

Homework 3

Deep Learning 2024 Spring

Due on 2024/4/13

1 True or False

Problem 1. Stochastic Gradient MCMC is designed to solve the optimization problem $\arg \max_{\theta} \mathbb{P}(\theta|\mathbf{X})$, where θ is the collection of parameters and \mathbf{X} represents data.

Problem 2. The convolution used in GLOW and WaveNet should be both invertible (i.e., the convolution kernels W should be invertible) since they are all flow models.

2 Q&A

Problem 3. (Importance Sampling) \mathbf{x} is a random variable. Given target distribution $p(\mathbf{x})$ and target random variable $y = f(\mathbf{x})$, importance sampling gives an estimator of $\mathbb{E}[y]$ from a proposal distribution $q(\mathbf{x})$:

$$\mathbb{E}_{\mathbf{x} \sim p} [f(\mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim q} \left[\frac{p(\mathbf{x})}{q(\mathbf{x})} f(\mathbf{x}) \right] \approx \frac{1}{N} \sum_{\mathbf{x} \sim q(\cdot)} \frac{p(\mathbf{x})}{q(\mathbf{x})} f(\mathbf{x}).$$

Prove that when q has the following form,

$$q^*(\mathbf{x}) \propto p(\mathbf{x}) |f(\mathbf{x})|$$

the variance of this estimator can be minimized.

Problem 4. (Markov Chain Monte Carlo)

1. Prove random-walk Metropolis-Hasting sampling (i.e., $\mathbf{s}' \leftarrow \mathbf{s} + \text{Gaussian noise}$) is a valid MCMC algorithm, i.e., it constructs a Markov chain which is ergodic and satisfies the detailed balance property.
2. Prove that Gibbs sampling is a special case of Metropolis-Hasting sampling, and that the acceptance rate of Gibbs sampling (i.e., $\alpha(\mathbf{s} \rightarrow \mathbf{s}')$) is 1.

Here we consider the following 2-step Gibbs proposal: (1) randomly sample a coordinate index i ; (2) sample coordinate \mathbf{s}_i from the coordinate proposal $q(\mathbf{s}_i \rightarrow \mathbf{s}'_i) = p(\mathbf{s}'_i | \mathbf{s}_{j \neq i})$.

3. **(Bonus Question)** In fact, Gibbs sampling is typically implemented in a *cyclic fashion*, i.e., running posterior sampling in a fixed order over all the dimensions. Prove that cyclic Gibbs sampling yields the same stationary distribution as random-order Gibbs sampling in the above question, as long as the Markov chain can access all states under the fixed ordering.

Problem 5. (Directed Probabilistic Model)

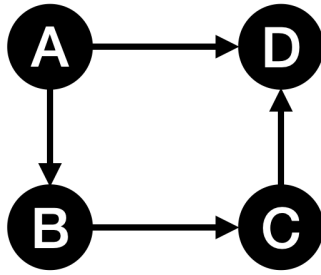


Figure 1: An example of directed graphical model.

Continuing Problem 6 in the previous homework, we consider **directed probabilistic model** in this problem. In a directed graphical model, an edge (arrow) implies dependency of the downstream random variable to the upstream one. Figure 1 shows an example, where the joint probability distribution can be factorized as

$$\mathbb{P}(A, B, C, D) = \mathbb{P}(A)\mathbb{P}(B|A)\mathbb{P}(C|B)\mathbb{P}(D|A, C).$$

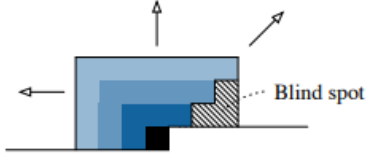
1. Judge whether the following statements regarding independence are true or not. $X \perp Y | Z$ denotes that X and Y are independent given Z .
 - (a) $A \perp C$
 - (b) $A \perp C | B$
 - (c) $A \perp C | D$
 - (d) $A \perp C | B, D$
 - (e) $B \perp D$
 - (f) $B \perp D | A$
 - (g) $B \perp D | C$
 - (h) $B \perp D | A, C$

2. Suppose the conditional distributions are given by

$$\mathbb{P}(B|A = a) = \mathcal{N}(a, 1)$$

$$\mathbb{P}(C|B = b) = \mathcal{N}(b, 1)$$

$$\mathbb{P}(D|A = a, C = c) = \mathcal{N}(c + a, 1).$$



(a) A sawtooth-shaped receptive field.

1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

(b) A PixelCNN causal mask.

The prior distribution of A is a Gaussian distribution with zero mean and unit variance. Derive the likelihood $\mathbb{P}(B, C, D|A)$ and posterior $\mathbb{P}(A|B, C, D)$ in closed form.

Problem 6. (PixelCNN) PixelCNN¹ is an auto-regressive generative model based on masked convolution kernels. Please answer the following questions:

1. Show that masked convolution kernels induce sawtooth-shaped receptive fields, as shown in Figure 2a.
2. An obvious flaw of PixelCNN is that the generated pixel can not condition on all the left and upper pixels due to the sawtooth-shaped receptive field, which is called the issue of *blind spot*. Denote the unmasked kernel as w , the corresponding mask as m (as shown in Figure 2b), and the image/feature map as x . The computation of a PixelCNN layer can be represented as

$$y = \text{conv}(m \cdot w, x)$$

Propose a minimal modification on PixelCNN to remove the blind spot while maintaining the auto-regressive property. Write down your per-layer computation process.

¹<https://arxiv.org/abs/1601.06759>