

Lecture 18: Langevin Dynamics and Diffusion models.

Last class: sampling over discrete domains.

This class: sampling over continuous domains.

of pixels \rightarrow

e.g. distributions over images are more naturally thought of as distributions over \mathbb{R}^d

(Super) crash course on continuous probability.

discrete distribution over $\{1, \dots, n\}$ is specified by probability mass function

$$p(i) = \Pr[X=i]. \quad \text{An event } E \subseteq \{1, \dots, n\} \text{ has probability}$$
$$p(i) \geq 0, \quad \sum_{i=1}^n p(i) = 1. \quad \Pr[X \in E] = \sum_{i \in E} p(i).$$

A continuous distribution over \mathbb{R} (or \mathbb{R}^d) is specified by a probability density function (pdf)

$$p: \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}.$$

$$p(x) \geq 0, \quad \int p(x) dx = 1.$$

For any event $E \subseteq \mathbb{R}^d$, we have

$$\Pr[X \in E] = \int_E p(x) dx.$$

It usually doesn't make much sense to talk about

$$\Pr[X=x] \leftarrow \text{will be } 0$$

For any function $f: [n] \rightarrow \mathbb{R}$,
its expectation is

$$\mathbb{E}[f(X)] = \sum_{i=1}^n f(i) p(i)$$

For any function $f: \mathbb{R}^d \rightarrow \mathbb{R}$,
its expectation is

$$\mathbb{E}[f(X)] = \int_{\mathbb{R}^d} f(x) p(x) dx$$

What are "nice" distributions over continuous domains?

Stereotypical: Gaussian distribution.

In 1-D:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$\begin{array}{ll} \text{mean } \mu & \mathbb{E}[X] = \mu \\ \text{variance } \sigma^2 & \mathbb{E}[(X-\mu)^2] = \sigma^2. \end{array}$$

In d dimensions:

$$p(x) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp(-(x-\mu)^T \Sigma^{-1} (x-\mu)).$$

$$\text{mean } \mu: \mathbb{E}[X] = \mu$$

$$\text{covariance matrix } \Sigma: \mathbb{E}[(X-\mu)(X-\mu)^T] = \Sigma.$$

If $\Sigma = I \rightarrow d$ independent univariate Gaussians.

General $\Sigma \rightarrow$ corresponds to $\Sigma^{1/2} \cdot X$, $X \sim \mathcal{N}(\mu, I)$.

Key point: the pdf has the form

$$p \propto \exp(f(x)), \quad p = \frac{1}{Z} \exp(f(x))$$

$$f(x) = -(x - \mu)^T \Sigma (x - \mu).$$

Key property of f : it is concave!

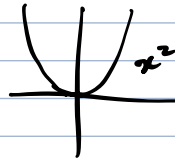
↑
partition function.



e.g. think about when $\mu = 0$, $\Sigma = I$.

$$f(x) = -\|x\|^2.$$

$\|x\|^2$ is convex.



so $-\|x\|^2$ is concave.

Def: We say a distribution $p \propto \exp(f(x))$ is log-concave if $f(x)$ is concave.

Note: Most distributions are not logconcave. The theory works best for logconcave, but in practice things work more generally (sometimes...).

Several questions:

1 Given access to $f(x)$ and/or $\nabla f(x)$, can you efficiently sample from

$$p(x) \propto \exp(f(x))?$$

2. Given samples $X_1, \dots, X_n \sim p(x)$, can you efficiently sample from p ?

generative modeling / distribution learning.

Langevin dynamics.

connection between
optimization \leftrightarrow sampling.

Let $\beta > 0$, and consider "Gibbs distribution with inverse temperature β ".

$$p_\beta \propto \exp(-\beta f(x))$$

Claim: As $\beta \rightarrow \infty$, almost all mass of p_β is at minimizer of f , if f is "nice"

easy to show in discrete case: $f: [n] \rightarrow \mathbb{R}$, we take

$$p_i \propto \exp(-\beta f(i)).$$

so $\Pr_{X \sim p_\beta}[X=i] = \frac{\exp(-\beta f(i))}{\sum_{i=1}^n \exp(-\beta f(i))}$

Let $B_\epsilon = \{i: f(i) \leq (1+\epsilon) f_{\min}\}$, $f_{\min} = \min_i f(i)$.

Then $\Pr_{X \sim p_\beta}[X \notin B_\epsilon] = \frac{\sum_{i \notin B_\epsilon} \exp(-\beta f(i))}{\sum_{i=1}^n \exp(-\beta f(i))} \quad (*)$

But we have:

$$\sum_{i \notin B_\epsilon} \exp(-\beta f(i)) \leq n \cdot \exp(-\beta(1+\epsilon)f_{\min})$$

$$\sum_{i=1}^n \exp(-\beta f(i)) \geq \exp(-\beta \cdot f_{\min}).$$

$$\text{so } (*) \leq \frac{n \cdot \exp(-\beta(1+\epsilon)f_{\min})}{\exp(-\beta f_{\min})}$$

$$= n \cdot \exp(-\epsilon \beta f_{\min}), \text{ so if } \beta \geq \frac{\log n + \log 1/\delta}{\epsilon f_{\min}}, \text{ then}$$

$$\leq n \cdot \exp(-(\log n + \log 1/\delta))$$

$$\leq \delta.$$

As $\beta \rightarrow \infty$, i.e. as temperature drops, our distribution puts more and more mass around minimizer.

For $\beta < \infty$, it's like minimizing, but also need to keep some mass overall.

For minimization / optimization, a standard method is gradient descent

$$x_{t+1} = x_t - \eta_t \nabla f(x_t)$$

\uparrow
Step size.

For "nice" f , this converges to minima.

Langevin dynamics:

$$x_{t+1} = x_t - \left(\frac{\epsilon}{2}\right) \nabla f(x_t) + \left(\sqrt{\epsilon}\right) z_t, \quad z_t \sim \mathcal{N}(0, I).$$

this scaling is important

A perspective from stochastic calculus: this is really a discretization of a continuous-time process.

Consider gradient descent, and reparametrize so that timesteps are small. x_1, \dots, x_T
 $x_0, x_{t+\Delta t}, x_{2\Delta t}, \dots, x_1$

$$x_{t+\Delta t} = x_t - \eta \nabla f(x_t)$$

$$\frac{x_{t+\Delta t} - x_t}{\Delta t} = -\eta \nabla f(x_t)$$

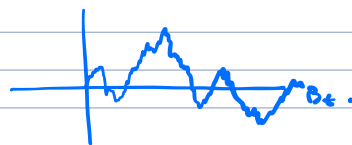
$$x_{t+\Delta t} - x_t = -\eta \nabla f(x_t) \Delta t$$

$$dx_t = -\eta \nabla f(x_t) dt \text{ as } \Delta t \rightarrow 0.$$

↖ PDE in x, t . This is called the "gradient flow"

Langevin: $dx_t = -\eta \nabla f(x_t) dt + \sqrt{2} dB_t$

↖ "Brownian motion"



This is called a stochastic differential equation

Thm: If f is (strongly) logconcave, then Langevin dynamics converge to p .

i.e. $d_{TV}(\text{law}(X_t), p) \rightarrow 0 \text{ as } t \rightarrow \infty.$

$$d_{TV}(a, \sigma) = \int |\pi(x) - \sigma(x)| dx.$$

One can also show that discretized Langevin converges.

Diffusion models: sampling via data.

given samples $x_1, \dots, x_n \sim p$, how can you generate fresh samples from p ?

Ornstein-Uhlenbeck (OU) process: Given a distribution $p = p_0$, the OU-process specifies a sequence of distributions p_t , where $p_t = \text{law}(X_t)$,

$$X_t = \underbrace{\exp(-t)}_{\alpha_t} X_0 + \underbrace{\sqrt{1 - \exp(-2t)}}_{\beta_t} Z_t, \quad X_0 \sim p_0, \quad Z_t \sim \mathcal{N}(0, I).$$

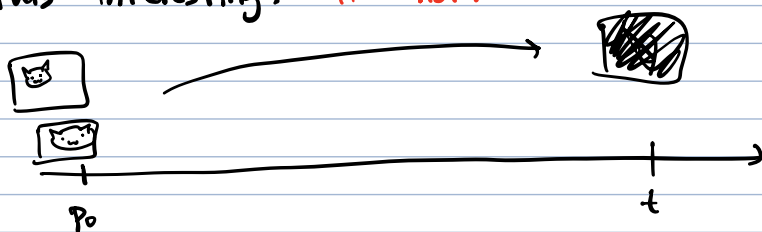
$$\alpha_t^2 + \beta_t^2 = 1.$$

The "right" way to think about OU is as an SDE:

$$dX_t = -X_t dt + \sqrt{2} dB_t.$$

Key point: OU takes a data distribution p_0 and turns it smoothly into noise.

Why is this interesting? *It's not!*



But: this process can be reversed

Fix some time T , define the reverse process $\bar{X}_t = X_{T-t}$ for $t \in [0, T]$.

Then this reverse process satisfies:

$$d\bar{X}_t = (\bar{X}_t + 2\nabla \log p_{T-t}(\bar{X}_t))dt + \sqrt{2} dB_t$$

this takes noise \rightarrow data. *This is very interesting!*

Two main problems:

1. This is still some continuous-time mumbo-jumbo.

Can discretize this SDE, just like how you can for Langevin.

discretize time $[0, T]$ into chunks of length Δt

each step is of the form:

$$\bar{X}_{t+\Delta t} = \alpha \bar{X}_t + \alpha (\beta \cdot \nabla \log p_{T-t}(\bar{X}_t) + \gamma Z_t),$$
$$Z_t \sim \mathcal{N}(0, I).$$

2. This depends on $\nabla \log p_{T-t} := s_{T-t}$, which we don't know.

Call this function the score function.

Score matching [Hyvärinen '05]

Suppose I want to find the best match to the score function

from some family of functions \mathcal{F} :

$$\arg \min_{s \in \mathcal{F}} \mathbb{E}_{X \sim p_t} [\|s(X) - \nabla \log p_t(X)\|^2]$$

the minimizer of this is the same as:

$$\arg \min_{s \in \mathcal{F}} \mathbb{E}_{\substack{X \sim p_0 \\ Z \sim \mathcal{N}(0, I)}} \left[\left\| s(X_t) + \frac{1}{\sqrt{1 - \exp(-2t)}} Z_t \right\|_2^2 \right], \quad X_t = \exp(-t) X_0 + \sqrt{1 - \exp(-2t)} Z_t$$

reparametrize: let $\hat{s}(X) = -s(X) \sqrt{1 - \exp(-2t)}$

$$\arg \min_{\hat{s}} \mathbb{E}_{\substack{x \sim p_0 \\ z \sim \mathcal{N}(0, I)}} \left[\|z - \hat{s}(x_t)\|_2^2 \right] \quad (*)$$

local denoising function: given noisy sample, predict what part of it is noise.

we can optimize this objective given samples from p_0 !

Given $x_1, \dots, x_n \sim p_0$, let $z_1, \dots, z_n \sim \mathcal{N}(0, I)$,

$$\text{form } y_i = \exp(-t)x_i + \sqrt{1 - \exp(-t)}z_i$$

$$\text{then } (*) \approx \frac{1}{n} \sum_{i=1}^n \|z_i - \hat{s}(y_i)\|_2^2$$

this is a regression problem!

In practice: take \hat{s} to be a large neural network.

→ 12k citations!

Denoising Diffusion Probabilistic Models (DDPM): [Ho-Abeel-Jain '20]

1. Learn score function from data.
2. plug it into discretized reverse process
3. ???
4. profit (literally, see e.g. stable diffusion).

The backbone of modern generative models!

Many variants now: DDIM, latent diffusion, consistency models, etc.

Thm [CCLSZ'23]: If the neural network learns the score effectively, then DDPM output is provably close to p_0 !