

Key points of this chapter are that:

- ① Knowing the motivation and concepts.
- ② Clarification of terms (Sampling, monte carlo, Gibbs Sampling, MCMC)
- ③ Sampling methods analysis.

Sampling and Monte Carlo Methods

- why sampling? { Calculating sum/integral efficiently (considering them as expectation)
We want to generate samples from training distribution.

- Basic idea of Monte Carlo Methods

Aim: we want to calculate sum/integral:

$$S = \sum_x p(x) f(x) = \mathbb{E}_p[f(x)] \quad \text{or} \quad S = \int_x p(x) f(x) = \mathbb{E}_p[f(x)]$$

we draw samples from $p(x)$

$$\hat{S}_n = \frac{1}{n} \sum_i f(x_i) \rightarrow \mathbb{E}[\hat{S}_n] = \frac{1}{n} \sum_i \mathbb{E}[f(x_i)] = S$$

if $n \rightarrow \infty$, $\hat{S}_n = S$. unbiased estimation

About the variance: $\text{Var}(S_n) = \frac{\text{Var}[f(x)]}{n}$ can be seen as an indicator of confidence

Importance Sampling

Motivation: Always, we cannot directly sampling from $p(x)$. we re-write:

$$p(x) f(x) = q(x) \frac{p(x) f(x)}{q(x)}$$

sampling from $q(x)$, and finally adding with importance weights

$$\mathbb{E}_p[f(x)] = \mathbb{E}_q\left(\frac{p}{q} f(x)\right)$$

$$\frac{p(x)}{q(x)}$$

About the variance: $\text{Var}(\hat{S}_q) = \text{Var}\left[\frac{p f}{q}\right] / n$, when $q^* = \frac{p f}{Z}$, we got minimum variance, which means that q is sufficiently close to p .

Bias importance sampling:

$$\hat{S}_{BIS} = \frac{\sum_{i=1}^n \frac{\tilde{p}(x_i)}{\tilde{q}(x_i)} f(x_i)}{\sum_{i=1}^n \frac{\tilde{p}(x_i)}{\tilde{q}(x_i)}}$$

Avoid to calculate partition function Z .

- Deep Learning use sampling: To compute some intractable calculations, making them a representation of expectation first and then do sampling.

Markov Chain Monte Carlo Methods

- Unlike directed graph, undirected graph cannot find an explicit dependency relation, we cannot do ancestral sampling.
- MCMC methods build a Markov chain to decide q , according to x^{t+1} :

we sampling $q(x^t) = \sum_x T(x'|x) \cdot q(x^t)$ \rightarrow 1-step probability of $q(x)$
 $q(x^t)$ at timestep t . \rightarrow Transition / Stochastic matrix

- Stationary Distribution of Markov Chain.

$$v^t = A v^{t+1} \rightarrow \text{Consists of } T(x'|x)$$

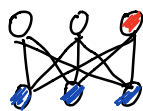
$v^t = A^t v^0$, then if $v = Av$, the distribution is stationary
 A has eigenvalue 1, and v is eigenvector

- MCMC operation: First burn into a stationary distribution, then sampling according to $q^t = A q^{t+1}$, sometimes fetching samples with n interval to ensure non-correlation between samples.

Gibbs Sampling

$q(x)$ is quite important in MCMC and importance sampling, how we decide $q(x)$ is related to how we choose $T(x'|x) / A$. means what kind of transition matrix is appropriate. two ways proposed $\left\{ \begin{array}{l} \text{Gibbs sampling} \\ \text{Directly parametrize } T. \end{array} \right.$

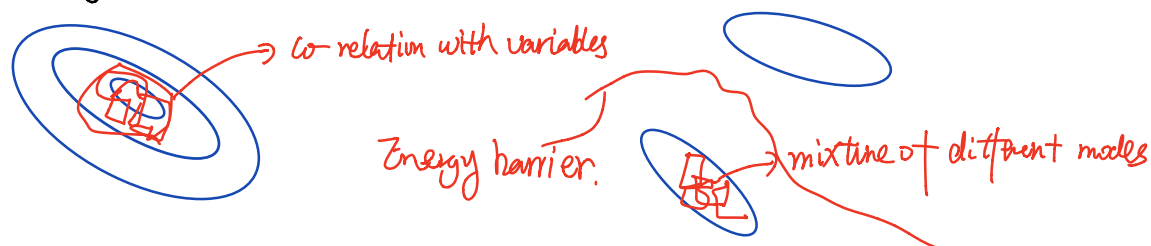
- Gibbs Sampling for $T(x'|x)$ means that $T(x'|x)$ with only one variable x_i is different in x' from x . such as RBM:



$T(x'|x)$ can be seen as shaded part \rightarrow sampling from $|||$, and we get $T(x'|x)$, if $||| = 1, 2, 3$, then $|||$ can be \dots .

The challenge of mixing between Separated Modes (models and Variables)

- MCMC: A significant pattern, we sample according to last time step, it will yield we are trapped within one mode.



- How we solve this issue:

- Tempering: $p(x) \propto \exp(-\beta \mathcal{E}(x))$ → Control to not be "so" peak!

- Depth may help: $\begin{matrix} \bigcirc & \bigcirc & \bigcirc & \rightarrow & x \\ \text{|||||} & & & \rightarrow & \text{Reeper structure} \\ \bigcirc & \bigcirc & \bigcirc & \rightarrow & h \end{matrix}$

observed → $p(x|h)$
latent → $p(h|x)$

means sample within huge "gap":