
Learning Stochastic Dynamical Systems via Bridge Sampling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and
2 right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points.
3 The word **Abstract** must be centered, bold, and in point size 12. Two line spaces
4 precede the abstract. The abstract must be limited to one paragraph.

5 The goal of this work is to enable automatic discovery of stochastic differential equations (SDE) from
6 time series data.

- 7 1. Literature review.
8 2. What is new and interesting about this work.

9 Points to cover:

- 10 • data specification
11 • Hermite polynomial and drift function representation
12 • Expectation and maximization formulas assuming data is filled in
13 • Filling data in with Brownian bridge
14 • MCMC iterations of brownian bridge using girsanov likelihood
15 • how synthetic data is generated
16 • results: 1D, 2D, 3D damped duffing, 3D lorenz
17 • plots: error of theta vs noise, error vs amount of data (number of data points) parametric
18 curves for noise levels, brownian bridge plots for illustration, ...
19 • Note: constant noise case, not inferring the gvec

20 1 Problem Setup

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

$$dX_t = f(X_t)dt + \Gamma dW_t. \tag{1}$$

24 For rigorous definitions of Brownian motion and SDE, see [Bhattacharya and Waymire 2009, Ok-
25 sendal 2007]. The nonlinear vector field $f : \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the *drift* function, and the $d \times d$ matrix
26 Γ is the *diffusion* matrix. To reduce the number of model parameters, we assume $\Gamma = \text{diag } \gamma$.

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

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

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

31 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, $\langle H_i, H_j \rangle_w =$
32 δ_{ij} , i.e., the Hermite polynomials are orthonormal with respect to the weighted inner product. In
33 fact, with respect to this inner product, the Hermite polynomials form an orthonormal basis of
34 $L_w^2(\mathbb{R}) = \{f : \langle f, f \rangle_w < \infty\}$.

35 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
36 $x^\alpha = \prod_j (x_j)^{\alpha_j}$ 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

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

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

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

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

42 There are $\binom{m+d-1}{d-1}$ possibilities for a d -dimensional multi-index α such that $|\alpha| = m$. Summing this
43 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
44 the i -th multi-index according to some ordering. Then we can write

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

45 **Data.** We consider our data $\mathbf{x} = \{x_j\}_{j=0}^L$ to be direct observations of X_t at discrete points in time
46 $\mathbf{t} = \{t_j\}_{j=0}^L$. Note that these time points do not need to be equispaced.

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

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

53 where $Z_{n+1} \sim \mathcal{N}(0, I)$ is a standard multivariate normal, independent of X_n . Note that $\widetilde{X}_{n+1}|\widetilde{X}_n =$
54 v has a $\mathcal{N}(v + f(v)h, h\Gamma^2)$ distribution. As h decreases, this Gaussian will converge to the true
55 transition density $X_{(n+1)h}|X_{nh} = v$, where X_t refers to the solution of (1).

56 **Diffusion Bridge.** To augment or complete the data, we employ diffusion bridge sampling, using
57 a Markov chain Monte Carlo (MCMC) method that goes back to [Roberts and Stramer, 2001;
58 Papaspiliopoulos, Roberts, and Stramer, 2013]. Let us describe this method here. We suppose our
59 current estimate of $\theta = (\beta, \gamma)$ is given. Then the goal is to generate a sample path of (1) conditioned
60 on both the initial value x_i at time t_i , and the final value x_{i+1} at time t_{i+1} . By a sample path, we
61 mean $F - 1$ new samples $\{z_{i,j}\}_{j=1}^{F-1}$ at times $t_i + jh$ with $h = (t_{i+1} - t_i)/F$.

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

$$d\widehat{X}_t = \Gamma dW_t \quad (7)$$

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

$$\hat{X}_t = \Gamma W_{t-t_i} + x_i - \frac{t-t_i}{t_{i+1}-t_i} (\Gamma W_{t_{i+1}-t_i} + x_i - x_{i+1}) \quad (8)$$

65 Let $\mathbf{z}^{(r)}$ denote the r^{th} diffusion bridge sample path:

$$z^{(r)} \sim z \mid x, \beta^{(k)} \quad (9)$$

The observed and sampled data can be interleaved together to create a time series (completed data)

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

66 of length $N = LF + 1$.

67 **EM.** The EM algorithm consists of two steps, computing the expectation of the log likelihood
68 function (on the completed data) and then maximizing it with respect to the parameters $\theta = (\beta, \gamma)$.

69 1. Start with an initial guess for the parameters, $\theta^{(0)}$.

70 2. For the expectation (or E) step,

$$Q(\theta, \theta^{(k)}) = \mathbb{E}_{\mathbf{z} \mid \mathbf{x}, \theta^{(k)}} [\log p(\mathbf{x}, \mathbf{z} \mid \theta)] \quad (10)$$

71 Our plan is to evaluate this expectation via bridge sampling. That is, we will sample from
72 diffusion bridges $\mathbf{z} \mid \mathbf{x}, \theta^{(k)}$. Then (\mathbf{x}, \mathbf{z}) will be a combination of the original data together
73 with sample paths.

74 3. For the maximization (or M) step, we start with the current iterate and a dummy variable θ
75 and define

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

76 It will turn out that we can maximize this quantity without numerical optimization. All we
77 will need to do is solve a least-squares problem.

78 4. Iterate Step 2 and 3 until convergence.

79 **Details.** With a fixed parameter vector $\theta^{(k)}$, the SDE (1) is specified completely, i.e., the drift and
80 diffusion terms have no further unknowns.

81 Suppose we form R such time series. The expected log likelihood can then be approximated by

$$\begin{aligned} Q(\theta, \theta^{(k)}) &= \mathbb{E}_{\mathbf{z} \mid \mathbf{x}, \theta^{(k)}} [\log p(\mathbf{x}, \mathbf{z} \mid \theta)] \\ &\approx \frac{1}{R} \sum_{r=1}^R \left[\sum_{j=1}^N \left[\sum_{i=1}^d -\frac{1}{2} \log(2\pi h \gamma_i^2) \right] \right. \\ &\quad \left. - \frac{1}{2h} (y_j^{(r)} - y_{j-1}^{(r)} - h \sum_{k=1}^M \beta_k \phi_k(y_{j-1}^{(r)}))^T \Gamma^{-2} (y_j^{(r)} - y_{j-1}^{(r)} - h \sum_{\ell=1}^M \beta_{\ell} \phi_{\ell}(y_{j-1}^{(r)})) \right] \end{aligned}$$

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

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

and the vector

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

We then solve the system

$$\mathcal{M}\beta = \rho$$

for β . Now that we have β , we maximize Q over γ . The solution can be obtained in closed form:

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

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

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

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

$$dX(t) = (\alpha X(t) + \beta X(t)^2 + \gamma) dt + g dW(t) \quad (12)$$

- 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:

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

2 Expectation Maximization Steps

The data provided is in the form of a time series, $X \in \mathbb{R}^d$ at regular time points $t_l, 0 \leq l \leq L$. Note: EM step in the sampling writeup.

3 Brownian bridge sampling

When the inter-observation time of the observed data X is large, the expectation step becomes less accurate. To mitigate this problem, we can fill in the observed data with a Brownian bridge. We generate many samples of the N -dimensional Brownian bridge and accept-reject samples using the Metropolis-Hastings algorithm. The approximation of the likelihood is obtained using the Girsanov likelihood function.

101 **3.1 Brownian bridge**

102 The \mathbb{R}^N dimensional Brownian bridge is defined by the integral:

$$I(t) = \int_0^t \frac{1-t}{1-T} dW(t) \quad (13)$$

103 **3.2 Metropolis Algorithm**