

Monte Carlo Methods

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics

1. Sampling and Monte Carlo Methods
2. Importance Sampling
3. Markov Chain Monte Carlo Methods
4. Gibbs Sampling
5. Mixing between separated modes

What are Monte Carlo Algorithms?

- We cannot obtain precise answers to many problems in machine learning
 - We can use either deterministic approximate algorithms or Monte Carlo approximations
 - Both approaches are common in machine learning
- Monte Carlo algorithms return answers with random amount of error
 - For a fixed computational budget, a Monte Carlo algorithm can provide an approximate answer
 - Las Vegas algorithms return exact answers but use random amount of resources

1. Why Sampling?

- Many ML algorithms are based on drawing samples from some probability distribution
 - and using these samples to form a Monte Carlo estimate of some desired quantity
- Sampling provides a flexible way to approximate many sums and integrals at reduced cost
 - Sometimes for speedup of a costly but tractable sum, e.g., subsample training cost with minibatches
 - In other cases, learning algorithms require us to approximate an intractable sum or integral
 - E.g., gradient of the log partition function of an undirected model

Basics of MC sampling

- When a sum or integral is intractable
 - E.g., has exponential no of terms and no exact simplification is known
 - It can often be approximated using MC sampling
- The idea is to view the sum or integral as an expectation under some distribution and to approximate the expectation by a corresponding average

Summation \rightarrow Expectation \rightarrow Average

- Sum or integral to approximate is

$$s = \sum_{\mathbf{x}} p(\mathbf{x})f(\mathbf{x}) \quad \text{or} \quad s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$$

- Rewriting expression as an expectation

$$s = E_p[f(\mathbf{x})] \quad \text{or} \quad s = E_p[f(\mathbf{x})]$$

– with the constraint that p is a probability distribution (for the sum) or a pdf (for the integral)

- Approximate s by drawing n samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ from p and then forming the empirical average

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)})$$

Justification of Approximation

- The sample average approximation is justified by a few different properties

1. The estimator \hat{s} is unbiased, since

$$\begin{aligned}
 E[\hat{s}_n] &= \frac{1}{n} \sum_{i=1}^n E[f(\mathbf{x}^{(i)})] && \text{since } \hat{s}_n \text{ is the sample average } \hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)}) \\
 &= \frac{1}{n} \sum_{i=1}^n s && \text{since } s = E_p[f(\mathbf{x})] \\
 &= s
 \end{aligned}$$

2. The law of large numbers states that if the samples $\mathbf{x}^{(i)}$ are i.i.d. then the average converges almost surely to the expected value

$$\lim_{n \rightarrow \infty} \hat{s}_n = s$$

Provided the variance of the individual terms $\text{Var}[f(\mathbf{x}^{(i)})]$ is bounded

Variance of estimate

- Consider variance of \hat{s}_n as n increases
 - $\text{Var}[\hat{s}_n]$ converges & decreases to 0 if $\text{Var}[f(\mathbf{x}^{(i)})] < \infty$:

$$\begin{aligned} \text{Var}[\hat{s}_n] &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}[f(\mathbf{x})] \\ &= \frac{\text{Var}[f(\mathbf{x})]}{n} \end{aligned}$$

- Result tells how to estimate error in MC average
 - Equivalently expected error of the approximation
 - Compute empirical average of the $f(\mathbf{x}^{(i)})$ and their empirical variance to determine estimator of $\text{Var}[\hat{s}_n]$
 - Central Limit Theorem tells us that distribution of the average \hat{s}_n has a normal distribution with mean s and variance $\text{Var}[\hat{s}_n]$. This allows us to estimate confidence intervals around the estimate \hat{s}_n using the Normal cdf

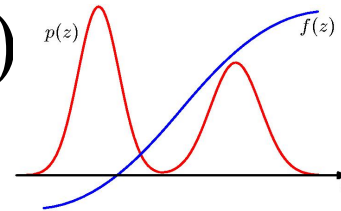
Sampling from base distribution $p(\boldsymbol{x})$

- Sampling from $p(\boldsymbol{x})$ is not always possible
- In such a case use importance sampling
- A more general approach is Monte Carlo Markov chains
 - To form a sequence of estimators that converge towards the distribution of interest

2. Importance Sampling

- Principal reason for sampling $p(\mathbf{x})$ is evaluating expectation of some $f(\mathbf{x})$

$$E[f] = \int f(\mathbf{x})p(\mathbf{x}) d\mathbf{x}$$



- Given samples $\mathbf{x}^{(i)}$, $i=1, \dots, n$, from $p(\mathbf{x})$, the finite sum approximation is

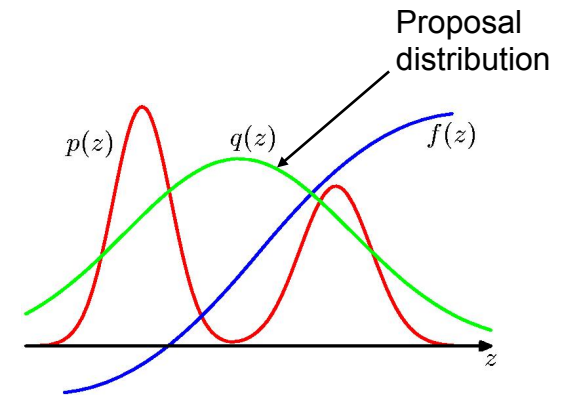
$$\hat{f} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)})$$

- But drawing samples $p(\mathbf{x})$ may be impractical
- Importance sampling uses:
 - a proposal distribution— like rejection sampling
 - But all samples are retained
 - Assumes that for any \mathbf{x} , $p(\mathbf{x})$ can be evaluated

Determining Importance weights

- Samples $\{\mathbf{x}^{(i)}\}$ are drawn from simpler dist. $p(\mathbf{x})$

$$\begin{aligned} E[f] &= \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \int f(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) d\mathbf{x} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)}) \end{aligned}$$



- Samples are weighted by ratios

$$r_i = p(\mathbf{x}^{(i)}) / q(\mathbf{x}^{(i)})$$

– Known as importance weights

- Which corrects the bias introduced by wrong distribution

Importance Sampling: Choice of $p(\mathbf{x})$

- More generally, wish to evaluate the expression using sampling:

$$s = \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x})$$

$$s = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

- Which part of the integrand has the role of the probability $p(\mathbf{x})$
 - [from which we sample $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$]
- And which part has role of $f(\mathbf{x})$ whose expected value (under the probability distribution) is estimated as

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)})$$

Decomposition of Integrand

- In the equation $s = \sum_x p(\mathbf{x})f(\mathbf{x})$ $s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$
- There is no unique decomposition because $p(\mathbf{x})f(\mathbf{x})$ can be rewritten as

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

– where we now sample from q and average

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

- In many cases the problem $s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$ is specified naturally as expectation of $f(\mathbf{x})$ given distribution $p(\mathbf{x})$
 - But it may not be optimal in no of samples required

Derivation of Optimal $q^*(\mathbf{x})$

- Any Monte Carlo estimator

$$\hat{s}_p = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim p}^n f(\mathbf{x}^{(i)})$$

- can be transformed into an importance sampling estimator

$$\hat{s}_q = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim q}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}$$

using

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

- It can be readily seen that the expected value of the estimator does not depend on q : $E_q[\hat{s}_q] = E_q[\hat{s}_p] = s$
- The variance of an importance sampling estimator is sensitive to the choice of q . The variance is

- The minimum variance occurs when q is

$$Var[\hat{s}_p] = Var\left[\frac{1}{n} \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]$$

$$q^*(\mathbf{x}) = \frac{p(\mathbf{x}) |f(\mathbf{x})|}{Z}$$

- where Z is the normalization constant chosen so that $q^*(\mathbf{x})$ sums or integrates to one

Choice of suboptimal $q(\mathbf{x})$

- Any choice of q is valid and q^* is the optimal one (yields minimum variance)
- Sampling from q^* is usually infeasible
- Other choices of q can be feasible, reducing variance somewhat

Biased Importance Sampling (BIS)

- Another approach is BIS
 - Has advantage of not requiring normalized p or q .
- For discrete variables, BIS estimator is

$$\begin{aligned}\hat{s}_{BIS} &= \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}} \\ &= \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \\ &= \frac{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}},\end{aligned}$$

where \hat{p} and \hat{q} are unnormalized forms of p and q and $\mathbf{x}^{(i)}$ are samples from q

Effect of choice of q

- Good $q \rightarrow$ efficient Monte Carlo estimation
- Poor choice of $q \rightarrow$ efficiency much worse
 - Looking at $Var[\hat{s}_p] = Var\left[\frac{1}{n} \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]$ if there are samples for which $\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}$ is large then the variance of the estimator can get large
 - Happens when $q(\mathbf{x})$ is tiny while neither $p(\mathbf{x})$ or $q(\mathbf{x})$ is small enough to cancel it
 - The q distribution is usually chosen to be a very simple distribution so that it is easy to sample from
 - When \mathbf{x} is high dimensional this simplicity causes it to match p of $p|f|$ poorly
 - Very small or very large ratios are possible when \mathbf{x} is high dimensional

Importance sampling in Deep learning

- Used in many ML algorithms incl. deep learning
- Examples
 - To accelerate training in neural language models with large vocabulary
 - Other neural nets with large no of outputs
 - Estimate partition function (normalize prob. distribution)
 - Estimate log-likelihood in deep directed models such as variational autoencoder
 - Estimate gradient in SGD where most of the cost comes from a small no of misclassified samples
 - Sampling more difficult examples can reduce variance of gradient

3. Markov Chain Monte Carlo Methods

- In many cases we wish to use a Monte Carlo technique but there is no tractable method for drawing exact samples from $p_{\text{model}}(\mathbf{x})$ or from a good (low variance) importance sampling distribution $q(\mathbf{x})$
- In deep learning this happens most often when $p_{\text{model}}(\mathbf{x})$ is represented by an undirected model
- In this case we use a mathematical tool called a Markov chain to sample from $p_{\text{model}}(\mathbf{x})$
- Algorithms that use Markov chains to perform Monte Carlo estimates are called MCMC

Def

MCMC and Energy-Based Models

- Guarantees for MCMC are when model does not assign zero probability to any state
- Thus convenient to present these techniques as sampling from an energy-based model

(EBM)

$$p(\mathbf{x}) \propto \exp(-E(\mathbf{x}))$$

$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x})$$

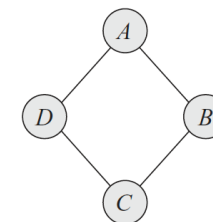
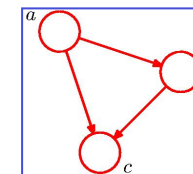
$$\tilde{p}(\mathbf{x}) = \prod_{C \in G} \phi(C)$$

$$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$$

$$\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x}))$$

Need more than ancestral sampling

- Ancestral Sampling for directed acyclic graphs:
 - Start with lowest numbered node
 - Draw a sample from the distribution $p(x_1)$ which we call \hat{x}_1
 - Work through each of the nodes in order
 - For node n we draw a sample from conditional distribution $p(x_n | pa_n)$
 - Defines an efficient single pass algorithm
- Not so simple in EBMs: chicken-egg problem
 - In order to sample from A we need to draw from $p(A | B, D)$
 - In order to sample from B we need to sample from $p(B | A, C)$
- Core idea of Markov chain
 - Have a state \mathbf{x} that begins with an arbitrary value
 - Over time we repeatedly update \mathbf{x}
 - Eventually \mathbf{x} becomes a fair sample from $p(\mathbf{x})$
 - Markov chain is defined by random state \mathbf{x} and transition distribution $T(\mathbf{x}' | \mathbf{x})$



Theoretical understanding of MCMC

- Reparameterize the problem
- Restrict attention to the case where r.v. x has countably many states
 - Represent the state as an integer x
- The different states are drawn from some distribution $q^{(t)}(x)$ where t is no of time steps
- Our goal is for $q^{(t)}(x)$ to converge to $p(x)$
 - If we have chosen T correctly then the stationary distribution q will be equal to the distribution p we wish to sample from
 - We describe how to choose T next

4. Gibbs Sampling

- The conceptually simplest approach for drawing samples from an undirected graph
- Suppose we have a graphical model over an n -dimensional vector of random variables \mathbf{x}
- We iteratively visit each variable x_i and draw a sample conditioned on all the other variables, i.e., from $p(x_i | \mathbf{x}_{-i})$
- Due to the separation properties of the graphical model, we can equivalently condition on only the neighbors of x_i

Gibbs Sampling with M variables

- Initialize first sample: $\{z_i, i=1, \dots, M\}$
- For $t=1, \dots, T$, $T = \text{no of samples}$
 - Sample $z_1^{(t+1)} \sim p(z_1 | z_2^{(t)}, z_3^{(t)}, \dots, z_M^{(t)})$
 - Sample $z_2^{(t+1)} \sim p(z_2 | z_1^{(t+1)}, z_3^{(t)}, \dots, z_M^{(t)})$
 -
 - Sample $z_j^{(t+1)} \sim p(z_j | z_1^{(t+1)}, \dots, z_{j-1}^{(t+1)}, z_{j+1}^{(t)}, \dots, z_M^{(t)})$
 -
 - Sample $z_M^{(t+1)} \sim p(z_M | z_1^{(t+1)}, z_2^{(t+1)}, \dots, z_{M-1}^{(t+1)})$
- $p(z_j | z_{-j})$ is called a *full conditional* for variable j