Probabilistic Artificial Intelligence

Problem Set 3

Oct 21, 2019

## 1. Markov chains and detailed balance

Assume that you are given a Markov chain with finite state space $\Omega$ and transition matrix $T$, which is defined for all $x, y \in \Omega$ and $t \geq 0$ as $T(x, y) := P(X_{t+1} = y \mid X_t = x)$. Furthermore, let $\pi$ be the stationary distribution of the chain.

(i) Show that, if for some $t$ the current state $X_t$ is distributed according to the stationary distribution and additionally the chain satisfies the detailed balance equations

$$\pi(x)T(x, y) = \pi(y)T(y, x), \text{ for all } x, y \in \Omega,$$

then the following holds for all $k \geq 0$ and $x_0, \ldots, x_k \in \Omega$:

$$P(X_t = x_0, \ldots, X_{t+k} = x_k) = P(X_t = x_k, \ldots, X_{t+k} = x_0).$$

(This is why a chain that satisfies detailed balance is called *reversible*.)

(ii) Show that, if $T$ is a symmetric matrix, then the chain satisfies detailed balance, and the uniform distribution on $\Omega$ is stationary for that chain.

## 2. Convergence of the Metropolis-Hastings algorithm

We use Markov Chain Monte Carlo (MCMC) methods to sample from a target distribution $P(x) = Q(x)/Z$ using a proposal distribution $R(x|x')$, without computing the normalization constant $Z$. A famous MCMC method is the METROPOLIS-HASTINGS algorithm, given below:

---
**Algorithm 1** METROPOLIS-HASTINGS

---
**Input:** Unnormalized target distribution $Q(x)$, proposal distribution $R(x|x')$
**Initialize:** $x_1$ arbitrary
**For** $t = 1, 2, \ldots, T$:

1. Sample proposal $x$ from the proposal distribution $R(x|x_t)$.
2. Compute the *acceptance probability* $\alpha = \min\left(1, \frac{Q(x)R(x_t|x)}{Q(x_t)R(x|x_t)}\right)$.
3. With probability $\alpha$, set $x_{t+1} = x$, else $x_{t+1} = x_t$.

---

The algorithm defines a Markov chain with transition kernel $T(x, x') = P(x_{t+1} = x'|x_t = x)$. In this exercise, we prove that the stationary distribution of this Markov chain is equal to the target distribution $P(x)$. *Remark:* While we show, that Metropolis-Hastings converges to the correct distribution, the proof doesn't tell us how fast it converges. In practice, we typically use samples only after a 'burn-in' period, which allows the chain to converge.

(i) Show that if a unnormalized distribution $Q$ on $\Omega$ satisfies the detailed balance equations,

$$Q(x)T(x,y) = Q(y)T(y,x), \text{ for all } x, y \in \Omega, \tag{1}$$

then $\pi(x) = \frac{1}{Z}Q(x)$ is the stationary distribution of the Markov chain defined by the transition kernel $T(x, x')$.

(ii) Show that if METROPOLIS-HASTINGS transitions to a new state, i.e. $x_{t+1} \neq x_t$, then the transition probability $T(x_t, x_{t+1})$ can be written as

$$T(x_t, x_{t+1}) = \frac{1}{Q(x_t)} \min \big( Q(x_t)R(x_{t+1}|x_t), Q(x_{t+1})R(x_t|x_{t+1}) \big). \tag{2}$$

Use this to show that the detailed balance equation for $Q$ is satisfied if $x_{t+1} \neq x_t$.

(iii) Finally, show that if $x_t = x_{t+1}$, the detailed balance condition is trivially satisfied. *Remark:* You can still compute the transition probability $T(x_t, x_{t+1})$ for this case, but the result follows independent of the exact transition probability.

## 3. Forward and Reverse KL for Variational Inference

In variational inference we seek to find a distribution $Q$ in a class of distributions $\mathcal{Q}$, that minimizes the KL-distance to a target distribution $P$, i.e.

$$Q \in \arg\min_{Q \in \mathcal{Q}} \text{KL}(Q\|P). \tag{3}$$

The KL-distance (for finite support) is defined as $\text{KL}(Q\|P) = \sum_x Q(x)\log(Q(x)/P(x))$. The KL-distance is not symmetric, so in general, $\text{KL}(Q\|P) \neq \text{KL}(P\|Q)$. We refer to $\text{KL}(P\|Q)$ as *forward KL* and $\text{KL}(Q\|P)$ as *reverse KL*. It is possible to use both, forward and reverse KL for variational inference, but the approximations will be different. In the plots below, we fit the true distribution $P(x)$ with a Gaussian $Q(x)$, using either forward or reverse KL. Explain which KL was used in either case!