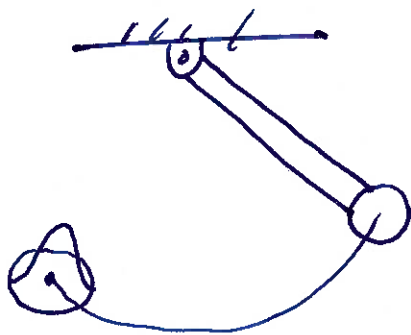


Today

1/12

- Integrating probability into learned models
- Stochastic dynamics
- Maximum likelihood
- Markov processes
- Euler - Maruyama

Motivation



- Tracking a point vs tracking a "blob" of probability

- Sensor noise, may have either
noisy observations } aleatoric uncertainty
or noisy physics

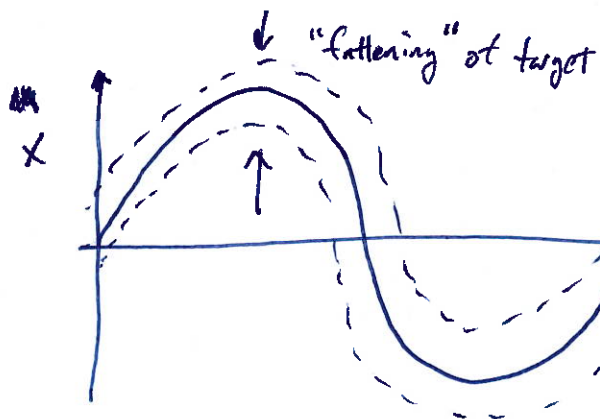
- or an incomplete description of physics (epistemic uncertainty)

- Even w/o noise, training a model in a probabilistic context is often easier

→ probabilistic regularization

$$\dot{X} = f(x) + g(x; \theta)$$

↑
random variable



References for probability

2

- 'Probability Essentials' Jacod + Protter
 - short paperback, rigorous but quick definitions
- 'Probability: theory + examples' Durrett
 - Measure theoretic → gnarly to get started from
- 'Machine learning, a probabilistic perspective' Kevin Murphy
 - Accessible, lots of background and follow-on refs
 - like wikipedia for ML (no ODEs/PDEs though)
- 'Introduction to stochastic integration' Kuo
 - where I'll pull refs from

Overall → I wouldn't recommend following a text for this course, because there isn't a good one for scientific computing + ML

Probability review

3

def a continuous random variable \underline{X} takes a random value $x \in \mathbb{R}$

def a cumulative distribution function (CDF) defines probability over a range of values

$$F_{\underline{X}}(x) = \mathbb{P}(\underline{X} \leq x)$$

def if CDF is differentiable, define probability density function

$$f_{\underline{X}}(x) = \frac{d}{dx} F_{\underline{X}}(x)$$

so $\mathbb{P}(a \leq X \leq b) = \int_a^b f(x) dx = F(b) - F(a)$

Drop subscripts
unless necessary

We'll talk about different RVs + events by different notions of pdfs:

Joint distribution $f(x_1, \dots, x_n) = \mathbb{P}(\underline{X}_1 = x_1, \dots, \underline{X}_n = x_n)$

Marginalization $f(\underline{X} = x) = \sum_y f(\underline{X} = x, \underline{Y} = y)$ Rule of total probability

conditional dist $f(\underline{Y} = y | \underline{X} = x) = \frac{f(\underline{X} = x, \underline{Y} = y)}{f(\underline{X} = x)}$

product rule $f(x, y) = f(x|y) f(y)$

prob. chain rule

$$\begin{aligned}
 f(x_1, \dots, x_n) &= f(x_2, \dots, x_n | x_1) f(x_1) \\
 &= f(x_3, \dots, x_n | x_1, x_2) f(x_2 | x_1) f(x_1) \\
 &\vdots \\
 &= f(x_n | x_1, \dots, x_{n-1}) \dots f(x_2 | x_1) f(x_1)
 \end{aligned}$$

Marginal indep

$$X \perp Y \quad \text{if} \quad f(x, y) = f(x) f(y)$$

$$\text{or in general} \quad f(x_1, \dots, x_n) = \prod_{i=1}^n f(x_i)$$

Conditional indep

$$X|Z \perp Y|Z \quad \text{if} \quad f(x, y | z) = f(x | z) f(y | z)$$

Fitting distributions to data w/ MLE

5

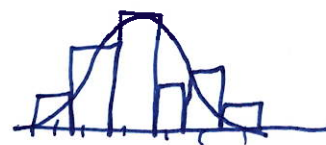
Given a dataset consisting of observations of a RV X

$$\mathcal{D} = \{x_i\}_{i=1}^{N_{\text{data}}}$$

Define the likelihood by evaluating the parameterized joint distribution ~~$f(x, \theta)$~~ $L = f(x_1, \dots, x_{N_{\text{data}}} | \theta)$

Fit the distribution to data by solving

$$\theta^* = \underbrace{-\log L(x_1, \dots, x_{N_{\text{data}}} | \theta)}_{\text{NLL (negative log likelihood)}}$$



Example

- Assume independent data

$$f(x_1, \dots, x_N | \theta) = \prod_{i=1}^N f(x_i | \theta)$$

- Choose marginal distribution

$$f(x_i | \theta) = N(x_i; \mu, \sigma^2) \rightarrow \theta = \{\mu, \sigma^2\}$$

$$N(x; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$

$$\log N = -C - \frac{1}{2} \log \sigma^2 - \frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2$$

$$\text{NLL} = \sum_i \frac{1}{2} \log \sigma^2 + \frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2$$

- solve for θ^*

$$0 = \nabla_{\mu} NLL = \sum_i \nabla_{\mu} \frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2$$

$$0 = - \sum_i (x_i - \mu)$$

$$\sum_i \mu = \sum_i x_i$$

$$N\mu = \sum_i x_i$$

$$\mu = \frac{\sum_i x_i}{N}$$

$$\text{similarly } \sigma^2 = \frac{\sum_i (x_i - \mu)^2}{N}$$

Stochastic differential equations

To account for random forcing/physics we'll expand our notion of ODEs

First, write

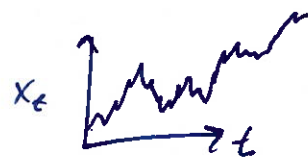
$$\frac{dx}{dt} = f(x, t)$$

Instead as

$$\int_0^t dx = \int_0^t f(x, t) dt$$

For short hand people omit the limits of integration

$$dx_t = f(x_t) dt, \quad x_t = x(t)$$



We will see that random forcing leads to solutions which aren't differentiable, so it's necessary to interpret in integral sense.

To account for stochastic terms, we will consider a SDE

$$dx_t = \underbrace{f(x_t, t) dt}_{\text{drift}} + \underbrace{g(x_t, t) dW_t}_{\text{diffusion}} \quad \swarrow \text{stochastic process}$$

Need to define dW_t and how to integrate it.

In general, we could spend a whole course studying stochastic processes - we'll give an accelerated version here.

We'll consider the Wiener process $W_t = \int_0^t dW_t$

defined via ① $W_0 = 0$

② W has independent increments for $t > 0, u \geq 0$

$W_{t+u} - W_t$ are independent of W_s , for any $s \leq t$

③ W has Gaussian increments w/ variance equal to time increment

$$W_{t+u} - W_t \sim N(0, u)$$

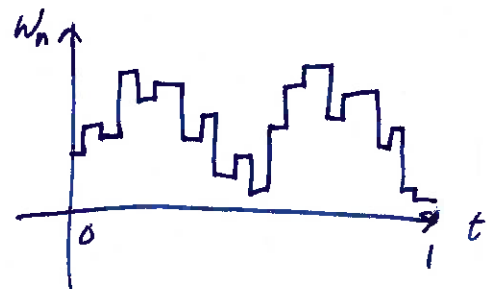
④ W is continuous

Many different constructions of W_t that satisfy these processes, e.g.

Let ξ_1, ξ_2, \dots be IID RVs w/ $E[\xi_i] = 0, \text{var}[\xi_i] = 1$

$$W_n(t) = \frac{1}{\sqrt{n}} \sum_{1 \leq k \leq \lfloor nt \rfloor} \xi_k$$

Then $\lim_{n \rightarrow \infty} W_n(t) = W_t$



To integrate against something like this, consider the
Riemann construction

$$\textcircled{*} \int_0^t g(x_t, t) dW_t = \lim_{\Delta t \rightarrow 0} \sum_{i=1}^{n-1} g(x_{t_i}, t_i) (dW_{t_{i+1}} - dW_{t_i})$$

It turns out that this leads to many pathological issues that violate the usual assumptions from standard calculus.

① We can prove that w/ probability 1, W_t is continuous

② Also w/ prob 1, W_t is nowhere differentiable

In a probability class we would get into details — but this is enough to pose a simple scheme for solving SDEs in our learning problems by discretizing the above $\textcircled{*}$

Euler-Maruyama Method

Given SDE

$$dx_t = f(x_t, t) dt + g(x_t, t) dW_t$$

Solve for $x_{tn} = x(t = nk)$

for $n = 0, 1, \dots$

$$x_{n+1} = x_n + k f(x_n, t_n) + \xi_n g(x_n, t_n)$$

$$\xi_n \sim N(0, k)$$

Some better integrators (see Milstein's method) but require a deeper dive

Using maximum likelihood we can use a Markov process 9
to fit an SDE to data

def For a K^{th} -order Markov process, given a discrete
time series $\vec{y} = \langle y_0, y_1, \dots, y_n \rangle$

$$P(y_i | y_0, \dots, y_{i-1}) = P(y_i | y_{i-1}, \dots, y_{i-K})$$

This means you only need to model the last
 K time steps \rightarrow this is like a multi-step integrator
w/ K steps.

Euler-Maruyama is Markov w/ $K=1$

$$P(x_{n+1} | x_n) = N(x_n + K f(x_n, t_n), K g(x_n, t_n)^2)$$

Pf If $x \sim N(\mu, \sigma^2)$
 $Ax + b \sim N(A\mu + b, A^2\sigma^2)$

Identify $b = x_n + K f_n$

$$A = g$$

$$\sigma^2 = K$$

Now we can go back and derive a NLL

10

swap g 's w/ f 's
on board...

$$\begin{aligned} -\log P(y_1, \dots, y_N) &= -\log P(y_N | y_{N-1}) P(y_{N-1} | y_{N-2}) \dots \\ &= -\log \prod_{i=1}^N P(y_i | y_{i-1}) \\ &= -\sum_{i=1}^N \log P(y_i | y_{i-1}) \\ &= -\sum_{i=1}^N \log N(x_i + \kappa f_i, \kappa g_i^2) \end{aligned}$$

$$= C + \sum_{i=1}^N \frac{1}{2} \log \kappa g_i^2 + \frac{1}{2} \frac{(x_{i+1} - x_i - \kappa f_i)^2}{\kappa g_i^2}$$

If we replace $f(x, t)$ w/ $f(x, t; \theta)$ i.e. make the
 $g(x, t)$ $g(x, t; \theta)$ SDE trainable

We can solve

$$\min_{\theta} \sum_{i=1}^N \log \kappa g_{i, \theta}^2 + \frac{(x_{i+1} - x_i - \kappa f_{i, \theta})^2}{\kappa g_{i, \theta}^2}$$

Example of structure preserving SDEs

11

Metriplectic / GENERIC formalism

$$dx_t = (L \partial_x E + M \partial_x S + K_B \partial_x \cdot M) dt + \sqrt{2K_B M} dW_t$$

where $E \rightarrow$ energy

$S \rightarrow$ entropy

$$L = -L^T$$

$$M = M^T \quad \text{SPD}$$

$K_B \rightarrow$ Boltzmann constant

Degeneracy Condition

$$* L \partial_x S = M \partial_x E = 0$$

- When "coarse-graining" a reversible system, the loss of information manifests as a entropy increase, or introduction of dissipation
- To counter this, stochastic forcing does work on the system.
- Exact fluctuation dissipation theorem guarantees this holds rigorously

Deterministic limit ($k_B \rightarrow 0$)

12

Thm Energy is conserved

Pf $\frac{dE}{dt} = \partial_x E^T \frac{dx}{dt}$ (chain rule)

$$= \partial_x E^T (L \partial_x E + M \partial_x S)$$

$$= \cancel{\partial_x E^T L \partial_x E} + \cancel{\partial_x S^T M \partial_x E}$$

By skew-symmetry By degen condition

Remark Metriplectic is therefore a generalization of Hamiltonian mechanics, with degen condition killing off cross terms between rev and irrev parts.

Thm Entropy is non-decreasing

Pf $\frac{dS}{dt} = \partial_x S^T \frac{dx}{dt}$

~~$$= \partial_x S^T L \partial_x E + \partial_x S^T M \partial_x S$$~~

$$= \underbrace{\partial_x E^T L \partial_x S}_{=0} + \underbrace{\partial_x S^T M \partial_x S}_{\geq 0} \geq 0$$