Monte Carlo Methods in Deep Learning

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

Monte Carlo vs Las Vegas Algorithms

Casino de Monte Carlo



- Monte Carlo algorithms return answers with random amount of error
 - For a fixed computational budget, MC algorithm provides an approximate answer

Las Vegas Strip



- Las Vegas algorithms return exact answers
 - But use random amount of resources

- Monte Carlo algorithms are fast but only probably correct
- Las Vegas algorithms are always correct but probably fast

Inference in Machine Learning

- Inference obtains answers from a model
 - 1. Learning

2. Inference



- Many ML problems are difficult
 - We cannot get precise answers
 - We can use either:
 - 1. Deterministic approximate algorithms or
 - 2. Monte Carlo approximations
- Both approaches are common

Role of Monte Carlo Methods

- 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

Deep Learning

Why Sampling?

- 1. Many ML algorithms are based on drawing samples from a probability distribution
 - To approximate sums and integrals during learning
 - Tractable sum:
 - Subsample for minibatches $\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \nabla_{\theta} L(\mathbf{x}^{(i)}, y^{(i)}, \theta)$
 - Gradient of log partition function of an RBM $\int_{x_{-p(x)} \nabla_{\theta} \log \tilde{p}(x) = \frac{1}{M} \sum_{i=1}^{M} \nabla_{\theta} \log \tilde{p}(x^{(m)}; \theta)}$
 - Intractable sum: $p(x, \theta) = \frac{1}{Z(\theta)} \tilde{p}(x, \theta)$ $Z(\theta) = \sum_{x} \tilde{p}(x, \theta)$
 - Gradient of log partition function of an MN
- 2. Sampling is actually our goal
 - Train a model so that we can sample from it
 - Style GAN



Preview of Sampling Methods Used

- MC: Importance Sampling Examples
 - 1. Neural language models with large vocabulary
 - Other neural nets with large no of outputs
 - 2. Log-likelihood in VAE
 - 3. Gradient in SGD where cost comes from few misclassified samples
 - Sampling difficult examples reduces variance of gradient
- MCMC: Gibbs Sampling Examples
 - 1. Inference in RBMs
 - 2. Deep Belief Nets

Basics of Monte Carlo sampling

Sum or integral is intractable

 Exponential no of terms and no exact simplification

 Approximate using Monte Carlo sampling

 View it as expectation under a distribution
 Approximate expectation by an average

Sampling for determining sums

Sampling allows approximating many sums and integrals at reduced cost

 $-\mathsf{E.g.}, p(\mathbf{v})=\Sigma_{\mathbf{h}} p(\mathbf{h},\mathbf{v})=\Sigma_{\mathbf{h}} p(\mathbf{v}/\mathbf{h})p(\mathbf{h})$



 When a sum/integral cannot be computed exactly, Monte Carlo methods view the sum as an expectation

Summation→Expectation→Average

1. Sum or integral to approximate is

 $s = \sum_{x} p(x) f(x)$ Or $s = \int p(x) f(x) dx$



2. Rewrite expression as an expectation

 $s = E_p[f(\boldsymbol{x})]$ Or $s = E_p[f(\boldsymbol{x})]$

- with the constraint that p is a

•probability distribution (for discrete case) or

•pdf (for the continuous case)

3. Approximate *s* by drawing *n* samples $x^{(1)}, ...x^{(n)}$ from *p* and then forming the empirical average

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\boldsymbol{x}^{(i)})$$
 10

Dee

Justification of Approximation

- Sample average $\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\boldsymbol{x}^{(i)})$ justified by 2 properties
 - 1. The estimator \hat{s}_n is unbiased, since

$$E[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n E\left[f\left(\boldsymbol{x}^{(i)}\right)\right]$$
$$= \frac{1}{n} \sum_{i=1}^n s$$
$$= s$$

2. Law of large numbers:

if $x^{(i)}$ are i.i.d. then average converges to expectation

 $\lim_{n \to \infty} \hat{s}_n = s \qquad \qquad s = E_p[f(\boldsymbol{x})]$

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

Deep

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

$$Var[\hat{s}_n] = \frac{1}{n^2} \sum_{i=1}^n Var[f(\boldsymbol{x})]$$
$$= \frac{Var[f(\boldsymbol{x})]}{n}$$

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

Sampling from base distribution *p*(*x*)

- Sampling from p(x) is not always possible
- In such a case use Importance Sampling
- A more general approach is Monte Carlo Markov chains
 - Form a sequence of estimators that converge towards the distribution of interest