

---

## Contents

---

<b>8</b>	<b>Variance reduction</b>	<b>3</b>
8.1	Overview of variance reduction . . . . .	3
8.2	Antithetics . . . . .	5
8.3	Example: expected log return . . . . .	8
8.4	Stratification . . . . .	10
8.5	Example: stratified compound Poisson . . . . .	14
8.6	Common random numbers . . . . .	17
8.7	Conditioning . . . . .	24
8.8	Example: maximum Dirichlet . . . . .	26
8.9	Control variates . . . . .	28
8.10	Moment matching and reweighting . . . . .	33
	End notes . . . . .	35
	Exercises . . . . .	38



---

## Variance reduction

---

Monte Carlo integration typically has an error variance of the form  $\sigma^2/n$ . We get a better answer by sampling with a larger value of  $n$ , but the computing time grows with  $n$ . Sometimes we can find a way to reduce  $\sigma$  instead. To do this, we construct a new Monte Carlo problem with the same answer as our original one but with a lower  $\sigma$ . Methods to do this are known as variance reduction techniques.

The techniques can be placed into groups, though no taxonomy is quite perfect. First we will look at antithetic sampling, stratification, and common random numbers. These methods all improve efficiency by sampling the input values more strategically. Next we will consider conditioning and control variates. These methods take advantage of closed form solutions to problems similar to the given one.

The last major method is importance sampling. Like some of the other methods, importance sampling also changes where we take the sample values, but rather than distributing them in more balanced ways it purposely oversamples from some regions and then corrects for this distortion by reweighting. It is thus a more radical reformulation of the problem and can be tricky to do well. We devote Chapter 9 to importance sampling. Some more advanced methods of variance reduction are given in Chapter 10.

### 8.1 Overview of variance reduction

Variance reductions are used to improve the efficiency of Monte Carlo methods. Before looking at individual methods, we discuss how to measure efficiency. Then we introduce some of the notation we need.

## Measuring efficiency

Methods of variance reduction can sometimes bring enormous improvements compared to plain Monte Carlo. It is not uncommon for the value  $\sigma^2$  to be reduced many thousand fold. It is also possible for a variance reduction technique to bring a very modest improvement, perhaps equivalent to reducing  $\sigma^2$  by only 10%. What is worse, some methods will raise  $\sigma^2$  in unfavorable circumstances.

The value of a variance reduction depends on more than the change in  $\sigma^2$ . It also depends on the computer's running time, possibly the memory consumed, and quite importantly, the human time taken to program and test the code.

Suppose for simplicity, that a baseline method is unbiased and estimates the desired quantity with variance  $\sigma_0^2/n$ , at a cost of  $nc_0$ , when  $n$  function evaluations are used. To get an error variance of  $\tau^2$  we need  $n = \sigma_0^2/\tau^2$  and this will cost  $c_0\sigma_0^2/\tau^2$ . Here we are assuming that cost is measured in time and that overhead cost is small.

If an alternative unbiased method has variance  $\sigma_1^2/n$  and cost  $nc_1$  under these conditions then it will cost us  $c_1\sigma_1^2/\tau^2$  to achieve the same error variance  $\tau^2$  that the baseline method achieved. The efficiency of the new method, relative to the standard method is

$$E = \frac{c_0\sigma_0^2}{c_1\sigma_1^2}. \quad (8.1)$$

At any fixed level of accuracy, the old method takes  $E$  times as much work as the new one.

The efficiency has two factors,  $\sigma_0^2/\sigma_1^2$  and  $c_0/c_1$ . The first is a mathematical property of the two methods that we can often handle theoretically. The second is more complicated. It can depend heavily on the algorithms used for each method. It can also depend on details of the computing environment, including the computer hardware, operating system, and implementation language. Numerical results for  $c_0/c_1$  obtained in one setting do not necessarily apply to another.

There is no fixed rule for how large an efficiency improvement must be to make it worth using. In some settings, such as rendering computer graphics for animated motion pictures, where thousands of CPUs are kept busy for months, a 10% improvement (i.e.,  $E = 1.1$ ) brings meaningful savings. In other settings, such as a one-off computation, a 60-fold gain (i.e.,  $E = 60$ ) which turns a one minute wait into a one second wait, may not justify the cost of programming a more complicated method.

Computation costs so much less than human effort that we ordinarily require large efficiency gains to offset the time spent programming up a variance reduction. The impetus to seek out an efficiency improvement may only come when we find ourselves waiting a very long time for a result, as for example, when we need to place our entire Monte Carlo calculation within a loop representing many variants of the problem. A very slow computation costs more than just the computer's time. It may waste time for those waiting for the answer. Also, slow computations reduce the number of alternatives that one can explore.

The efficiency gain necessary to justify using a method is less if the programming effort can be amortized over many applications. The threshold is high for a one time program, lower for something that we are adding to our personal library, lower still for code to share with a few coworkers and even lower for code to be put into a library or simulation tool for general use.

In the numerical examples in this chapter, some of the methods achieve quite large efficiency gains, while others are more modest. These results should not be taken as inherent to the methods. All of the methods are capable of a great range of efficiency improvements.

## Notation

Monte Carlo problems can be formulated through expectations or integrals or for discrete random variables, as sums. Generally, we will pick whichever format makes a given problem easiest to work with.

We suppose that the original Monte Carlo problem is to find  $\mu = \mathbb{E}(f(\mathbf{X}))$  where  $\mathbf{X}$  is a random variable from the set  $\mathcal{D} \subset \mathbb{R}^d$  with distribution  $p$ . When  $p$  is a probability density function we may write  $\mu = \int_{\mathcal{D}} f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$ . Most of the time we just write  $\mu = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$  with the understanding that  $p(\mathbf{x}) = 0$  for  $\mathbf{x} \notin \mathcal{D}$ . The integral version is convenient when we are reparameterizing the problem. Then, following the rules for integration is the best way to be sure of getting the right answer.

Monte Carlo sampling of  $\mathbf{X} \sim p$  is often based on  $s$  uniform random variables through a transformation  $\mathbf{X} = \psi(\mathbf{U})$ , for  $\mathbf{U} \sim \mathbf{U}(0,1)^s$ . Some variance reductions (e.g., antithetic sampling and stratification) are easier to apply directly to  $\mathbf{U}$  rather than to  $\mathbf{X}$ . For this case we write  $\mu = \int_{(0,1)^s} f(\psi(\mathbf{u}))d\mathbf{u}$ , or  $\mu = \int_{(0,1)^s} f^*(\mathbf{u})d\mathbf{u}$ , where  $f^*(\mathbf{u}) = f(\psi(\mathbf{u}))$ . When we don't have to keep track of both transformed and untransformed versions, then we just write  $\mu = \int_{(0,1)^d} f(\mathbf{u})d\mathbf{u}$ , subsuming  $\psi$  into  $f$ . This expression may be abbreviated to  $\mu = \int f(\mathbf{u})d\mathbf{u}$  when the domain of  $\mathbf{u}$  is clear from context.

Similar expressions hold for discrete random variables. Also some of the methods extend readily to  $d = \infty$ .

## 8.2 Antithetics

When we are using Monte Carlo averages of quantities  $f(\mathbf{X}_i)$  then the randomness in the algorithm leads to some error cancellation. In antithetic sampling we try to get even more cancellation. An antithetic sample is one that somehow gives the opposite value of  $f(\mathbf{x})$ , being low when  $f(\mathbf{x})$  is high and vice versa. Ordinarily we get an opposite  $f$  by sampling at a point  $\tilde{\mathbf{x}}$  that is somehow opposite to  $\mathbf{x}$ .

Let  $\mu = \mathbb{E}(\mathbf{X})$  for  $\mathbf{X} \sim p$ , where  $p$  is a symmetric density on the symmetric set  $\mathcal{D}$ . Here, symmetry is with respect to reflection through the center point  $\mathbf{c}$  of  $\mathcal{D}$ . If we reflect  $\mathbf{x} \in \mathcal{D}$  through  $\mathbf{c}$  we get the point  $\tilde{\mathbf{x}}$  with  $\tilde{\mathbf{x}} - \mathbf{c} = -(\mathbf{x} - \mathbf{c})$ , that is  $\tilde{\mathbf{x}} = 2\mathbf{c} - \mathbf{x}$ . Symmetry means that  $p(\tilde{\mathbf{x}}) = p(\mathbf{x})$  including the constraint

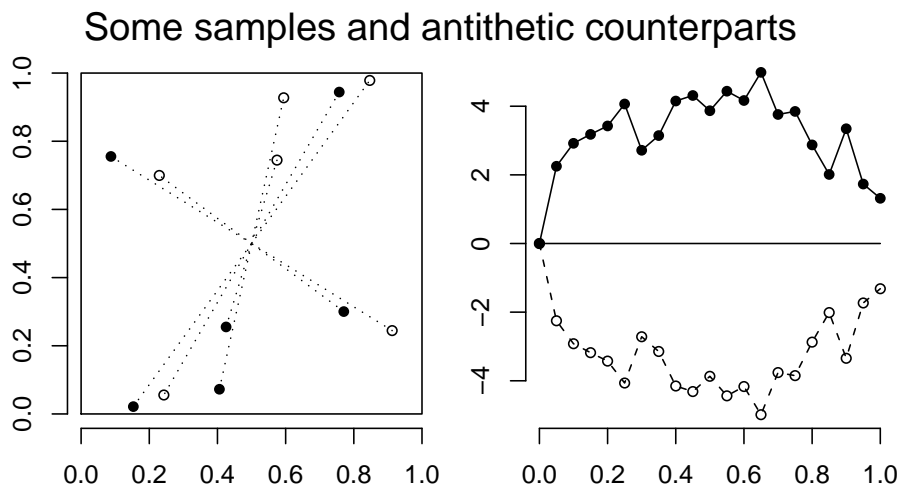


Figure 8.1: The left panel shows 6 points  $\mathbf{u}_i \in [0, 1]^2$  as solid points, connected to their antithetic counterparts  $\tilde{\mathbf{u}}_i = 1 - \mathbf{u}_i$ , shown as open circles. The right panel shows one random trajectory of 20 points joined by solid lines and connected to the origin, along with its antithetic mirror image in open points.

that  $\mathbf{x} \in \mathcal{D}$  if and only if  $\tilde{\mathbf{x}} \in \mathcal{D}$ . For basic examples, when  $p$  is  $\mathcal{N}(0, \Sigma)$  then  $\tilde{\mathbf{x}} = -\mathbf{x}$ , and when  $p$  is  $\mathbf{U}(0, 1)^d$  we have  $\tilde{\mathbf{x}} = 1 - \mathbf{x}$  componentwise. The antithetic counterpart of a random curve could be its reflection in the horizontal axis. See Figure 8.1 for examples. From the symmetry it follows that  $\tilde{\tilde{\mathbf{x}}} = \mathbf{x}$ .

The **antithetic sampling** estimate of  $\mu$  is

$$\hat{\mu}_{\text{anti}} = \frac{1}{n} \sum_{i=1}^{n/2} (f(\mathbf{X}_i) + f(\tilde{\mathbf{X}}_i)), \quad (8.2)$$

where  $\mathbf{X}_i \stackrel{\text{iid}}{\sim} p$ , and  $n$  is an even number.

The rationale for antithetic sampling is that each value of  $\mathbf{x}$  is balanced by its opposite  $\tilde{\mathbf{x}}$  satisfying  $(\mathbf{x} + \tilde{\mathbf{x}})/2 = \mathbf{c}$ . Whether this balance is helpful depends on  $f$ . Clearly if  $f$  is nearly linear we could obtain a large improvement. Suppose that  $\sigma^2 = \mathbb{E}((f(\mathbf{X}) - \mu)^2) < \infty$ . Then the variance in antithetic sampling is

$$\begin{aligned} \text{Var}(\hat{\mu}_{\text{anti}}) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^{n/2} f(\mathbf{X}_i) + f(\tilde{\mathbf{X}}_i)\right) \\ &= \frac{n/2}{n^2} \text{Var}(f(\mathbf{X}) + f(\tilde{\mathbf{X}})) \\ &= \frac{1}{2n} \left( \text{Var}(f(\mathbf{X})) + \text{Var}(f(\tilde{\mathbf{X}})) + 2\text{Cov}(f(\mathbf{X}), f(\tilde{\mathbf{X}})) \right) \end{aligned}$$

$$= \frac{\sigma^2}{n}(1 + \rho) \quad (8.3)$$

where  $\rho = \text{Corr}(f(\mathbf{X}), f(\tilde{\mathbf{X}}))$ .

From  $-1 \leq \rho \leq 1$  we obtain  $0 \leq \sigma^2(1 + \rho) \leq 2\sigma^2$ . In the best case, antithetic sampling gives the exact answer from just one pair of function evaluations. In the worst case it doubles the variance. Both cases do arise.

It is clear that a negative correlation is favorable. If  $f$  happens to be monotone in all  $d$  components of  $\mathbf{x}$ , then it is known that  $\rho < 0$ . Monotonicity of  $f$  is a safe harbor: if  $f$  is monotone then we're sure antithetic sampling will reduce the variance. We can often establish monotonicity theoretically, for example by differentiating  $f$ . But  $\rho < 0$  can hold without  $f$  being monotone in any of its inputs. Conversely  $\rho$  can be just barely negative when  $f$  is monotone. As a result, monotonicity alone is not a good guide to whether antithetic sampling will bring a large gain. See Exercise 8.1.

To get a qualitative understanding of antithetic sampling, break  $f$  into even and odd parts via

$$\begin{aligned} f(\mathbf{x}) &= \frac{f(\mathbf{x}) + f(\tilde{\mathbf{x}})}{2} + \frac{f(\mathbf{x}) - f(\tilde{\mathbf{x}})}{2} \\ &\equiv f_{\text{E}}(\mathbf{x}) + f_{\text{O}}(\mathbf{x}). \end{aligned}$$

The even part satisfies  $f_{\text{E}}(\mathbf{x}) = f_{\text{E}}(\tilde{\mathbf{x}})$  and  $\int_{\mathcal{D}} f_{\text{E}}(\mathbf{x})p(\mathbf{x}) \, d\mathbf{x} = \mu$ . The odd part satisfies  $f_{\text{O}}(\mathbf{x}) = -f_{\text{O}}(\tilde{\mathbf{x}})$  and  $\int_{\mathcal{D}} f_{\text{O}}(\mathbf{x})p(\mathbf{x}) \, d\mathbf{x} = 0$ .

The even and odd parts of  $f$  are orthogonal. This is not a surprise, because the product  $f_{\text{O}}(\mathbf{x})f_{\text{E}}(\mathbf{x})$  is an odd function. But to be careful and rule out  $\mathbb{E}(|f_{\text{O}}(\mathbf{X})f_{\text{E}}(\mathbf{X})|) = \infty$ , we compute directly that

$$\begin{aligned} \int_{\mathcal{D}} f_{\text{E}}(\mathbf{x})f_{\text{O}}(\mathbf{x})p(\mathbf{x}) \, d\mathbf{x} &= \int_{\mathcal{D}} \left( \frac{f(\mathbf{x}) + f(\tilde{\mathbf{x}})}{2} \right) \left( \frac{f(\mathbf{x}) - f(\tilde{\mathbf{x}})}{2} \right) p(\mathbf{x}) \, d\mathbf{x} \\ &= \frac{1}{4} \int_{\mathcal{D}} (f(\mathbf{x})^2 - f(\tilde{\mathbf{x}})^2) p(\mathbf{x}) \, d\mathbf{x} = 0. \end{aligned}$$

Now it follows easily that  $\sigma^2 = \sigma_{\text{E}}^2 + \sigma_{\text{O}}^2$  where  $\sigma_{\text{E}}^2 = \int_{\mathcal{D}} (f_{\text{E}}(\mathbf{x}) - \mu)^2 p(\mathbf{x}) \, d\mathbf{x}$  and  $\sigma_{\text{O}}^2 = \int_{\mathcal{D}} f_{\text{O}}(\mathbf{x})^2 p(\mathbf{x}) \, d\mathbf{x}$ .

Reworking equation (8.3) yields  $\hat{\mu}_{\text{anti}} = (2/n) \sum_{i=1}^{n/2} f_{\text{E}}(\mathbf{X}_i)$ . Therefore  $\text{Var}(\hat{\mu}_{\text{anti}}) = 2\sigma_{\text{E}}^2/n$  and we can combine this with the variance of ordinary Monte Carlo sampling as follows:

$$\begin{pmatrix} V(\hat{\mu}) \\ V(\hat{\mu}_{\text{anti}}) \end{pmatrix} = \frac{1}{n} \begin{pmatrix} 1 & 1 \\ 2 & 0 \end{pmatrix} \begin{pmatrix} \sigma_{\text{E}}^2 \\ \sigma_{\text{O}}^2 \end{pmatrix}. \quad (8.4)$$

We see from (8.4) that antithetic sampling eliminates the variance contribution of  $f_{\text{O}}$  but doubles the contribution from  $f_{\text{E}}$ . Antithetic sampling is extremely beneficial for integrands that are primarily odd functions of their inputs, having  $\sigma_{\text{O}}^2 \gg \sigma_{\text{E}}^2$ . The connection to correlation is via  $\rho = (\sigma_{\text{E}}^2 - \sigma_{\text{O}}^2)/(\sigma_{\text{E}}^2 + \sigma_{\text{O}}^2)$  (Exercise 8.3).

The analysis above shows that antithetic sampling reduces variance if  $\rho = \text{Corr}(f(\mathbf{X}), f(\widetilde{\mathbf{X}})) < 0$ , or equivalently, if  $\sigma_{\mathbf{O}}^2 > \sigma_{\mathbf{E}}^2$ . That analysis is appropriate when the most of the computation is in evaluating  $f$  and there is no economy in evaluating both  $f(\mathbf{X})$  and  $f(\widetilde{\mathbf{X}})$ .

Variance reduction is only part of the story because the cost of antithetic sampling using  $n$  points could well be smaller than the cost of plain Monte Carlo with  $n$  points. That will happen if it is expensive to generate  $\mathbf{X}$ , compared to the cost of computing  $f$ , but inexpensive to generate  $\widetilde{\mathbf{X}}$ . For example,  $\mathbf{X}$  might be a carefully constructed and expensive sample path from a Gaussian process while  $\widetilde{\mathbf{X}} = -\mathbf{X}$ .

We can explore this effect by letting  $c_x$  be the cost of generating  $\mathbf{X}$ , and  $c_f$  be the cost of computing  $f(\mathbf{X})$  once we have  $\mathbf{X}$ . We also let  $\tilde{c}_x$  and  $\tilde{c}_f$  be the corresponding costs for the antithetic sample. For illustration, suppose that to a reasonable approximation  $\tilde{c}_x = 0$  and  $\tilde{c}_f = c_f$ . In special circumstances  $\tilde{c}_f < c_f$  because it may be possible to reuse some computation.

Under the assumptions we are exploring, the efficiency of antithetic sampling relative to plain Monte Carlo is

$$E_{\text{anti}} = \frac{2c_x + 2c_f}{c_x + 2c_f} \times \frac{\sigma_{\mathbf{O}}^2 + \sigma_{\mathbf{E}}^2}{2\sigma_{\mathbf{E}}^2}.$$

Then antithetic sampling is more efficient than plain Monte Carlo if

$$\frac{\sigma_{\mathbf{O}}^2}{\sigma_{\mathbf{E}}^2} > \frac{c_f}{c_x + c_f}.$$

If generating  $\mathbf{x}$  costs ten times as much as computing  $f$  then antithetic sampling pays off when  $\sigma_{\mathbf{O}}^2/\sigma_{\mathbf{E}}^2 > 1/11$ .

Because antithetic samples have dependent values within pairs, the usual variance estimate must be modified. The most straightforward approach is to analyze the data as a sample of size  $m = n/2$  values of  $f_{\mathbf{E}}(\mathbf{X})$ . Let  $Y_i = f_{\mathbf{E}}(\mathbf{X}_i) = (f(\mathbf{X}_i) + f(\widetilde{\mathbf{X}}_i))/2$  for  $i = 1, \dots, m = n/2$ . Then take

$$\hat{\mu}_{\text{anti}} = \frac{1}{m} \sum_{i=1}^m Y_i, \quad \text{and}$$

$$s_{\text{anti}}^2 = \frac{1}{m-1} \sum_{i=1}^m (Y_i - \hat{\mu}_{\text{anti}})^2,$$

and use  $s_{\text{anti}}^2/m$  as the estimate of  $\text{Var}(\hat{\mu}_{\text{anti}})$ .

### 8.3 Example: expected log return

As an example of antithetic sampling we consider the expected logarithmic return of a portfolio. There are  $K$  stocks and the portfolio has proportion



$\lambda_k \geq 0$  invested in stock  $k$  for  $k = 1, \dots, K$ , with  $\sum_{k=1}^K \lambda_k = 1$ . The expected logarithmic return is

$$\mu(\lambda) = \mathbb{E}(\log(\sum_{k=1}^K \lambda_k e^{X_k})) \quad (8.5)$$

where  $\mathbf{X} \in \mathbb{R}^K$  is the vector of returns. At the end of the time period, the allocations are proportional to  $\lambda_k e^{X_k}$ . By selling some of the stocks with the largest  $X_k$  and buying some with the smallest  $X_k$ , it is possible to rebalance the portfolio so that the fraction of value in stock  $k$  is once again  $\lambda_k$ .

The expected logarithmic return is interesting because if one keeps reinvesting and rebalancing the portfolio at  $N$  regular time intervals then, by the law of large numbers, one's fortune grows as  $\exp(N\mu + o_p(N))$ , assuming of course that vectors  $\mathbf{X}$  for each time period are independent and identically distributed. See Luenberger (1998, Chapter 15). The log-optimal choice  $\lambda$  is the allocation that maximizes  $\mu$ . Log-optimal portfolios are of interest to very long term investors. Luenberger (1998) describes other criteria as well.

Finding a model for the distribution of  $\mathbf{X}$  and then choosing  $\lambda$  are challenging problems, but to illustrate antithetic sampling, simplified choices serve as well as elaborate ones. We focus on the problem of evaluating  $\mu(\lambda)$  for a given  $\lambda$ . We probably have to solve that problem en route to finding the best  $\lambda$  and definitely need to solve it once we have chosen  $\lambda$ . Here we take  $\lambda_k = 1/K$  for  $k = 1, \dots, K$  with  $K = 500$ . We also suppose that each marginal distribution is  $X_k \sim \mathcal{N}(\delta, \sigma^2)$  but that  $\mathbf{X}$  has the  $t(0, \nu, \Sigma)$  copula. Here  $\delta = 0.001$ ,  $\sigma = 0.03$ ,  $\nu = 4$  and  $\Sigma = \rho \mathbf{1}_K \mathbf{1}_K^T + (1 - \rho)I_K$  for  $\rho = 0.3$ . These values of  $\delta$  and  $\sigma$  are chosen to reflect roughly a one week time frame.

Letting  $f(\mathbf{X}) = \log(\sum_{k=1}^K e^{X_k}/K)$ , the plain Monte Carlo estimate of  $\mu$  is  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i)$ . The antithetic counterpart to  $\mathbf{X}_i$  has  $\tilde{\mathbf{X}}_{ik} = 2\delta - X_{ik}$ . Using  $n = 10,000$  sample values we find  $\hat{\rho}(f(\mathbf{X}), f(\tilde{\mathbf{X}})) \doteq -0.999508$  and so the variance reduction factor from antithetic sampling is  $(1 + \rho)^{-1} \doteq 2030.0$ .

For those  $n = 10,000$  pairs we let  $Y_i = (f(\mathbf{X}_i) + f(\tilde{\mathbf{X}}_i))/2 = f_E(\mathbf{X}_i)$  and get the estimate

$$\hat{\mu}_{\text{anti}} = \frac{1}{n} \sum_{i=1}^n Y_i \doteq 0.00132.$$

The standard deviation is  $s = ((n-1)^{-1} \sum_{i=1}^n (Y_i - \hat{\mu}_{\text{anti}})^2)^{1/2} \doteq 0.000252$ . The 99% confidence interval for  $\mu$  is

$$\hat{\mu}_{\text{anti}} \pm 2.58sn^{-1/2} \doteq 0.00132 \pm 6.49 \times 10^{-6}.$$

Antithetic sampling worked so well here because the function is nearly linear. The exponentials in (8.5) operate on a random variable that is usually near 0 and the logarithm operates on an argument that is usually near 1, and as a result the random variable whose expectation we take is nearly linear in  $\mathbf{X}$ . This near linearity is not limited to the particular  $\lambda$  and  $\Sigma$  we have used.

When  $\mathbf{X}$  varies more widely, then the curvature of the exponential and logarithmic functions makes more of a difference and antithetic sampling will

Stocks	Period	Correlation	Reduction	Estimate	Uncertainty
20	week	-0.99957	2320.0	0.00130	$6.35 \times 10^{-6}$
500	week	-0.99951	2030.0	0.00132	$6.49 \times 10^{-6}$
20	year	-0.97813	45.7	0.06752	$3.27 \times 10^{-4}$
500	year	-0.99512	40.2	0.06850	$3.33 \times 10^{-4}$

Table 8.1: This table summarizes the results of the antithetic sampling to estimate the expected log return of a portfolio, as described in the text. The first column has the number  $K$  of stocks. The second column indicates whether the return was for a week or a year. The third column is the correlation between log returns and their antithetic counterpart. The fourth column turns this correlation into a variance reduction factor. Then comes the estimate of expected log return and the half width of a 99% confidence interval.

lose some effectiveness. Let's consider for example, annual rebalancing, and take  $\delta = 52 \times 0.01$  and  $\sigma = \sqrt{52} \times 0.03$ . The annualized  $\mathbf{X}$  has the same mean and variance as the sum of 52 IID copies of the weekly random variable. It does not have quite the same copula. We ignore that small difference and simulate using the same  $t$  copula as before. In this case, we find a reduced but still substantial variance reduction of about 40 fold. Conversely, running an example with  $K = 20$  instead of 500 leads to a slightly bigger advantage for antithetic sampling. Four cases are summarized in Table 8.1.

## 8.4 Stratification

The idea in stratified sampling is to split up the domain  $\mathcal{D}$  of  $\mathbf{X}$  into separate regions, take a sample of points from each such region, and combine the results to estimate  $\mathbb{E}(f(\mathbf{X}))$ . Intuitively, if each region gets its fair share of points then we should get a better answer. Figure 8.2 shows two small examples of stratified domains. We might be able to do better still by oversampling within the important strata and undersampling those in which  $f$  is nearly constant.

We begin with the notation for stratified sampling. Then we show that stratified sampling is unbiased, find the variance of stratified sampling and show how to estimate that variance.

Our goal is to estimate  $\mu = \int_{\mathcal{D}} f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$ . We partition  $\mathcal{D}$  into mutually exclusive and exhaustive regions  $\mathcal{D}_j$ , for  $j = 1, \dots, J$ . These regions are the strata. We write  $\omega_j = \mathbb{P}(\mathbf{X} \in \mathcal{D}_j)$  and to avoid trivial issues, we assume  $\omega_j > 0$ . Next let  $p_j(\mathbf{x}) = \omega_j^{-1}p(\mathbf{x})\mathbb{1}_{\mathbf{x} \in \mathcal{D}_j}$ , the conditional density of  $\mathbf{X}$  given that  $\mathbf{X} \in \mathcal{D}_j$ .

To use stratified sampling, we must know the sizes  $\omega_j$  of the strata, and we must also know how to sample  $\mathbf{X} \sim p_j$  for  $j = 1, \dots, J$ . These conditions are quite reasonable. When we are defining strata, we naturally prefer ones we can sample from. If however, we know  $\omega_j$  but are unable to sample from  $p_j$ , then the method of post-stratification described on page 12 is available.

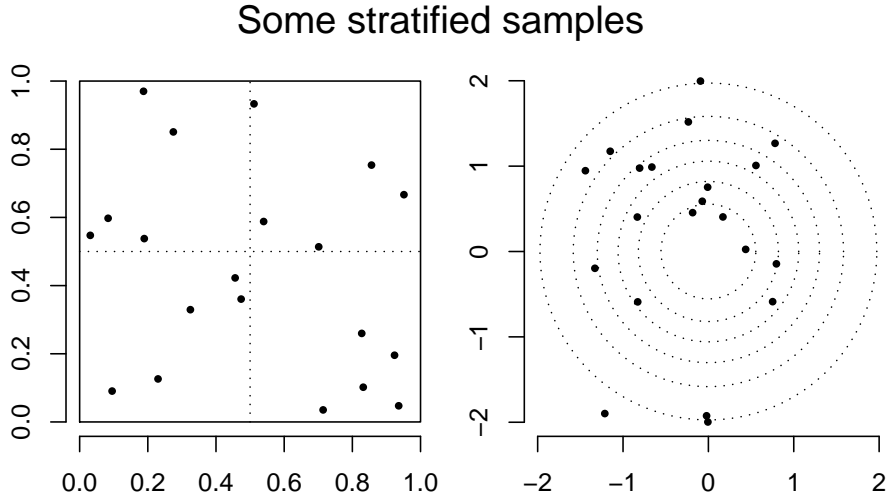


Figure 8.2: The left panel shows 20 points  $\mathbf{x}_i \in [0, 1]^2$  of which 5 are sampled uniformly from within each of four quadrants. The right panel shows 21 points from the  $\mathcal{N}(0, I_2)$  distribution. There are 6 concentric rings separating the distribution into 7 equally probable strata. Each stratum has 3 points sampled from within it.

Let  $\mathbf{X}_{ij} \sim p_j$  for  $i = 1, \dots, n_j$  and  $j = 1, \dots, J$  be sampled independently. The stratified sampling estimate of  $\mu$  is

$$\hat{\mu}_{\text{strat}} = \sum_{j=1}^J \frac{\omega_j}{n_j} \sum_{i=1}^{n_j} f(\mathbf{X}_{ij}). \quad (8.6)$$

We choose  $n_j > 0$  so that  $\hat{\mu}_{\text{strat}}$  is properly defined. Unless otherwise specified, we make sure that  $n_j \geq 2$ , which will allow the variance estimate (8.10) below to be applied.

Now

$$\begin{aligned} \mathbb{E}(\hat{\mu}_{\text{strat}}) &= \sum_{j=1}^J \omega_j \mathbb{E}\left(\frac{1}{n_j} \sum_{i=1}^{n_j} f(\mathbf{X}_{ij})\right) = \sum_{j=1}^J \omega_j \int_{\mathcal{D}_j} f(\mathbf{x}) p_j(\mathbf{x}) \, d\mathbf{x} \\ &= \sum_{j=1}^J \int_{\mathcal{D}_j} f(\mathbf{x}) p(\mathbf{x}) \, d\mathbf{x} = \int_{\mathcal{D}} f(\mathbf{x}) p(\mathbf{x}) \, d\mathbf{x} = \mu, \end{aligned} \quad (8.7)$$

and so stratified sampling is unbiased.

We study the variance of  $\hat{\mu}_{\text{strat}}$  to determine when stratification is advantageous, and to see how to design an effective stratification. Let  $\mu_j = \int_{\mathcal{D}_j} f(\mathbf{x}) p_j(\mathbf{x}) \, d\mathbf{x}$  and  $\sigma_j^2 = \int_{\mathcal{D}_j} (f(\mathbf{x}) - \mu_j)^2 p_j(\mathbf{x}) \, d\mathbf{x}$  be the  $j$ 'th stratum mean

and variance, respectively. The variance of the stratified sampling estimate is

$$\text{Var}(\hat{\mu}_{\text{strat}}) = \sum_{j=1}^J \omega_j^2 \frac{\sigma_j^2}{n_j}. \quad (8.8)$$

An immediate consequence of (8.8) is that  $\text{Var}(\hat{\mu}_{\text{strat}}) = 0$  for integrands  $f$  that are constant within strata  $\mathcal{D}_j$ . The variance of  $f(\mathbf{X})$  can be decomposed into within- and between-stratum components as follows

$$\sigma^2 = \sum_{j=1}^J \omega_j \sigma_j^2 + \sum_{j=1}^J \omega_j (\mu_j - \mu)^2. \quad (8.9)$$

Equation (8.9) is simply  $\text{Var}(f(\mathbf{X})) = \mathbb{E}(\text{Var}(f(\mathbf{X} | Z))) + \text{Var}(\mathbb{E}(f(\mathbf{X} | Z)))$  where  $Z \in \{1, \dots, J\}$  is the stratum containing the random point  $\mathbf{X}$ .

For error estimation, we write

$$\begin{aligned} \hat{\mu}_j &= \frac{1}{n_j} \sum_{i=1}^{n_j} Y_{ij}, & s_j^2 &= \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (Y_{ij} - \hat{\mu}_j)^2, \quad \text{and} \\ \widehat{\text{Var}}(\hat{\mu}_{\text{strat}}) &= \sum_{j=1}^J \omega_j^2 \frac{s_j^2}{n_j}. \end{aligned} \quad (8.10)$$

Clearly  $\mathbb{E}(s_j^2) = \sigma_j^2$  and so  $\mathbb{E}(\widehat{\text{Var}}(\hat{\mu}_{\text{strat}})) = \text{Var}(\hat{\mu}_{\text{strat}})$ . A central limit theorem based 99% confidence interval for  $\mu$  is

$$\hat{\mu}_{\text{strat}} \pm 2.58 \sqrt{\widehat{\text{Var}}(\hat{\mu}_{\text{strat}})}. \quad (8.11)$$

The CLT-based interval (8.11) is reasonable if all the  $n_j$  are large enough that each  $\hat{\mu}_j$  is nearly normally distributed. This condition is sufficient but not necessary. The estimate  $\hat{\mu}_{\text{strat}}$  is a sum of  $J$  terms  $\omega_j \hat{\mu}_j$ . Even if every  $n_j = 2$ , it might be reasonable to apply a central limit theorem holding as  $J \rightarrow \infty$  as described in Karr (1993, Chapter 7).

If we know  $\omega_j$  but prefer not to sample  $\mathbf{X} \sim p_j$  (or if we cannot do that), then we may still use the strata. In **post-stratification** we sample  $\mathbf{X}_i \sim p$  and assign  $\mathbf{X}_i$  to their strata after the fact. We let  $n_j$  be the number of sample points  $\mathbf{X}_i \in \mathcal{D}_j$ , let  $\hat{\mu}_j$  be the average of  $f(\mathbf{X}_i)$  for those points and  $s_j^2$  be their sample variance. Then we estimate  $\mu$  by the same  $\hat{\mu}_{\text{strat}}$  in (8.8) and use the same confidence interval (8.11) as before.

The main difference is that  $n_j$  are now random. There is also a risk of getting some  $n_j = 0$  in which case we cannot actually compute  $\hat{\mu}_{\text{strat}}$  by (8.7). However  $\mathbb{P}(\min_j n_j = 0) \leq \sum_{j=1}^J (1 - \omega_j)^n$  which we can make negligible by choosing  $n$  and the strata appropriately. Similarly, a sound choice for  $n$  and the strata  $\mathcal{D}_j$  will make  $n_j < 2$  very improbable.

Post-stratified sampling is a special case of the method of control variates. We will see this in Example 8.4 of §8.9.

A natural choice for stratum sample sizes is **proportional allocation**,  $n_j = n\omega_j$ . In our analysis, we'll suppose that all the  $n_j$  are integers. We can usually choose  $n$  and  $\mathcal{D}_j$  to make this so, or else accept small non-proportionalities due to rounding.

For proportional allocation, equation (8.6) for  $\hat{\mu}_{\text{strat}}$  reduces to the ordinary sample mean

$$\hat{\mu}_{\text{prop}} = \frac{1}{n} \sum_{j=1}^J \sum_{i=1}^{n_j} f(\mathbf{X}_{ij}). \quad (8.12)$$

Also, with proportional allocation, equation (8.8) for  $\text{Var}(\hat{\mu}_{\text{strat}})$  becomes

$$\sum_{j=1}^J \omega_j^2 \frac{\sigma_j^2}{n\omega_j} = \frac{1}{n} \sum_{j=1}^J \omega_j \sigma_j^2. \quad (8.13)$$

Equation (8.13) allows us to show that stratified sampling with proportional allocation cannot have larger variance than ordinary MC sampling. Let  $\sigma_{\text{W}}^2 = \sum_{j=1}^J \omega_j \sigma_j^2$  and  $\sigma_{\text{B}}^2 = \sum_{j=1}^J \omega_j (\mu_j - \mu)^2$  be the within- and between-stratum variances. We can compare IID and proportional stratification in one equation:

$$\begin{pmatrix} \text{Var}(\hat{\mu}) \\ \text{Var}(\hat{\mu}_{\text{prop}}) \end{pmatrix} = \frac{1}{n} \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \sigma_{\text{B}}^2 \\ \sigma_{\text{W}}^2 \end{pmatrix}. \quad (8.14)$$

A good stratification scheme is one that reduces the within-stratum variance, ideally leaving  $\sigma_{\text{B}}^2 \gg \sigma_{\text{W}}^2$ . If sampling from  $p_j$  is slower than sampling from  $p$ , then that reduces any efficiency gain from stratification.

Another way to look at proportional allocation is to construct the piece-wise constant function  $h(\mathbf{x})$  with  $h(\mathbf{x}) = \mu_j$  when  $\mathbf{x} \in \mathcal{D}_j$ . Then (Exercise 8.5),

$$\text{Var}(\hat{\mu}_{\text{prop}}) = (1 - \rho^2) \text{Var}(\hat{\mu}), \quad (8.15)$$

where  $\rho$  is the correlation between  $f(\mathbf{X})$  and  $h(\mathbf{X})$  for  $\mathbf{X} \sim p$ .

A proportional allocation is not necessarily the most efficient. For instance, given two strata with equal  $\omega_j$  but unequal  $\sigma_j^2$ , we benefit by taking fewer points from the less variable stratum. In the extreme, if  $\sigma_j = 0$  then  $n_j = 1$  is enough to tell us  $\mu_j$ .

The problem of optimal sample allocation to strata has been solved in the survey sampling literature. The result is known as the Neyman allocation, and the formulation allows for unequal sampling costs from the different strata. Suppose that for unit costs  $c_j > 0$  the stratified sampling costs  $C + \sum_{j=1}^J n_j c_j$  to generate random variables and evaluate  $f$ . Here  $C \geq 0$  is an overhead cost and  $c_j$  is the (expected) cost to generate  $\mathbf{X}$  from  $p_j$  and then compute  $f(\mathbf{X})$ . To minimize variance subject to an upper bound on cost, take

$$n_j \propto \frac{\omega_j \sigma_j}{\sqrt{c_j}}. \quad (8.16)$$

The solution (8.16) also minimizes cost subject to a lower bound on variance. Equation (8.16) can be established by the method of Lagrange multipliers. These optimal values  $n_j$  usually need to be rounded to integers and some may have to be raised, if for other reasons we insist that all  $n_j$  be above some minimum such as 2.

When the sampling cost  $c_j$  is the same in every stratum then the optimal allocation has

$$n_j = \frac{n\omega_j\sigma_j}{\sum_{k=1}^J \omega_k\sigma_k}. \quad (8.17)$$

Let  $\hat{\mu}_{n\text{-opt}}$  be the stratified sampling estimate (8.6) with optimal  $n_j$  from (8.17). By substituting (8.17) into the stratified sampling variance (8.8) we find that

$$\text{Var}(\hat{\mu}_{n\text{-opt}}) = \frac{1}{n} \left( \sum_{j=1}^J \omega_j\sigma_j \right)^2 \leq \frac{1}{n} \sum_{j=1}^J \omega_j\sigma_j^2 = \text{Var}(\hat{\mu}_{\text{prop}}). \quad (8.18)$$

Equality holds in (8.18), only when  $\sigma_j$  is constant in  $j$ .

In typical applications, the values of  $\sigma_j$  are not known. We might make an educated guess  $\hat{\sigma}_j$  and then employ  $n_j \propto \omega_j\hat{\sigma}_j$ . The optimal allocation only depends on  $\sigma_1, \dots, \sigma_J$  through ratios  $\sigma_j/\sigma_k$  for  $j \neq k$ , and so only the ratios  $\hat{\sigma}_j/\hat{\sigma}_k$  need to be accurate. Non-proportional allocations carry some risk. The optimal allocation assuming  $\sigma_j = \hat{\sigma}_j$  can be worse than proportional allocation if it should turn out that  $\sigma_j$  are not proportional to  $\hat{\sigma}_j$ . It can even give higher variance than ordinary Monte Carlo sampling, completely defeating the effort put into stratification.

From results in survey sampling (Cochran, 1977), it is known how to construct theoretically optimal strata. The variance minimizing strata take the form  $\mathcal{D}_k = \{\mathbf{x} \mid a_{k-1} \leq f(\mathbf{x}) < a_k\}$  for some constants  $a_0 < a_1 < \dots < a_J$ . There are also guidelines for choosing the  $a_j$ . In practice we cannot usually locate the contours of  $f$  and even when we can it will usually be very hard to sample between them. But the intuition is still valuable: we want strata within which  $f$  is as flat as possible.

## 8.5 Example: stratified compound Poisson

Compound Poisson models are commonly used for rainfall. Here we will look at stratifying such a model.

In our model setting, the number of rainfall events (storms) in the coming month is  $S \sim \text{Poi}(\lambda)$  with  $\lambda = 2.9$ . The depth of rainfall in storm  $i$  is  $D_i \sim \text{Weib}(k, \sigma)$  with shape  $k = 0.8$  and scale  $\sigma = 3$  (centimeters) and the storms are independent. If the total rainfall is below 5 centimeters then an emergency water allocation will be imposed.

The total rainfall is thus  $X = \sum_{s=1}^S D_s$  taking the value 0 when  $S = 0$ . It is easy to get the mean and variance of  $X$ , but here we want  $\mathbb{P}(X < 5)$ , that is  $\mathbb{E}(f(\mathbf{X}))$  where  $f(\mathbf{X}) = \mathbb{1}_{X < 5}$ . In a direct simulation, depicted in Figure 8.3,

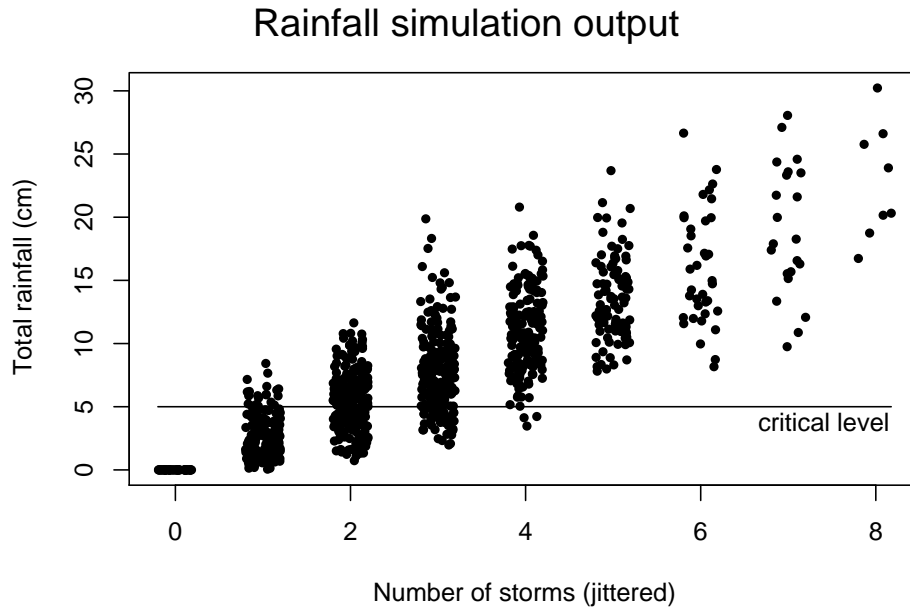


Figure 8.3: This figure depicts 1000 simulations of the compound Poisson model for rainfall described in the text.

the rainfall was below the critical level 353 times out of 1000. Thus the estimate of  $\mathbb{P}(\mathbf{X} < 5)$  is  $\hat{\mu} = 0.353$ . Because this probability is not near 0 or 1 a simple 99% confidence interval of  $\hat{\mu} \pm 2.58\sqrt{\hat{\mu}(1-\hat{\mu})/n}$  is adequate, and it yields the confidence interval  $0.314 \leq \mathbb{P}(\mathbf{X} < 5) \leq 0.392$ .

From simple Monte Carlo, we learn that the probability of a critically low total rainfall is roughly 30 to 40 percent. From Figure 8.3 we see that this probability depends strongly on the number of rainfall events.

Consider stratifying  $S$  according to a proportional allocation. The number of times  $S = s$  in 1000 trials should be  $n_s = 1000 e^{-\lambda} \lambda^s / s!$  where  $\lambda = 2.9$ . Two issues come up immediately. First, the sample sizes  $n_s$  are not integers. That is not a serious problem. We can use rounded sample sizes, in an approximately proportional allocation and still obtain an unbiased estimate of  $\mathbb{P}(\mathbf{X} < 5)$  and a workable variance estimate. The second issue to come up is that when  $S = 0$  we don't really need to simulate at all. In that case we are sure that  $\mathbf{X} < 5$ . For the second issue we will take  $n_0 = 2$ . That way we can use the plain stratified sampling formulas (8.6) and (8.10), and we only waste 2 of 1000 simulations on the foregone conclusion that with no storms there will be a water shortage.

Taking  $n_0 = 2$  samples with  $S = 0$  and allocating the remaining 998 in proportion to  $\mathbb{P}(S = s)/(1 - \mathbb{P}(S = 0))$  we get the counts

s	0	1	2	3	4	5	$\geq 6$
$n_s$	2	169	244	236	171	99	79

$s$	$\omega_s$	$n_s$	$T_s$	$\hat{\mu}_s$	$\hat{\sigma}_s^2$
0	0.055	2	2	1.000	0.000
1	0.160	169	152	0.899	0.091
2	0.231	244	111	0.455	0.249
3	0.224	236	33	0.140	0.121
4	0.162	171	3	0.018	0.017
5	0.094	99	1	0.010	0.010
6+	0.074	79	1	0.013	0.013

Table 8.2: This table shows the results of a stratified simulation of the compound Poisson rainfall model from the text. Here  $s$  is the number of storms. The last stratum is for  $s \geq 6$ . Continuing,  $\omega_s$  is  $\mathbb{P}(S = s)$  under a Poisson model, and  $n_s$  is the number of simulations allocated to  $S = s$ . Of  $n_s$  trials, there were  $T_s$  below the critical level. Then  $\hat{\mu}_s$  and  $\hat{\sigma}_s^2$  are estimated within stratum means and variances.

where the values from 6 on up have been merged into one stratum.

The  $S \geq 6$  stratum is more complicated to sample from than the others. One way is to first find  $q_6 = \sum_{s=0}^6 e^{-\lambda} \lambda^s / s!$ . Then draw  $S = F_\lambda^{-1}(q_6 + (1 - q_6)U)$  where  $U \sim \mathbf{U}(0, 1)$  and  $F_\lambda$  is the  $\text{Poi}(\lambda)$  CDF.

The results of this simulation are shown in Table 8.2. Using those values, the estimated probability of a shortage is  $\hat{\mu}_{\text{strat}} = \sum_s \omega_s \hat{\mu}_s \doteq 0.334$ . Using equation (8.10),  $\widehat{\text{Var}}(\hat{\mu}_{\text{strat}}) = \sum_s \omega_s^2 \hat{\sigma}_s^2 / n_s \doteq 9.84 \times 10^{-5}$ . The plain Monte Carlo simulation has an estimated variance of  $\hat{p}(1 - \hat{p})/n \doteq 0.353 \times 0.643/1000 \doteq 2.28 \times 10^{-4}$ , about 2.3 times as large as the estimated variance for stratified sampling.

This value 2.3 is only an estimate, but it turns out to be close to correct. In 10,000 independent replications of both methods the sample variance of the 10,000 plain Monte Carlo simulation answers was 2.24 times as large as that of the 10,000 stratified sampling answers.

A variance reduction of just over 2-fold is helpful but not enormous. Such a variance reduction would only justify the extra complexity of stratified sampling, if we needed to run many simulations of this sort.

The estimated factor of 2.24 does not take into account running time. Stratification has the possibility of being slightly faster here because most of the samples are deterministic: instead of sampling 1000 Poisson random variables, we generate 79 variables from the right tail of the Poisson distribution and use pre-chosen values for the other 921 Poisson random variables.

A further modest variance reduction can be obtained by reducing the number of observations with  $s \geq 5$ , increasing the number with  $s = 2$  or 3, and replacing the estimate from  $s = 1$  by  $\mathbb{P}(\text{Weib}(k, \sigma) \leq 5)$ . None of these steps can bring a dramatic increase in accuracy because the strata  $s = 2$  and 3 have high variance. Stratifying on  $S$  cannot help with the variance of  $f(\mathbf{X})$  given  $S = s$ .



## 8.6 Common random numbers

Suppose that  $f$  and  $g$  are closely related functions and that we want to find  $\mathbb{E}(f(\mathbf{X}) - g(\mathbf{X}))$  for  $\mathbf{X} \sim p$ . Perhaps  $f(\mathbf{x}) = h(\mathbf{x}, \theta)$  for a parameter  $\theta \in \mathbb{R}^p$ , and then to study the effect of  $\theta$  we look at  $g(\mathbf{x}) = h(\mathbf{x}, \tilde{\theta})$  for some  $\tilde{\theta} \neq \theta$ . We assume at first that neither  $f$  nor  $g$  (nor  $h$ ) makes any use of random numbers other than  $\mathbf{X}$ . Later we relax that assumption.

Because  $\mathbb{E}(f(\mathbf{X}) - g(\mathbf{X})) = \mathbb{E}(f(\mathbf{X})) - \mathbb{E}(g(\mathbf{X}))$  we clearly have two different ways to go. We could estimate the difference by

$$\widehat{D}_{\text{com}} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i) - g(\mathbf{X}_i), \quad (8.19)$$

for  $\mathbf{X}_i \stackrel{\text{iid}}{\sim} p$ , or by differencing averages

$$\widehat{D}_{\text{ind}} = \frac{1}{n_1} \sum_{i=1}^{n_1} f(\mathbf{X}_{i1}) - \frac{1}{n_2} \sum_{i=1}^{n_2} g(\mathbf{X}_{i2}) \quad (8.20)$$

for  $\mathbf{X}_{ij} \stackrel{\text{iid}}{\sim} p$ . Taking  $n = n_1 = n_2$  makes the computing costs in (8.19) and (8.20) comparable, assuming that costs of computing  $f$  and  $g$  dominate those of generating  $\mathbf{X}$ .

The sampling variances of these methods are

$$\begin{aligned} \text{Var}(\widehat{D}_{\text{com}}) &= \frac{1}{n} (\sigma_f^2 + \sigma_g^2 - 2\rho\sigma_f\sigma_g) \\ \text{Var}(\widehat{D}_{\text{ind}}) &= \frac{1}{n} (\sigma_f^2 + \sigma_g^2), \end{aligned} \quad (8.21)$$

where  $\sigma_f^2$  and  $\sigma_g^2$  are individual function variances and  $\rho = \text{Corr}(f(\mathbf{X}), g(\mathbf{X}))$ . When  $\rho > 0$  we are better off using common random numbers. There is no guarantee that  $\rho > 0$ . When  $f$  and  $g$  compute similar quantities then we anticipate that  $\rho > 0$ , and if so, then  $\widehat{D}_{\text{com}}$  is more effective than  $\widehat{D}_{\text{ind}}$ .

Most people would instinctively use the common variates. So at first sight, the method looks more like avoiding a variance increase than engineering a variance decrease. Later, when we relax the rule forbidding  $f$ ,  $g$ , and  $h$  to use other sources of randomness, we will find that retaining some common random numbers requires considerable care in synchronization. The added complexity might well tip the balance against using common random numbers.

Much the same problem arises if we are comparing  $\mathbb{E}(f(\mathbf{X}))$  for  $\mathbf{X} \sim p$  and  $\mathbb{E}(f(\widetilde{\mathbf{X}}))$  for  $\widetilde{\mathbf{X}} \sim \tilde{p}$ . Sometimes we can rewrite that problem in terms of common random variables that get transformed to a different distribution before  $f$  is applied. For instance, if the first simulation has  $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$  and the second has  $\widetilde{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2)$  then we can sample  $Z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$  and use

$$\widehat{D}_{\text{com}} = \frac{1}{n} \sum_{i=1}^n f(\mu + \sigma Z_i) - f(\tilde{\mu} + \tilde{\sigma} Z_i).$$

More generally, when  $\mathbf{X}_i$  is generated via a transformation  $\Psi(\mathbf{U}_i; \theta)$  of  $\mathbf{U}_i \sim \mathbf{U}(0, 1)^s$  then we can average  $f(\Psi(\mathbf{U}_i; \theta)) - f(\Psi(\mathbf{U}_i; \tilde{\theta}))$ .

Acceptance-rejection sampling of  $\mathbf{X}$  does not fit cleanly into this framework, because the number  $s$  of needed uniform random variables is not fixed and may vary with  $\theta$ .

The construction above is a **coupling** of the random vectors  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$ . Any joint distribution on  $(\mathbf{X}, \tilde{\mathbf{X}})$  with  $\mathbf{X} \sim p$  and  $\tilde{\mathbf{X}} \sim \tilde{p}$  is a coupling. Common random numbers provide a particularly close coupling between  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$ .

### Example: dosage content uniformity

Medicines are typically sold with a label claim giving the amount of active ingredient that should be in each dose. The actual amount fluctuates but should be close to the claim. Sampling schemes are used to determine whether a given lot has high enough quality. The average dose should be close to the target and the standard deviation should not be too large.

There are many different types of test, depending on the product (tablet, capsule, aerosol, skin patch, etc.). Here is one, based on the US Pharmacopeial Convention content uniformity test. We will measure the dose as a percentage of the label claim, and assume that the target value is 100% of label claim. In some instances targets over 100% are considered, perhaps to compensate for declining dosage in storage.

To describe the test, we need to introduce the function

$$M(x) = \begin{cases} 98.5, & x < 98.5 \\ x, & 98.5 \leq x \leq 101.5 \\ 101.5, & x > 101.5. \end{cases} \quad (8.22)$$

This function will be used to make the test less sensitive to tiny fluctuations in the average dose. Exercise 8.22 looks at whether using  $M(x)$  makes any difference to the acceptance probability.

The test first samples 10 units, getting measured values  $x_1, \dots, x_{10}$ . Then the values

$$\bar{x}_1 = \frac{1}{10} \sum_{j=1}^{10} x_j, \quad s_1^2 = \frac{1}{9} \sum_{j=1}^{10} (x_j - \bar{x}_1)^2, \quad \text{and} \quad M_1 = M(\bar{x}_1)$$

are computed. The lot passes if  $|\bar{x}_1 - M_1| + 2.4s_1 \leq 15$ . Otherwise, 20 more units are sampled giving  $x_{11}, \dots, x_{30}$ . Then the values

$$\bar{x}_2 = \frac{1}{30} \sum_{j=1}^{30} x_j, \quad s_2^2 = \frac{1}{29} \sum_{j=1}^{30} (x_j - \bar{x}_2)^2, \quad \text{and} \quad M_2 = M(\bar{x}_2)$$

are computed. The lot passes if  $|\bar{x}_2 - M_2| + 2.0s_2 \leq 15$  and  $\min_{1 \leq j \leq 30} x_j \geq 0.75M_2$  and  $\max_{1 \leq j \leq 30} x_j \leq 1.25M_2$ . Otherwise it fails.

### Estimated probability to pass content uniformity test

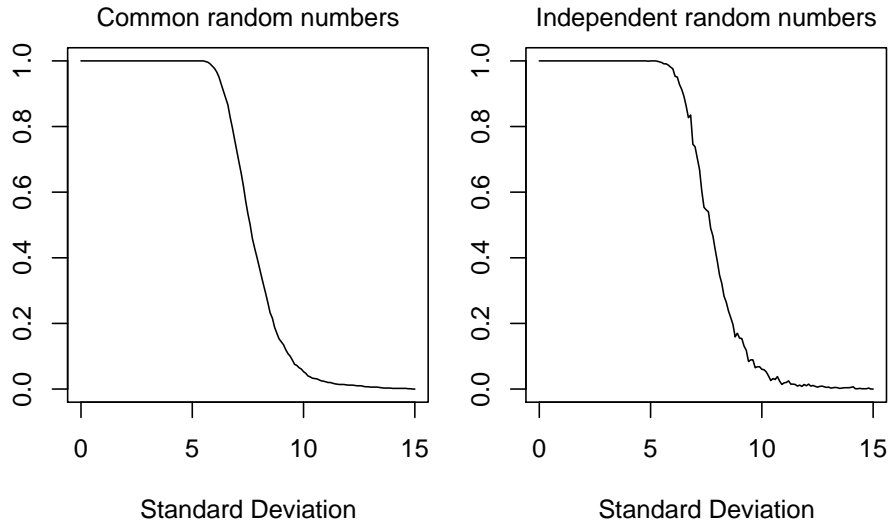


Figure 8.4: Each panel shows the estimated probability of passing the content uniformity test for  $X_i \sim \mathcal{N}(100, \sigma^2)$  as the standard deviation increases from 0 to 15 units. The smooth curve on the left is based on common random numbers. The rougher curve on the right uses independent random numbers. Both were based on  $n = 1000$  replications.

When the quality is high, the product usually passes at the first stage, and then the two stage test saves time and expense. But the two stage test is not amenable to closed form analysis even when  $x_j \sim \mathcal{N}(\mu, \sigma^2)$ . Monte Carlo methods are well suited to studying the probability of passing the test.

A direct simulation of the process is easy to do. But suppose that we want to compare the effects of varying  $\mu$  and  $\sigma$  on the passage probability. Then it makes sense to use a common random number scheme with  $Z_1, \dots, Z_{30}$  sampled independently from  $\mathcal{N}(0, 1)$  and  $X_j = \mu + \sigma Z_j$  for  $j = 1, \dots, 30$ . To keep the simulation synchronized, we always reserve the values  $Z_{11}, \dots, Z_{30}$  for the second stage, even when the test is accepted at stage 1.

When  $\mu = 100$ , the test will tend to fail if  $\sigma$  is high enough. Figure 8.4 shows Monte Carlo estimates of the probability of passing the uniformity test for  $X_j \sim \mathcal{N}(100, \sigma^2)$  with  $0 \leq \sigma \leq 15$ . The probability of passing is very high for  $\sigma \leq 5.5$  or so, but then it starts to drop quickly. When common random numbers are used, the estimated probability is very smooth, and also monotone, in  $\sigma$ . When independent random numbers are used, the estimated probability is non-monotone and the non-smoothness is even visible to the eye.

For large enough  $n$ , the non-smoothness would not be visible, but it would still result in less accurate estimation of differences in acceptance probability.

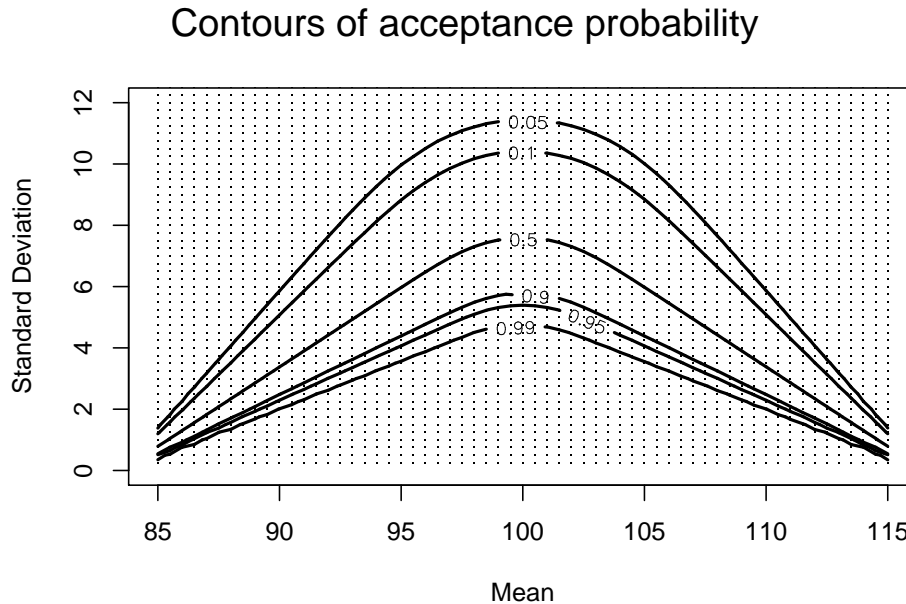


Figure 8.5: This plot shows contours of the acceptance probability of the content uniformity test when the data are  $\mathcal{N}(\mu, \sigma^2)$ . The horizontal axis is  $\mu$  and the vertical axis is  $\sigma$ . Monte Carlo sampling with  $n = 100,000$  points was run at each point of the grid shown in light dots. Values of  $\mu$  run from 85 to 115 in steps of 0.5, while  $\sigma$  runs from 0.25 to 12.0 in steps of 0.25. Common random numbers were used.

The acceptance probability is mapped out as a function of  $\mu$  and  $\sigma$  in Figure 8.5. That figure was created by using a common random numbers Monte Carlo sample on a grid of  $(\mu, \sigma)$  pairs. There is a roughly triangular region in the  $(\mu, \sigma)$  plane where the success probability is over 99%. Because the probability is between 99 and 100 percent there, and is monotone in  $\sigma$  and  $|\mu - 100|$ , the surface is very flat within this triangle. The region with 99.9% success probability (not shown) is just barely smaller than the one with 99% probability. There is a tiny bit of wiggle in some of the contours partly because the grid spacing is wide and partly because those contours go through a region where failures are rare events.

### Implementing common random numbers

We want to estimate  $\mu_j = \mathbb{E}(h(\mathbf{X}, \theta_j))$  for  $j = 1, \dots, m$  using  $n$  random inputs  $\mathbf{X}_i$ , for  $i = 1, \dots, n$ . The content uniformity example had a large value of  $m$  but in the simplest case,  $m = 2$  and we're interested in  $\mu_1 - \mu_2$ . We still assume that  $h$  really is a function of  $\mathbf{X}$  and  $\theta$  and in particular our implementation of  $h$  does not cause more random numbers to be generated.

**Algorithm 8.1** Common random numbers algorithm I

---

```

setseed(seed)
 $\hat{\mu}_j \leftarrow 0, \quad 1 \leq j \leq m$ 
for  $i = 1$  to  $n$  do
   $\mathbf{X}_i \sim p$ 
   $\hat{\mu}_j \leftarrow \hat{\mu}_j + h(\mathbf{X}_i, \theta_j), \quad 1 \leq j \leq m$ 
 $\hat{\mu}_j \leftarrow \hat{\mu}_j/n, \quad 1 \leq j \leq m$ 
deliver  $\hat{\mu}_1, \dots, \hat{\mu}_m$ 

```

---

This algorithm shows the method of common random numbers with the outer loop over random samples. The only random numbers used in  $h$  are from  $\mathbf{X}_i$ . Setting the seed keeps the  $\mathbf{X}_i$  reproducible if we change our list of  $\theta_j$ . The vectorized approach of equation (8.23) may be convenient.

---

We can run a nested loop over samples indexed by  $i$  and parameter values indexed by  $j$ . There are two main approaches that we can take, depending on which is the outer loop.

Algorithm 8.1 shows common random numbers with the outer loop over  $\mathbf{X}_i$  for  $i = 1, \dots, n$ . When  $\mathbf{X}_i$  is multi-dimensional we have to make sure that every component of  $\mathbf{X}$  needed for any value of  $\theta_j$  is provided. In the content uniformity problem (page 18) we generate  $Z_{11}$  for every simulated batch even though some only use  $Z_1, \dots, Z_{10}$ .

A vectorized implementation of Algorithm 8.1 is advantageous. It uses a function  $H$  that takes  $\mathbf{X}$  and a list  $\Theta = (\theta_1, \dots, \theta_m)$  of parameter values. This  $H$  returns a list  $(h(\mathbf{X}, \theta_1), \dots, h(\mathbf{X}, \theta_m))$  and the simulation computes

$$(\hat{\mu}_1, \dots, \hat{\mu}_m) = \frac{1}{n} \sum_{i=1}^n H(\mathbf{X}_i, \Theta). \quad (8.23)$$

This vectorized  $H$  makes it easier to separate the code that creates  $\Theta$  from that which evaluates  $h$ .

Algorithm 8.2 shows common random numbers with the outer loop over the parameters. It regenerates all  $n$  vectors  $\mathbf{X}_i$  for each  $j$ . To keep these vectors synchronized it keeps resetting the random seed. If we look at the output from Algorithm 8.1 partway through the computation, we will see incomplete estimates for all of the  $\theta_j$ . If we do that for Algorithm 8.2 we will see completed estimates for a subset of the  $\theta_j$ .

Now suppose that we relax our constraint on  $h$  and allow it to sample random numbers. That creates some messy synchronization issues described on page 36 of the end notes. Algorithm 8.1 is more robust to this change than Algorithm 8.2, but both could bring unpleasant surprises. Such a relaxation leaves us with only partially common numbers that we look at next.

**Algorithm 8.2** Common random numbers algorithm II

---

```

for  $j = 1$  to  $m$  do
  setseed(seed),  $\hat{\mu}_j \leftarrow 0$ 
  for  $i = 1$  to  $n$  do
     $\mathbf{X}_i \sim p$ 
     $\hat{\mu}_j \leftarrow \hat{\mu}_j + h(\mathbf{X}_i, \theta_j)$ 
   $\hat{\mu}_j \leftarrow \hat{\mu}_j/n$ 
deliver  $\hat{\mu}_1, \dots, \hat{\mu}_m$ 

```

---

This algorithm shows the method of common random numbers with the outer loop over the parameter list. It keeps resetting the seed and regenerating the data. The only random numbers used in  $h$  are from  $\mathbf{X}_i$ .

---

**Partial common random numbers**

Sometimes we can take some but not all of the random variables in two simulations to be common. For instance, suppose that we want to simulate how a coffee shop operates. There is a process by which customers arrive and choose what to order. Then another process defines how quickly their order is fulfilled. We might want to compare two or more service processes. Perhaps the shop adds one more barista at peak hours, or changes how the customers line up, or buys new equipment. Under any of these changes we should be able to run the same sequence of simulated customers through the shop. But there may be no practical way to implement any form of common service times.

In general, we may be trying to find  $\mu = \mathbb{E}(f(\mathbf{X}, \mathbf{Y}) - g(\mathbf{X}, \mathbf{Z}))$  for independent inputs  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{Z}$ . In the coffee shop example,  $\mathbf{X}$  drives the customer arrivals while  $\mathbf{Y}$  (or  $\mathbf{Z}$ ) determines their service times conditionally on the set of arrival times. We can use

$$\hat{\mu}_{\text{ind}} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i, \mathbf{Y}_i) - g(\widetilde{\mathbf{X}}_i, \mathbf{Z}_i)$$

where  $\mathbf{X}_i$ ,  $\widetilde{\mathbf{X}}_i$ ,  $\mathbf{Y}_i$  and  $\mathbf{Z}_i$  are mutually independent. To make a more accurate comparison we would rather have  $\widetilde{\mathbf{X}}_i = \mathbf{X}_i$ . Then we use

$$\hat{\mu}_{\text{com}} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i, \mathbf{Y}_i) - g(\mathbf{X}_i, \mathbf{Z}_i)$$

This is only a ‘partial common random numbers’ algorithm because some but not all of the inputs are common.

**Example 8.1** (Coupling Poisson variables and processes). Suppose that  $X \sim \text{Poi}(\mu)$  and  $Y \sim \text{Poi}(\eta)$  with  $0 < \mu < \eta$ . We can sample  $X$  and  $Y$  by inversion from a common random variable  $U \sim \mathbf{U}(0, 1)$  and they will be closely coupled. We can also simulate  $X \sim \text{Poi}(\mu)$ ,  $Z \sim \text{Poi}(\eta - \mu)$ , and take  $Y = X + Z$ . This second approach does not generate quite as close a connection between  $X$  and  $Y$  but it underlies a useful generalization to Poisson processes.

Let  $\lambda_j \geq 0$  for  $j = 1, 2$  be two intensity functions on  $[0, T]$  with corresponding cumulative intensity functions  $\Lambda_j(t) = \int_0^t \lambda_j(t) dt$ . We can sample these two processes via  $T_{i,j} = \Lambda_j^{-1}(\Lambda_j(T_{i-1,j}) + E_i)$ ,  $j = 1, 2$ , using the common random numbers  $E_i \stackrel{\text{iid}}{\sim} \text{Exp}(1)$ .

The processes  $T_{i,1}$  and  $T_{i,2}$  are simulated from common random numbers but they won't have any common event times. When common event times are desired, we can proceed as follows. We define  $\underline{\lambda}(t) = \min(\lambda_1(t), \lambda_2(t))$  and  $\lambda_j^*(t) = \lambda_j(t) - \underline{\lambda}(t)$  for  $j = 1, 2$ . These have cumulative intensities  $\underline{\Lambda}$  and  $\Lambda_j^*$ , respectively, and they generate Poisson process realizations  $\underline{T}_i$  for  $i = 1, \dots, \underline{N}$  and  $T_j^*$  for  $i = 1, \dots, N_j^*$ . Now we take

$$\begin{aligned} \{T_{1,1}, \dots, T_{N_1,1}\} &= \{\underline{T}_1, \dots, \underline{T}_{\underline{N}}\} \cup \{T_{1,1}^*, \dots, T_{N_1^*,1}^*\}, \quad \text{and} \\ \{T_{1,2}, \dots, T_{N_2,2}\} &= \{\underline{T}_1, \dots, \underline{T}_{\underline{N}}\} \cup \{T_{1,2}^*, \dots, T_{N_2^*,2}^*\}. \end{aligned}$$

If necessary, we sort the points of each process. These processes share  $\underline{N}$  common event times while having  $N_j^*$  unshared event times each. Anderson and Higham (2012) use this method to couple multilevel simulations of continuous time Markov chains.

## Derivative estimation

An extreme instance of the value of common random numbers arises in estimating a derivative. Suppose that  $\mu(\theta) = \mathbb{E}(h(\mathbf{X}, \theta))$  and that we want to estimate  $\mu'(\theta_0) = d\mu/d\theta|_{\theta=\theta_0}$ . We assume that  $h(\mathbf{x}, \theta)$  is well behaved enough to satisfy

$$\frac{d}{d\theta} \int h(\mathbf{x}, \theta) p(\mathbf{x}) d\mathbf{x} = \int \frac{\partial}{\partial \theta} h(\mathbf{x}, \theta) p(\mathbf{x}) d\mathbf{x}$$

at  $\theta = \theta_0$ . If we can compute the needed partial derivative, then we can take

$$\hat{\mu}'(\theta_0) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} h(\mathbf{X}_i, \theta)$$

for  $\mathbf{X}_i \sim p$ . Otherwise, we may need to use divided differences, such as the forward or centered estimators,

$$\begin{aligned} \hat{\mu}'_F(\theta_0) &= \frac{1}{n} \sum_{i=1}^n \frac{h(\mathbf{X}_i, \theta_0 + \epsilon) - h(\mathbf{X}_i, \theta_0)}{\epsilon}, \quad \text{or} \\ \hat{\mu}'_C(\theta_0) &= \frac{1}{n} \sum_{i=1}^n \frac{h(\mathbf{X}_i, \theta_0 + \epsilon) - h(\mathbf{X}_i, \theta_0 - \epsilon)}{2\epsilon}, \end{aligned}$$

respectively, for some small  $\epsilon > 0$ .

Using common random variables we can take a very small  $\epsilon > 0$ , limited only by numerical stability of the required differences. By contrast, with independent random variables, the variance would be

$$\frac{\text{Var}(h(\mathbf{X}, \theta_0)) + \text{Var}(h(\mathbf{X}, \theta_0 + \epsilon))}{n\epsilon^2}$$

leading to certain failure as  $\epsilon \rightarrow 0$ .

If we cannot use common random numbers then there is a bias-variance tradeoff in choosing the optimal  $\epsilon$  given the sample size  $n$ . We can sketch the result using Taylor series centered at  $\theta_0$  for each of  $\theta_0 + \epsilon$  and  $\theta_0 - \epsilon$ . If  $h$  has three partial derivatives with respect to  $\theta$  then

$$h(\mathbf{X}, \theta_0 \pm \epsilon) = h(\mathbf{X}, \theta_0) \pm \epsilon \frac{\partial}{\partial \theta} h(\mathbf{X}, \theta_0) + \frac{\epsilon^2}{2} \frac{\partial^2}{\partial \theta^2} h(\mathbf{X}, \theta_0) \pm \frac{\epsilon^3}{6} \frac{\partial^3}{\partial \theta^3} h(\mathbf{X}, \theta_{\pm})$$

where  $\theta_{\pm}$  is between  $\theta_0$  and  $\theta_0 \pm \epsilon$  and may depend on  $\mathbf{X}$ . Therefore

$$\frac{h(\mathbf{X}, \theta_0 + \epsilon) - h(\mathbf{X}, \theta_0 - \epsilon)}{2\epsilon} = \frac{\partial}{\partial \theta} h(\mathbf{X}, \theta_0) + O_p(\epsilon^2).$$

The result is that the bias in  $\hat{\mu}'_C(\theta_0)$  is  $O_p(\epsilon^2)$  while the variance is  $O(1/(n\epsilon^2))$ . The optimal tradeoff has  $\epsilon \propto n^{-1/6}$  with a mean squared error of  $O(n^{-1/3})$ . Some references on page 36 of the end notes give more information on estimating derivatives.

## 8.7 Conditioning

Sometimes we can do part of the problem in closed form, and then do the rest of it by Monte Carlo or some other numerical method. Suppose for example that we want to find  $\mu = \int_0^1 \int_0^1 f(x, y) dx dy$  where  $f(x, y) = e^{g(x)y}$ . It is easy to integrate out  $y$  for fixed  $x$ , yielding  $h(x) = (e^{g(x)} - 1)/g(x)$ . Then we have a one dimensional problem, which may be simpler to handle. If  $g$  is complicated, such as  $g(x) = \sqrt{5/4 + \cos(2\pi x)}$ , then we cannot easily integrate  $x$  out of  $h(x)$ . Nor, it seems, can we integrate  $f(x, y)$  over  $x$  for fixed  $y$  in closed form.

In general, suppose that  $\mathbf{X} \in \mathbb{R}^k$  and  $\mathbf{Y} \in \mathbb{R}^{d-k}$  are random vectors and that we want to estimate  $\mathbb{E}(f(\mathbf{X}, \mathbf{Y}))$ . The natural estimate is  $\hat{\mu} = (1/n) \sum_{i=1}^n f(\mathbf{X}_i, \mathbf{Y}_i)$  where  $(\mathbf{X}_i, \mathbf{Y}_i) \in \mathbb{R}^d$  are independent samples from the joint distribution of  $(\mathbf{X}, \mathbf{Y})$ . Now let  $h(\mathbf{x}) = \mathbb{E}(f(\mathbf{X}, \mathbf{Y}) | \mathbf{X} = \mathbf{x})$ . We might also estimate  $\mu$  by

$$\hat{\mu}_{\text{cond}} = \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i) \tag{8.24}$$

where  $\mathbf{X}_i$  are independently sampled from the distribution of  $\mathbf{X}$ . The justification for the method is that  $\mathbb{E}(f(\mathbf{X}, \mathbf{Y})) = \mathbb{E}(\mathbb{E}(f(\mathbf{X}, \mathbf{Y}) | \mathbf{X})) = \mathbb{E}(h(\mathbf{X}))$ . The function  $h(\cdot)$  gives the conditional mean of  $\mathbf{Y}$  in closed form and then we complete the job by Monte Carlo sampling. The method is called **conditioning**, or **conditional Monte Carlo**, for obvious reasons. The main requirement for conditioning is that we must be able to compute  $h(\cdot)$ . We also need a method for sampling  $\mathbf{X}$ , but we have that already if we can sample  $(\mathbf{X}, \mathbf{Y})$  jointly.

We easily find that

$$\text{Var}(\hat{\mu}_{\text{cond}}) = \frac{1}{n} \text{Var}(h(\mathbf{X})) = \frac{1}{n} \text{Var}(\mathbb{E}(f(\mathbf{X}, \mathbf{Y}) | \mathbf{X})).$$



Recalling the elementary expression

$$\text{Var}(f(\mathbf{X}, \mathbf{Y})) = \mathbb{E}(\text{Var}(f(\mathbf{X}, \mathbf{Y}) \mid \mathbf{X})) + \text{Var}(\mathbb{E}(f(\mathbf{X}, \mathbf{Y}) \mid \mathbf{X}))$$

it is immediately clear that conditional Monte Carlo cannot have higher variance than ordinary Monte Carlo sampling of  $f$  has and will typically have strictly smaller variance. We summarize that finding as follows:

**Theorem 8.1.** *Let  $(\mathbf{X}, \mathbf{Y})$  have joint distribution  $F$  and let  $f(\mathbf{x}, \mathbf{y})$  satisfy  $\text{Var}(f(\mathbf{X}, \mathbf{Y})) = \sigma^2 < \infty$ . Define  $h(\mathbf{x}) = \mathbb{E}(f(\mathbf{X}, \mathbf{Y}) \mid \mathbf{X} = \mathbf{x})$  for  $(\mathbf{X}, \mathbf{Y}) \sim F$ . Suppose that  $(\mathbf{X}_i, \mathbf{Y}_i) \sim F$ . Then*

$$\text{Var}\left(\frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i)\right) \leq \text{Var}\left(\frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i, \mathbf{Y}_i)\right).$$

Conditioning is a special case of **derandomization**. The function  $f(\mathbf{X}, \mathbf{Y})$  has two sources of randomness,  $\mathbf{X}$  and  $\mathbf{Y}$ . For any given  $\mathbf{x}$  and random  $\mathbf{Y}$  we replace the random value  $f(\mathbf{x}, \mathbf{Y})$  by its expectation  $h(\mathbf{x})$ , removing one of the two sources of randomness. For the function  $f(x, y) = e^{g(x)y}$  at the beginning of this section, derandomization brings a nice, but not overwhelming, variance reduction. See Exercise 8.9.

Conditioning is sometimes called **Rao-Blackwellization** in reference to the Rao-Blackwell theorem in theoretical statistics. In that theorem, the quantity being conditioned on has to obey quite stringent conditions. Those conditions usually don't hold in Monte Carlo applications and, from Theorem 8.1, we don't need them. As a result, the term Rao-Blackwellization is not really descriptive of the way conditioning is used in Monte Carlo sampling.

Even though derandomization by conditioning always reduces variance, it is not always worth doing. We could find our estimate is less efficient if computing  $h$  costs much more than computing  $f$  does. For instance, to average

$$f(\mathbf{x}) = \cos\left(g(x_1) + \sum_{j=1}^d a_j x_j\right)$$

over  $\mathbf{x} \in (0, 1)^d$ , we can derandomize and average

$$\frac{1}{a_d} \left( \sin\left(g(x_1) + \sum_{j=1}^{d-1} a_j x_j + a_d\right) - \sin\left(g(x_1) + \sum_{j=1}^{d-1} a_j x_j\right) \right)$$

over  $(0, 1)^{d-1}$  instead. We have reduced the variance but will have nearly doubled the cost, if evaluating  $\sin(\cdot)$  is the most expensive part of computing  $f$ . Derandomizing  $d - 1$  times would leave us with a one dimensional integrand that requires  $2^{d-1}$  sinusoids to evaluate.

**Example 8.2** (Hit or miss). Let  $C = \{(x, y) \mid a \leq x \leq b, 0 \leq y \leq f(x)\}$ . Suppose that  $f(x) \leq c$  holds for  $a \leq x \leq b$ . Then the hit or miss Monte Carlo

estimate of  $\mathbf{vol}(C)$  is

$$\widehat{\mathbf{vol}}(C) = \frac{c(b-a)}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \leq f(X_i)}$$

where  $(X_i, Y_i) \sim \mathbf{U}([a, b] \times [0, c])$  are independent for  $i = 1, \dots, n$ . Now  $h(x) = \mathbb{E}(\mathbb{1}_{Y \leq f(X)} \mid X = x) = f(x)/c$ . Derandomizing hit or miss Monte Carlo by conditioning, yields the estimate

$$\frac{c(b-a)}{n} \sum_{i=1}^n \frac{f(X_i)}{c} = \frac{b-a}{n} \sum_{i=1}^n f(X_i).$$

The result is perhaps the most obvious way to estimate  $\mathbf{vol}(C)$  by Monte Carlo, and it has lower variance than hit or miss. A case could be made for hit or miss when the average cost of determining whether  $Y \leq f(X)$  holds is quite small compared to the cost of precisely evaluating  $f(X)$  itself. But outside of such special circumstances, there is little reason to use hit or miss MC for finding the area under a curve.

Conditioning can be used in combination with other variance reduction methods. The most straightforward way is to apply those other methods to the problem of estimating  $\mathbb{E}(h(\mathbf{X}))$ . The combination of conditioning with stratified and/or antithetic sampling of  $\mathbf{X}$  is thus simple, provided that the distribution of  $\mathbf{X}$  is amenable to stratification or has some natural symmetry that we can exploit in antithetic sampling.

Conditioning brings a dimension reduction in addition to the variance reduction, because the dimension  $k$  of  $\mathbf{X}$  is smaller than the dimension  $d$ , of  $(\mathbf{X}, \mathbf{Y})$ . When  $k$  is very small, then stratification methods or even quadrature can be used to compound the gain from conditioning. The example in §8.8 has  $d = 38$  and  $k = 1$ .

## 8.8 Example: maximum Dirichlet

The gambler Allan Wilson once tabulated the results of 79,800 plays at a roulette table. Those values are given in the column labeled ‘Wheel 1’ in Table 8.3. The wheel on that table had 38 slots, numbered 1 through 36 along with 0 and 00, which we’ll denote by 37 and 38 respectively. The wheel seemed to be imperfect, either due to manufacture or maintenance. The number 19 came up more often than any other.

Suppose that the counts  $\mathbf{C} = (C_1, \dots, C_{38})$  for wheel 1 follow a  $\text{Mult}(N, \mathbf{p})$  distribution with  $N = 79,800$  and  $\mathbf{p} = (p_1, \dots, p_{38})$ . If we adopt a prior distribution with  $\mathbf{p} \sim \text{Dir}(1, \dots, 1)$  then the posterior distribution of  $\mathbf{p}$  given that  $\mathbf{C} = \mathbf{c}$  is  $\text{Dir}(\alpha_1, \dots, \alpha_{38})$  where  $\alpha_j = c_j + 1$ . For this posterior distribution, we would like to know  $\mathbb{P}(p_{19} = \max_{1 \leq j \leq 38} p_j)$ , the probability that number 19 really does come up most often.

Number	Wheel 1	Wheel 2
00	2127	1288
1	2082	1234
13	2110	1261
36	2221	1251
24	2192	1164 <sub>w</sub>
3	2008	1438 <sub>b</sub>
15	2035	1264
34	2113	1335
22	2099	1342
5	2199	1232
17	2044	1326
32	2133	1302
20	1912 <sub>w</sub>	1227
7	1999	1192
11	1974	1278
30	2051	1336
26	1984	1296
9	2053	1298
28	2019	1205
0	2046	1189
2	1999	1171
14	2168	1279
35	2150	1315
23	2041	1296
4	2047	1256
16	2091	1304
33	2142	1304
21	2196	1351
6	2153	1281
18	2191	1392
31	2192	1306
19	2284 <sub>b</sub>	1330
8	2136	1266
12	2110	1224
29	2032	1190
25	2188	1229
10	2121	1320
27	2158	1336
Avg	2100	1279.16

Table 8.3: This table gives counts from two roulette wheels described in Wilson (1965, Appendix E). The best and worst holes, for the customer, are marked with **b** and **w** respectively.

In §5.4 we represented the Dirichlet distribution as normalized independent Gamma random variables. Here we can define  $\mathbf{X} = (X_1, \dots, X_{38})$  where  $X_j \sim \text{Gam}(\alpha_j)$  are independent, and  $p_j = X_j / \sum_{k=1}^{38} X_k$ . Clearly  $p_{19}$  is the largest  $p_j$  if and only if  $X_{19}$  is the largest  $X_j$ . Therefore, we want to find  $\mu = \mathbb{E}(f(\mathbf{X}))$  where

$$f(\mathbf{X}) = \begin{cases} 1, & X_{19} = \max_{1 \leq j \leq 38} X_j \\ 0, & X_{19} < \max_{1 \leq j \leq 38} X_j. \end{cases}$$

A direct Monte Carlo estimate of  $\mu$  proceeds by repeatedly sampling  $\mathbf{X} \in [0, \infty)^{38}$  and averaging  $f(\mathbf{X})$ . Here we condition on  $X_{19}$ . Given that  $X_{19} = x_{19}$ , the probability that  $X_{19}$  is largest is

$$h(x_{19}) = \prod_{j=1, j \neq 19}^{38} G_{\alpha_j}(x_{19}) \quad (8.25)$$

where  $G_\alpha(x) = \int_0^x e^{-y} y^{\alpha-1} dy / \Gamma(\alpha)$  is the CDF of the  $\text{Gam}(\alpha)$  distribution. To find the answer for this roulette wheel, do Exercise 8.10.

By conditioning, we replace  $(1/n) \sum_{i=1}^n f(\mathbf{X}_i)$  where  $X_{ij} \sim \text{Gam}(\alpha_j)$  are independent by  $(1/n) \sum_{i=1}^n h(Y_i)$  where  $Y_i \sim \text{Gam}(\alpha_{19})$  are independent.

Computations for the function  $h(y)$  could, in some instances, underflow. That does not happen for the roulette example, but if we want to get the probability that the apparent worst number is actually the best, the values of  $h$  become very small. Similarly for problems with higher dimensional Dirichlet distributions and more unequal counts, underflow is more likely. Underflow can be mitigated by working with software that computes  $\log(G_{\alpha_j})$  directly. To find the probability that component  $j_0$  is the largest of  $J$  components, we can define  $\tilde{h}(y) = \sum_{j=1, j \neq j_0}^J \log(G_{\alpha_j}(y))$  find  $h^* = \max_{1 \leq i \leq n} \tilde{h}(y_i)$  for the sampled  $y_i$  values and report the answer as  $\exp(h^*)$  times  $(1/n) \sum_{i=1}^n \exp(\tilde{h}_i - h^*)$ .

## 8.9 Control variates

We saw in §8.7 on conditioning how to get a better estimate by doing part of the problem in closed form. Control variates provide another way to exploit closed form results. With control variates we use some other problem, quite similar to our given one, but for which an exact answer is known. The precise meaning of 'similar' depends on how we will use this other problem, and more than one method is given below. As for 'exact', we will mean it literally, but in practice it may just mean known with an error negligible compared to Monte Carlo errors.

Suppose first that we want to find  $\mu = \mathbb{E}(f(\mathbf{X}))$  and that we know the value  $\theta = \mathbb{E}(h(\mathbf{X}))$  where  $h(\mathbf{x}) \approx f(\mathbf{x})$ . Letting  $\hat{\mu} = (1/n) \sum_{i=1}^n f(\mathbf{X}_i)$  and  $\hat{\theta} = (1/n) \sum_{i=1}^n h(\mathbf{X}_i)$  we can estimate  $\mu$  by the **difference estimator**

$$\hat{\mu}_{\text{diff}} = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{X}_i) - h(\mathbf{X}_i)) + \theta = \hat{\mu} - \hat{\theta} + \theta. \quad (8.26)$$

The expected value of  $\hat{\mu}_{\text{diff}}$  is  $\mu$  because  $\mathbb{E}(\hat{\theta}) = \theta$ . The variance of  $\hat{\mu}_{\text{diff}}$  is

$$\text{Var}(\hat{\mu}_{\text{diff}}) = \frac{1}{n} \text{Var}(f(\mathbf{X}) - h(\mathbf{X})).$$

So if  $h$  is similar to  $f$  in the sense that the difference  $f(\mathbf{X}) - h(\mathbf{X})$  has smaller variance than  $f(\mathbf{X})$  has, we will get reduced variance by using  $\hat{\mu}_{\text{diff}}$ .

In this setting  $h(\mathbf{X})$ , the random variable whose mean is known, is the **control variate**. The difference estimator is not the only way to use a control variate. The ratio and product estimators

$$\hat{\mu}_{\text{ratio}} = \hat{\mu} \theta / \hat{\theta}, \quad \text{and} \quad (8.27)$$

$$\hat{\mu}_{\text{prod}} = \hat{\mu} \hat{\theta} / \theta \quad (8.28)$$

respectively, are also used. These estimators are undefined when  $\theta = 0$ , but otherwise they generally converge to  $\mu$  as  $n \rightarrow \infty$ . See Exercise 8.18 for the product estimator. The ratio and product estimators are usually biased because  $\mathbb{E}(\hat{\theta}/\hat{\mu}) \neq \theta/\mu$  and  $\mathbb{E}(\hat{\theta}\hat{\mu}) \neq \theta\mu$  in general. It is possible to generalize the control variate method in very complicated ways. Maybe we could use  $\hat{\mu} \cos(\hat{\theta} - \theta)$  or some more imaginative quantity. But we don't. By far the most common way of using a control variate is through the regression estimator, considered next.

For a value  $\beta \in \mathbb{R}$ , the **regression estimator** of  $\mu$  is

$$\hat{\mu}_{\beta} = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{X}_i) - \beta h(\mathbf{X}_i)) + \beta \theta = \hat{\mu} - \beta(\hat{\theta} - \theta). \quad (8.29)$$

Taking  $\beta = 0$  yields the simple MC estimator  $\hat{\mu}$  and  $\beta = 1$  gives us the difference estimator. The regression estimator is unbiased:  $\mathbb{E}(\hat{\mu}_{\beta}) = \mu$  for all  $\beta$  because  $\mathbb{E}(\hat{\theta}) = \theta$ .

The variance of the regression estimator is

$$\text{Var}(\hat{\mu}_{\beta}) = \frac{1}{n} (\text{Var}(f(\mathbf{X})) - 2\beta \text{Cov}(f(\mathbf{X}), h(\mathbf{X})) + \beta^2 \text{Var}(h(\mathbf{X}))).$$

By differentiating, we find that the best value of  $\beta$  is

$$\beta_{\text{opt}} = \frac{\text{Cov}(f(\mathbf{X}), h(\mathbf{X}))}{\text{Var}(h(\mathbf{X}))} = \frac{\mathbb{E}((h(\mathbf{X}) - \theta)f(\mathbf{X}))}{\mathbb{E}((h(\mathbf{X}) - \theta)^2)},$$

and after some algebra, the resulting minimal variance is

$$\text{Var}(\hat{\mu}_{\beta_{\text{opt}}}) = \frac{\sigma^2}{n} (1 - \rho^2),$$

where  $\rho = \text{Corr}(f(\mathbf{X}), h(\mathbf{X}))$ . In the regression estimator, any control variate that correlates with  $f$  is helpful, even one that correlates negatively.

In practice we don't know  $\beta_{\text{opt}}$  and so we estimate it by

$$\hat{\beta} = \frac{\sum_{i=1}^n (f(\mathbf{X}_i) - \bar{f})(h(\mathbf{X}_i) - \bar{h})}{\sum_{i=1}^n (h(\mathbf{X}_i) - \bar{h})^2},$$

where  $\bar{f} = (1/n) \sum_{i=1}^n f(\mathbf{X}_i)$  and  $\bar{h} = (1/n) \sum_{i=1}^n h(\mathbf{X}_i)$ . Then the regression estimator of  $\mu$  is  $\hat{\mu}_{\hat{\beta}}$ . In general  $\mathbb{E}(\hat{\mu}_{\hat{\beta}}) \neq \mu$ , but this bias is usually small. We postpone study of the bias until later (equation (8.34)) when we consider multiple control variates. The estimated variance of  $\hat{\mu}_{\hat{\beta}}$  is

$$\widehat{\text{Var}}(\hat{\mu}_{\hat{\beta}}) = \frac{1}{n^2} \sum_{i=1}^n (f(\mathbf{X}_i) - \hat{\mu}_{\hat{\beta}} - \hat{\beta}(h(\mathbf{X}_i) - \bar{h}))^2,$$

and a 99% confidence interval is  $\hat{\mu}_{\hat{\beta}} \pm 2.58\sqrt{\widehat{\text{Var}}(\hat{\mu}_{\hat{\beta}})}$ .

The variance with a control variate is  $\sigma^2(1 - \rho^2)/n$  which is never worse than  $\sigma^2/n$  and usually better. Whether the control variate is helpful ultimately depends on how much it costs to use it. Suppose that the total cost of generating  $\mathbf{X}_i$  and then computing  $f(\mathbf{X}_i)$  is, on average,  $c_f$ . Let  $c_h$  be the extra cost incurred by the control variate on average. That includes the cost to evaluate  $h(\mathbf{X}_i)$  but not the cost of sampling  $\mathbf{X}_i$ . We will suppose that the cost to compute  $\hat{\beta}$  is small. If not then  $c_h$  should be increased to reflect it. Control variates improve efficiency when  $(1 - \rho^2)(c_f + c_h) < c_f$ , that is when  $|\rho| > \sqrt{c_h/(c_f + c_h)}$ . For illustration, if  $c_h = c_f$  then we need  $|\rho| > \sqrt{1/2} \doteq 0.71$  in order to benefit from the control variate.

**Example 8.3** (Arithmetic and geometric Asian option). A well known and very effective control variate arises in finance. Let  $f(\mathbf{X}) = \max(0, (1/m) \sum_{k=1}^m S(t_k) - K)$  be the value of an Asian call option, from §6.4, in terms of a geometric Brownian motion  $S(t)$  generated from  $\mathbf{X} \sim \mathbf{U}(0, 1)^d$ . Now let  $h(\mathbf{X}) = \max(0, \prod_{k=1}^m S(t_k)^{1/m} - K)$ , be the same option except that the arithmetic average has been replaced by a geometric average. The geometric average has a lognormal distribution. Thus  $\theta$  can be computed by a one dimensional integral with respect to the normal probability density function. The result is the Black-Scholes formula.

A significant advantage of the regression estimator is that it generalizes easily to handle multiple control variates. The potential value is greatest when  $f$  is expensive but is approximately equal to a linear combination of inexpensive control variates.

Suppose that  $\mathbb{E}(h_j(\mathbf{X})) = \theta_j$  are known values for  $j = 1, \dots, J$ . Let  $h(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_J(\mathbf{x}))^\top$  be a vector of functions with  $\mathbb{E}(h(\mathbf{X})) = \theta = (\theta_1, \dots, \theta_J)^\top$ , and let  $\beta = (\beta_1, \dots, \beta_J)^\top \in \mathbb{R}^J$ . The regression estimator for  $J \geq 1$  is

$$\hat{\mu}_{\beta} = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{X}_i) - \beta^\top h(\mathbf{X}_i)) + \beta^\top \theta = \hat{\mu} - \beta^\top \bar{H} + \beta^\top \theta \quad (8.30)$$

where  $\bar{H} = (1/n) \sum_{i=1}^n h(\mathbf{X}_i)$ . As before,  $\mathbb{E}(\hat{\mu}_{\beta}) = \mu$ .

The variance of  $\hat{\mu}_{\beta}$  is  $\sigma_{\beta}^2/n$  where

$$\sigma_{\beta}^2 = \mathbb{E}((f(\mathbf{X}) - \mu - \beta^\top (h(\mathbf{X}) - \theta))^2). \quad (8.31)$$

**Algorithm 8.3** Control variates by regression

---

**given**  $f(\mathbf{x}_i)$ ,  $h_j(\mathbf{x}_i)$ ,  $\theta_j = \mathbb{E}(h_j(\mathbf{X}))$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, J$   
 $Y_i \leftarrow f(\mathbf{x}_i)$ ,  $i = 1, \dots, n$   
 $Z_{ij} \leftarrow h_j(\mathbf{x}_i) - \theta_j$   $i = 1, \dots, n$ ,  $j = 1, \dots, J$  // centering  
MLR  $\leftarrow$  multiple linear regression of  $Y_i$  on  $Z_{ij}$   
 $\hat{\mu}_{\text{reg}} \leftarrow$  estimated intercept from MLR  
 $\text{se} \leftarrow$  intercept standard error from MLR  
**deliver**  $\hat{\mu}$ ,  $\text{se}$

---

This algorithm shows how to use linear regression software to do control variate computation. It is **essential** to center the control variates. It may be necessary to drop one or more control variates, if they are linearly dependent in the sample.

---

To minimize (8.31) with respect to  $\beta$  is a least squares problem and the solution vector  $\beta$  satisfies  $\text{Var}(h(\mathbf{X}))\beta = \text{Cov}(h(\mathbf{X}), f(\mathbf{X}))$ . If the  $J$  by  $J$  matrix  $\text{Var}(h(\mathbf{X}))$  is singular, then one of the  $h_j$  is a linear combination of the other  $J - 1$  control variates. There is no harm in deleting that redundant variate. As a result we can assume that the matrix  $\text{Var}(h(\mathbf{X}))$  is not singular. Then the optimal value of  $\beta$  is

$$\begin{aligned} \beta_{\text{opt}} &= \text{Var}(h(\mathbf{X}))^{-1} \text{Cov}(h(\mathbf{X}), f(\mathbf{X})) \\ &= (\mathbb{E}([h(\mathbf{X}) - \theta][h(\mathbf{X}) - \theta]^\top))^{-1} \mathbb{E}([h(\mathbf{X}) - \theta]f(\mathbf{X})). \end{aligned} \quad (8.32)$$

In applications we ordinarily do not know  $\beta_{\text{opt}}$ . The usual way to estimate it is by replacing expectations by sample averages:

$$\hat{\beta} = \left( \frac{1}{n} \sum_{i=1}^n (h(\mathbf{X}_i) - \bar{H})(h(\mathbf{X}_i) - \bar{H})^\top \right)^{-1} \frac{1}{n} \sum_{i=1}^n (h(\mathbf{X}_i) - \bar{H})f(\mathbf{X}_i). \quad (8.33)$$

Equation (8.33) is the least squares estimate of  $\beta_{\text{opt}}$ .

The usual estimate of  $\mu$  with control variates is  $\hat{\mu}_{\hat{\beta}}$ . The estimated variance is

$$\widehat{\text{Var}}(\hat{\mu}_{\hat{\beta}}) = \frac{1}{n^2} \sum_{i=1}^n (f(\mathbf{x}_i) - \hat{\mu}_{\hat{\beta}} - \hat{\beta}^\top (h(\mathbf{x}_i) - \bar{h}))^2.$$

Both the estimate, and its standard error  $\sqrt{\widehat{\text{Var}}(\hat{\mu}_{\hat{\beta}})}$ , can be computed using standard multiple linear regression software. See Algorithm 8.3. The key insight is to treat  $\mu$  as the intercept in a multiple linear regression relating  $f(\mathbf{X})$  to predictors  $h_j(\mathbf{X}) - \theta_j$ . The regression formula is  $f(\mathbf{X}) \approx \mu + (h(\mathbf{X}) - \theta)^\top \beta$ . It is crucial to subtract  $\theta_j$  from the control variates in order to make  $\mu = \mathbb{E}(f(\mathbf{X}))$  match the regression intercept.

The error of the regression estimator using  $\beta = \hat{\beta}$  is

$$\begin{aligned} \hat{\mu}_{\hat{\beta}} - \mu &= \hat{\mu}_{\hat{\beta}} - \hat{\mu}_{\beta_{\text{opt}}} + \hat{\mu}_{\beta_{\text{opt}}} - \mu \\ &= (\hat{\mu} - \hat{\beta}^\top \bar{H} + \hat{\beta}^\top \theta) - (\hat{\mu} - \beta_{\text{opt}}^\top \bar{H} + \beta_{\text{opt}}^\top \theta) + \hat{\mu}_{\beta_{\text{opt}}} - \mu \end{aligned}$$

$$= (\hat{\beta} - \beta_{\text{opt}})^\top (\theta - \bar{H}) + \hat{\mu}_{\beta_{\text{opt}}} - \mu. \quad (8.34)$$

The first term in (8.34) is the product of two components of mean zero, while the second term is the error in the unknown optimal regression estimator. The second term has mean zero, but the first does not in general, because the expected value of a product is not necessarily the same as the product of the expected values. As a result, the control variate estimator is usually biased.

Although estimating  $\beta$  from the sample data brings a bias, that bias is ordinarily negligible. Each of the factors  $\hat{\beta} - \beta_{\text{opt}}$  and  $\bar{H} - \theta$  is  $O_p(n^{-1/2})$  so their product is  $O_p(n^{-1})$ . The second term  $\hat{\mu}_{\beta_{\text{opt}}} - \mu$  in (8.34) is of larger magnitude  $O_p(n^{-1/2})$ . For large  $n$ , the first term is negligible while the second term is unbiased. On closer inspection, the first term in (8.34) is the sum of  $J$  contributions, so the bias might be regarded as a  $J/n$  term. Ordinarily  $J$  is not large enough to cause us to change our mind about whether the sum of  $J$  terms of size  $O_p(n^{-1})$  is negligible compared to a single  $O_p(n^{-1/2})$  term. Thus, for applications with  $J \ll \sqrt{n}$ , it is common to neglect the bias from using estimated control variate coefficients.

When an unbiased estimator is required, then we can get one by using an estimate of  $\beta_{\text{opt}}$  that is independent of the  $\mathbf{X}_i$  used in  $\hat{\mu}_{\tilde{\beta}}$ . For example  $\tilde{\beta}$  can be computed from (8.33) using only a pilot sample  $\widetilde{\mathbf{X}}_1, \dots, \widetilde{\mathbf{X}}_m \stackrel{\text{iid}}{\sim} p$  independent of the  $\mathbf{X}_i$ . Then  $\hat{\mu}_{\tilde{\beta}}$  can be computed by (8.30) using  $\mathbf{X}_1, \dots, \mathbf{X}_n$  and taking  $\beta = \tilde{\beta}$ . Now  $\mathbb{E}(\hat{\mu}_{\tilde{\beta}}) = \mu$  and

$$\text{Var}(\hat{\mu}_{\tilde{\beta}}) = \mathbb{E}(\text{Var}(\hat{\mu}_{\tilde{\beta}} \mid \widetilde{\mathbf{X}}_1, \dots, \widetilde{\mathbf{X}}_m)) = \frac{1}{n} \mathbb{E}(\sigma_{\tilde{\beta}}^2).$$

If  $f(\mathbf{X}_i)$  and  $h(\mathbf{X}_i)$  have finite fourth moments then  $\tilde{\beta} = \beta_{\text{opt}} + O_p(1/\sqrt{m})$ . Since  $\sigma_{\tilde{\beta}}^2$  is differentiable with respect to  $\beta$  and takes its minimum at  $\beta_{\text{opt}}$  we have  $\sigma_{\tilde{\beta}}^2 = \sigma_{\beta_{\text{opt}}}^2 + O_p(1/m)$ . Exercise 8.20 asks you to allocate computation between the  $m$  pilot observations and the  $n$  followup observations. See page 35 of the end notes for more sophisticated bias removal.

**Example 8.4** (Post-stratification). Suppose that we have strata  $\mathcal{D}_1, \dots, \mathcal{D}_J$  as in §8.4, but instead of a stratified sample, we take  $\mathbf{X}_i \stackrel{\text{iid}}{\sim} p$  for  $i = 1, \dots, n$ . Let  $h_j(\mathbf{x}) = \mathbb{1}\{\mathbf{x} \in \mathcal{D}_j\}$  for  $j = 1, \dots, J$ . The stratum probabilities  $\omega_j \equiv \mathbb{P}(\mathbf{X} \in \mathcal{D}_j)$  are known. Therefore we can use  $h_j(\mathbf{x})$  as a control variate with mean  $\theta_j = \omega_j$ . If we use the stratum indicators as control variates, we get the same estimate  $\hat{\mu}_{\text{strat}}$  as in post-stratification. The corresponding variance estimate is slightly different.

Using control variates multiplies the asymptotic variance of  $\hat{\mu}$  by a factor  $1 - R^2$  where the  $R^2$  is the familiar proportion of variance explained coefficient from linear regression. If  $J = 1$  then  $R^2 = \rho^2$  where  $\rho$  is the correlation of  $f(\mathbf{X})$  and  $h(\mathbf{X})$ .

If the cost of computing  $h$  is high, then the variance reduction from control variates may need to be large in order to make it worthwhile. Let  $c_f$  be the cost



of computing  $f(\mathbf{X})$  including the cost of computing  $\mathbf{X}$ . Let  $c_h$  be the additional cost of computing the vector  $h(\mathbf{X})$  given that we are already committed to computing  $\mathbf{X}$  and  $f(\mathbf{X})$ . If some parts of the  $f$  computation can be saved and reused in computing  $h$ , then the related costs should be included in  $c_f$  but not in  $c_h$ . The cost of computing  $\hat{\beta}$  has an  $O(J^3)$  term and an  $O(nJ^2)$  term. We suppose that the part that grows proportionally to  $n$  is included in  $c_h$  unless somehow it was needed for computing  $f(\mathbf{X})$ . We also suppose that  $J \ll n$ , so that the  $O(J^3)$  cost may be neglected.

Under the assumptions above, using control variates multiplies the variance by  $1 - R^2$  but multiplies the cost per observation by  $(c_f + c_h)/c_f$ . It improves efficiency if

$$(1 - R^2) \times \frac{c_f + c_h}{c_f} < 1.$$

As a simple special case, suppose that  $c_h = Jc_f$ . After some rearrangement, we find efficiency is improved if  $R^2 > J/(J + 1)$ .

When  $J$  is large it will be very hard to have  $R^2 > J/(J + 1)$ . Multiple control variates may still be worthwhile if they are much less expensive than  $f$ . Suitable control variates include low order polynomials in the components of  $\mathbf{X}$ . These are either inexpensive to compute, or nearly free if we already had to compute them in order to compute  $f(\mathbf{X})$ . When the control variates cost on average  $\epsilon$  times as much as  $f$ , then they improve efficiency if  $R^2 > J\epsilon/(J\epsilon + 1)$ .

## 8.10 Moment matching and reweighting

When we know the value of  $\mathbb{E}(\mathbf{X}) \equiv \theta$  we can use it to improve our estimate of  $\mu = \mathbb{E}(f(\mathbf{X}))$  via control variates as described in §8.9. A simple and very direct alternative approach is to adjust the sample values, setting

$$\widetilde{\mathbf{X}}_i = \mathbf{X}_i + \theta - \bar{\mathbf{X}} \tag{8.35}$$

where  $\bar{\mathbf{X}} = (1/n) \sum_{i=1}^n \mathbf{X}_i$ , and then estimate  $\mu$  by the **moment matching** estimator

$$\hat{\mu}_{\text{mm}} = \frac{1}{n} \sum_{i=1}^n f(\widetilde{\mathbf{X}}_i). \tag{8.36}$$

Moment matching can also be applied to the variance of  $\mathbf{X}$ . Suppose that we know  $\mathbb{E}((\mathbf{X} - \theta)(\mathbf{X} - \theta)^\top) \equiv \Sigma$ , as we would for a simulation based on  $\mathbf{X}_i \sim \mathcal{N}(\theta, \Sigma)$ . Let  $\widehat{\Sigma} = (1/n) \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top$  be the sample variance matrix, and suppose that  $\widehat{\Sigma}$  has full rank, as it will for large enough  $n$ , if  $\Sigma$  has full rank. We can then set

$$\widetilde{\mathbf{X}}_i = \theta + \Sigma^{1/2} \widehat{\Sigma}^{-1/2} (\mathbf{X}_i - \bar{\mathbf{X}})$$

and use (8.36).

In financial applications a multiplicative form of moment matching is commonly used replacing geometric Brownian motion sample paths  $X_i(t)$  by

$$\tilde{X}_i(t) = X_i(t) \times \frac{\mathbb{E}(X_i(t))}{\bar{X}(t)}, \quad \text{where} \quad \bar{X}(t) = \frac{1}{n} \sum_{i=1}^n X_i(t).$$

An analysis in Boyle et al. (1997) shows that moment matching is asymptotically like using the known moments in control variates but with a non-optimal value for the coefficient  $\beta$ .

It is harder to get confidence intervals for moment matching estimators. The  $n$  values  $\tilde{X}_i$  are no longer independent. To get a variance estimate we can repeat the computation  $K$  times independently getting  $\hat{\mu}_{\text{mm},1}, \dots, \hat{\mu}_{\text{mm},K}$ , and then use

$$\hat{\mu}_{\text{mm}} = \frac{1}{K} \sum_{k=1}^K \hat{\mu}_{\text{mm},k}, \quad \text{and}$$

$$\widehat{\text{Var}}(\hat{\mu}_{\text{mm}}) = \frac{1}{K(K-1)} \sum_{k=1}^K (\hat{\mu}_{\text{mm},k} - \hat{\mu}_{\text{mm}})^2.$$

The pooled estimate  $\hat{\mu}_{\text{mm}}$  ordinarily has a small bias.

Despite their lesser accuracy and greater complexity, a motivation to use moment matching arises in financial valuation, where the expectations correspond to various prices. There one reasons that the Monte Carlo must reproduce certain known prices, in order to be credible. If one decides to buy (or sell) securities at a price determined by a Monte Carlo model that is higher (respectively lower) than the market price, then an adversarial trader could exploit that difference.

Another way to meet the goal of moment matching is to reweight the sample. We can replace the equal weight estimator by

$$\sum_{i=1}^n w_i f(\mathbf{X}_i) \tag{8.37}$$

using the same carefully chosen weights  $w_i$  for each function  $f$ . The weights should satisfy  $\sum_{i=1}^n w_i \mathbf{X}_i = \theta$  in the case of (8.35) above. They should also satisfy  $\sum_{i=1}^n w_i = 1$ .

It turns out that control variate estimates of  $\mu$  already take the form (8.37). Suppose that the vector  $h$  of control variates has  $\mathbb{E}(h(\mathbf{X})) = \theta \in \mathbb{R}^J$ . The case (8.35) simply has  $h(\mathbf{X}) = \mathbf{X}$ . Then the estimator (8.33) of  $\beta$  takes the form

$$\hat{\beta} = \sum_{i=1}^n S_{HH}^{-1} (h(\mathbf{X}_i) - \bar{H}) f(\mathbf{X}_i)$$

for  $S_{HH}^{-1} = \sum_{i=1}^n (h(\mathbf{X}_i) - \bar{H})(h(\mathbf{X}_i) - \bar{H})^\top$ . As a result the control variate estimator is

$$\hat{\mu}_\beta = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i) - \hat{\beta}^\top (\bar{H} - \theta) = \sum_{i=1}^n w_i f(\mathbf{X}_i), \quad \text{for}$$

$$w_i = \frac{1}{n} - (h(\mathbf{X}_i) - \bar{H})^\top S_{HH}^{-1}(\bar{H} - \theta).$$

One slim advantage of moment matching over control variates is that it will automatically obey some natural constraints. For example, if  $f(\mathbf{x}) = \exp(\mathbf{x})$  then we know that  $\mathbb{E}(f(\mathbf{X}))$  cannot be negative. It is possible for control variates to supply a negative estimate for such a quantity that must be positive. By contrast, we can be sure that  $\hat{\mu}_{\text{mm}}$  is not negative when  $f(\mathbf{x}) \geq 0$  always holds. Some methods to find non-negative weights with  $\sum_i w_i h(\mathbf{X}_i) = \theta$  and  $\sum_i w_i = 1$  (when they exist) are describe on page 38 of the end notes.

Moment matching and related methods allow one to bake in certain desirable properties of the sample points  $\tilde{\mathbf{X}}_i$ . Their main attraction arises when those properties are important enough to give up on some estimation accuracy and simplicity of forming confidence intervals.

## Chapter end notes

There is a large literature on variance reduction methods. For surveys, see Wilson (1984) and L'Ecuyer (1994).

Antithetic sampling was introduced by Hammersley and Morton (1956). Some generalizations of antithetic sampling are considered in Chapter 10.

Stratification is a classic survey sampling method. See Cochran (1977), for issues of variance estimation and also for design of strata. It is not just stratification. Antithetics, control variates and importance sampling (Chapter 9) have direct antecedents in the survey sampling literature.

The difference estimator is also commonly used in classical quadrature methods. Suppose that both  $h(\mathbf{x})$  and  $f(\mathbf{x})$  are unbounded, but  $f(\mathbf{x}) - h(\mathbf{x})$  is bounded, and  $\int h(\mathbf{x}) d\mathbf{x}$  is known. Then it often pays to use numerical quadrature on  $f - h$  and add in the known integral of  $h$ . For Monte Carlo sampling it will ordinarily be better to use regression estimator. However for quasi-Monte Carlo and randomized quasi-Monte Carlo (Chapters 15 through 17) we may prefer the difference estimator if  $f - h$  is then of bounded variation.

The ratio and product estimators are not available when  $\theta = \mathbb{E}(h(\mathbf{X})) = 0$ . Their typical applications are in problems where  $h(\mathbf{x}) > 0$ . The reason that complicated nonlinear control variates are seldom used is that, in large samples, they are almost equivalent to the regression estimator, which is simple to use. See Glynn and Whitt (1989).

The regression estimator for control variates has a mildly annoying bias. Avramidis and Wilson (1993) describe a way to get rid of it. They split the sample into  $m \geq 2$  subsets of equal size and arrange that each coefficient estimate  $\hat{\beta}$  is always applied to points independent of it. The result is an unbiased estimate of  $\mu$  using control variates. When  $m \geq 3$  they are also able to get an unbiased estimate of  $\text{Var}(\hat{\mu})$ .

Kahn and Marshall (1953) make an early mention of the method of common random numbers, referring to it as correlation of samples. They liken it to

pairing and blocking which had long been an important part of the design of physical experiments.

Lunney and Anderson (2009) use Monte Carlo methods to measure the power of the content uniformity test under some alternatives with non-normally distributed data.

Asmussen and Glynn (2007, Chapter VII) cover Monte Carlo estimation of derivatives. They include many algorithms of varying complexity for the case where  $\mathbf{X}$  is a process and  $\theta$  is a parameter of that process. Burgos and Giles (2012) look at multilevel Monte Carlo for estimation of derivatives.

Hesterberg and Nelson (1998) explore the use of control variates for quantile estimation. For random pairs  $(X_i, Y_i)$  one or more known quantiles of the  $X$  distribution can be used as control variates for  $\alpha$  quantile of the  $Y$  distribution. The most direct approach is to estimate  $\mathbb{E}(\mathbb{1}_{Y \leq y})$  using  $\mathbb{1}_{X \leq x_1}, \dots, \mathbb{1}_{X \leq x_s}$  as control variates, and estimate the  $\alpha$  quantile of  $Y$  to be the value  $y$  for which  $\widehat{\mathbb{E}}(\mathbb{1}_{Y \leq y}) = \alpha$ . They consider using a small number of values  $x_j$  at or near the  $\alpha$  quantile of  $X$ . Extreme variance reductions are hard to come by because it is hard to find regression variables that are extremely predictive of a binary value like  $\mathbb{1}_{Y \leq y}$ .

Barraquand (1995) and Duan and Simonato (1998) use some moment matching methods on sample paths of geometric Brownian motion. Cheng (1985) gives an algorithm to generate  $n$  random vectors from the distribution  $\mathcal{N}(0, I_p)$  conditionally on their sample mean being  $\mu$  and sample covariance being  $\Sigma$ . In that approach the constraints are built in to the sample generation rather than imposed by transformation afterwards. Pullin (1979) had earlier done this for samples from  $\mathcal{N}(0, 1)$ .

### Common random numbers with randomness in $h$

Here we allow the function  $h(\mathbf{X}, \theta)$  in common random numbers to generate further random numbers. We assume that the number of further random numbers  $h(\mathbf{X}, \theta)$  uses depends on both  $\mathbf{X}$  and  $\theta$ . If instead  $h()$  always takes the same number of uniform random numbers we can include them in  $\mathbf{X}$  and proceed as if  $h$  does not generate random variables.

We begin with Algorithm 8.1 and we assume that all the dependence we wanted to incorporate comes through the shared  $\mathbf{X}_i$  and so  $h(\mathbf{X}_i, \theta_j)$  for  $1 \leq i \leq n$  and  $1 \leq j \leq m$  are conditionally independent given  $\mathbf{X}_1, \dots, \mathbf{X}_n$ .

In Algorithm 8.1, an  $h$  that consumes random numbers would advance the random number stream by some number of positions and thereby change  $\mathbf{X}_2, \dots, \mathbf{X}_n$ . The differences  $\hat{\mu}_j - \hat{\mu}_k$  would still be unbiased estimates of  $\mu_j - \mu_k$ . The additional randomness in  $h$  would increase the variance of  $\hat{\mu}_j - \hat{\mu}_k$ , reducing the gain from common random numbers. Because the  $n$  sample differences  $h(\mathbf{X}_i, \theta_j) - h(\mathbf{X}_i, \theta_k)$  going into that estimate are still statistically independent, our confidence intervals remain reliable.

The challenge with Algorithm 8.1 starts when we consider changing our parameter list  $\theta_1, \dots, \theta_m$ , perhaps by adding  $\theta_{m+1}, \dots, \theta_{m+k}$ . To account for changing parameters it is less ambiguous to write  $\hat{\mu}(\theta_j)$  instead of  $\hat{\mu}_j$ . When

$h$  consumes random numbers, then changing the parameter list  $\theta_1, \dots, \theta_m$ , can change  $\mathbf{X}_2$  and all subsequent  $\mathbf{X}_i$  that Algorithm 8.1 uses.

If we add new parameters  $\theta_{m+1}, \dots, \theta_{m+k}$  to our list and rerun Algorithm 8.1 for all  $m+k$  parameter values, then it is likely that all of our old estimates  $\hat{\mu}(\theta_j)$  for  $j \leq m$  will have changed. The estimates still reflect common random numbers. But we might have preferred those old values to remain fixed.

A very serious problem (i.e., an error) arises when we store the values  $\hat{\mu}(\theta_j)$  for  $j = 1, \dots, m$ , and then instead of re-running Algorithm 8.1 on the whole list, we just run it on the list of  $k$  new parameter values. Then the new estimates  $\hat{\mu}(\theta_{m+1}), \dots, \hat{\mu}(\theta_{m+k})$  will not have been computed with the same  $\mathbf{X}_i$  that the old ones used. Even though that algorithm starts by setting the seed, synchronization will already be lost for  $\mathbf{X}_2$  because  $h$  generated random numbers. We would have lost the accuracy advantage of common random numbers for comparisons involving one of the first  $m$  parameters and one of the last  $k$ . Also, some of the random numbers used to generate  $\mathbf{X}_i$  for the first set of parameters may end up incorporated into both  $\mathbf{X}_i$  and  $\mathbf{X}_{i+1}$  (or some other set of variables) for the second set. The differences  $h(\mathbf{X}_i, \theta_r) - h(\mathbf{X}_i, \theta_s)$ ,  $i = 1, \dots, n$  would not be independent if  $r \leq m < s$ . So we would get unreliable standard deviations for those comparisons.

To be sure that  $\mathbf{X}_i$  is the same for all sets  $\Theta$ , we should not let  $h$  use the same stream of random numbers that  $\mathbf{X}_i$  are generated from. Even giving  $h$  its own stream of random numbers leaves us with synchronization problems. Computing  $h(\mathbf{X}_1, \theta_{m+1})$  would affect the random numbers that  $h(\mathbf{X}_2, \theta_1)$  sees.

If we want  $\hat{\mu}(\theta_j)$  to be unaffected by the other  $\theta_k \in \Theta$ , then the solution is to give  $h$  a different random number stream for each value of  $\theta$  that we use. One approach is to maintain a lookup table of  $\theta$ 's and their corresponding seeds. Another is to hash the value of  $\theta_j$  into a seed (or a stream identifier) for  $h$  to use. If each  $\theta_j$  gets its own stream, as in L'Ecuyer et al. (2002) then the common seed for all of those streams gets set at the beginning of the algorithm. If each  $\theta_j$  is hashed into its own seed for a random number generator like the Mersenne Twister (Matsumoto and Nishimura, 1998), then seeded copies of that generator should be created at the beginning of the algorithm. Now each  $\hat{\mu}(\theta)$  is a reproducible function of  $\theta$  and  $n$  and the seeds used.

Now consider Algorithm 8.2 where  $h$  consumes random numbers. For each  $\theta_j$  it sets the seed then does a Monte Carlo sample. It is more fragile than Algorithm 8.1. That algorithm still works if we run all  $\theta_j$  at once and do not mind having  $\hat{\mu}(\theta_j)$  depend on the set of other  $\theta$  values in  $\Theta$ . For Algorithm 8.2, if  $h$  generates random numbers then  $\mathbf{X}_i$  for  $i \geq 2$  will vary with  $\theta_j$  and we lose synchronization. To ensure that  $\mathbf{X}_i$  are really common we should not let  $h$  use the stream that we use to generate  $\mathbf{X}_i$ . To keep each  $\hat{\mu}(\theta)$  unaffected by changes to the set  $\Theta$ , we should once again give every value of  $\theta$  its own stream, and set the seed for that stream at the same time the  $\mathbf{X}$  stream's seed is set.

### Alternative reweightings

As described in §8.10, control variates reweight the sample values but might include some negative weights. We would prefer to have weights  $w_i$  that satisfy

$$w_i \geq 0, \quad \sum_{i=1}^n w_i = 1, \quad \text{and} \quad \sum_{i=1}^n w_i h(\mathbf{X}_i) = \theta. \quad (8.38)$$

Ideally, the weights  $w_i$  should be as close to  $1/n$  as possible, subject to the constraints in (8.38). Then we may estimate  $\mu$  by

$$\hat{\mu}_w = \sum_{i=1}^n w_i f(\mathbf{X}_i).$$

Constraints (8.38) cannot always be satisfied. If  $\min_{1 \leq i \leq n} h_j(\mathbf{X}_i) > \theta_j$  then there is no way to satisfy (8.38). More generally, if  $\theta$  is outside the convex hull of  $\{h(\mathbf{X}_1), \dots, h(\mathbf{X}_n)\}$ , so that there exists a hyperplane with  $\theta \in \mathbb{R}^J$  on one side and all of  $h(\mathbf{X}_i)$  on the other, then (8.38) cannot be satisfied. If  $\theta$  is outside the convex hull of  $h(\mathbf{X}_i)$  then maybe  $n$  is too small, or  $J$  is too large, or the functions  $h_j$  are poorly chosen.

If a solution to (8.38) exists then there is an  $n - J - 1$  dimensional family of solutions. To choose weights in this family we need to choose a measure of their distance from  $(1/n, \dots, 1/n)$ . One such way is to maximize the log empirical likelihood  $-\sum_{i=1}^n \log(nw_i)$  subject to (8.38). A second way is to maximize the entropy  $-\sum_{i=1}^n w_i \log(w_i)$  subject to (8.38). Both of these criteria favor  $w_i$  that are nearly equal. Each of them leads to weighted Monte Carlo estimates with the same asymptotic variance that  $\hat{\mu}_{\hat{\beta}}$  has.

If we maximize the empirical likelihood, then a Lagrange multipliers argument yields

$$w_i^{\text{EL}} = \frac{1}{n} \frac{1}{1 + \lambda^\top (h(\mathbf{X}_i) - \theta)}$$

where the Lagrange multiplier  $\lambda \in \mathbb{R}^J$  satisfies

$$\sum_{i=1}^n \frac{h(\mathbf{X}_i) - \theta}{1 + \lambda^\top (h(\mathbf{X}_i) - \theta)} = 0.$$

(Owen, 2001, Chapter 3) gives details including computation of  $\lambda$ . Empirical likelihood and entropy are two members in a family of non-negative weighting methods. For Monte Carlo applications where non-negativity is not needed, regression based control variates are simpler to use.

## Exercises

### Antithetics

**8.1.** Given  $\epsilon > 0$ , construct an increasing function  $f(x)$  on  $0 \leq x \leq 1$  such that

$$0 \geq \text{Corr}(f(X), f(1 - X)) > -\epsilon$$

for  $X \sim \mathbf{U}(0, 1)$ .

**8.2.** Find an example for the following set of conditions, or prove that it is impossible to do so:  $0 < \text{Var}(\hat{\mu}_{\text{anti}}) < \text{Var}(\hat{\mu}) = \infty$ . Here  $\hat{\mu}$  is ordinary Monte Carlo sampling with a finite even number  $n \geq 2$  of function values and  $\hat{\mu}_{\text{anti}}$  is antithetic sampling with  $n/2$  pairs. If this is possible, then  $\hat{\mu}_{\text{anti}}$  has an infinite efficiency relative to ordinary Monte Carlo without having 0 variance.

**8.3.** Show that the correlation in antithetic sampling is

$$\rho = \frac{\sigma_{\mathbf{E}}^2 - \sigma_{\mathbf{O}}^2}{\sigma_{\mathbf{E}}^2 + \sigma_{\mathbf{O}}^2},$$

in the notation of §10.2.

**8.4** (Antithetic sampling and spiky integrands). Here we investigate what happens with antithetic sampling and a spiky function. We will use

$$f(x) = \begin{cases} 0, & 0 < x \leq 0.9 \\ 100, & 0.9 < x \leq 0.91 \\ 0, & 0.91 < x < 1 \end{cases}$$

for  $X \sim \mathbf{U}(0, 1)$  as a prototypical spiky function.

- Determine whether antithetic sampling is helpful, harmful, or neutral for the example  $f$ . You may do this by finding the variance of  $\hat{\mu}$  under IID and under antithetic sampling using the same sample size. You may find the variances either theoretically or from a large enough simulation.
- Explain your findings from the part above, in terms of the even and odd parts of  $f$ .
- Construct a spiky function for which you would have reached a very different conclusion about the effectiveness of antithetic sampling.

## Stratification

**8.5.** Prove equation (8.15), which represents the variance reduction from proportional allocation in terms of a correlation between  $f$  and the within stratum mean of  $f$ .

**8.6.** Equation (8.14) expresses the sampling variance of the stratified estimator and the ordinary MC estimator in terms of between and within variances  $\sigma_{\mathbf{B}}^2$  and  $\sigma_{\mathbf{W}}^2$ . Given  $f$  with  $\int f(\mathbf{x})^2 d\mathbf{x} < \infty$  show how to construct functions  $f_{\mathbf{B}}(\mathbf{x})$  and  $f_{\mathbf{W}}(\mathbf{x})$  such that  $f(\mathbf{x}) = f_{\mathbf{B}}(\mathbf{x}) + f_{\mathbf{W}}(\mathbf{x})$  with  $\int f_{\mathbf{W}}(\mathbf{x}) d\mathbf{x} = \int f_{\mathbf{B}}(\mathbf{x}) f_{\mathbf{W}}(\mathbf{x}) d\mathbf{x} = 0$  and  $\int f_{\mathbf{B}}(\mathbf{x}) d\mathbf{x} = \int f(\mathbf{x}) d\mathbf{x} = \mu$ ,  $\int f_{\mathbf{W}}(\mathbf{x})^2 d\mathbf{x} = \sigma_{\mathbf{W}}^2$ ,  $\int (f_{\mathbf{B}}(\mathbf{x}) - \mu)^2 d\mathbf{x} = \sigma_{\mathbf{B}}^2$ , and for which the stratified sampling estimate of the mean of  $f_{\mathbf{B}}$  has variance zero.

## Stratified Brownian motion

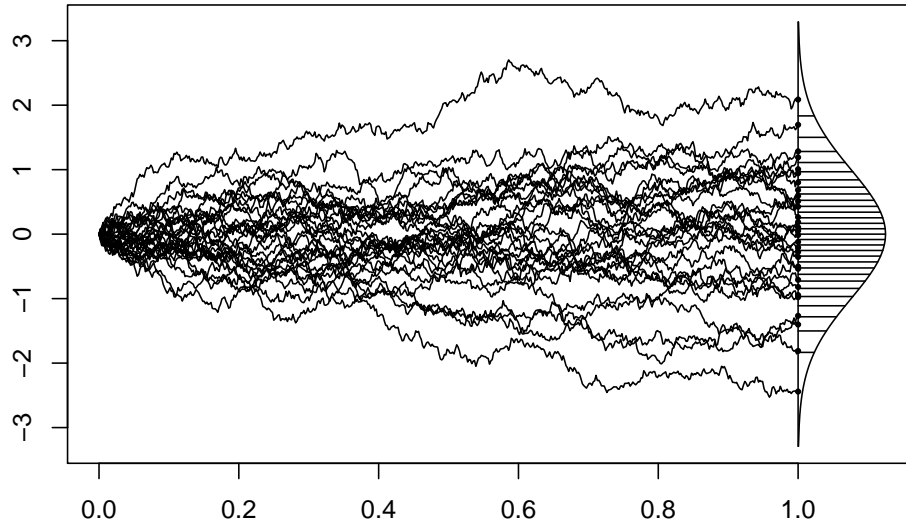


Figure 8.6: This figure shows 30 sample paths of standard Brownian motion  $B(\cdot) \sim \text{BM}(0,1)$  on  $\mathcal{T} = [0,1]$ . They are stratified on  $B(1) \sim \mathcal{N}(0,1)$ . See Exercise 8.7. Also shown is the  $\mathcal{N}(0,1)$  density function partitioned into 30 equi-probable intervals.

**8.7** (Stratified Brownian motion). Here we investigate stratified Brownian motion, as shown in Figure 8.6. Let path  $i$  at time  $t$  take the value  $B_i(t)$  for  $i = 1, \dots, N$  and  $t \in \{1/M, 2/M, \dots, 1\}$ . To stratify standard Brownian motion on its endpoint, we take  $B_i(1) = \Phi^{-1}((i - U_i)/N)$  for independent  $U_1, \dots, U_N \sim \mathbf{U}(0,1)$ . Points  $B_i(j/M)$ , for  $j = 1, \dots, M - 1$  are then sampled conditionally on  $B_i(1)$ . See §xxx.

- a) Write a function to generate stratified standard Brownian motion. It should take arguments  $M, N \in \mathbb{N}$ , and  $i \in \{1, \dots, N\}$ . It should produce the sample path of stratified  $B_i(t)$  at  $t = j/M$  for  $j = 1, \dots, M$ . Describe how you sampled the path  $B_i(\cdot)$ , conditionally on  $B_i(1)$ , with enough detail to make it clear that your method is correct. Turn in your code with comments. [Note: if you prefer, you may instead write the function to generate and return all  $N$  paths  $i = 1, \dots, N$  at once.]
- b) Generalize your function to generate stratified Brownian motion with drift  $\delta \in \mathbb{R}$  and volatility  $\sigma \geq 0$  on the interval  $\mathcal{T} = [0, T]$  for  $T > 0$ . As before the value of  $B(T)$  is stratified. Explain how your generalization works, and turn in your code. You may either pass the new arguments  $\delta$ ,  $\sigma$ , and  $T$  into a generalized version of your previous function, or you may write a wrapper function that calls your previous function and modifies its output



to take account of  $\delta$ ,  $\sigma$  and  $T$ .

- c) Let  $S(\cdot) \sim \text{GBM}(S_0, \delta, \sigma)$  be geometric Brownian motion (§6.4). For  $M = 100$ , let

$$f(S(\cdot)) = \max_{0 \leq j \leq M} S(j/M) - \min_{0 \leq j \leq M} S(j/M).$$

We want  $\mu = \mathbb{E}(f(S(\cdot)))$  for  $\delta = 0.05$ ,  $\sigma = 0.3$ , and  $T = 1$ . For  $N = 1000$  and  $M = 100$  generate two independent stratified Geometric Brownian motions with these parameters. Estimate  $\mu$  and give a 99% confidence interval. [Hint: the two independent stratified samples can be pooled into one stratified sample of  $n = 2N$  paths, with  $J = N$  strata having  $n_j = 2$  for  $j = 1, \dots, N$ .]

The function  $f$  is related to the value of a lookback option whose payoff is equivalent to buying at the minimum and selling at the maximum price in the time interval  $[0, T]$ . As given,  $f$  omits the discount factor  $e^{-\delta T}$  that compensates for waiting until time  $T$  to collect the payoff.

- d) Estimate the variance reduction obtained from stratification. Use  $R$  independent replications of the stratified sampling method on  $n = 2N$  paths, where  $R \geq 300$ . The variance should be compared to that obtained by plain Monte Carlo with  $2N$  paths.
- e) Compare the time required to compute  $2N = 2000$  sample paths of length  $M = 100$  by stratification to that required to compute  $2N$  sample paths of length  $M$  without stratification. Report the details of the hardware, operating system, and the software in which you made the comparison.

**8.8** (Stratification with  $n_j = 1$ ). Consider proportional allocation (see §8.4) in the special case where all the strata have equal probability. Then  $\omega_j = 1/J$  and  $n_j = m$  for  $j = 1, \dots, J$  where the sample size is  $n = mJ$ .

- a) Suppose first that  $m \geq 2$  and let  $s_j^2$  be as given in (8.10). Define  $\bar{s}^2 = (1/J) \sum_{j=1}^J s_j^2$ . Show that the formula for  $\widehat{\text{Var}}(\hat{\mu}_{\text{strat}})$  in (8.10) reduces to  $\bar{s}^2/n$ .
- b) Now suppose that  $m = 1$  and that  $n = J$  is an even number. We saw in §xxx that the stratified sampling estimate  $\hat{\mu}_{\text{strat}}$  is  $\bar{Y} = (1/n) \sum_{i=1}^n Y_i$ , in this setting where  $Y_i = f(X_i)$ . For  $m = 1$  we cannot use equation (8.10) for  $\widehat{\text{Var}}(\hat{\mu}_{\text{strat}})$ . For  $j = 1, \dots, n/2$  let  $\tilde{s}_j^2 = (f(X_{2j-1}) - f(X_{2j}))^2$ , put  $\bar{\tilde{s}}^2 = (2/n) \sum_{j=1}^{n/2} \tilde{s}_j^2$  and let  $\tilde{V} = \bar{\tilde{s}}^2/n$ . Prove that  $\mathbb{E}(\tilde{V}) \geq \text{Var}(\bar{Y})$ .
- c) Suppose now that stratum  $i$  is  $[(i-1)/n, i/n)$ , that  $n$  is very large, and  $f$  has two derivatives on  $[0, 1]$ . Roughly how large will  $\mathbb{E}(\tilde{V})/\text{Var}(\bar{Y})$  be?

## Conditioning

**8.9.** Let  $(X, Y) \sim \mathbf{U}(0, 1)^2$  and put  $f(x, y) = e^{g(x)y}$  for  $g(x) = \sqrt{5/4 + \cos(2\pi x)}$ . Let  $h(x) = (e^{g(x)} - 1)/g(x)$ .

- a) Using  $n = 10^6$  samples estimate the variance of  $f(x, y)$ . Similarly, estimate the variance of  $h(x)$ .
- b) Report the efficiency gain from conditioning assuming that  $f$  and  $h$  cost the same amount of computer time. Then report the efficiency gain taking account of the time it takes to compute both  $f$  and  $h$ . In this case give details of the computing environment that you obtained the results for. Also hand in your source code.
- c) Repeat the two steps above for  $g(x) = \sqrt{1 + \cos(2\pi x)}$  taking special care near  $x = 1/2$ . (Hint: you may need a Taylor expansion.)

Exercises 8.10 through 8.13 require a function that computes the CDF of the Gamma distribution.

**8.10.** Here we find the answer to the roulette problem of §8.8, using conditional Monte Carlo, but no other variance reductions.

- a) What is the numerical value of  $\alpha_{19}$  for wheel 1?
- b) Use conditional Monte Carlo to find the probability that number 19 has the highest probability of coming up on wheel 1 of §8.8. Give a 99% confidence interval.
- c) Estimate the probability that 3 is the highest probability number for wheel 2 of Table 8.3 and give a 99% confidence interval.
- d) Give a 99% confidence interval for  $p_{19}$  of wheel 1 and  $p_3$  of wheel 2. A gambler will make money in the long run by betting on a wheel with  $p > 1/36$ , and lose if  $p < 1/36$ , while the game is fair if  $p = 1/36$ . Do these confidence intervals include  $1/36$ ? You don't need to do a Monte Carlo for this part, the Monte Carlo you need is reported in Table 8.3.
- e) On wheel 1, the second most common number was 36. Estimate the probability that number 36 is the most probable, and give a 99% confidence interval.

**8.11.** Devise a strategy to find the probability that number 19 is the **second best** number for wheel 1 based on the data in Table 8.3. Give a formula for your method, and implement it, reporting the answer and a 99% confidence interval.

**8.12.** For the simulation in Exercise 8.10b estimate how much the variance was reduced by conditioning.

**8.13.** For the simulation in Exercise 8.10b sample  $p_{19}$  by stratified sampling, with 2 observations per stratum and the same sample size you used there (plus one if your sample size was odd). Report the ratio of the estimated variance of  $\hat{p}_{19}$  using ordinary IID sampling to that using stratified sampling. Both Monte Carlos in the ratio should employ conditioning.

**8.14.** In introductory probability exercises we might imagine a perfect roulette wheel with  $p_j = 1/38$  exactly. In Exercise 8.10 we considered  $\mathbf{p}$  uniformly distributed over all possible probability vectors. Neither of these models is reasonable. A more plausible model is that  $\mathbf{p} \sim \text{Dir}(A, A, \dots, A)$  for some value of  $A$  with  $1 < A < \infty$ . Then  $\mathbf{p} \mid \mathbf{X} \sim \text{Dir}(A + C_1, \dots, A + C_{38})$ .

a) For what value of  $A$  does

$$\mathbb{E}\left(\sum_{j=1}^{38}(p_j - 1/38)^2\right) = \sum_{j=1}^{38}(C_j/N - 1/38)^2$$

hold, where the counts  $C_j$  come from wheel  $j$ , and  $N = \sum_{j=1}^{38} C_j$ ?

b) Consider the following empirical Bayes analysis. Taking the number  $A$  obtained from part **a** replace the prior  $\text{Dir}(1, \dots, 1)$  by  $\text{Dir}(A, \dots, A)$ . This empirical Bayes analysis will change the estimated probability that 19 is really the best hole for wheel 1. Assuming that  $A > 1$ , the prior for  $\mathbf{p}$  will concentrate closer to the center of the simplex, and we anticipate a lower probability that wheel 19 is best.

How much does  $\mathbb{P}(p_{19} \geq \max_{1 \leq j \leq 38} p_j)$  change when we replace  $\alpha_j = 1$  by  $\alpha_j = A$  in the prior distribution? Use conditional Monte Carlo and common random numbers with  $n = 10,000$  sample points to estimate the difference in these probabilities.

## Control variates

**8.15.** Let  $f$  and  $h$  be two functions of the random variable  $\mathbf{X} \sim p$ . Define  $\mu = \mathbb{E}(f(\mathbf{X}))$ ,  $\theta = \mathbb{E}(h(\mathbf{X}))$ , and  $\Delta = \mu - \theta$ . Assume that we know  $\theta$  and that our goal is to estimate  $\Delta$ . Two estimators come to mind. The first estimator is  $\hat{\Delta}_1 = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - \theta)$ . The second estimator,  $\hat{\Delta}_2$  is obtained by estimating the mean of  $f(\mathbf{X}) - h(\mathbf{X})$ , using  $h(\mathbf{X})$  as a control variate.

For which values of  $\rho = \text{Corr}(f(\mathbf{X}), h(\mathbf{X}))$  is  $\hat{\Delta}_2$  more efficient than  $\hat{\Delta}_1$ ? You may use the following simplifying assumptions:

- i)  $\text{Var}(f(\mathbf{X})) = \text{Var}(h(\mathbf{X})) = \sigma^2 \in (0, \infty)$ .
- ii) The cost to evaluate  $h$  is the same as that for  $f$ .
- iii) The cost to sample  $\mathbf{X}$  is negligible.
- iv)  $n$  is large enough that the delta method approximation to the variance of the regression estimator is accurate enough.

**8.16.** In quadrature problems it is common to subtract a singularity that we can handle analytically. Here we look at what might happen if we used control variates instead.

Let  $f(x) = x^{-1/2} + x$  for  $x \in (0, 1)$ . Let  $h(x) = x^{-1/2}$ . We know that  $\theta \equiv \int_0^1 h(x) dx = 2$ , and of course  $\mu \equiv \int_0^1 f(x) dx = 5/2$ . Suppose that  $X \sim \mathbf{U}(0, 1)$ . Here we estimate  $\mathbb{E}(f(X))$  by Monte Carlo using  $h$  as a control variate, and forgetting for the moment that we know  $\mu$ . That is, we use  $\hat{\mu}_{\hat{\beta}}$  instead of  $\hat{\mu}_1 = (1/n) \sum_{i=1}^n (f(x_i) - 1(h(x_i) - 2))$ .

- a) Show that  $\text{Var}(f(X) - \beta h(X)) < \infty$  if and only if  $\beta = 1$ . State the variance of  $f(X) - h(X)$ .
- b) Let  $\hat{\mu}_{\hat{\beta}}$  be the usual control variate estimate of  $\mu$ . Suppose that  $n = 1000$ . Do a nested Monte Carlo analysis that repeats the size  $n$  simulation  $R =$

10,000 times. Report the sample mean, sample variance and histogram of  $\hat{\beta}$  over the  $R$  replicates. Does  $\hat{\beta}$  look like it is roughly normally distributed around the true value  $\beta = 1$ ?

- c) Show the sample mean, sample variance and histogram of  $\hat{\mu}_{\hat{\beta}}$  over the  $R$  estimates. Compare  $\hat{\mu}_{\hat{\beta}}$  to  $\hat{\mu}_2$ , by judging their sample squared errors. For practical purposes, do they appear to have very similar or sharply different accuracy? Either way, which one came out better than the other, in your simulations?
- d) Repeat the previous two parts with  $R = 10,000$  and  $n = 50$ .
- e) Inspect the histogram of  $\hat{\beta}$  values from part b. Find an apparent upper bound  $\hat{\beta} \leq A$  and then prove it holds. [Hint: Chebyshev's sum inequalities may be useful. If  $a_1 \geq a_2 \geq \dots \geq a_n$  and  $b_1 \geq b_2 \geq \dots \geq b_n$  and  $c_1 \leq c_2 \leq \dots \leq c_n$  then  $n \sum_i a_i b_i \geq \sum_i a_i \sum_i b_i$  and  $n \sum_i a_i c_i \leq \sum_i a_i \sum_i c_i$ .]

**8.17.** Suppose that  $\mathbb{E}(f(\mathbf{X})^2) < \infty$  and  $\mathbb{E}(h(\mathbf{X})^2) < \infty$  and  $\theta = \mathbb{E}(h(\mathbf{X})) \neq 0$ . Consider the ratio estimator  $\hat{\mu}_R = \theta \sum_{i=1}^n f(\mathbf{X}_i) / \sum_{i=1}^n h(\mathbf{X}_i)$ . Show that  $\mathbb{P}(|\hat{\mu}_R - \mu| > \epsilon) \rightarrow 0$  holds for any  $\epsilon > 0$ , and  $\mu = \mathbb{E}(f(\mathbf{X}))$ .

**8.18.** Under the conditions of Exercise 8.17, show that  $\mathbb{P}(|\hat{\mu}_P - \mu| > \epsilon) \rightarrow 0$ , where  $\hat{\mu}_P = \left(\frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i)\right) \left(\frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i)\right) / \theta$ .

**8.19.** Suppose that a control variate  $g(\mathbf{X})$  has a correlation of 0.1 with the variable  $f(\mathbf{X})$  of interest. By how much does its use reduce the variance of  $\mathbb{E}(f(\mathbf{X}))$ ? How much faster than  $f$  does the control variate function have to be for its use to improve the efficiency measure (8.1)?

**8.20.** For the unbiased control variate problem suppose that we will take  $N = n + m$  observations. The fraction of the sample allocated to finding the pilot estimate  $\hat{\beta}$  is  $f = m/N$ . Then a fraction  $1 - f$  is used for the final estimate. Suppose that the mean squared error takes the form  $(1/n)(A + \sigma_0^2/m)$  for constants  $A > 0$  and  $\sigma_0^2 > 0$ .

- a) Find the value of  $f$  that minimizes the mean squared error for fixed  $N > 0$  over the interval  $0 \leq f \leq 1$ . Let  $f$  vary continuously, even though  $fN$  must really be an integer.
- b) Let  $m(N)$  be the optimal solution from part a. For what  $r$ , if any, does  $m(N)/N^r$  approach a limit as  $N \rightarrow \infty$ ?

If the answer in part b is  $r = 1$  then the pilot sample should be a fixed fraction of the total data set. For  $r = 0$  we get a fixed number of pilot samples.

**8.21.** If  $\hat{\mu}$  and  $\hat{\theta}$  are positively correlated then  $\hat{\mu}/\hat{\theta}$  should be more stable because fluctuations in the numerator and denominator will offset each other. If they are negatively correlated we would expect  $\hat{\mu}\hat{\theta}$  to be more stable. Investigate this intuition by finding the delta method approximation to the variance of  $\hat{\mu}_R$  and  $\hat{\mu}_P$ . Assume that  $0 < \text{Var}(f(\mathbf{X})) = \sigma^2 < \infty$ ,  $0 < \text{Var}(\hat{\theta}) = \tau^2 < \infty$ ,  $\text{cor}(f(\mathbf{X}), h(\mathbf{X})) = \rho \in (-1, 1)$ , and that  $\theta \neq 0$ . By comparing the variances, decide whether  $\rho > 0$  favors the product estimator, or the ratio estimator, or neither as  $n \rightarrow \infty$ .

### Common random numbers

**8.22.** The content uniformity test on page 18 involved a small shift of the target value from 100 towards  $\bar{x}$ , but not going more than a distance of 1.5 units. This was implemented by the target shifting function  $M(x)$  in equation (8.22). It is natural to wonder whether target shifting makes much difference to the acceptance probability. We can turn off that feature by replacing  $M(x)$  with  $\tilde{M}(x) = 100$  for all  $x$ . Assume throughout that  $X_j \sim \mathcal{N}(\mu, \sigma^2)$  for  $j = 1, \dots, 30$  are independent.

- a) Suppose that  $\mu = 102$  and  $\sigma = 3$ . Estimate the amount (and direction) of the change in acceptance probability that arises from the use of target shifting.
- b) Now suppose that  $\mu = 100$ . Is there any  $\sigma$  for which target shifting changes the acceptance probability by more than 5%?
- c) Are there any  $(\mu, \sigma)$  pairs for which the acceptance probability changes by more than 50% due to target shifting? If so, describe the region where this happens. If not, what is the greatest change one can find? In either case, indicate which  $(\mu, \sigma)$  pairs result in the greatest change in acceptance probability.

Make a reasonable choice of Monte Carlo method for this problem, explaining the reasons for your choice. State the sample size you used. There will necessarily be numerical uncertainty because you cannot sample all configurations and  $n$  must be bounded.

**8.23.** In the content uniformity test, a really good product will pass at the first level, while a very bad one will not pass at all. Which combinations of  $\mu$  and  $\sigma$  lead to the greatest probability that the test will have to carry on to the second level, but will then pass?

**8.24.** Figure 8.5 was made with  $n = 100,000$  simulated cases, which may have been more than necessary. How could one determine whether a given sample size  $n$  is large enough for such a contour plot?

**8.25.** In financial applications one often needs the partial derivatives of an option value with respect to parameters like  $\delta$  and  $\sigma$ . These derivatives, termed ‘Greeks’ are needed for hedging. For the lookback option function  $f$  of Exercise 8.7c define  $g(\delta, \sigma, T) = \mathbb{E}(f(S(\cdot)))$  for the given values of  $\delta$ ,  $\sigma$ , and  $T$ . Using plain Monte Carlo, without stratification, estimate the following:

- a)  $g(0.051, 0.3, 1) - g(0.05, 0.3, 1)$ ,
- b)  $g(0.05, 0.31, 1) - g(0.05, 0.3, 1)$ , and
- c)  $g(0.05, 0.3, 1.01) - g(0.05, 0.3, 1)$ .

Give a confidence interval in each case. Make a reasonable choice for  $n$ .

**8.26.** Give an example where common random numbers increases variance. That is, find a distribution  $p$  and functions  $f$  and  $g$  and prove that  $\text{Var}(\hat{D}_{\text{com}}) > \text{Var}(\hat{D}_{\text{ind}})$  holds with your  $p$ ,  $f$  and  $g$ .

**8.27.** Exercise 5.13 is about sampling a bivariate distribution with Gaussian margins and the same copula that the Marshall-Olkin bivariate exponential distribution has.

In the notation of that exercise, suppose that  $\lambda_1 = \lambda_2 = 1$  and that we want  $\text{Corr}(Y_1, Y_2) = 0.7$ .

- a) What value of  $\lambda_3$  should we use? Devise a way to solve this problem using common random numbers and a fixed  $n \times 3$  matrix with independent components that were sampled from the  $\mathbf{U}(0, 1)$  distribution. Report the value of  $\lambda_3$  that you get.
- b) Repeat the previous part 10 times independently and report the 10 values you get.

---

## Bibliography

---

- Anderson, D. F. and Higham, D. J. (2012). Multilevel Monte Carlo for continuous time Markov chains, with applications in biochemical kinetics. *Multiscale Modeling & Simulation*, 10(1):146–179.
- Asmussen, S. and Glynn, P. W. (2007). *Stochastic simulation*. Springer, New York.
- Avramidis, A. N. and Wilson, J. R. (1993). A splitting scheme for control variates. *Operations Research Letters*, 14:187–198.
- Barraquand, J. (1995). Numerical valuation of high dimensional multivariate European securities. *Management Science*, 41(12):1882–1891.
- Boyle, P. P., Broadie, M., and Glasserman, P. (1997). Monte Carlo methods for security pricing. *Journal of economic dynamics and control*, 21(8):1267–1321.
- Burgos, S. and Giles, M. B. (2012). Computing Greeks using multilevel path simulation. In *Monte Carlo and Quasi-Monte Carlo Methods 2010*, pages 281–296. Springer.
- Cheng, R. C. H. (1985). Generation of multivariate normal samples with given sample mean and covariance matrix. *Journal of Statistical Computation and Simulation*, 21(1):39–49.
- Cochran, W. G. (1977). *Sampling Techniques (3rd Ed)*. John Wiley & Sons, New York.
- Duan, J.-C. and Simonato, J. (1998). Empirical martingale simulation for asset prices. *Management Science*, 44(9):1218–1233.
- Glynn, P. W. and Whitt, W. (1989). Indirect estimation via  $l = \lambda w$ . *Operations Research*, 37(1):82–103.

- Hammersley, J. M. and Morton, K. W. (1956). A new Monte Carlo technique: antithetic variates. *Mathematical proceedings of the Cambridge philosophical society*, 52(3):449–475.
- Hesterberg, T. C. and Nelson, B. L. (1998). Control variates for probability and quantile estimation. *Management Science*, 44(9):1295–1312.
- Kahn, H. and Marshall, A. (1953). Methods of reducing sample size in Monte Carlo computations. *Journal of the Operations Research Society of America*, 1(5):263–278.
- Karr, A. F. (1993). *Probability*. Springer, New York.
- L'Ecuyer, P. (1994). Efficiency improvement and variance reduction. In *Proceedings of the 1994 Winter Simulation Conference*, pages 122–132.
- L'Ecuyer, P., Simard, R., Chen, E. J., and Kelton, W. D. (2002). An object-oriented random number package with many long streams and substreams. *Operations research*, 50(6):131–137.
- Luenberger, D. G. (1998). *Investment Science*. Oxford University Press, New York.
- Lunney, P. D. and Anderson, C. A. (2009). Investigation of the statistical power of the content uniformity tests using simulation studies. *Journal of Pharmaceutical Innovation*, 4(1):24–35.
- Matsumoto, M. and Nishimura, T. (1998). Mersenne twister: a 623-dimensionally equidistributed uniform pseudo-random number generator. *ACM transactions on modeling and computer simulation*, 8(1):3–30.
- Owen, A. B. (2001). *Empirical Likelihood*. Chapman & Hall/CRC, Boca Raton, FL.
- Pullin, D. I. (1979). Generation of normal variates with given sample mean and variance. *Journal of Statistical Computation and Simulation*, 9(4):303–309.
- Wilson, A. (1965). *The Casino Gambler's Guide*. Harper & Row, New York.
- Wilson, J. R. (1984). Variance reduction techniques for digital simulation. *American Journal of Mathematical and Management Sciences*, 4(3):277–312.