

Sampling (ch. 14 from Bishop's "Deep Learning")

Tuesday, November 12, 2024 2:18 PM

Last part of our class will focus on **generative models** and I plan to follow Christopher Bishop's "Deep Learning" book from 2024

- ch 14 Sampling (rejection, MCMC, Langevin) background
- ch 17 GANs
- ch 18 Normalizing Flows
- ch 19 (variational) Autoencoders (VAE)
- ch 20 diffusion models

Goal: Sample $z \sim p(z)$ ("distribution" / PDF...)

Ch 14: Background on sampling

Given: a source that samples $z \sim \text{Unif}(0,1)$ or $z \sim N(0,1)$

- Inverse transform sampling for simple distributions

(details: see §14.1.2 or [wikipedia](#))

Ex: you can sample $z \sim \text{Unif}(0,1)$]source

you want samples from $\text{Unif}(-1,1)$]target

Solution: $y = g(z), g(z) = 2z - 1$

motivation: ① p is a function we know and in fact our end goal is to find Z_p !
② $p(z|y) = \frac{p(y|z)p(z)}{p(y)}$ ← unknown

• Rejection Sampling

Want to sample from pdf $p(z)$

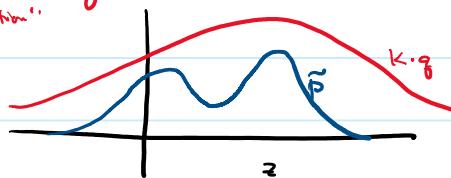
Assume we can sample from an unnormalized version \tilde{p}

i.e. $p(z) = \frac{1}{Z_p} \cdot \tilde{p}(z) \quad Z_p = \int \tilde{p} = \text{"partition function"}$

i.e. one we already know how to sample from
Use a simple "proposal distribution" g or "envelope function"

$$\tilde{p}(z) \leq K \cdot g(z)$$

comparison function



To generate a sample:

Until "ACCEPT"

1) draw $z \sim g(z)$

2) with probability $\frac{\tilde{p}(z)}{K \cdot g(z)}$, ACCEPT (else: REJECT)

i.e. draw $U_0 \in \text{Unif}[0, K \cdot g(z)]$ and accept if

$$U_0 \leq \tilde{p}(z)$$

Inefficient if it takes a while to accept... want $\tilde{p} \approx K \cdot g$.

Variants: adaptively refine g

Sampling (p. 2)

Tuesday, November 12, 2024 2:42 PM

- Importance Sampling For related problem of estimating

$$\mathbb{E}_{z \sim p}[f] := \int f(z) p(z) dz$$

Idea: compute $\int f(z) \frac{p(z)}{g(z)} \cdot g(z) dz =: \mathbb{E}_{z \sim g}[f \cdot \frac{p}{g}]$

Error depends on $\text{Var}[f \cdot \frac{p}{g}]$, so want $g \approx f \cdot p$ so that
(i.e., # of samples) it's almost constant
(i.e., low variance)

Variants: Sampling-importance-resampling (§14.1.6)

- Markov Chain Monte Carlo (MCMC) family of methods

History: Metropolis + Ulam 1949

→ Stanislaw Ulam, 1909-1984. Manhattan project

CU math dept. '61-'62, '65-'67, chair '67-'75

Target: $p(z)$

like rejection sampling, we have a **proposal distribution** g

but now it's conditioned on the current state, z^t ,

so our samples $z^{(1)}, z^{(2)}, z^{(3)}, \dots$ form a Markov Chain.

also, as before, assume we can sample from \tilde{p} , where $\tilde{p}(z) = \frac{1}{Z_p} \tilde{p}(z)$

- Metropolis MCMC (1953) pick a simple $g(\cdot | \cdot)$ and

for Metropolis, assume it's symmetric: $g(z|\tilde{z}) = g(\tilde{z}|z)$

Given z^t ,

- sample $z^* \sim g(z|z^t)$

- ACCEPT w/ probability $A(z^*, z^t) := \min(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^t)})$

If ACCEPTED, $z^{t+1} = z^*$

else, discard z^* and set $z^{t+1} = z^t$ (unlike pure rejection sampling which loops)

Under some conditions,

$$\lim_{t \rightarrow \infty} z^t = p(z)$$

in distribution

... but we don't have independent samples!
 $M \gg 1$
(in practice, take every M^{th} sample, pretend it's independent)

Does a given transition matrix T have a stationary/invariant/distribution? equilibrium

eigenvalue problem, $T_p = p$

Yes, via Perron-Frobenius thm (1907, 1912)

Sampling (p. 3)

Tuesday, November 12, 2024 5:30 PM

- Metropolis-Hastings MCMC ... if q isn't symmetric.

$$\text{ACCEPT w.p. } A_k(z^*, z^t) := \min\left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^t)} \cdot \frac{q_k(z^t | z^*)}{q_k(z^* | z^t)}\right)$$

- Gibbs Sampling MCMC 1984, special case of Metropolis-Hastings

see § 14.2.4 or wikipedia

- Ancestral Sampling MCMC 1990

see § 14.2.5

more efficient MCMC, two big directions

Hamiltonian MC (aka Hybrid MC)

Langevin Sampling
(metropolis-adjusted Langevin algo)

Background: Energy-based models

we're looking for a generative model, sometimes unconditional ($p(x)$),

other times conditional ($p(x|w)$)
eg. $w = Ax + v$ inverse problem

often write $p(x|w) = \frac{1}{Z(w)} e^{-E(x,w)}$

* energy function (e.g. neural net)

$$* \text{partition function} = \int e^{-E(x,w)} dw \quad \begin{cases} \text{difficult, and for} \\ \text{conditional models} \\ \text{you'd have to} \\ \text{recompute for every } w \end{cases}$$

for MLE, neg. likelihood (of $p(x|w)$),

plugging in observations $\{x_1, \dots, x_N\}$, is $\min_w \left(-\sum_{n=1}^N E(x_n, w) - N \cdot \ln(Z(w)) \right)$

Sometimes we could ignore $-N \cdot \log Z$ term

ex: $p(x|w) \sim N(w, \sigma^2 I)$ The mean w doesn't affect Z so we can treat it as a constant

other times we can't

ex: $p(x|w) \sim N(0, w^2 I)$, pdf is

$$\frac{1}{w \sqrt{2\pi}} e^{-\frac{1}{2} \frac{x^2}{w^2}}$$

* partition function depends on w

we'll talk later about **Score matching** for diffusion models

for now, MCMC perspective for both:

1) training (so find $\nabla_w (-\log(\text{likelihood}))$)

2) inference (once w trained, generate sample $x \sim p(x|w)$)

Sampling (p. 4)

Tuesday, November 12, 2024 8:02 PM

$$\nabla_w \ln p(x|w) = -\nabla_w E(x, w) - \nabla_w \ln Z(w) \quad \text{for a single point}$$

or for a dataset $D = \{x_1, \dots, x_N\}$

$$\nabla_w \mathbb{E}_{x \sim p_D} \ln p(x|w) = -\mathbb{E}_{x \sim p_D} \nabla_w E(x, w) - \nabla_w \cdot \ln Z(w)$$

fancy way of saying

$$\mathbb{E}_{x \sim p_D} f(x) := \frac{1}{N} \sum_{n=1}^N f(x_n), \quad \text{"empirical" law/distribution}$$

$$-\nabla_w \ln Z(w) = -\frac{1}{Z(w)} \cdot \nabla_w Z(w) \quad \text{chain rule}$$

$$= -\frac{1}{Z(w)} \cdot \nabla_w \int e^{-E(x, w)} dx \quad \text{def'n of } Z$$

$$= -\frac{1}{Z(w)} \int \nabla_w e^{-E(x, w)} dx \quad \begin{matrix} \text{wished thinking} \\ (\text{or expand via DCT}) \\ \text{if } e^{-x} \text{ converges...} \\ \text{which it should} \end{matrix}$$

$$= -\frac{1}{Z(w)} \int e^{-E(x, w)} \cdot \nabla_w (-E(x, w)) dx \quad \text{chain rule again}$$

$$= \int p(x|w) \cdot \nabla_w E(x, w) dx \quad p(x|w) := \frac{1}{Z(w)} e^{-E(x, w)}$$

$$= \mathbb{E}_{x \sim p(x|w)} \nabla_w E(x, w)$$

So ...

$$\nabla_w \mathbb{E}_{x \sim p_D} \ln p(x|w) = -\mathbb{E}_{x \sim p_D} \nabla_w E(x, w) + \mathbb{E}_{x \sim p(x|w)} \nabla_w E(x, w)$$

* easy: discrete

how to sample?
Langevin Sampling

uses score function

$$S(x, w) := \nabla_x \ln p(x|w)$$

not ∇_w

$$= -\nabla_x E(x, w) \quad \leftarrow Z \text{ goes away since } \nabla_x Z(w) = 0$$

Algo: Langevin Dynamics

$$x^{t+1} \leftarrow x^t + \eta \underbrace{\nabla_x p(x^t, w)}_{= S(x^t, w)} + \sqrt{2\eta} \cdot \varepsilon^t \quad \varepsilon^t \sim N(0, I)$$

as $\eta \rightarrow 0$, and $t \rightarrow \infty$, $x^t \sim p(x)$ } in our case, it's really $x^t \sim p(x|w)$

So, do this to draw several (M) samples $x_m \sim p(x|w)$

$$\text{Then approximate (*) via } \mathbb{E}_{x \sim p(x|w)} \nabla_w E(x, w) \approx \frac{1}{M} \sum_{m=1}^M \nabla_w E(x_m, w)$$