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

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

33 In this paper, we develop an algorithm to learn SDE models from high-dimensional time series. To
34 our knowledge, this is the most general expectation maximization (EM) approach to learning an
35 SDE with multidimensional drift vector field and diagonal diffusion matrix. Prior EM approaches
36 were restricted to one-dimensional SDE [8], or used a Gaussian process approximation, linear drift
37 approximation, and approximate maximization [22]. To develop our method, we use diffusion bridge
38 sampling as in [12, 29], which focused on Bayesian nonparametric methods for SDE in \mathbb{R}^1 . After

39 augmenting the data using bridge sampling, we are left with a least-squares problem, generalizing
40 the work of [6] from the ODE to the SDE context.

41 In the literature, variational Bayesian methods are the only other SDE learning methods that have
42 been tested on high-dimensional problems [31]. These methods use approximations consisting of
43 linear SDE with time-varying coefficients [1], kernel density estimates [2], or Gaussian processes
44 [3]. In contrast, we parameterize the drift vector field using tensor products of Hermite polynomials;
45 as mentioned above, the resulting SDE has much higher capacity than linear and/or Gaussian process
46 models.

47 Many other techniques explored in the statistical literature focus on scalar SDE [4, 13, 14, 30].

48 As mentioned, differential equation discovery problems have attracted considerable recent interest.
49 A variety of methods have been developed to learn ODE [6, 7, 17, 23, 25, 26, 28] as well as PDE [18,
50 19, 21, 24]. Unlike many of these works, we do not focus on model selection and/or regularization;
51 if needed, our methods can be combined with model selection procedures developed in the ODE
52 context [10, 11].

- 53 • Expectation and maximization formulas assuming data is filled in, CLOSE2DONE
- 54 • how synthetic data is generated, EDIT
- 55 • results: 1D, 2D, 3D damped duffing, 3D lorenz
- 56 • plots: error of theta vs noise, error vs amount of data (number of data points) parametric
- 57 curves for noise levels, brownian bridge plots for illustration, ...

58 1 Problem Setup

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

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

62 For rigorous definitions of Brownian motion and SDE, see [5, 32]. The nonlinear vector field $f :$
63 $\Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the *drift* function, and the $d \times d$ matrix Γ is the *diffusion* matrix. To reduce the
64 number of model parameters, we assume $\Gamma = \text{diag } \gamma$.

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

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

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

69 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,
70 $\langle H_i, H_j \rangle_w = \delta_{ij}$, i.e., the Hermite polynomials are orthonormal with respect to the weighted inner
71 product. In fact, with respect to this inner product, the Hermite polynomials form an orthonormal
72 basis of $L_w^2(\mathbb{R}) = \{f : \langle f, f \rangle_w < \infty\}$.

73 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
74 $x^\alpha = \prod_j (x_j)^{\alpha_j}$ 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

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

75 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.

There are $\binom{m+d-1}{d-1}$ possibilities for a d -dimensional multi-index α such that $|\alpha| = m$. Summing 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) denote the i -th multi-index according to some ordering. Then we can write

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

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

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

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

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

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

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

Diffusion Bridge. To augment or complete the data, we employ diffusion bridge sampling, using a Markov chain Monte Carlo (MCMC) method that goes back to [16, 20]. Let us describe our version here. We suppose our current estimate of $\theta = (\beta, \gamma)$ is given. Define the diffusion bridge process to 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 goal is to generate sample paths of this diffusion bridge. By a sample path, we mean $F - 1$ new samples $\{z_{i,j}\}_{j=1}^{F-1}$ at times $t_i + jh$ with $h = (t_{i+1} - t_i)/F$.

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

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

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

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

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

Let \mathbb{P} denote the law of the diffusion bridge process, and let \mathbb{Q} denote the law of the Brownian bridge (9). Using Girsanov's theorem [15], we can show that

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

where the constant C depends only on x_i and x_{i+1} . The left-hand side is a Radon-Nikodym derivative, equivalent to a density or likelihood; the ratio of two such likelihoods is the accept/reject ratio in the Metropolis algorithm [27].

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

- 114 1. Use (9) to generate samples of \widehat{X}_t at times $t_i + jh$, for $j = 1, 2, \dots, F - 1$ and $h =$
 115 $(t_{i+1} - t_i)/F$. This is the proposal $\mathbf{z}^* = \{z_{i,j}^*\}_{j=1}^{F-1}$.
 116 2. Numerically approximate the integrals in (10) to compute the likelihood of the proposal.
 117 Specifically, we compute

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

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

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

122 We initialize this algorithm with a Brownian bridge path, run for 10 burn-in steps, and then use
 123 subsequent steps as the diffusion bridge samples we seek.

124 **Expectation Maximization (EM).** Let us now give details to justify the intuition expressed above,
 125 that employing the diffusion bridge to augment the data on a fine scale will enable estimation. Let
 126 $\mathbf{z}^{(r)} = \{z_{i,j}^{(r)}\}_{j=1}^{F-1}$ be the r -th diffusion bridge sample path. We interleave this sampled data together
 127 with the observed data \mathbf{x} to create the completed time series

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

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

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

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(k)}) = \mathbb{E}_{\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^{(k)}} [\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})] \\ \approx \frac{1}{R} \sum_{r=1}^R \log p(\mathbf{y}^{(r)} | \boldsymbol{\theta}) \\ = \frac{1}{R} \sum_{r=1}^R \sum_{j=1}^{N-1} \log p(y_{j+1}^{(r)} | y_j^{(r)}, \boldsymbol{\theta}) \\ = -\frac{1}{R} \sum_{r=1}^R \sum_{j=1}^{N-1} \left[\sum_{i=1}^d \frac{1}{2} \log(2\pi h_j \gamma_i^2) \right. \\ \left. + \frac{1}{2h_j} \left\| \Gamma^{-1} (y_{j+1}^{(r)} - y_j^{(r)} - h_j \sum_{\ell=1}^{\widetilde{M}} \beta_{(\ell)} H_{(\ell)}(y_j^{(r)})) \right\|^2 \right].$$

- 138 2. Now for the M step:

$$\boldsymbol{\theta}^{(k+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(k)}) \quad (11)$$

139 To maximize Q over θ , we first assume $\Gamma = \text{diag } \gamma$ is known and maximize over β . This
 140 is a least squares problem. The solution is given by forming the matrix

$$\mathcal{M}_{k,\ell} = \frac{1}{R} \sum_{r=1}^R \sum_{j=1}^N h_j \phi_k^T(y_{j-1}^{(r)}) \Gamma^{-2} \phi_\ell^T(y_{j-1}^{(r)})$$

141 and the vector

$$\rho_k = \frac{1}{R} \sum_{r=1}^R \sum_{j=1}^N \phi_k^T(y_{j-1}^{(r)}) \Gamma^{-2} (y_j^{(r)} - y_{j-1}^{(r)}).$$

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

$$\gamma_i^2 = \frac{1}{RNh} \sum_{r=1}^R \sum_{j=1}^N ((y_j^{(r)} - y_{j-1}^{(r)} - h \sum_{\ell=1}^M \beta_\ell \phi_\ell(y_{j-1}^{(r)})) \cdot e_i)^2$$

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

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

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

152 2 Experiments.

153 1d case - The observed data is created using a known dynamical system and additive Gaussian noise.
 154 Euler-Maruyama time stepping is used from a randomly chosen initial condition to march forward
 155 in time. The initial time is 0 and the final time point is 10. While the time stepping happens with
 156 $h = 0.0001$ and 100,000 internal time steps are taken, only 1000 of these time steps are saved.
 157 A smaller time scale for data generation is used so that the error in the observed data is mostly
 158 accredited to the Gaussian noise and not the error of the numerical stepping method. The system
 159 described above has a stable equilibrium at 1.6 and an unstable equilibrium at -0.6 . The initial
 160 conditions were thus specified to be > -0.6 to avoid the unstable equilibrium.

161 For the 1d case, the dynamical system equation is given as:

$$dX_t = (1 + x - x^2)dt + g dW_t \quad (12)$$

162 In the nonparametric form it can be written as:

$$f_1(x) = \sum \Phi(x) \beta + g dW_t \quad (13)$$

163 where Φ is the matrix created by the hermite basis functions and beta are the parameters to be
 164 inferred. With Hermite functions upto degree 3, Φ can be represented as:

$$\Phi(x) = [h_0(x) \quad h_1(x) \quad h_2(x^2 - 1) \quad h_3(x^3 - 3x)] \quad (14)$$

$$\beta = [\beta_0 \quad \beta_1 \quad \beta_2 \quad \beta_3] \quad (15)$$

165 For inference, 4 sets of experiments were performed to monitor the error in the estimated parameter
 166 β with changes in the Expectation Maximization parameters. The initial guess of β is randomly
 167 generated with mean 0 and variance 0.5 as we assume the β to be sparse. While maximizing the
 168 likelihood of β , noise in the data is known and kept constant at 0.05. The tolerance for relative error
 169 in the estimated parameters is set to be 0.01. The MCMC iteration includes 10 burnin steps and 100
 170 recorded iterations. Different variations of the number of time series, number of data points, number
 171 of sub intervals and amount of noise in the data is used in the following experiments:

- Varying amount of data by providing more time series with the same number of data points in each. We varied the number of time series from 1 to 10 with 101 data points each.
- Varying amount of data by providing random time points ranging from 11 to 101 points from the same 10 time series.
- Varying amount of noise. 10 time series were created using varying amount of noise, ranging from 0.5 to 0.0001 while keeping other parameters constant
- Varying number of subintervals were used for the Brownian bridge sampling ranging from 1 to 9. With 1 sub interval, no Brownian bridge is created and the expectation is merely the sum of the observed data. With 2 intervals, 1 additional latent point is introduced into the system and Brownian bridges are created to compute the likelihood of the observed and the latent data points combined.

For error computation, the Frobenius norm of the difference between estimated $\tilde{\beta}$ and true β in the Hermite space is considered:

$$e = \sqrt{\sum |\beta_i - \tilde{\beta}_i|^2} \quad (16)$$

We demonstrate the method for 1, 2 and 3 dimensional systems.

- For the 1-dimensional system, we use the ? oscillator:

$$dX(t) = (\alpha X(t) + \beta X(t)^2 + \gamma) dt + g dW(t) \quad (17)$$

- For the 2-dimensional system, we use the undamped Duffing oscillator:

$$\begin{aligned} dX_1(t) &= X_2(t)dt + g_1 dW_1(t) \\ dX_2(t) &= (-X_1(t) - X_1^3(t))dt + g_2 dW_2(t) \end{aligned}$$

- For the 3-dimensional case, we consider 2 different form of equations. The first one is the damped Duffing oscillator, a general form of the damped oscillator considered in the 2-dimensional case:

$$\begin{aligned} dX_1(t) &= X_2(t) dt + g_1 dW_1(t) \\ dX_2(t) &= (\alpha X_1(t) - \beta X_1(t) - \delta X_2(t) + \gamma \cos(X_3(t))) dt + g_2 dW_2(t) \\ dX_3(t) &= \omega dt + g_3 dW_3(t) \end{aligned}$$

- Another example considered for the 3-dimensional case is the Lorenz oscillator:

$$\begin{aligned} dX_1(t) &= \sigma(X_2(t) - X_1(t)) dt + g_1 dW_1(t) \\ dX_2(t) &= (X_1(t)(\rho - X_3(t)))dt + g_2 dW_2(t) \\ dX_3(t) &= (X_1(t)X_2(t) - \beta X_3(t)) dt + g_3 dW_3(t) \end{aligned}$$

For simplicity, consider the example where the $X \in \mathbb{R}^2$ and the highest degree of the Hermite polynomial is three, including four Hermite polynomials:

$$\begin{aligned} f(x_1, x_2) &= \sum_{m=0}^2 \sum_{i+j=m} \zeta_{i,j} \psi_{i,j} \\ &= \sum_{d=0}^3 \sum_{i+j=d} \zeta_{i,j} H_i(x_1) H_j(x_2) \\ &= \sum_{i+j=0} \zeta_{i,j} H_i(x_1) H_j(x_2) + \sum_{i+j=1} \zeta_{i,j} H_i(x_1) H_j(x_2) + \sum_{i+j=2} \zeta_{i,j} H_i(x_1) H_j(x_2) + \sum_{i+j=3} \zeta_{i,j} H_i(x_1) H_j(x_2) \\ &= \zeta_{0,0} H_0(x_1) H_0(x_2) + \zeta_{0,1} H_0(x_1) H_1(x_2) + \zeta_{1,0} H_1(x_1) H_0(x_2) + \zeta_{0,2} H_0(x_1) H_2(x_2) \\ &\quad + \zeta_{2,0} H_2(x_1) H_0(x_2) + \zeta_{1,1} H_1(x_1) H_1(x_2) + \zeta_{0,3} H_0(x_1) H_3(x_2) + \zeta_{3,0} H_3(x_1) H_0(x_2) \\ &\quad + \zeta_{2,1} H_2(x_1) H_1(x_2) + \zeta_{1,2} H_1(x_1) H_2(x_2) \end{aligned}$$

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. ISSN 1742-5468. doi: 10.1088/1742-5468/2016/08/083404. URL <http://stacks.iop.org/1742-5468/2016/i=8/a=083404?key=crossref.4719755181cc98b942ea066bb4d58264>.
- [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. Google-Books-ID: 89dZjIid6XcC.
- [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. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1517384113. URL <http://www.pnas.org/lookup/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. ISSN 0162-1459. doi: 10.1080/01621459.2016.1229197. URL <https://doi.org/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. URL <https://papers.nips.cc/paper/1594-learning-nonlinear-dynamical-systems-using-an-em-algorithm.pdf>.
- [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. Google-Books-ID: BCvtssomlCMC.
- [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. ISSN 1364-5021, 1471-2946. doi: 10.1098/rspa.2017.0009. URL <http://rspa.royalsocietypublishing.org/content/473/2204/20170009>.
- [12] F. v. d. Meulen, M. Schauer, and J. v. Waaij. Adaptive nonparametric drift estimation for diffusion processes using Faber–Schauder expansions. *Statistical Inference for Stochastic Processes*, pages 1–26, June 2017. ISSN 1387-0874, 1572-9311. doi: 10.1007/s11203-017-9163-7. URL <https://link.springer.com/article/10.1007/s11203-017-9163-7>.
- [13] 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>.
- [14] J. Nicolau. NONPARAMETRIC ESTIMATION OF SECOND-ORDER STOCHASTIC DIFFERENTIAL EQUATIONS. *Econometric Theory*, 23(05):880, Oct. 2007. ISSN 0266-4666, 1469-4360. doi: 10.1017/S0266466607070375. URL http://www.journals.cambridge.org/abstract_S0266466607070375.
- [15] O. Papaspiliopoulos and G. Roberts. Importance sampling techniques for estimation of diffusion models. *Statistical methods for stochastic differential equations*, 124:311–340, 2012.
- [16] O. Papaspiliopoulos, G. O. Roberts, and O. Stramer. Data Augmentation for Diffusions. *Journal of Computational and Graphical Statistics*, 22(3):665–688, July 2013. ISSN 1061-8600, 1537-2715. doi: 10.1080/10618600.2013.783484. URL <http://www.tandfonline.com/doi/abs/10.1080/10618600.2013.783484>.
- [17] M. Quade, M. Abel, J. N. Kutz, and S. L. Brunton. Sparse Identification of Nonlinear Dynamics for Rapid Model Recovery. *arXiv:1803.00894 [nlin, physics:physics]*, Mar. 2018. URL <http://arxiv.org/abs/1803.00894>. arXiv: 1803.00894.
- [18] 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. ISSN 0021-9991. doi: 10.1016/j.jcp.2017.11.039. URL <http://www.sciencedirect.com/science/article/pii/S0021999117309014>.

- [19] 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. ISSN 0021-9991. doi: 10.1016/j.jcp.2017.07.050. URL <http://www.sciencedirect.com/science/article/pii/S0021999117305582>.
- [20] 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. ISSN 0006-3444. doi: 10.1093/biomet/88.3.603. URL <https://academic.oup.com/biomet/article/88/3/603/340094>.
- [21] 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. ISSN 2375-2548. doi: 10.1126/sciadv.1602614. URL <http://advances.sciencemag.org/content/3/4/e1602614>.
- [22] A. Rutter, 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. URL <http://papers.nips.cc/paper/4967-approximate-gaussian-process-inference-for-the-drift-function-in-stoc>
- [23] 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. ISSN 1364-5021, 1471-2946. doi: 10.1098/rspa.2016.0446. URL <http://rspa.royalsocietypublishing.org/lookup/doi/10.1098/rspa.2016.0446>.
- [24] H. Schaeffer, R. Caflisch, C. D. Hauck, and S. Osher. Sparse dynamics for partial differential equations. *Proceedings of the National Academy of Sciences*, 110(17):6634–6639, Apr. 2013. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1302752110. URL <http://www.pnas.org/cgi/doi/10.1073/pnas.1302752110>.
- [25] H. Schaeffer, G. Tran, and R. Ward. Extracting Sparse High-Dimensional Dynamics from Limited Data. *arXiv:1707.08528 [math]*, July 2017. URL <http://arxiv.org/abs/1707.08528>. arXiv: 1707.08528.
- [26] T. B. Schön, A. Svensson, L. Murray, and F. Lindsten. Probabilistic learning of nonlinear dynamical systems using sequential Monte Carlo. *arXiv preprint arXiv:1703.02419*, 2017. URL <https://arxiv.org/abs/1703.02419>.
- [27] A. M. Stuart. Inverse problems: A Bayesian perspective. *Acta Numerica*, 19:451–559, May 2010. ISSN 0962-4929, 1474-0508. doi: 10.1017/S0962492910000061. URL http://www.journals.cambridge.org/abstract_S0962492910000061.
- [28] G. Tran and R. Ward. Exact Recovery of Chaotic Systems from Highly Corrupted Data. *Multiscale Modeling & Simulation*, 15(3):1108–1129, Jan. 2017. ISSN 1540-3459. doi: 10.1137/16M1086637. URL <https://epubs.siam.org/doi/abs/10.1137/16M1086637>.
- [29] F. van der Meulen, M. Schauer, and H. van Zanten. Reversible jump MCMC for non-parametric drift estimation for diffusion processes. *Computational Statistics & Data Analysis*, 71:615–632, Mar. 2014. ISSN 0167-9473. doi: 10.1016/j.csda.2013.03.002. URL <http://www.sciencedirect.com/science/article/pii/S016794731300090X>.
- [30] N. Verzelen, W. Tao, H.-G. Müller, and others. Inferring stochastic dynamics from functional data. *Biometrika*, 99(3):533–550, 2012. URL <http://nicolas.verzelen.free.fr/pdf/2012-Ver-Biometrika.pdf>.
- [31] M. D. Vrettas, M. Opper, and D. Cornford. Variational mean-field algorithm for efficient inference in large systems of stochastic differential equations. *Physical Review E*, 91(1):012148, Jan. 2015. doi: 10.1103/PhysRevE.91.012148. URL <https://link.aps.org/doi/10.1103/PhysRevE.91.012148>.
- [32] B. Øksendal. *Stochastic Differential Equations: An Introduction with Applications*. Universitext. Springer-Verlag, Berlin Heidelberg, 6 edition, 2003. ISBN 978-3-540-04758-2. URL <http://www.springer.com/us/book/9783540047582>.