

SAMPLING.

Monday, March 23, 2020 4:11 PM

Yuhu

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Sample $x \in \mathbb{R}^n$ with prob $\propto e^{-f(x)}$.

e.g. $f(x) = \|x\|^2$ Gaussian

$$f(x) = \begin{cases} 1 & x \in K \\ \infty & x \notin K. \end{cases}$$

Intractable in general.

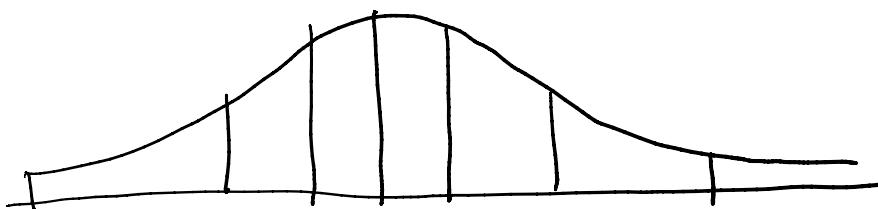
What classes of functions have polytime algorithms?

P1 how to sample from $v(x) \propto e^{-\frac{x^T Ax}{2}} \cdot \mathbf{1}_{\{x \geq 0\}}$?

$e^{-\frac{x^T Ax}{2}}$ is easy: get $y \sim N(0, I)$

apply $x = A^{-\frac{1}{2}}y$ then $x \sim N(0, A^{-1})$

$N(0, I)$: generate each $x_i \sim N(0, 1)$



many numerically efficient methods

many numerically efficient ...

What about PI?

Idea: use a randomized process

$$dX_t = -\nabla f(X_t)dt \rightarrow \text{to go to a local min of } f.$$

$dX_t = -\nabla f(X_t)dt + \sqrt{2} dW_t$

Langevin Dynamics
a Diffusion

(discrete time)

$$X^{k+1} = X^k - \nabla f(X_k) \cdot h + \sqrt{2h} Z \quad \begin{matrix} \uparrow \\ \text{infinitesimal Brownian motion.} \end{matrix}$$
$$Z \sim N(0,1).$$

Theorem: For any smooth f , the density $\propto e^{-f}$ is stationary for Langevin Diffusion.

This theorem follows from a very general analysis of stochastic processes of this form.

SDE Stochastic Differential Equation

$$dX_t = \mu(t)dt + \sigma(t)dW_t$$

..

m

$$dX_t = \mu(t) dt + \sigma(t) dW_t$$

$$X_t, \mu_t \in \mathbb{R}^n \quad \sigma(t) \in \mathbb{R}^{n \times m} \quad dW_t \in \mathbb{R}^m.$$

Thm 2 [Fokker-Planck] $X \sim p_0$

$$\frac{dp_t(x)}{dt} = - \sum_i \frac{\partial}{\partial x_i} (N_i(t) p_t(x)) + \frac{1}{2} \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} (D_{ij}(t) p_t(x))$$

$$- \nabla \cdot (\mu p) + \frac{1}{2} \nabla \cdot (\nabla \cdot (p D))$$

$$D(t) = \sigma(t) \sigma(t)^T$$

Pf (of Thm 1). For stationary p_t , with $D = 2 I$

$$0 = \frac{dp_t(x)}{dt} = - \sum_i \frac{\partial}{\partial x_i} (-p_t(x) \frac{\partial f(x)}{\partial x_i}) + \sum_i \frac{\partial^2}{\partial x_i^2} (p_t(x))$$

$$0 = \sum_i p_t(x) \frac{\partial^2 f(x)}{\partial x_i^2} + \frac{\partial p_t(x)}{\partial x_i} \frac{\partial f(x)}{\partial x_i} + \frac{\partial^2}{\partial x_i^2} p_t(x)$$

OR

$$0 = \sum_i \frac{\partial}{\partial x_i} \left(p_t(x) \frac{\partial f(x)}{\partial x_i} + \frac{\partial}{\partial x_i} p_t(x) \right)$$

$$= \sum_i \frac{\partial}{\partial x_i} \left(p_t(x) \left(\frac{\partial f(x)}{\partial x_i} + \frac{\partial}{\partial x_i} \log p_t(x) \right) \right)$$

$$= \sum_i \frac{\partial}{\partial x_i} \left(p_t(x) \left(\frac{\partial f^m}{\partial x_i} + \frac{\partial}{\partial x_i} \cdot \alpha \cdot \dots \right) \right)$$

$$0 = \sum_i \frac{\partial}{\partial x_i} \left(p_t(x) \frac{\partial}{\partial x_i} \log \left(\frac{p_t(x)}{e^{-f(x)}} \right) \right)$$

$p_t(x) = C \cdot e^{-f(x)}$ is a solution.

We will review the proof of F-P at the end.

What is the rate of convergence of LD?

First in continuous time. How to compare distributions?

$$p_0 \sim p_t \sim \phi \propto e^{-f}$$

$d_w(p, q)$ Wasserstein distance metric $\|\cdot\|$ on domain of p .

$$= \min_{\Pi: \text{coupling of } X \sim p, Y \sim q} \mathbb{E}(\|X - Y\|^2)$$

Π : coupling of

$$X \sim p, Y \sim q$$

marginal of Π in X is p , in Y is q .

$$\text{More familiar } d_{TV}(p, q) = \frac{1}{2} \int |p(x) - q(x)| dx.$$

and KL-divergence (asymmetric)

$$d_{KL}(p\|q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

Thm 3. Assume f is μ -strongly convex.

Then LD starting at X_0, Y_0 can be coupled

st. $E(\|X_t - Y_t\|^2) \leq e^{-2\mu t} \|X_0 - Y_0\|^2$

(or. LD converges to \bar{e}^f .

Pf. Take $Y_0 \sim \bar{e}^{-f}$ then $Y_t \sim \bar{e}^{-f}$.

Pf (of Thm 3). Use identity coupling.

$$dX_t = -\nabla f(X_t)dt + \sqrt{2} dW_t \quad) \text{equal.}$$

$$dY_t = -\nabla f(Y_t)dt + \sqrt{2} dW_t$$

$$\frac{d}{dt}(X_t - Y_t) = -(\nabla f(X_t) - \nabla f(Y_t))$$

so $\frac{d}{dt} \|X_t - Y_t\|^2 = 2 \langle X_t - Y_t, \frac{d}{dt}(X_t - Y_t) \rangle$

$\dots \langle x - y, x - y \rangle >$

$$\frac{d}{dt} = -2 \langle \nabla f(x_t) - \nabla f(y_t), x_t - y_t \rangle$$

Using strong convexity,

$$f(x_t) - f(y_t) \geq \langle \nabla f(y_t), x_t - y_t \rangle + \frac{\mu}{2} \|x_t - y_t\|^2$$

$$f(y_t) - f(x_t) \geq \langle \nabla f(x_t), y_t - x_t \rangle + \frac{\mu}{2} \|x_t - y_t\|^2$$

$$\langle \nabla f(x_t) - \nabla f(y_t), x_t - y_t \rangle \geq \mu \|x_t - y_t\|^2$$

$$\Rightarrow \frac{d}{dt} \|x_t - y_t\|^2 \leq -2\mu \|x_t - y_t\|^2$$

i.e.

$$\frac{d}{dt} \log \|x_t - y_t\|^2 \leq -2\mu$$

$$\Rightarrow \|x_t - y_t\|^2 \leq e^{-2\mu t} \|x_0 - y_0\|^2.$$

back to Fokker-Planck and SDE

Given $dX_t = \mu(X_t) dt + \sigma(t) dW_t$

Given $dX_t = f(X_t) dt + \sigma dW_t$

What can we say about $df(X_t)$?

In usual calculus, this is the chain rule,

$$df(X_t) = \nabla f(X_t) dX_t$$

Which comes from Taylor expansion of f

$$f(X_t + h) = f(X_t) + \nabla f(X_t) h + \frac{1}{2} h^T \nabla^2 f(X_t) h + \dots$$

$h \rightarrow 0$ only second term matters in

$$\frac{f(X_t + h) - f(X_t)}{h} \rightarrow \nabla f(X_t).$$

But with an SDE, we have both dt and dW_t .

Note that $W_t - W_0 \sim N(0, t)$

$$\mathbb{E}(W_t^2) = t$$

$$dW_t = \omega dt$$

$$\mathbb{E}((dW_t)^2) = \mathbb{E}(\omega^2 dt) = dt.$$

So,

$$df(X_t) = \nabla f(X_t) dX_t + \frac{1}{2} dX_t^T \nabla^2 f(X_t) dX_t$$

$$\begin{aligned}
 d f(x_t) &= \nabla f(x_t) dX_t + \frac{1}{2} dX_t \nabla^2 f(x_t) dX_t \\
 &= \nabla f(x_t) P(x_t) dt + \nabla f(x_t) \sigma(t) dW_t \\
 &\quad + \frac{1}{2} P(x_t)^T \nabla^2 f(x_t) P(x_t) (dt)^2 \rightarrow 0 \\
 &\quad + P(x_t)^T \nabla^2 f(x_t) \sigma(t) dW_t dt \rightarrow 0 \\
 &\quad + \frac{1}{2} \underline{dW_t^T \sigma(t)^T \nabla^2 f(x_t) \sigma(t)} \underline{dW_t}
 \end{aligned}$$

$$\begin{aligned}
 d f(x_t) &= \nabla f(x_t)^T P(x_t) dt + \frac{1}{2} \nabla f(x_t)^T \sigma(t) dW_t \\
 &\quad + \frac{1}{2} \langle \nabla^2 f(x_t), \sigma(t) \sigma(t)^T \rangle dt
 \end{aligned}$$

Itô's rule : chain rule for stochastic calculus.

Now let's prove F-P.

Pf. (of Thm 2) ϕ smooth function

$$E_p(\phi(x)) = E_p(\phi(x_t))$$

... + ... L + ... left sides.

Take derivatives wrt t on both sides.

$$\int \phi(x) d\hat{p}_t(x) dx = \mathbb{E} \left(\nabla \phi(x_t) \cdot dX_t + \frac{1}{2} \langle \nabla^2 \phi(x_t), \sigma(x_t) \sigma(x_t)^T \rangle dt \right)$$

$$= \mathbb{E} \left(\nabla \phi(x_t)^T \nu(x_t) dt + \underbrace{\nabla \phi(x_t) \sigma(x_t) dW_t}_{\textcircled{1}} + \frac{1}{2} \langle \quad \rangle dt \right) \textcircled{2}$$

$$\mathbb{E}(dW_t) = 0$$

Note that $x_t \sim \hat{p}_t$. So,

$$\textcircled{1} = \mathbb{E}_{\hat{p}_t} \left(\nabla \phi(x_t)^T \nu(x_t) dt \right) = \mathbb{E}_{\hat{p}_t} \left(\nabla \phi(x) \nu(x) dt \right)$$

$$= \int \nabla \phi(x)^T \nu(x) \hat{p}_t(x) dt \cdot dx$$

integrate by parts

$$= \left[\phi(x) \nu(x) \hat{p}_t(x) \right] - \int \phi(x) \sum_i \frac{\partial (\nu_i(x) \hat{p}_t(x))}{\partial x_i} dt \cdot dx$$

$$= 0 .$$

$$\textcircled{2} = \frac{1}{2} \int \langle \nabla^2 \phi(x), \sigma(x) \sigma(x)^T \rangle p_t(x) dt$$

integrate by parts

$$= \frac{1}{2} \int \langle \nabla \phi(x), \sum_i \frac{\partial}{\partial x_i} (p_t(x) (\sigma(x) \sigma(x)^T)_{ii}) \rangle dt dx$$

$$= \frac{1}{2} \int \phi(x) \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} (p_t(x) (\sigma(x) \sigma(x)^T)_{ij}) dt$$

$$\therefore \int \phi(x) \left[\frac{dp_t(x)}{dt} + \sum_i \frac{\partial}{\partial x_i} (p_t(x) \nabla_i(x)) - \frac{1}{2} \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} (p_t(x) (\sigma(x) \sigma(x)^T)_{ij}) \right] dx = 0$$

+ smooth ϕ

\Rightarrow

$= 0 \quad \checkmark$

There is a more direct connection to gradient descent.

We start at p_0

$$\dots + n \propto \bar{f}$$

we want to end at $q \propto e^{-t}$.

Can we phrase this as an OPT problem?

YES.

$$\text{Min } d(p_t, q)$$

What distance to use?

is the gradient implementable?

One nice answer.

use $d_{KL}(p_t \| q)$ minimized when $p_t = q$

with Wasserstein metric as distance between measures.

Then G_t is LD!