

# Chapter 17

## Monte Carlo Methods

Randomized algorithms fall into two rough categories: Las Vegas algorithms and Monte Carlo algorithms. Las Vegas algorithms always return precisely the correct answer (or report that they failed). These algorithms consume a random amount of resources, usually memory or time. In contrast, Monte Carlo algorithms return answers with a random amount of error. The amount of error can typically be reduced by expending more resources (usually running time and memory). For any fixed computational budget, a Monte Carlo algorithm can provide an approximate answer.

Many problems in machine learning are so difficult that we can never expect to obtain precise answers to them. This excludes precise deterministic algorithms and Las Vegas algorithms. Instead, we must use deterministic approximate algorithms or Monte Carlo approximations. Both approaches are ubiquitous in machine learning. In this chapter, we focus on Monte Carlo methods.

### 17.1 Sampling and Monte Carlo Methods

Many important technologies used to accomplish machine learning goals are based on drawing samples from some probability distribution and using these samples to form a Monte Carlo estimate of some desired quantity.

#### 17.1.1 Why Sampling?

There are many reasons that we may wish to draw samples from a probability distribution. Sampling provides a flexible way to approximate many sums and

integrals at reduced cost. Sometimes we use this to provide a significant speedup to a costly but tractable sum, as in the case when we subsample the full training cost with minibatches. In other cases, our learning algorithm requires us to approximate an intractable sum or integral, such as the gradient of the log partition function of an undirected model. In many other cases, sampling is actually our goal, in the sense that we want to train a model that can sample from the training distribution.

### 17.1.2 Basics of Monte Carlo Sampling

When a sum or an integral cannot be computed exactly (for example the sum has an exponential number of terms and no exact simplification is known) it is often possible to approximate it using Monte Carlo sampling. The idea is to view the sum or integral as if it was an expectation under some distribution and to *approximate the expectation by a corresponding average*. Let

$$s = \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x}) = E_p[f(\mathbf{x})] \quad (17.1)$$

or

$$s = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = E_p[f(\mathbf{x})] \quad (17.2)$$

be the sum or integral to estimate, rewritten as an expectation, with the constraint that  $p$  is a probability distribution (for the sum) or a probability density (for the integral) over random variable  $\mathbf{x}$ .

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

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)}). \quad (17.3)$$

This approximation is justified by a few different properties. The first trivial observation is that the estimator  $\hat{s}$  is unbiased, since

$$\mathbb{E}[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(\mathbf{x}^{(i)})] = \frac{1}{n} \sum_{i=1}^n s = s. \quad (17.4)$$

But in addition, 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:

$$\lim_{n \rightarrow \infty} \hat{s}_n = s, \quad (17.5)$$

provided that the variance of the individual terms,  $\text{Var}[f(\mathbf{x}^{(i)})]$ , is bounded. To see this more clearly, consider the variance of  $\hat{s}_n$  as  $n$  increases. The variance  $\text{Var}[\hat{s}_n]$  decreases and converges to 0, so long as  $\text{Var}[f(\mathbf{x}^{(i)})] < \infty$ :

$$\text{Var}[\hat{s}_n] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[f(\mathbf{x})] \quad (17.6)$$

$$= \frac{\text{Var}[f(\mathbf{x})]}{n}. \quad (17.7)$$

This convenient result also tells us how to estimate the uncertainty in a Monte Carlo average or equivalently the amount of expected error of the Monte Carlo approximation. We compute both the empirical average of the  $f(\mathbf{x}^{(i)})$  and their empirical variance,<sup>1</sup> and then divide the estimated variance by the number of samples  $n$  to obtain an estimator of  $\text{Var}[\hat{s}_n]$ . The **central limit theorem** tells us that the distribution of the average,  $\hat{s}_n$ , converges to a normal distribution with mean  $s$  and variance  $\frac{\text{Var}[f(\mathbf{x})]}{n}$ . This allows us to estimate confidence intervals around the estimate  $\hat{s}_n$ , using the cumulative distribution of the normal density.

However, all this relies on our ability to easily sample from the base distribution  $p(\mathbf{x})$ , but doing so is not always possible. When it is not feasible to sample from  $p$ , an alternative is to use importance sampling, presented in section 17.2. A more general approach is to form a sequence of estimators that converge towards the distribution of interest. That is the approach of Monte Carlo Markov chains (section 17.3).

## 17.2 Importance Sampling

An important step in the decomposition of the integrand (or summand) used by the Monte Carlo method in equation 17.2 is deciding which part of the integrand should play the role the probability  $p(\mathbf{x})$  and which part of the integrand should play the role of the quantity  $f(\mathbf{x})$  whose expected value (under that probability distribution) is to be estimated. There is no unique decomposition because  $p(\mathbf{x})f(\mathbf{x})$  can always be rewritten as

$$p(\mathbf{x})f(\mathbf{x}) = q(\mathbf{x}) \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}, \quad (17.8)$$

where we now sample from  $q$  and average  $\frac{pf}{q}$ . In many cases, we wish to compute an expectation for a given  $p$  and an  $f$ , and the fact that the problem is specified

---

<sup>1</sup>The unbiased estimator of the variance is often preferred, in which the sum of squared differences is divided by  $n - 1$  instead of  $n$ .

from the start as an expectation suggests that this  $p$  and  $f$  would be a natural choice of decomposition. However, the original specification of the problem may not be the optimal choice in terms of the number of samples required to obtain a given level of accuracy. Fortunately, the form of the optimal choice  $q^*$  can be derived easily. The optimal  $q^*$  corresponds to what is called optimal importance sampling.

Because of the identity shown in equation 17.8, any Monte Carlo estimator

$$\hat{s}_p = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim p}^n f(\mathbf{x}^{(i)}) \quad (17.9)$$

can be transformed into an importance sampling estimator

$$\hat{s}_q = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim q}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}. \quad (17.10)$$

We see readily that the expected value of the estimator does not depend on  $q$ :

$$\mathbb{E}_q[\hat{s}_q] = \mathbb{E}_q[\hat{s}_p] = s. \quad (17.11)$$

However, the variance of an importance sampling estimator can be greatly sensitive to the choice of  $q$ . The variance is given by

$$\text{Var}[\hat{s}_q] = \text{Var}\left[\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]/n. \quad (17.12)$$

The minimum variance occurs when  $q$  is

$$q^*(\mathbf{x}) = \frac{p(\mathbf{x})|f(\mathbf{x})|}{Z}, \quad (17.13)$$

where  $Z$  is the normalization constant, chosen so that  $q^*(\mathbf{x})$  sums or integrates to 1 as appropriate. Better importance sampling distributions put more weight where the integrand is larger. In fact, when  $f(\mathbf{x})$  does not change sign,  $\text{Var}[\hat{s}_{q^*}] = 0$ , meaning that *a single sample is sufficient* when the optimal distribution is used. Of course, this is only because the computation of  $q^*$  has essentially solved the original problem, so it is usually not practical to use this approach of drawing a single sample from the optimal distribution.

Any choice of sampling distribution  $q$  is valid (in the sense of yielding the correct expected value) and  $q^*$  is the optimal one (in the sense of yielding minimum variance). Sampling from  $q^*$  is usually infeasible, but other choices of  $q$  can be feasible while still reducing the variance somewhat.

Another approach is to use **biased importance sampling**, which has the advantage of not requiring normalized  $p$  or  $q$ . In the case of discrete variables, the biased importance sampling estimator is given by

$$\hat{s}_{BIS} = \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}} \quad (17.14)$$

$$= \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \quad (17.15)$$

$$= \frac{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}}, \quad (17.16)$$

where  $\tilde{p}$  and  $\tilde{q}$  are the unnormalized forms of  $p$  and  $q$  and the  $\mathbf{x}^{(i)}$  are the samples from  $q$ . This estimator is biased because  $\mathbb{E}[\hat{s}_{BIS}] \neq s$ , except asymptotically when  $n \rightarrow \infty$  and the denominator of equation 17.14 converges to 1. Hence this estimator is called asymptotically unbiased.

Although a good choice of  $q$  can greatly improve the efficiency of Monte Carlo estimation, a poor choice of  $q$  can make the efficiency much worse. Going back to equation 17.12, we see that if there are samples of  $q$  for which  $\frac{p(\mathbf{x})|f(\mathbf{x})|}{q(\mathbf{x})}$  is large, then the variance of the estimator can get very large. This may happen when  $q(\mathbf{x})$  is tiny while neither  $p(\mathbf{x})$  nor  $f(\mathbf{x})$  are 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 in  $q$  causes it to match  $p$  or  $p|f|$  poorly. When  $q(\mathbf{x}^{(i)}) \gg p(\mathbf{x}^{(i)})|f(\mathbf{x}^{(i)})|$ , importance sampling collects useless samples (summing tiny numbers or zeros). On the other hand, when  $q(\mathbf{x}^{(i)}) \ll p(\mathbf{x}^{(i)})|f(\mathbf{x}^{(i)})|$ , which will happen more rarely, the ratio can be huge. Because these latter events are rare, they may not show up in a typical sample, yielding typical underestimation of  $s$ , compensated rarely by gross overestimation. Such very large or very small numbers are typical when  $\mathbf{x}$  is high dimensional, because in high dimension the dynamic range of joint probabilities can be very large.

In spite of this danger, importance sampling and its variants have been found very useful in many machine learning algorithms, including deep learning algorithms. For example, see the use of importance sampling to accelerate training in neural language models with a large vocabulary (section 12.4.3.3) or other neural nets with a large number of outputs. See also how importance sampling has been used to estimate a partition function (the normalization constant of a probability

distribution) in section 18.7, and to estimate the log-likelihood in deep directed models such as the variational autoencoder, in section 20.10.3. Importance sampling may also be used to improve the estimate of the gradient of the cost function used to train model parameters with stochastic gradient descent, particularly for models such as classifiers where most of the total value of the cost function comes from a small number of misclassified examples. Sampling more difficult examples more frequently can reduce the variance of the gradient in such cases (Hinton, 2006).

## 17.3 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 the distribution  $p_{\text{model}}(\mathbf{x})$  or from a good (low variance) importance sampling distribution  $q(\mathbf{x})$ . In the context of deep learning, this most often happens when  $p_{\text{model}}(\mathbf{x})$  is represented by an undirected model. In these cases, we introduce a mathematical tool called a **Markov chain** to approximately sample from  $p_{\text{model}}(\mathbf{x})$ . The family of algorithms that use Markov chains to perform Monte Carlo estimates is called **Markov chain Monte Carlo methods** (MCMC). Markov chain Monte Carlo methods for machine learning are described at greater length in Koller and Friedman (2009). The most standard, generic guarantees for MCMC techniques are only applicable when the model does not assign zero probability to any state. Therefore, it is most convenient to present these techniques as sampling from an energy-based model (EBM)  $p(\mathbf{x}) \propto \exp(-E(\mathbf{x}))$  as described in section 16.2.4. In the EBM formulation, every state is guaranteed to have non-zero probability. MCMC methods are in fact more broadly applicable and can be used with many probability distributions that contain zero probability states. However, the theoretical guarantees concerning the behavior of MCMC methods must be proven on a case-by-case basis for different families of such distributions. In the context of deep learning, it is most common to rely on the most general theoretical guarantees that naturally apply to all energy-based models.

To understand why drawing samples from an energy-based model is difficult, consider an EBM over just two variables, defining a distribution  $p(a, b)$ . In order to sample  $a$ , we must draw  $a$  from  $p(a | b)$ , and in order to sample  $b$ , we must draw it from  $p(b | a)$ . It seems to be an intractable chicken-and-egg problem. Directed models avoid this because their graph is directed and acyclic. To perform **ancestral sampling** one simply samples each of the variables in topological order, conditioning on each variable's parents, which are guaranteed to have already been sampled (section 16.3). Ancestral sampling defines an efficient, single-pass method

of obtaining a sample.

In an EBM, we can avoid this chicken and egg problem by sampling using a Markov chain. The core idea of a Markov chain is to have a state  $\mathbf{x}$  that begins as an arbitrary value. Over time, we randomly update  $\mathbf{x}$  repeatedly. Eventually  $\mathbf{x}$  becomes (very nearly) a fair sample from  $p(\mathbf{x})$ . Formally, a Markov chain is defined by a random state  $\mathbf{x}$  and a transition distribution  $T(\mathbf{x}' | \mathbf{x})$  specifying the probability that a random update will go to state  $\mathbf{x}'$  if it starts in state  $\mathbf{x}$ . Running the Markov chain means repeatedly updating the state  $\mathbf{x}$  to a value  $\mathbf{x}'$  sampled from  $T(\mathbf{x}' | \mathbf{x})$ .

To gain some theoretical understanding of how MCMC methods work, it is useful to reparametrize the problem. First, we restrict our attention to the case where the random variable  $\mathbf{x}$  has countably many states. We can then represent the state as just a positive integer  $x$ . Different integer values of  $x$  map back to different states  $\mathbf{x}$  in the original problem.

Consider what happens when we run infinitely many Markov chains in parallel. All of the states of the different Markov chains are drawn from some distribution  $q^{(t)}(x)$ , where  $t$  indicates the number of time steps that have elapsed. At the beginning,  $q^{(0)}$  is some distribution that we used to arbitrarily initialize  $x$  for each Markov chain. Later,  $q^{(t)}$  is influenced by all of the Markov chain steps that have run so far. Our goal is for  $q^{(t)}(x)$  to converge to  $p(x)$ .

Because we have reparametrized the problem in terms of positive integer  $x$ , we can describe the probability distribution  $q$  using a vector  $\mathbf{v}$ , with

$$q(x = i) = v_i. \quad (17.17)$$

Consider what happens when we update a single Markov chain's state  $x$  to a new state  $x'$ . The probability of a single state landing in state  $x'$  is given by

$$q^{(t+1)}(x') = \sum_x q^{(t)}(x) T(x' | x). \quad (17.18)$$

Using our integer parametrization, we can represent the effect of the transition operator  $T$  using a matrix  $\mathbf{A}$ . We define  $\mathbf{A}$  so that

$$A_{i,j} = T(\mathbf{x}' = i | \mathbf{x} = j). \quad (17.19)$$

Using this definition, we can now rewrite equation 17.18. Rather than writing it in terms of  $q$  and  $T$  to understand how a single state is updated, we may now use  $\mathbf{v}$  and  $\mathbf{A}$  to describe how the entire distribution over all the different Markov chains (running in parallel) shifts as we apply an update:

$$\mathbf{v}^{(t)} = \mathbf{A} \mathbf{v}^{(t-1)}. \quad (17.20)$$

















