

Lecture 16, Apr 9, 2024

Stochastic Variational Inference (SVI)

- *Stochastic variational inference* is the technique of approximating the true conditional $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}$ by a simpler distribution, $q(\mathbf{z}|\boldsymbol{\theta})$
 - We want $q(\mathbf{z}|\boldsymbol{\theta})$ to be “close to” $p(\mathbf{z}|\mathbf{x})$; to do this we need to define “closeness” of distributions
 - We can choose $q(\mathbf{z}|\boldsymbol{\theta})$ to come from a known family of distributions, e.g. Gaussians

Definition

The *Kullback-Leibler (KL) divergence* of two distributions $p(\mathbf{z})$ and $q(\mathbf{z})$ is

$$KL(q \parallel p) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\log \frac{q(\mathbf{z})}{p(\mathbf{z})} \right]$$

with the following properties:

- $KL(q \parallel p) \geq 0$
 - $KL(q \parallel p) = 0 \iff q = p$
 - $KL(q \parallel p) \neq KL(p \parallel q)$
- KL divergence is always positive and zero when distributions are equal, however it is not symmetric!
 - For *reverse-KL* (aka *information projection*), we take $KL(q \parallel p)$, which penalizes q having mass where p has none
 - * When p is large where q is small, the KL divergence is small
 - * When p is small where q is large, the KL divergence is large
 - * This will compress q so it fits to one of the peaks of p
 - For *forward-KL* (aka *moment projection*), we take $KL(p \parallel q)$, which penalizes q missing mass where p has some
 - * When p is large where q is small, the KL divergence is large
 - * When p is small where q is large, the KL divergence is small
 - * This will stretch out q to cover all the peaks of p
 - The choice of which KL divergence to optimize leads to different fits
 - * In practice however we normally use reverse KL for computational reasons

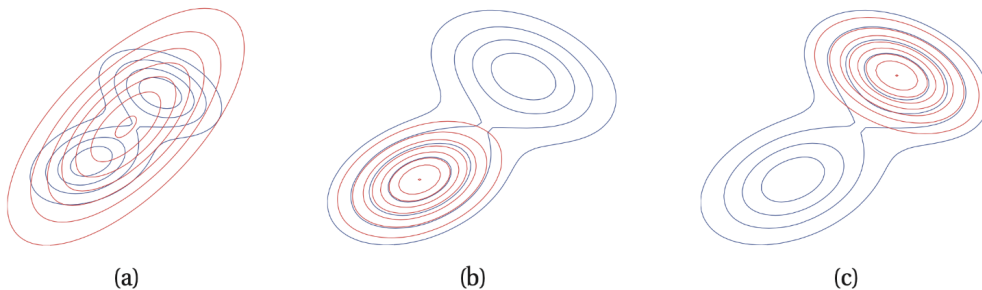


Figure 1: Approximating a bimodal distribution by a unimodal distribution; (a) minimizes forward KL, (b) and (c) minimize reverse KL.

- SVI tries to minimize the KL divergence of p and q

- $KL(q(\mathbf{z}|\boldsymbol{\theta}) \parallel p(\mathbf{z}|\mathbf{x})) = \mathbb{E}_{\mathbf{z} \sim q} \left[\log \frac{q(\mathbf{z}|\boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x})} \right]$

$$= \mathbb{E}_{\mathbf{z} \sim q} \left[\log \left(q(\mathbf{z}|\boldsymbol{\theta}) \frac{p(\mathbf{x})}{p(\mathbf{z}, \mathbf{x})} \right) \right]$$

$$= \mathbb{E}_{\mathbf{z} \sim q} \left[\log \left(\frac{q(\mathbf{z}|\boldsymbol{\theta})}{p(\mathbf{z}, \mathbf{x})} \right) \right] + \log p(\mathbf{x})$$

$$= -\mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) + \log p(\mathbf{x})$$
- $\mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) = -\mathbb{E}_{\mathbf{z} \sim q} \left[\log \frac{q(\mathbf{z}|\boldsymbol{\theta})}{p(\mathbf{z}, \mathbf{x})} \right]$ is the *evidence lower bound* (ELBO)
- Since $-\mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) + \log p(\mathbf{x}) \geq 0$ (since KL is positive), the ELBO is a lower bound for $\log p(\mathbf{x})$
- As $\log p(\mathbf{x})$ is constant, to minimize the KL divergence we have to maximize the ELBO; therefore we do not have to compute the normalization, which is infeasible to do
- The ELBO gradient is $\vec{\nabla}_{\boldsymbol{\theta}} = \vec{\nabla}_{\boldsymbol{\theta}} \int q(\mathbf{z}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} d\mathbf{z}$, which must be estimated since we cannot compute this high-dimension integral
 - The *score function* (aka *REINFORCE*) gradient estimator
 - * $\vec{\nabla}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q} \left[\vec{\nabla}_{\boldsymbol{\theta}} \log q(\mathbf{z}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} \right]$
 - * Using Monte Carlo, $\vec{\nabla}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) \approx \frac{1}{B} \sum_{i=1}^B \vec{\nabla}_{\boldsymbol{\theta}} \log q(\mathbf{z}^{(i)}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}, \mathbf{z}^{(i)})}{q(\mathbf{z}^{(i)}|\boldsymbol{\theta})}$
 - B is the number of samples
 - This is an unbiased estimator and easy to compute
 - * In practice, this has higher variance than the pathwise gradient estimator
 - * Use in specific domains such as reinforcement learning
 - The *pathwise* (aka *reparametrization*) gradient estimator factors out all the randomness of the distribution into a parameterless fixed source of noise, $p(\boldsymbol{\varepsilon})$
 - * Find $T(\boldsymbol{\varepsilon}, \boldsymbol{\theta})$ such that for $\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})$, then $\mathbf{z} = T(\boldsymbol{\varepsilon}, \boldsymbol{\theta}) \implies \mathbf{z} \sim q(\mathbf{z}|\boldsymbol{\theta})$
 - e.g. for a Gaussian, $\boldsymbol{\theta} = \{\mu, \sigma\}$, let $\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\varepsilon}|0, 1)$ and $T(\boldsymbol{\varepsilon}, \boldsymbol{\theta}) = \sigma\boldsymbol{\varepsilon} + \mu$, then $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\mu, \sigma)$
 - * Using the above, $\vec{\nabla}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) = \mathbb{E}_{\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})} \left[\vec{\nabla}_{\boldsymbol{\theta}} \log \frac{p(\mathbf{x}, T(\boldsymbol{\varepsilon}, \boldsymbol{\theta}))}{q(T(\boldsymbol{\varepsilon}, \boldsymbol{\theta})|\boldsymbol{\theta})} \right]$
 - * This can then be estimated using Monte Carlo
- The main drawback of SVI is the challenge of determining how good the approximation is after the optimization terminates

Monte Carlo and Importance Sampling

- So far we've examined methods of estimating the full distribution $p(\mathbf{x})$, but sometimes we're only interested in the expectation of some function $\phi(\mathbf{x})$ under the distribution, i.e. $I = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\phi(\mathbf{x})]$
- The Monte Carlo approximation is given by $I = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\phi(\mathbf{x})] \approx \hat{I} = \frac{1}{R} \sum_{i=1}^R \phi(\mathbf{x}^{(i)})$
 - This is unbiased, with a standard deviation proportional to $\frac{1}{\sqrt{R}}$, independent of the dimension of \mathbf{x}
- If we only need the expectation, we only need to be able to sample from the distribution, and apply Monte Carlo to find the expectation
 - However, sampling is hard because we typically only have the unnormalized distribution, $\tilde{p}(\mathbf{x}) = Zp(\mathbf{x})$; even if we did have the full distribution, sampling from a high-dimension distribution is hard
- *Importance sampling* is a method for approximating the expectation when we only have the unnormalized distribution
 - A notable example is the particle filter for state estimation in robotics
- Let $q(\mathbf{x})$ be the *sampler density*, a simpler density function that we can easily sample from

$$\begin{aligned}
- I &= \int p(\mathbf{x})\phi(\mathbf{x}) \, d\mathbf{x} \\
&= \int \phi(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \, d\mathbf{x} \\
&= \frac{\int \frac{\phi(\mathbf{x})p(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \, d\mathbf{x}}{\int \frac{p(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \, d\mathbf{x}} \\
&= \frac{\int \frac{\phi(\mathbf{x})\frac{1}{2}\tilde{p}(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \, d\mathbf{x}}{\int \frac{\frac{1}{2}\tilde{p}(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \, d\mathbf{x}} \\
&= \frac{\int \frac{\phi(\mathbf{x})p(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \, d\mathbf{x}}{\int \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \, d\mathbf{x}} \\
&= \frac{\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\frac{\phi(\mathbf{x})p(\mathbf{x})}{q(\mathbf{x})} \right]}{\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \right]}
\end{aligned}$$

- Now we can use Monte Carlo to approximate the expectations

$$- \hat{I} = \frac{\frac{1}{R} \sum_{r=1}^R \frac{\phi(\mathbf{x}^{(r)})\tilde{p}(\mathbf{x}^{(r)})}{q(\mathbf{x}^{(r)})}}{\frac{1}{R} \sum_{r=1}^R \frac{\tilde{p}(\mathbf{x}^{(r)})}{q(\mathbf{x}^{(r)})}} = \frac{\sum_r w_r \phi(\mathbf{x}^{(r)})}{\sum_r w_r}$$

- Each $w_r = \frac{\tilde{p}(\mathbf{x}^{(r)})}{q(\mathbf{x}^{(r)})}$ is referred to as the *importance weight*

- * Intuitively, if at a point $p(\mathbf{x}^{(r)}) > q(\mathbf{x}^{(r)})$, then sampling from q will under-represent this point; therefore the points are weighted more in the sum, since w_r will be larger
 - * Conversely $p(\mathbf{x}^{(r)}) < q(\mathbf{x}^{(r)})$ means q over-represents the point, so in this case w_r will be small and less weight is applied to it
 - * When $p(\mathbf{x}^{(r)}) = q(\mathbf{x}^{(r)})$ we can show that \hat{I} applies no reweighing to samples
- The sampler density should have heavy tails (e.g. a Cauchy distribution instead of a Gaussian), since we need to compensate for the difference between distribution
 - If the sampler is chosen improperly, the variance of the result can be extremely high
 - In high dimensions, if the sampler distribution is not a near-perfect approximation of the target, then the entire sum will likely be dominated by a few samples with a huge weight, leading to a very bad estimate