

CS 886 Winter 2024:

Lecture Notes

1	Introduction	2
1.1	Probability basics	2
1.2	Brownian motion	4
2	The stochastic integral	7
2.1	Integration, writ large	7
2.2	Itô's integral	8
2.3	Martingales	9
2.4	Differentials	10
3	Stochastic differential equations	11
4	Fokker–Planck–Kolmogorov equation	12
4.1	Deriving FPK	12
4.2	Applying to the SDE	14
4.3	Score matching	14

Lecture notes taken, unless otherwise specified, by myself during the Winter 2024 offering of CS 886, taught by Yaoliang Yu.

I am not a graduate student. This is going to be a disaster.

Lectures		Lecture 2	Jan 16	6		
		Lecture 3	Jan 23	11		
Lecture 1	Jan 9	2	Lecture 4	Jan 30	12

Chapter 1

Introduction

A diffusion model is basically a limit of infinite auto-regressive models. We can construct this $\mathbf{x}_{t+1} \approx \mathbf{x}_t + \eta_t \cdot \mathbf{f}_t(\mathbf{x}_t)$ as an ODE and take the limit:

$$d\mathbf{x}_{t+1} = \mathbf{f}_t(\mathbf{x}_t)$$

Make a stochastic differential equation: add a perturbation $G_t(\mathbf{x}_t)$.

$$d\mathbf{x}_{t+1} = \mathbf{f}_t(\mathbf{x}_t) dt + G_t(\mathbf{x}_t) dB_t$$

Clearly, an ODE is an SDE with $G_t \equiv \mathbf{0}$. But an SDE can be integrated to recover an ODE using the score function.

Given an SDE, we can reverse it. That is, if the forward SDE goes from data to noise, a backwards SDE can generate data from noise. The key is estimating that score function.

If we have some normalized distribution $p(x) = \frac{\exp(E(x))}{\int \exp(E(x))}$ for some energy function E , then $\log p(x) = E(x) - \log \int \dots$ which gives us an ignorable integral constant. Then, we can define a loss function.

This is all very vague.

1.1 Probability basics

Definition 1.1.1 (random variable)

Fix a sample space Ω equipped with a σ -algebra \mathcal{F} where $\mu : \mathcal{F} \rightarrow [0, 1]$ assigns probability.

A random variable is a function $\mathbf{X} : (\Omega, \mathcal{F}, \mu) \rightarrow (\mathbb{S}, \mathcal{B}, \mathbf{X}_\# \mu)$ to the state space \mathbb{S} (in this course, always \mathbb{R}).

The distribution of \mathbf{X} , notated $\mathbf{X}_\# \mu$, is a probability measure on $\mathcal{B} \subseteq 2^{\mathbb{S}}$, i.e.,

$$(\mathbf{X}_\# \mu)(S) := \mu(\{\omega : \mathbf{X}(\omega) \in S\}) = \mu(\mathbf{X}^{-1}(S))$$

so long as $\mathbf{X}^{-1}(\mathcal{B}) \subseteq \mathcal{F}$. If $\omega \simeq \mu$, then $\mathbf{X}(\omega) \simeq \mathbf{X}_\# \mu$.

*Lecture 1
Jan 9*

Example 1.1.2. Say $X \simeq \mathcal{N}(0, 1)$.

That formally means $X : (\Omega, \mathcal{F}, \mu) \rightarrow (\mathbb{R}, \mathcal{B})$ and $X_{\#}\mu = \mathcal{N}(0, 1)$.

Example 1.1.3. Say $Y \simeq \chi_1^2$.

That formally means $Y : (\Omega, \mathcal{F}, \mu) \rightarrow (\mathbb{R}, \mathcal{B})$ and $Y_{\#}\mu = \chi_1^2$.

Consider $f : (\mathbb{R}, \mathcal{B}) \rightarrow (\mathbb{R}, \mathcal{B})$ where $x \mapsto x^2$. Compose $f(X) : (\Omega, \mathcal{F}, \mu) \rightarrow (\mathbb{R}, \mathcal{B})$. Then, $f(X) \simeq \chi_1^2$ by definition of the χ^2 distribution.

Observe that the distribution of $f(X)$ is $(f \circ X)_{\#}\mu = f_{\#}[X_{\#}\mu]$.

We want to go the other direction, to solve the inverse problem going from distributions to functions.

Problem 1.1.4

Given distributions P and Q , find f such that $f_{\#}P = Q$.

In a generative model, $P = \mathcal{N}(\mathbf{0}, I)$ is noise, and Q is the data distribution. We want to find f such that if we draw $X \simeq P$, then we can apply $f(X) \simeq Q$.

Theorem 1.1.5 (Representation through Push-forward)

Let P be any continuous distribution on \mathbb{R}^m . For any distribution Q on \mathbb{R}^d , there exist push-forward maps $f : \mathbb{R}^m \rightarrow \mathbb{R}^d$ such that $Z \simeq P \implies f(Z) \simeq Q$ (equivalently $f_{\#}P = Q$).

We need P to be continuous so that we can send it to anything.

In practice, we never have a continuous Q , since our data is discrete. Instead, we only have an approximated empirical \hat{Q} .

Example 1.1.6. Let $X \simeq \mathcal{N}(0, 1)$ and $Y \simeq \chi_1^2$.

Solution. By definition, $f(x) = x^2$ works.

But we can also consider the CDF Φ of $\mathcal{N}(0, 1)$. If we apply $\Phi(X)$:

$$\Pr[\Phi(X) \leq u] = \Pr[X \leq \Phi^{-1}(u)] = \Phi(\Phi^{-1}(u)) = u$$

we get a uniform distribution.

Then, apply the inverse CDF Ψ^{-1} of χ^2 :

$$\Pr[\Psi^{-1}(\Phi(X)) \leq t] = \Pr[\Phi(X) \leq \Psi(t)] = \Psi(t)$$

since $\Psi(X)$ is uniform, which means that $\Psi^{-1}(\Phi(X))$ has CDF Ψ .

That is, $f = \Psi^{-1} \circ \Phi$ works as well. In fact, we know this is a distinct solution, since this function is always increasing, but x^2 is not. \square

Remark 1.1.7. If we add the condition that f is monotonically increasing, the only f is $\Psi^{-1} \circ \Phi$. We call this the optimal transformation.

We can consider $X \xrightarrow{g} \mu$ as a composition of functions $X \xrightarrow{g_1} X_1 \xrightarrow{g_2} X_2 \xrightarrow{g_3} \dots \xrightarrow{g_n} X_n \approx \mu$.

Then, to go backwards, $Y \xleftarrow{h} \mu \dots$

This does not work because [reason]

Definition 1.1.8 (stochastic process)

Consider “time” $t \in \mathbb{T}$. Equivalently:

- A collection of random variables $X : \mathbb{T} \rightarrow \mathbb{R}^\Omega : t \mapsto X(t, \cdot)$
- A random function to paths $X : \Omega \rightarrow \mathbb{R}^\mathbb{T} : \omega \mapsto X(\cdot, \omega)$
- A bivariate function $X(t, \omega) : \mathbb{T} \times \Omega \rightarrow \mathbb{R}$

depending on if we fix the time or state.

We can equivalently write $X(t)$, $X(t, \omega)$, or X_t .

1.2 Brownian motion

Definition 1.2.1 (Brownian motion)

A stochastic process $\{B_t : t \geq 0\}$ such that:

- $B_0 \equiv 0$
- Each increment $B_{t_1} - B_{t_0}$ is independent of all others
- All increments $B_t - B_s \simeq B_{t-s}$
- $B_t \simeq \mathcal{N}(0, t)$
- The function $t \mapsto B_t(\omega)$ is continuous for all ω

Brownian motion is a Gaussian process with covariance kernel $\kappa(s, t) := \mathbb{E}[B_s B_t] = s \wedge t$ (i.e., the minimum of s and t). We can simulate this (discretely) now, since we can simulate a Gaussian process.

Side note: “white noise” is typically defined as the derivative B'_t of a Brownian motion, since $\kappa' := \partial_{12}\kappa$ is also a kernel.

However, we still need the continuity condition. Luckily, we cheat and cite some theorem:

Theorem 1.2.2 (Continuity condition of stochastic processes)

Let \mathbf{X}_t be a stochastic process with index $t \in \mathbb{R}^m$. If for some $\alpha, \beta, L > 0$,

$$\mathbb{E}[\|\mathbf{X}_s - \mathbf{X}_t\|^\alpha] \leq L \|t - s\|^{m+\beta}$$

for all s and t , then there exists a modification $\tilde{\mathbf{X}}_t$ that is locally Hölder continuous of order $\gamma < \beta/\alpha$.

To be Hölder continuous at s of order γ means that for all t around s , $\|\mathbf{X}_s - \mathbf{X}_t\| \leq c \cdot \|s - t\|^\gamma$.

We can show this condition holds for Brownian motion.

Kolmogorov's construction For finitely many t_1, \dots, t_n , define $\mathbf{B}_{1:n} := (B_{t_i}) \simeq \mathcal{N}(\mathbf{0}, K_n)$ where $K_n(t_i, t_j) = t_i \wedge t_j$. Then, by the Kolmogorov extension theorem, a continuous process \mathbf{B}_t exists that satisfies all conditions except continuity.

We know that

$$\begin{aligned} \mathbb{E}|\mathbf{B}_s - \mathbf{B}_t|^{2k} &= \mathbb{E}|\mathbf{B}_{s-t}|^{2k} && \text{stationary} \\ &= \mathbb{E}|\sqrt{t-s} \cdot \mathbf{B}_1|^{2k} && \text{Gaussian} \\ &= |t-s|^k \cdot \mathbb{E}|\mathbf{B}_1|^{2k} \end{aligned}$$

by properties of Brownian motion. With $\alpha = 2k$, $m = 1$, $\beta = k - m = k - 1$, we have continuity of order $\gamma < \frac{k-1}{2k}$ which maxes out to $\frac{1}{2}$.

Theorem 1.2.3 (Irregularity)

Brownian motion is nowhere Hölder continuous of order $\gamma > \frac{1}{2}$.

That is, it is differentiable nowhere.

We can make a heuristic argument as well. At some time τ , Brownian motion is independent of the information up to τ . That is, it is a memoryless Markov process. But then the left and right derivatives can never agree, since they are independently generated.

Definition 1.2.4 (Brownian bridge)

A stochastic process $\{\mathbf{B}_t^\circ : t \in [0, 1]\}$ where:

- $\mathbf{B}_0^\circ \equiv \mathbf{B}_1^\circ \equiv 0$
- Each increment $\mathbf{B}_{t_1}^\circ - \mathbf{B}_{t_0}^\circ$ is independent of all others
- All increments $\mathbf{B}_t^\circ - \mathbf{B}_s^\circ \simeq \mathbf{B}_{t-s}^\circ$
- $\mathbf{B}_t^\circ \simeq \mathcal{N}(0, t(1-t))$
- The function $t \mapsto \mathbf{B}_t^\circ(\omega)$ is continuous for all ω

We can go back and forth between Brownian motions and bridges. If t is restricted to $[0, 1]$, $\mathbf{B}_t^\circ \simeq \mathbf{B}_t - t\mathbf{B}_1$ and $\mathbf{B}_t \simeq \mathbf{B}_t^\circ + t\mathbf{Z}$ with perturbation $\mathbf{Z} \simeq \mathcal{N}(0, 1) \perp \mathbf{B}_t^\circ$.

Since t is in the variance, $\frac{1}{\sqrt{c}}B_{ct}$ is also Brownian, as is $tB_{1/t}$.

[aside: Poisson and Lévy processes]

Now, consider a bunch more ways to construct Brownian motion.

Wiener's construction Let G_n be i.i.d. Gaussian variables. Then, based on a Fourier process, take

$$B_t = tG_0 + \sum_{n=1}^{\infty} \frac{\sin(n\pi t)}{nt} G_n$$

We can apply this by truncating the series to get a Brownian path.

Ciecielski's construction For Haar wavelets (square waves) $\phi_{k/2^n}(t)$, notice that we can write:

$$\int_0^1 B'_t \cdot \phi_{k/2^n}(t) dt \simeq G_{k/2^n}$$

which means we can go the other way to get

$$B_t = \int B'_t dt = tG_0 + \sum G_{k/2^n} \int \phi_{k/2^n}(s) ds$$

which is calculable, I guess.

Lévy's construction (which is actually simple-ish) Initialize points $B_0 = 0$, $B_1 \simeq \mathcal{N}(0, 1)$. Recursively interpolate between points, adding a Gaussian perturbation to each point before saving it. In fact, it's possible to prove this is the same as Ciecielski's method, but without the weird shenanigans.

Donsker's construction Take n i.i.d. random variables $\xi_i \simeq F$ each with mean 0 and unit variance from literally any distribution. Let $S_n = \sum_{i=1}^n \xi_i$ be the cumulative sum. Then, $X_t^n := \frac{1}{\sqrt{n}}S_{[nt]}$ will eventually converge to Brownian motion even though there's no Gaussian. In other words, Brownian motion is a sort of limiting behaviour of a random walk.

In fact, this is a weird stronger statement of the central limit theorem. The CLT says that at S_n , the sum converges to a Gaussian. Donsker says the whole path to get to S_n converges.

*Lecture 2
Jan 16*

Chapter 2

The stochastic integral

Recall: Brownian motion can be considered as a function of ω, t but it's easier to think of it as either:

1. A distribution $\Omega \rightarrow C(\mathbb{R}_+)$ from the sample space to the space of continuous functions of t
2. A collection of random variables $B(\cdot, t)$ for each time point t

It is hard to define Brownian motion if we just say

$$dX_t = f_t(X_t) dt + G_t(X_t) dB_t$$

since Brownian motion is nowhere differentiable, i.e., dB_t does not exist. Instead, write the integral

$$X_T = X_0 + \int_0^T f_t(X_t) dt + \int_0^T G_t(X_t) dB_t$$

where the first term is pretty normal but the second is funny.

How can we define this?

2.1 Integration, writ large

An integral has an integrand $G_t(X_t)$ living in some space and integrator B_t living in some other space. The integral is just a pair $\langle G_t(X_t); B_t \rangle$ that respects some properties.

An integral should have:

- a bilinear form: $\int (\alpha f + \beta g) dB = \alpha \int f dB + \beta \int g dB$ and $\int f d(\alpha B + \beta M) = \alpha \int f dB + \beta \int f dM$
- continuity with respect to integrand and integrator
- generality to a big class of integrators/integrands
- some sort of change of variable formula
- computability

Definition 2.1.1 (Weiner integral)

Let $g : [0, T] \rightarrow \mathbb{R}$ be of bounded variation where $g(0) = g(T) = 0$. We define the integral

$$\int_0^T g(t) dX_t = \cancel{g(t)B_t|_0^T} - \int_0^T X_t dg(t)$$

as the integration by parts.

This works nicely since X_t is continuous with respect to t and if we have sufficiently nice $g(t)$, this works. But diffusion models won't have nice $g(t)$.

What about $\int_0^T B_t dB_t$, which is just...cursed.

Let's try a Riemann sum from MATH 137.

Definition 2.1.2 (Riemann-Stieltjes approximation)

Let $0 = t_0 < t_1 < \dots < t_n < t_{n+1} = 1$ with $\delta_n := \max_{0 \leq k \leq n} |t_{k+1} - t_k| \rightarrow 0$.

Also pick $\tau_k := (1 - \lambda)t_k + \lambda t_{k+1}$ for some $\lambda \in [0, 1]$.

Then, define the approximation R of $\int_0^1 B_t dB_t$

$$\begin{aligned} R &:= R(\delta_n, \lambda) = \sum_{k=0}^n B_{\tau_k} [B_{t_{k+1}} - B_{t_k}] \\ &= \frac{B_1^2}{2} - \underbrace{\frac{1}{2} \sum_{k=0}^n [B_{t_{k+1}} - B_{t_k}]^2}_{S_n(1)} + \underbrace{\sum_{k=0}^n [B_{\tau_k} - B_{t_k}]^2}_{S_n(\lambda)} + \underbrace{\sum_{k=0}^n [B_{t_{k+1}} - B_{\tau_k}][B_{\tau_k} - B_{t_k}]}_{\text{some Gaussians with mean 0}} \\ &= \frac{B_1^2}{2} - \frac{1}{2} + \lambda \end{aligned}$$

Claim that $S_n(\lambda) \rightarrow \lambda$.

Now, depending on what we choose for λ , i.e., what point we pick in the Riemann slice, we get different integrals. Itô picks $\lambda = 0 \implies R = (B_1^2 - 1)/2$ and Stratonovich picks $\lambda = \frac{1}{2} \implies R = (B_1^2)/2$. For this course, we pick Itô's integral.

2.2 Itô's integral

We construct Itô's integral by just wanting simple properties, e.g., the integral over a union should be the sum, integration should preserve convergence (i.e., continuity), etc.

Recall X_t is actually $X(t, \omega)$ a function of two variables. Define an indicator $X(t, \omega) = \mathbf{1}_{(\varsigma, \tau]}(t) \cdot \mathbf{1}_A(\omega)$ for some $A \subseteq \Omega$ that depends only on information up to time ς .

Now, define $\int X_t dB_t := [B_\tau - B_\varsigma] \cdot \mathbf{1}_A$. This gives us linearity (good!) and integrates out t (better!).

For more complex $X(t, \omega) = \sum_k c_k \mathbf{1}_{(\varsigma_k, \tau_k]} \cdot A_k$, define $\int X_t dB_t := \sum_k c_k [B_{\tau_k} - B_{\varsigma_k}]$ to get more linearity.

Suppose $X_t \in \mathcal{F}_t = \sigma(\{B_s : 0 \leq s \leq t\})$ is an arbitrary non-anticipating function depending on information up to t . Also suppose $t \mapsto X_t$ is left continuous. Then, we can discretize X_t by taking rectangles from the left point:

$$X_t^n := X_{([t2^n]-1)/2^n} = \sum_{k=0}^{\infty} X_{k/2^n} \mathbb{I}[k/2^n < t \leq (k+1)/2^n]$$

We can approximate this as $X_{k/2^n}(\omega) \approx \sum_i \frac{i}{2^m} \mathbf{1}_{A_i}(\omega)$ using indicator functions on the sets $A_i := \{(i-1)/2^m < X_{k/2^n} \leq i/2^m\}$.

Finally, define $\int X_t dB_t = \lim_n \int X_t^n dB_t$. This is the same as defining a Lebesgue integral.

We now try to define an isometry. Recall $X(t, \omega) : \mathbb{R}_+ \times \Omega \rightarrow \mathbb{R}$.

Definition 2.2.1 (Doléans measure)

$\lambda = \lambda_{B^2}$ as follows:

$$\begin{aligned} \lambda\{(s, t] \times \mathbf{1}_A\} &:= \mathbb{E}[(B_t^2 - B_s^2) \mathbf{1}_A] = \mathbb{E}[(B_t - B_s)^2 \mathbf{1}_A] \\ &= \mathbb{E}[\mathbb{E}[(B_t^2 - 2(B_t - B_s)B_s - B_s^2) \mathbf{1}_A \mid \mathcal{F}_s]] \\ &= \mathbb{E}[\mathbb{E}[(B_t^2 - 0) \mid \mathcal{F}_s] \mathbf{1}_A - B_s^2 \mathbf{1}_A] \\ &= \mathbb{E}[B_t^2 - B_s^2] \mu(A) \\ &= (t - s) \cdot \mu(A) \end{aligned}$$

We can now treat $X(t, \omega)$ as a random variable from $\mathcal{L}^2(\mathbb{R}_+ \times \Omega, \mathcal{P}, \lambda)$ to \mathbb{R} . Then, the integral $\int X_t dB_t$ is a linear map that sends $X \in \mathcal{L}^2(\mathbb{R}_+ \times \Omega, \mathcal{P}, \lambda)$ to $\mathcal{L}^2(\Omega, \mu)$.

Also, we have the isometry

$$\|X\|_{L^2(\mathbb{R}_+ \times \Omega, \mathcal{P}, \lambda)}^2 = \left\| \int X_t dB_t \right\|_{\mathcal{L}^2(\Omega, \mu)}^2$$

This is nice because if we know $\|X^n(t, \omega) - X(t, \omega)\|^2 \rightarrow 0$, then $\|X^n(t, \omega) - X^m(t, \omega)\|^2 \rightarrow 0$. That is, given convergence of X^n , we get Cauchy. Then, by the definition, the integral exists.

2.3 Martingales

Define $M_t = \int_0^t X_\tau dB_\tau = \int \mathbf{1}_{(0, t]} \cdot X_\tau dB_\tau$. We get these results:

- If $t \mapsto X_t$ is continuous, so is $t \mapsto M_t$
- Defining $M_0 := 0$, we get $\mathbb{E}[M_t] = 0$.
- The usual $\int_s^t X_\tau dB_\tau = \int_0^t X_\tau dB_\tau - \int_0^s X_\tau dB_\tau$
- $\mathbb{E}[M_t \mid \mathcal{F}_s] = \mathbb{E}[(M_t - M_s) + M_s \mid \mathcal{F}_s] = M_s$
- $\mathbb{E}[M_t^2 - M_s^2] = \mathbb{E}[M_t^2 - 2(M_t - M_s)M_s + M_s^2] = \mathbb{E}[M_t - M_s]^2$

Notice that we never used any properties of Brownian motion other than the martingale property. That is, we have defined

$$\int Z_t dM_t$$

for *any* $t \mapsto Z_t$ left-continuous/non-anticipating and M_t continuous martingale. Then, we can define

$$\int_0^t X Z dB = \int_0^t Z dM \quad \text{where} \quad M = \int_0^t X dB$$

as the integral with respect to the integral since the first integral was a martingale.

2.4 Differentials

Define the quadratic variation $[M]_t = S_t^n := \sum_{k=0}^n [M_{t_{k+1}} - M_{t_k}]^2$. As the sum of non-negative terms, it is increasing. Also, for Brownian motion $[B]_t = t$.

Theorem 2.4.1 (Itô's formula)

For all functions $f(M_t, V_t)$ where M_t is a continuous martingale and V_t is continuously differentiable in t ,

$$df(M_t, V_t) = f_x(M_t, V_t) dM_t + f_y(M_t, V_t) dV_t + \frac{1}{2} f_{xx}(M_t, V_t) d[M]_t$$

Then, we can calculate $[\int_0^t X dM] = \int_0^t X^2 d[M]$

Chapter 3

Stochastic differential equations

Lecture 3
Jan 23

Chapter 4

Fokker–Planck–Kolmogorov equation

4.1 Deriving FPK

*Lecture 4
Jan 30*

Let $p(s, \mathbf{x}, t, y)$ be a Markov kernel such that $p(s, \mathbf{x}, t, y) \geq 0$ and $\int p(s, \mathbf{x}, t, \mathbf{y}) d\mathbf{y} = 1$. We also enforce the Chapman-Kolmogorov equation

$$p(s, \mathbf{x}, t, y) = \int p(s, \mathbf{x}, \tau, z) p(\tau, \mathbf{z}, t, \mathbf{y}) d\mathbf{z}$$

for all $s < \tau < t$. Also, define an operator $T_{s,t}$ such that

$$T_{s,t}f(x) = [T_{s,t}(f)](x) = \int p(s, \mathbf{x}, t, \mathbf{y}) f(\mathbf{y}) d\mathbf{y}$$

to integrate out \mathbf{y} .

For example, let $f(\mathbf{y}) = \mathbb{1}_B(\mathbf{y})$. Then, $T_{s,t}f(\mathbf{x}) = \int p(s, \mathbf{x}, t, \mathbf{y}) \mathbb{1}_B(\mathbf{y}) d\mathbf{y} = \Pr[s, \mathbf{x}, t, B]$.

Now, using T notation, we can write

$$\begin{aligned} T_{s,t} &= T_{s,\tau} T_{\tau,t} \\ &= \iint p(s, \mathbf{x}, \tau, \mathbf{z}) p(\tau, \mathbf{z}, t, \mathbf{y}) f(\mathbf{y}) d\mathbf{z} d\mathbf{y} \\ &= \int p(s, \mathbf{x}, \tau, \mathbf{z}) T_{\tau,t}f(\mathbf{z}) d\mathbf{z} \\ &= T_{s,\tau} T_{\tau,t}f(\mathbf{x}) \end{aligned}$$

Can we use this to define a derivative? What happens if $s \rightarrow t$? Define the generator of the Markov process

$$L_t = \lim_{h \downarrow 0} \frac{T_{t,t+h} - T_{t,t}}{h} = \lim_{h \downarrow 0} \frac{T_{t,t+h} - \text{id}}{h}$$

and claim it exists because proving it rigorously is painful.

Consider a generator that can be written as

$$L_t f(\mathbf{x}) = \frac{1}{2} A(t, \mathbf{x}) \cdot \nabla^2 f(\mathbf{x}) + \mathbf{b}(t, \mathbf{x}) \cdot \nabla f(\mathbf{x})$$

which all generators from SDEs can be. This holds when

1. $\lim_{h \downarrow 0} p(t, \mathbf{x}, t+h, B_\varepsilon(\mathbf{x})) = 0$. That is, for a small timestep h , the chance of leaving a small ball ε around \mathbf{x} is also very small. Intuitively, we have “continuity”.
2. $\lim_{h \downarrow 0} \int_{B_\varepsilon(\mathbf{x})} (\mathbf{y} - \mathbf{x}) p(t, \mathbf{x}, t+h, \mathbf{y}) d\mathbf{y} = \mathbf{b}(t, \mathbf{x})$. That is, $\mathbf{b}(t, \mathbf{x})$ is the expected drift (average movement).
3. $\lim_{h \downarrow 0} \int_{B_\varepsilon(\mathbf{x})} (\mathbf{y} - \mathbf{x})(\mathbf{y} - \mathbf{x})^\top p(t, \mathbf{x}, t+h, \mathbf{y}) d\mathbf{y} = A(t, \mathbf{x})$. That is, $A(t, \mathbf{x})$ is the expected covariance. Since the matrix $(\mathbf{y} - \mathbf{x})(\mathbf{y} - \mathbf{x})^\top$ is rank 1 and positive semi-definite, $A(t, \mathbf{x})$ must also be positive semi-definite.

Equivalently, we can express this by saying

1. $\Pr[\|\mathbf{x}_{t+h} - \mathbf{x}_t\| \geq \varepsilon \mid \mathbf{x}_t] \rightarrow 0$ as $h \downarrow 0$.
2. $\frac{\mathbb{E}[\mathbf{x}_{t+h} - \mathbf{x}_t \mid \mathbf{x}_t = \mathbf{x}]}{h} \rightarrow \mathbf{b}(t, \mathbf{x})$
3. $\frac{\mathbb{E}[(\mathbf{x}_{t+h} - \mathbf{x}_t)(\mathbf{x}_{t+h} - \mathbf{x}_t)^\top \mid \mathbf{x}_t = \mathbf{x}]}{h} \rightarrow A(t, \mathbf{x})$

Now, let us write $\langle L_t f; p \rangle$ for the integral $\int L_t f \cdot p d\mathbf{y}$. We have a useful property that

$$\langle L_t f; p \rangle = \langle f; L_t^* p \rangle$$

where L_t^* is the adjoint of L_t . For example, with our generator,

$$\begin{aligned} \left\langle \frac{1}{2} A \cdot \nabla^2 f + \mathbf{b} \nabla \cdot f; p \right\rangle &= \left\langle \nabla^2 f; \frac{1}{2} p A \right\rangle + \langle \nabla f; p \mathbf{b} \rangle \\ &= \left\langle f; \nabla^2 \frac{1}{2} p A \right\rangle - \langle f; \nabla p \mathbf{b} \rangle \\ &= \left\langle f; \nabla^2 \cdot \frac{1}{2} p A - \nabla \cdot p \mathbf{b} \right\rangle \\ L_t^* p &= \nabla^2 \frac{1}{2} A - \nabla \mathbf{b} \end{aligned}$$

using the integration by parts formula $\langle f'; g \rangle = -\langle f; g' \rangle$.

We can now write the (forward) FPK equation

$$\partial_t p(s, \mathbf{x}, t, \mathbf{y}) = L_t^* p(s, \mathbf{x}, t, \mathbf{y})$$

and the backward FPK equation

$$\partial_s p(s, \mathbf{x}, t, \mathbf{y}) = -L_s^* p(s, \mathbf{x}, t, \mathbf{y})$$

Proof. To derive this, take a function f and consider the derivative of the integral

$$\begin{aligned} \partial_t \langle f; p(s, \mathbf{x}, t, \mathbf{y}) \rangle &= \partial_t T_{s,t} f(\mathbf{x}) \\ &= \lim_{h \downarrow 0} \frac{T_{s,t+h} - T_{s,t}}{h} f(\mathbf{x}) \\ &= \lim_{h \downarrow 0} T_{s,t} \frac{T_{t,t+h} - T_{t,t}}{h} f(\mathbf{x}) \\ &= T_{s,t} L_t f(\mathbf{x}) \\ &= \langle L_t f; p(s, \mathbf{x}, t, \mathbf{y}) \rangle \\ &= \langle f; L_t^* p(s, \mathbf{x}, t, \mathbf{y}) \rangle \end{aligned}$$

Then, if $\langle f; \partial_t p \rangle = \langle f; L_t^* p \rangle$ for all f , we must have $\partial_t p = L_t^* p$. □

4.2 Applying to the SDE

Recall the SDE $dX_t = \mathbf{b}_t(X_t) dt + G_t(X_t) dB_t$.

Claim that it has a generator $L_t f = \mathbf{b}_t \cdot \nabla f + \frac{1}{2} [G_t G_t^\top] \cdot \nabla^2 f$.

Then, the forward FPK equation is $\partial_t p = L_t^* p = -\nabla \cdot (p \mathbf{b}_t) + \frac{1}{2} \nabla^2 \cdot (p G_t G_t^\top)$

Proof. Take a function f . By Itô's formula,

$$\begin{aligned} df(X_t) &= \nabla f(X_t) dX_t + \frac{1}{2} \nabla^2 f(X_t) \cdot G_t G_t^\top dt \\ &= \left[\mathbf{b}_t \cdot \nabla f + \frac{1}{2} \nabla^2 f + G_t G_t^\top \right] dt + \nabla f \cdot G_t dB_t \\ d \langle f; p(s, \mathbf{x}, t, \cdot) \rangle &= \langle L_t f; p(s, \mathbf{x}, t, \cdot) \rangle dt + 0 \\ \langle f; \partial_t p(s, \mathbf{x}, t, \cdot) \rangle &= \langle f; L_t^* p(s, \mathbf{x}, t, \cdot) \rangle dt \end{aligned}$$

and as above, since the formula must hold for all f , we have what we want. \square

Now, suppose $G = 0$ so that $dX_t = \mathbf{b}_t(X_t) dt$. Then, $\partial_t p = -\nabla \cdot (p \mathbf{b})$. This is the continuity equation.

Recall the heat equation $\partial_t p = \Delta p = \sum \partial_i^2 p$. We recover the heat equation when $G = \sqrt{2}I$. In other words, the heat equation is the SDE $dX_t = \sqrt{2} dB_t$.

To solve this, try transforming the PDE into an ODE. Let $p(t, x) = p(t, \cdot)$ and consider $p : t \mapsto p(t, \cdot)$ as a function of just t . Then, we have $\frac{dp_t}{dt} = -\nabla f(p_t)$ with the energy function $f(p) = \frac{1}{2} \|\nabla p\|_2^2$.

Running gradient descent on f will find us the solution. Consider the derivative

$$\begin{aligned} \left. \frac{df(p + \varepsilon q)}{d\varepsilon} \right|_{\varepsilon=0} &= \frac{d \frac{1}{2} \|\nabla(p + \varepsilon q)\|}{d\varepsilon} \\ &= \frac{d \frac{1}{2} \int \langle \nabla p + \varepsilon \nabla q, \nabla p + \varepsilon \nabla q \rangle}{d\varepsilon} \\ &= \int \langle \nabla p, \nabla q \rangle dx \\ &= \langle \nabla p; \nabla q \rangle \\ &= -\langle \Delta p; q \rangle \end{aligned}$$

Then, $\frac{dp_t}{dt} = \Delta p_t$.

4.3 Score matching

Consider two densities p and q . We define the Fisher divergence

$$F(p \parallel q) = \mathbb{E}_{X \sim q} \|\nabla \log p(x) - \nabla \log q(x)\|_2^2 = \mathbb{E}_{X \sim q} \|\mathbf{s}_p(x) - \mathbf{s}_q(x)\|^2$$

Also, define noisy versions $p_t = p * \mathcal{N}(0, 2t)$ and $q_t = q * \mathcal{N}(0, 2t)$. Then, the KL divergence

$$\frac{d}{dt} \text{KL}(p_t \parallel q_t) = \frac{d}{dt} \int p_t \log \frac{p_t}{q_t} dx = -F(p_t \parallel q_t)$$

Proof. Recall that the standard Gaussian $\mathcal{N}(0, 2t)$ will satisfy the heat equation $\partial_t p = \Delta p$. Then, since it is a convolution, p_t also satisfies the heat equation. That is,

$$\partial_t(p * \mathcal{N}(0, 2t)) = p * \partial_t \mathcal{N}(0, 2t) = p * \Delta \mathcal{N}(0, 2t) = \Delta(p * \mathcal{N}(0, 2t))$$

Now,

$$\begin{aligned} & \frac{d}{dt} \text{KL}(p_t \parallel q_t) \\ &= \int \frac{dp_t}{dt} \cdot \log \frac{p_t}{q_t} dx + \int \cancel{p_t \cdot \frac{\partial_t p_t}{p_t}} dx - \int p_t \frac{\partial_t q_t}{q_t} dx \\ &= \int \Delta p_t \cdot \log p_t dx - \int \Delta p_t \log q_t dx - \int p_t \frac{\Delta q_t}{q_t} dx \end{aligned}$$

The first term

$$\int \Delta p_t \cdot \log p_t dx = \langle \Delta p_t; \log p_t \rangle = \left\langle p_t \frac{\Delta p_t}{p_t}; \log p_t \right\rangle = - \int p_t \|\mathbf{s}_{p_t}\|^2 dx$$

and the second term

$$- \int \Delta p_t \log q_t dx = \int p_t \langle \mathbf{s}_{p_t}, \mathbf{s}_{q_t} \rangle dx$$

and the third term

$$\begin{aligned} - \int p_t \frac{\Delta q_t}{q_t} dx &= - \left\langle \frac{p_t}{q_t}; \Delta q_t \right\rangle \\ &= \left\langle \nabla \frac{p_t}{q_t}; \nabla q_t \right\rangle \\ &= \int \frac{\nabla p_t \cdot q_t - p_t \cdot \nabla q_t}{q_t^2} \cdot \nabla q_t dx \\ &= \int \frac{\nabla p_t}{p_t} \cdot \frac{\nabla q_t}{q_t} p_t - p_t \|\mathbf{s}_{q_t}\|^2 dx \\ &= \int \langle \mathbf{s}_{p_t}, \mathbf{s}_{q_t} \rangle \cdot p_t - p_t \|\mathbf{s}_{q_t}\|^2 dx \end{aligned}$$

Then, we have

$$- \int p_t \|\mathbf{s}_{p_t}\|^2 dx + 2 \int p_t \langle \mathbf{s}_{p_t}, \mathbf{s}_{q_t} \rangle dx - \int p_t \|\mathbf{s}_{q_t}\|^2 dx = - \int p_t \|\mathbf{s}_{p_t} - \mathbf{s}_{q_t}\|^2 dx = -F(p_t \parallel q_t)$$

as desired. \square

Now, given some distributions p and q , we want to find T such that $q = T_{\#}^{-1}p$ (i.e., $p = T_{\#}q$).

For example, if x is noise, we want to have $T^{-1}(x)$ be data. Mathematically, this is the same as taking data and pushing it to noise. If we can do one, we can do the other.

We build T over time. Let $p_t := (T_t^{-1})_{\#}p$ and $q_t := (T_t)_{\#}q$.

We can build p_t as the change in variables formula $p_t(x) = p(T_t x) \cdot |\det T'_t(x)|$.

To train, we want to minimize KL-divergence

$$\min_{T_t} \text{KL}(q \parallel p_t) = \min_{T_t} \int q \log \frac{q}{p_t} = \min_{T_t} - \int q \cdot \log p_t \approx \max_{T_t} \sum_i \log p_t \mathbf{x}_i$$

we end up with a maximum log-likelihood estimator.