drifts—see (16)—decreases significantly. From left to right, we have plotted results for the 1D, 2D, and 3D systems.

187 **Experiment 3: Varying Noise Strength.** Here we vary the noise strength γ , stepping from 0.5 to
 188 0.0001 while keeping other parameters constant. Specifically, we take $S = 10$ time series each of
 189 length $L + 1 = 101$. In Figure 5, we have plotted Frobenius errors for all three systems. Though
 190 the error in the estimated coefficients for the 3D system may seem large, the estimated and true drift
 191 functions are close—see Figure 6.

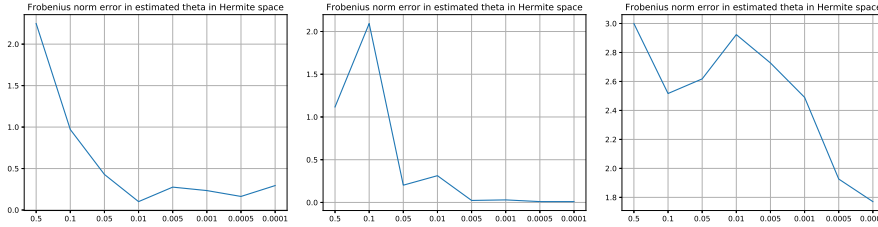


Figure 5: Varying the strength of the noise in the simulated data alters the quality of estimated drift coefficients, quantified using the Frobenius error (16). We proceed from left to right. For the 1D and 2D systems, the maximum noise strength of 0.5 remains below the magnitude of the drift field coefficients. For these systems, as the noise strength decreases, the error drops close to zero. For the 3D system, the maximum noise strength of 0.5 is greater than or equal to two of the drift field coefficients, leading to apparently decreased performance—however, see Figure 6.

192 **Experiment 4: Varying Data Augmentation.** We start with $S = 10$ time series with $L + 1 = 51$
 193 points each. Here we vary the number of interleaved diffusion bridge samples: $F = 1, \dots, 10$. For
 194 $F = 1$, no diffusion bridge is created; the likelihood is computed by applying the Gaussian transition
 195 density directly to the observed data. The results, plotted in Figures 7 and 8, show that increased
 196 data augmentation dramatically improves the quality of estimated drifts. Though the Frobenius error
 197 for the 3D system exceeds 2.6, Figure 8 shows that EM’s estimates are still accurate.

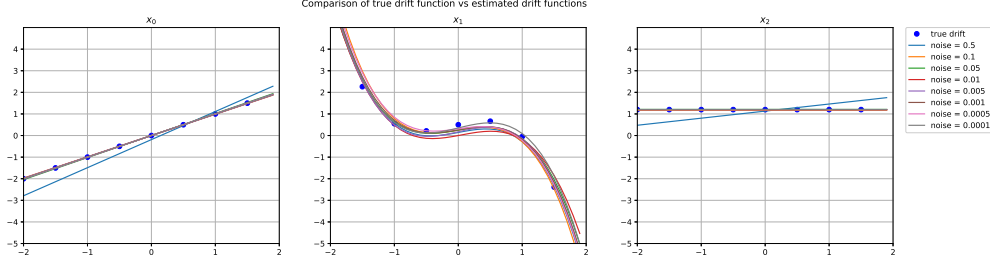


Figure 6: Though Figure 5 shows a Frobenius norm error for the 3D system greater than ≈ 1.8 at all noise levels, when plotted, the estimated drift functions lie close to the true drift function. The three components of the vector field are plotted as in the third row of Figure 1.

198 We have not plotted results for the scarce data regime where we have $S = 10$ time series with $L = 11$
 199 points each. In this regime, data augmentation enables highly accurate estimation for the 2D and 3D
 200 systems. For the 1D system, the observations do not explore phase space properly, leading to poor
 201 estimation of the drift.

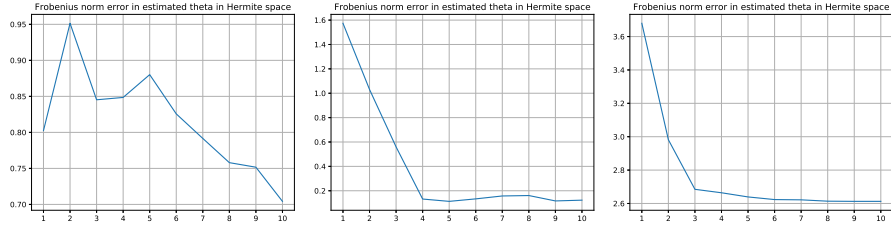


Figure 7: As we increase the length F of the diffusion bridge interleaving observed data points, the quality of estimated drifts improves considerably. From left to right, we have plotted Frobenius errors (16) between true and estimated coefficients, for the 1D, 2D, and 3D systems.

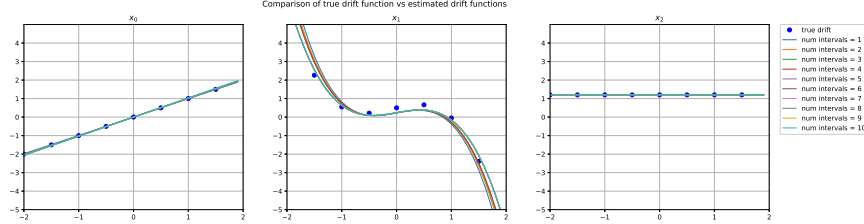


Figure 8: Though Figure 7 shows a Frobenius norm error for the 3D system greater than ≈ 2.6 at all noise levels, when plotted, the estimated drift functions lie close to the true drift function. The three components of the vector field are plotted as in the third row of Figure 1.

202 4 Conclusion

203 We have developed an EM algorithm for estimation of drift functions and diffusion matrices for
 204 SDE. We have demonstrated the conditions under which the algorithm succeeds in estimating SDE.
 205 Specifically, our tests show that with enough data volume and data augmentation, the EM algo-
 206 rithm produces highly accurate results. In future work, we seek to further test our method on high-
 207 dimensional, nonlinear problems, problems with non-constant diffusion matrices, and real experi-
 208 mental data. As we move to higher-dimensional problems, we will also explore regularization and
 209 model selection techniques.

References

- [1] C. Archambeau, M. Opper, Y. Shen, D. Cornford, and J. S. Shawe-Taylor. Variational inference for diffusion processes. In *Advances in Neural Information Processing Systems*, pages 17–24, 2008.
- [2] P. Batz, A. Ruttor, and M. Opper. Variational estimation of the drift for stochastic differential equations from the empirical density. *Journal of Statistical Mechanics: Theory and Experiment*, 2016(8):083404, Aug. 2016. doi: 10.1088/1742-5468/2016/08/083404.
- [3] P. Batz, A. Ruttor, and M. Opper. Approximate Bayes learning of stochastic differential equations. *arXiv preprint arXiv:1702.05390*, 2017. URL <https://arxiv.org/abs/1702.05390>.
- [4] H. S. Bhat and R. W. M. A. Madushani. Nonparametric Adjoint-Based Inference for Stochastic Differential Equations. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 798–807, Oct. 2016. doi: 10.1109/DSAA.2016.69.
- [5] R. N. Bhattacharya and E. C. Waymire. *Stochastic Processes with Applications*. SIAM, Aug. 2009. ISBN 978-0-89871-689-4.
- [6] S. L. Brunton, J. L. Proctor, and J. N. Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, Apr. 2016. doi: 10.1073/pnas.1517384113.
- [7] S. Chen, A. Shojaie, and D. M. Witten. Network Reconstruction From High-Dimensional Ordinary Differential Equations. *Journal of the American Statistical Association*, 112(520):1697–1707, Oct. 2017. doi: 10.1080/01621459.2016.1229197.
- [8] Z. Ghahramani and S. T. Roweis. Learning nonlinear dynamical systems using an EM algorithm. *Advances in Neural Information Processing Systems (NIPS)*, pages 431–437, 1999.
- [9] P. E. Kloeden and E. Platen. *Numerical Solution of Stochastic Differential Equations*. Springer Science & Business Media, June 2011. ISBN 978-3-540-54062-5.
- [10] N. M. Mangan, S. L. Brunton, J. L. Proctor, and J. N. Kutz. Inferring Biological Networks by Sparse Identification of Nonlinear Dynamics. *IEEE Transactions on Molecular, Biological and Multi-Scale Communications*, 2(1):52–63, June 2016. doi: 10.1109/TMBMC.2016.2633265.
- [11] N. M. Mangan, J. N. Kutz, S. L. Brunton, and J. L. Proctor. Model selection for dynamical systems via sparse regression and information criteria. *Proc. R. Soc. A*, 473(2204):20170009, Aug. 2017. doi: 10.1098/rspa.2017.0009.
- [12] H.-G. Müller, F. Yao, and others. Empirical dynamics for longitudinal data. *The Annals of Statistics*, 38(6):3458–3486, 2010. URL <http://projecteuclid.org/euclid.aos/1291126964>.
- [13] J. Nicolau. Nonparametric estimation of second-order stochastic differential equations. *Econometric Theory*, 23(05):880, Oct. 2007. doi: 10.1017/S0266466607070375.
- [14] O. Papaspiliopoulos and G. O. Roberts. Importance sampling techniques for estimation of diffusion models. *Statistical methods for stochastic differential equations*, 124:311–340, 2012.
- [15] O. Papaspiliopoulos, G. O. Roberts, and O. Stramer. Data Augmentation for Diffusions. *Journal of Computational and Graphical Statistics*, 22(3):665–688, July 2013. doi: 10.1080/10618600.2013.783484.
- [16] M. Quade, M. Abel, J. N. Kutz, and S. L. Brunton. Sparse Identification of Nonlinear Dynamics for Rapid Model Recovery. Mar. 2018. URL <http://arxiv.org/abs/1803.00894>. arXiv: 1803.00894.
- [17] M. Raissi and G. E. Karniadakis. Hidden physics models: Machine learning of nonlinear partial differential equations. *Journal of Computational Physics*, 357:125–141, Mar. 2018. doi: 10.1016/j.jcp.2017.11.039.
- [18] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Machine learning of linear differential equations using Gaussian processes. *Journal of Computational Physics*, 348:683–693, Nov. 2017. doi: 10.1016/j.jcp.2017.07.050.
- [19] G. O. Roberts and O. Stramer. On inference for partially observed nonlinear diffusion models using the Metropolis–Hastings algorithm. *Biometrika*, 88(3):603–621, Oct. 2001. doi: 10.1093/biomet/88.3.603.
- [20] S. H. Rudy, S. L. Brunton, J. L. Proctor, and J. N. Kutz. Data-driven discovery of partial differential equations. *Science Advances*, 3(4):e1602614, Apr. 2017. doi: 10.1126/sciadv.1602614.
- [21] A. Ruttor, P. Batz, and M. Opper. Approximate Gaussian process inference for the drift function in stochastic differential equations. In *Advances in Neural Information Processing Systems*, pages 2040–2048, 2013.
- [22] H. Schaeffer. Learning partial differential equations via data discovery and sparse optimization. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Science*, 473(2197):20160446, Jan. 2017. doi: 10.1098/rspa.2016.0446.

- 265 [23] H. Schaeffer, R. Caflisch, C. D. Hauck, and S. Osher. Sparse dynamics for partial differential equations.
 266 *Proceedings of the National Academy of Sciences*, 110(17):6634–6639, Apr. 2013. doi: 10.1073/pnas.
 267 1302752110.
- 268 [24] H. Schaeffer, G. Tran, and R. Ward. Extracting Sparse High-Dimensional Dynamics from Limited Data.
 269 *arXiv:1707.08528 [math]*, July 2017. arXiv: 1707.08528.
- 270 [25] T. B. Schön, A. Svensson, L. Murray, and F. Lindsten. Probabilistic learning of nonlinear dynamical
 271 systems using sequential Monte Carlo. 2017. URL <https://arxiv.org/abs/1703.02419>.
- 272 [26] A. M. Stuart. Inverse problems: A Bayesian perspective. *Acta Numerica*, 19:451–559, May 2010. doi:
 273 10.1017/S0962492910000061.
- 274 [27] G. Tran and R. Ward. Exact Recovery of Chaotic Systems from Highly Corrupted Data. *Multiscale*
 275 *Modeling & Simulation*, 15(3):1108–1129, Jan. 2017. doi: 10.1137/16M1086637.
- 276 [28] F. van der Meulen, M. Schauer, and H. van Zanten. Reversible jump MCMC for nonparametric drift
 277 estimation for diffusion processes. *Computational Statistics & Data Analysis*, 71:615–632, Mar. 2014.
 278 doi: 10.1016/j.csda.2013.03.002.
- 279 [29] F. van der Meulen, M. Schauer, and J. van Waaij. Adaptive nonparametric drift estimation for diffusion
 280 processes using Faber–Schauder expansions. *Statistical Inference for Stochastic Processes*, pages 1–26,
 281 June 2017. doi: 10.1007/s11203-017-9163-7.
- 282 [30] N. Verzelen, W. Tao, H.-G. Müller, and others. Inferring stochastic dynamics from functional data.
 283 *Biometrika*, 99(3):533–550, 2012.
- 284 [31] M. D. Vrettas, M. Opper, and D. Cornford. Variational mean-field algorithm for efficient inference in
 285 large systems of stochastic differential equations. *Physical Review E*, 91(1):012148, Jan. 2015. doi:
 286 10.1103/PhysRevE.91.012148.
- 287 [32] B. Øksendal. *Stochastic Differential Equations: An Introduction with Applications*. Universitext. Springer-
 288 Verlag, Berlin Heidelberg, 6 edition, 2003. ISBN 978-3-540-04758-2.