

Monte Carlo

Galin L. Jones

School of Statistics

University of Minnesota

Draft: August 18, 2024

Contents

1	Introduction	2
1.1	Motivating Examples	2
2	Monte Carlo	6
2.1	Producing a random sample	7
2.1.1	Inversion	7
2.1.2	Accept-Reject	8
2.1.3	Ratio of Uniforms	9
2.1.4	Linchpin Variables	11
2.2	Estimation Theory for Monte Carlo	13
2.2.1	Monte Carlo estimation	13
2.2.2	Monte Carlo Error	14
2.2.3	Generalized Monte Carlo Error	15

1 Introduction

The “Monte Carlo method” means using a computer to simulate data in order estimate fixed unknown quantities (i.e. features or parameters) from a specified distribution. These quantities can be expectations, probabilities, density functions, quantiles, and so on. This likely sounds familiar since it is basically a description of much of classical statistics. The main difference is that it is based on data produced by a computer, rather than collected from an external experiment. There are two fundamental issues in implementing the Monte Carlo method: (1) designing experiments that produce useful observations and (2) using these observations to estimate features of the given distribution.

The simulated data can be a random sample as in classical Monte Carlo¹ or a realization of a Markov chain as in Markov chain Monte Carlo (MCMC). An independent and identically distributed (iid) sequence is a trivial Markov chain so it is not too surprising that Monte Carlo and MCMC have much in common. However, simulating a Markov chain introduces complications beyond those encountered when iid samples are available. Thus this chapter begins at the beginning and focuses on the simpler setting where simulation of an iid sample is possible.

1.1 Motivating Examples

Suppose distribution F has support \mathbf{X} and there is either an associated probability density function (pdf) or probability mass function (pmf), either of which will be denoted by f and referred to as a *probability function* (pf). If, at various points, something more general is needed, it will be carefully stated. The goal is to estimate $\theta \in \mathbb{R}^p$, $p \geq 1$ which is a vector of fixed, unknown features of F . For example, components of θ often include the following.

¹“Classical Monte Carlo” is cumbersome so “Monte Carlo” will be used instead. Hopefully, this will not cause confusion as “Monte Carlo” is often used as shorthand for “Monte Carlo method”.

1. Expectations. Let $h : \mathbf{X} \rightarrow \mathbb{R}$ and

$$\mu_h = E_F[h(X)] = \int_{\mathbf{X}} h(x) F(dx) .$$

This notation is used throughout so as to avoid having separate formulas for the continuous case where it denotes $\mu_h = \int_{\mathbf{X}} h(x) f(x) dx$ and the discrete case where it denotes $\mu_h = \sum_{x \in \mathbf{X}} h(x) f(x)$.

2. Quantiles. Set $V = h(X)$ and let F_V denote the distribution of V . If $0 < q < 1$, the q th quantile is

$$\xi_q = F_V^{-1}(q) = \inf\{v : F_V(v) \geq q\} .$$

Monte Carlo can be used for more than estimating expectations and quantiles, but these applications are common. A few examples follow which are intended to illustrate more concretely the types of settings where Monte Carlo might be useful. Although it might not be obvious at first glance, this includes settings that have no inherent probabilistic component.

Example 1.1. Consider

$$\int_0^1 \frac{1}{(x+1)^{2.3} [\log(x+3)]^2} dx ,$$

which does not appear easy to solve, but can be expressed as an expectation. If f is a pdf on $(0, 1)$, then

$$\begin{aligned} \int_0^1 \frac{1}{(x+1)^{2.3} [\log(x+3)]^2} dx &= \int_0^1 \frac{1}{(x+1)^{2.3} [\log(x+3)]^2} \frac{f(x)}{f(x)} dx \\ &= E_F \left[\frac{1}{f(x)(x+1)^{2.3} [\log(x+3)]^2} \right] . \end{aligned}$$

The example also makes the point that an expectation can be converted to an expectation with respect to a different distribution, say G having pf g and support containing \mathbf{X} :

$$\mu_h = \int_{\mathbf{X}} h(x) f(x) dx = \int_{\mathbf{X}} \frac{h(x) f(x)}{g(x)} g(x) dx = \mu_{hf/g} .$$

Example 1.2 (Bayesian Logistic Regression). Let X be a known $n \times p$ matrix with rows x_i and let β be a p -vector of parameters. Set

$$h(x_i) = \frac{\exp(x_i \beta)}{1 + \exp(x_i \beta)}$$

and assume $Y_i \sim \text{Bernoulli}(h(x_i))$, independently for $i = 1, \dots, m$. Let y denote all of the observed data and assume $\beta \sim N_p(0, I_p)$ so that the posterior is characterized by

$$q(\beta \mid y) \propto \left[\prod_{i=1}^n \frac{e^{-y_i x_i^T \beta}}{1 + e^{-x_i^T \beta}} \right] e^{-\frac{1}{2} \beta^T \beta} .$$

The normalizing constant or marginal density is

$$m(y) = \int \left[\prod_{i=1}^n \frac{e^{-y_i x_i \beta}}{1 + e^{-x_i \beta}} \right] e^{-\frac{1}{2} \beta^T \beta} d\beta ,$$

which is analytically intractable. A typical goal is to calculate the posterior mean of β :

$$\mu_\beta = E_q[\beta \mid y] = \int_{\mathbb{R}^p} \beta q(\beta \mid y) d\beta .$$

Posterior inference also often requires expectations of other functions $h(\beta)$, such as second moments, so the general goal is calculation of

$$\mu_h = E_q[h(\beta) \mid y] = \int_{\mathbb{R}^p} h(\beta) q(\beta \mid y) d\beta .$$

Posterior credible intervals can also be based on quantiles. For example, suppose the analysis requires a .95 credible interval $(\xi_{.025}, \xi_{.975})$ for the first component of β , that is β_1 . If $q(\beta_1 \mid y)$ is the marginal posterior density of β_1 , finding $(\xi_{.025}, \xi_{.975})$ requires

$$\int_{-\infty}^{\xi_{.025}} q(\beta_1 \mid y) d\beta_1 = .025 \quad \text{and} \quad \int_{-\infty}^{\xi_{.975}} q(\beta_1 \mid y) d\beta_1 = .975 .$$

Example 1.3 (Bayesian Linear Model). Suppose that, for $i = 1, \dots, k$ and $j = 1, \dots, m_i$,

$$Y_{ij} \mid \tau_i, \lambda_e \sim N(\tau_i, \lambda_e^{-1})$$

$$\tau_i \mid \mu, \lambda_e \sim N(\mu, \lambda_e^{-1})$$

$$\mu \sim N(m_0, s_0^{-1})$$

$$\lambda_e \sim \text{Gamma}(a_1, b_1)$$

$$\lambda_t \sim \text{Gamma}(a_2, b_2) .$$

Letting y denote all of the observed data and τ denote all of the τ_i , the posterior distribution is characterized by

$$q(\tau, \mu, \lambda_e, \lambda_t \mid y) \propto f(y \mid \tau, \lambda_e) f(\tau \mid \mu, \lambda_t) f(\mu) f(\lambda_e) f(\lambda_t) .$$

The normalizing constant or marginal density of y is

$$m(y) = \int f(y \mid \tau, \lambda_e) f(\tau \mid \mu, \lambda_t) f(\mu) f(\lambda_e) f(\lambda_t) d\tau d\mu d\lambda_e d\lambda_t ,$$

and is analytically intractable. Interest often centers on posterior expectations and quantiles. For example,

$$\mu_\tau = E_q[\tau \mid y] = \int \tau q(\tau \mid y) d\tau = \int \tau q(\tau, \mu, \lambda_e, \lambda_t \mid y) d\mu d\lambda_e d\lambda_t d\tau .$$

Because $m(y)$ is unavailable to us, analytical evaluation of posterior expectations or quantiles is unavailable.

While Monte Carlo methods have had a profound impact on the implementation of Bayesian inference, they are also important to the implementation of classical inference. Here is a simple example to consider.

Example 1.4 (Logit-Normal Generalized Linear Mixed Model). Let

$$p(\beta, u) = \frac{e^{\beta+u}}{1 + e^{\beta+u}} .$$

Suppose $Y \mid u, \beta \sim \text{Bernoulli}(p(\beta, u))$ and $U \mid \lambda \sim N(0, \lambda)$.

Then the likelihood is

$$L(\beta, \lambda) = \int_{-\infty}^{\infty} f(y \mid u, \beta) f(u \mid \lambda) du = \frac{1}{\sqrt{2\pi\lambda}} \int_{-\infty}^{\infty} \frac{e^{\beta+u}}{(1 + e^{\beta+u})^2} e^{-\frac{1}{2\lambda}u^2} du .$$

Now $L(\beta, \lambda)$ can be written as an expectation: let G be a distribution having density g on \mathbb{R} so that

$$\begin{aligned} L(\beta, \lambda) &= \int_{-\infty}^{\infty} \frac{f(y \mid u, \beta) f(u \mid \lambda)}{g(u)} g(u) du \\ &= E_G \left[\frac{f(y \mid u, \beta) f(u \mid \lambda)}{g(u)} \right] . \end{aligned}$$

Since the likelihood can be expressed as an expectation Monte Carlo can be used to approximate the function [4].

2 Monte Carlo

Settings where Monte Carlo is appropriate often begin with a given distribution F and the goal is to estimate θ , a vector of features of F . In Monte Carlo experiments, observations X_1, \dots, X_m are simulated and are used to construct an estimator θ_m in such a way that $\theta_m \approx \theta$ for large m . Calculation of the estimator θ_m alone is an incomplete solution to the problem. No matter how large m is, there will be an unknown *Monte Carlo error*, $\theta_m - \theta$ and hence θ_m will be more valuable if a measure of the Monte Carlo error is included.

Monte Carlo sample size is used to mean the size of the simulation effort and will be denoted m while n will be used to denote the sample size associated with the original statistical setting. A couple of examples may help illuminate the difference.

Example 2.1. Suppose Y_1, Y_2, Y_3 are iid $\text{Poisson}(\lambda)$ and $\lambda \sim \text{Gamma}(2, 3)$. Then the posterior is $\lambda \mid y_1, y_2, y_3 \sim \text{Gamma}(3\bar{y} + 2, 6)$, where \bar{y} is the sample mean of the three observations. There is no simple closed form for the median of the posterior, but simulation can be used to estimate it. Simulate a large number, 10000 say, observations from the posterior distribution. The sample median is an estimate of the true posterior median. Here $n = 3$ and $m = 10000$.

Example 2.2. Suppose Y_1, \dots, Y_{100} and $x_i = i/50$ for $i = 1, \dots, 100$. Let's use linear regression to model the observations as $Y_i = \beta x_i + \varepsilon_i$ and $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$. The goal is to test the hypotheses $H_0 : \beta = 0$ versus $H_1 : \beta \neq 0$ with a type 1 error rate of $\alpha = .05$. Then β is estimated using least squares and a standard t-test may be used for testing the hypotheses.

How robust is this procedure to departures from the assumption $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$? If, in fact, $\varepsilon_i \stackrel{iid}{\sim} \text{Cauchy}(0, \sigma)$, but the t-test is used, what happens to the type 1 error rate? One way to find out is via simulation. Do the following 1000 times: fix $\beta = 0$, simulate $\varepsilon_i \stackrel{iid}{\sim} \text{Cauchy}(0, \sigma)$ for $i = 1, \dots, 100$ and conduct the t-test. The proportion out of 1000 that reject is an estimate of the type 1 error. Here $n = 100$ and $m = 1000$.

While the theory of Monte Carlo is classical large sample theory, the asymptotics are as the Monte Carlo sample size m increases. Typically, the observed data sample size n is treated as fixed and known, but there are situations where it makes more sense to let n increase to infinity while

fixing m or letting m and n increase simultaneously. These last two settings will not be addressed further in this chapter, but will be considered later.

2.1 Producing a random sample

Monte Carlo methods are based on the ability to have the computer generate independent $\text{Uniform}(0, 1)$ observations. Of course, the observations are not random since they are produced by deterministic methods. However, good pseudorandom number generators produce sequences that effectively mimic independent $\text{Uniform}(0, 1)$ observations and hence are known as pseudorandom sequences. It is not clear that truly random sequences would be desirable since repeatability would be problematic, making debugging much more challenging. For the most part, this issue will not be considered further since the distinction between pseudorandom and random often will not be useful here.

This section considers some basic ways of obtaining a random sample from (non-uniform) F , including inversion, ratio of uniforms, the accept-reject algorithm, and linchpin variable sampling. This presentation is not intended to be comprehensive and the interested reader may consult many other texts for a more thorough treatment [e.g. 2, 3, 6].

2.1.1 Inversion

Define $F^{-1} : (0, 1) \rightarrow \mathbb{R}$ by

$$F^{-1}(y) = \inf\{x : F(x) \geq y\}.$$

The *quantile function theorem* [e.g. 1] is the foundation for simulating random variates from an arbitrary distribution given the ability to simulate from a Uniform distribution.

Theorem 2.1. *If $U \sim \text{Uniform}(0, 1)$ and $X = F^{-1}(U)$, then $X \sim F$.*

Example 2.3. Suppose $\beta > 0$. Inversion can be used to construct a draw from an $\text{Exp}(\beta)$ distribution. Then $F^{-1}(y) = -\beta \log(1 - y)$ so if $U \sim \text{Uniform}(0, 1)$, then $F^{-1}(U) = -\beta \log(1 - U) \sim \text{Exp}(\beta)$.

When F^{-1} is explicitly available and calculation of it is fast, inversion is practical, but these limitations often prevent its use.

2.1.2 Accept-Reject

The accept-reject algorithm uses draws from a convenient distribution G , having pf g , say, and converts them into draws from F . Suppose the support of G contains X and that

$$M = \sup_{x \in \mathsf{X}} \frac{f(x)}{g(x)} < \infty .$$

Algorithm 1 Accept-Reject

1: Draw $Y \sim G$

2: Draw $U \sim \text{Uniform}(0, 1)$

3: If

$$u \leq \frac{f(y)}{Mg(y)}$$

accept y as a draw from F ; otherwise return to step 1.

The accept-reject algorithm is a stochastic algorithm in that the accept-reject step is random. The probability of acceptance on a given step is

$$\begin{aligned} P(U \leq f(y)/Mg(y)) &= E [P(U \leq f(y)/Mg(y)) \mid Y] \\ &= E \left[\frac{f(Y)}{Mg(Y)} \right] \\ &= \frac{1}{M} . \end{aligned}$$

Theorem 2.2. *Algorithm 1 produces $X \sim F$.*

Proof. Notice that

$$\begin{aligned} P(X \leq x) &= P(Y \leq x \mid U \leq f(y)/Mg(y)) \\ &= \frac{P(Y \leq x, U \leq f(y)/Mg(y))}{P(U \leq f(y)/Mg(y))} . \end{aligned}$$

Now consider numerator and denominator separately:

$$\begin{aligned}
P(Y \leq x, U \leq f(y)/Mg(y)) &= E[P(Y \leq x, U \leq f(y)/Mg(y)) \mid Y] \\
&= E\left[I(Y \leq x) \frac{f(Y)}{Mg(Y)}\right] \\
&= \frac{1}{M} F(x)
\end{aligned}$$

and

$$\begin{aligned}
P(U \leq f(y)/Mg(y)) &= E[P(U \leq f(y)/Mg(y)) \mid Y] \\
&= E\left[\frac{f(Y)}{Mg(Y)}\right] \\
&= \frac{1}{M} .
\end{aligned}$$

Putting these together yields $P(X \leq x) = F(x)$, which proves the claim. \square

Clearly, the choice of proposal density g is crucial to the success of the algorithm. Note that M is the expected number of proposals required before a draw is obtained. To make M smaller and the algorithm more efficient requires a proposal g that mimics f in its tails.

Also notice that accept-reject can be used when the normalizing constant for f is unknown. Let $f(x) = c_1 h(x)$, $g(x) = c_2 l(x)$, and

$$\sup_{x \in \mathbf{X}} \frac{h(x)}{l(x)} = K .$$

If $U \sim \text{Uniform}(0, 1)$ and $Y \sim G$, then

$$U \leq \frac{h(y)}{Kl(y)} = \frac{f(y)}{\frac{c_1}{c_2} K g(y)} = \frac{f(y)}{Mg(y)}$$

and hence yields an equivalent accept-reject algorithm.

2.1.3 Ratio of Uniforms

Let h be a positive integrable function on (a, b) where a and b are not necessarily finite. Define

$$\mathbf{A}_h = \{(x, y) : 0 \leq x \leq h^{1/2}(y/x), \ a < y/x < b\} . \quad (1)$$

Example 2.4. Suppose $X \sim \text{Cauchy}(0, 1)$, then $h(x) = [1 + x^2]^{-1}$ and

$$\mathbf{A}_h = \{(x, y) : x > 0 \text{ and } x^2 + y^2 \leq 1\}.$$

Theorem 2.3. *If (U, V) is uniformly distributed on \mathbf{A}_h , then $X = V/U$ has pdf $f(x) \propto h(x)$.*

The proof of Theorem 2.3 is left as an exercise. The theorem suggests a simple algorithm for generating from a non-uniform distribution having pf f .

Algorithm 2 Ratio of Uniforms

- 1: Draw (U, V) uniformly on \mathbf{A}_h
 - 2: Set $X = V/U$
-

Algorithm 2 can be efficient, but generating uniformly on \mathbf{A}_h can be challenging. Fortunately, there is a special case of the accept-reject algorithm for avoiding this bottleneck. Set $a = 0$,

$$b = \sup_x \sqrt{h(x)} < \infty, \quad c = \sup_x x \sqrt{h(x)} < \infty, \text{ and } d = \inf_x x \sqrt{h(x)} > -\infty.$$

Then $\mathbf{A}_h \subseteq A = [a, b] \times [c, d]$ which suggests the following ratio of uniforms algorithm using accept-reject.

Algorithm 3 Ratio of Uniforms using Accept-Reject

- 1: Draw $U \sim \text{Uniform}(a, b)$
 - 2: Draw $V \sim \text{Uniform}(c, d)$
 - 3: If $U \leq h^{1/2}(V/U)$, set $X = V/U$; otherwise, repeat step 1.
-

The probability of acceptance is

$$\frac{\text{area}(\mathbf{A}_h)}{\text{area}(A)} = \frac{\text{area}(\mathbf{A}_h)}{b(d - c)} \tag{2}$$

and hence that the mean number of proposals until success is finite.

Example 2.5. This is a continuation of example 2.4. Notice that

$$A = [0, 1] \times [-1, 1]$$

and the acceptance probability (2) is $\pi/4$.

2.1.4 Linchpin Variables

The accept-reject and ratio of uniforms algorithms can be difficult to apply in multivariate settings, that is, when $d > 1$. However, a complicated multivariate simulation setting often can be converted to simulating from a simpler distribution. Suppose f can be expressed as a product of a conditional pf $f_{X|Y}$ and a marginal pf f_Y so that

$$f(x, y) = f_{X|Y}(x | y)f_Y(y)$$

If sampling from $f_{X|Y}$ is straightforward, then say Y is a *linchpin variable*.

Algorithm 4 Linchpin Variable Algorithm

- 1: Draw $Y \sim F_Y$
 - 2: Draw $X \sim F_{X|Y}(\cdot | Y)$
-

It should be obvious that the linchpin variable algorithm will be useful only when sampling from the marginal F_Y is easier than sampling from the joint F .

Example 2.6. Consider the Rosenbrock (or banana) density on \mathbb{R}^2

$$f(x, y) \propto \exp \left\{ -\frac{1}{20} [100(x - y^2)^2 + (1 - y)^2] \right\}.$$

This has become a benchmark example for illustrating the performance of sampling methods in highly correlated settings. In particular, because the contour plots resemble the shape of a banana, it can be a challenge to implement an effective sampling algorithm. Notice that by inspection of the joint density, $X|Y = y \sim N(y^2, 10^{-1})$ and integrating $f(x, y)$ with respect to x yields that $Y \sim N(1, 10)$. Hence Y is a linchpin variable and it is simple to implement conditional sampling.

Example 2.7. For $i = 1, \dots, n$ assume $t_i > 0$ is known and let $Y_i | \lambda_i \stackrel{ind}{\sim} \text{Poisson}(t_i \lambda_i)$. Assume priors $\lambda_i \stackrel{ind}{\sim} \text{Gamma}(a, \beta)$ and $\beta \sim \text{Gamma}(c, d)$ with a , c , and d known positive constants. The posterior is characterized by

$$f(\beta, \lambda | y) \propto \left(\prod_{i=1}^n \lambda_i^{a+y_i-1} e^{-(\beta+t_i)\lambda_i} \right) \beta^{an+c-1} e^{-\beta d}.$$

By inspection the conditionals are $\lambda_i \mid \beta \stackrel{ind}{\sim} \text{Gamma}(a + y_i, \beta + t_i)$ and the marginal pf for β is characterized by

$$f(\beta \mid y) \propto \beta^{an+c-1} e^{-\beta d} (\beta + t_i)^{-(a+y_i)} .$$

Thus β is a linchpin variable. It is left as an exercise to show that there is an efficient accept-reject algorithm using a $\text{Gamma}(an + c, d)$ proposal distribution to sample from the $\beta \mid y$ marginal. How does the acceptance rate change as d changes?

Example 2.8. For $i = 1, \dots, K$ suppose that for known $a, b, c > 0$,

$$\begin{aligned} Y_i \mid \theta_i &\sim \text{N}(\theta_i, a) & \theta_i \mid \mu, \lambda &\sim \text{N}(\mu, \lambda) \\ \lambda &\sim \text{IG}(b, c) & f(\mu) &\propto 1 . \end{aligned}$$

Then the hierarchy yields a proper posterior $f(\theta, \mu, \lambda \mid y)$ with $\theta = (\theta_1, \dots, \theta_K)^T$ and $y = (y_1, \dots, y_K)^T$. Consider the factorization [see 5]

$$f(\theta, \mu, \lambda \mid y) = f(\theta \mid \mu, \lambda, y) f(\mu \mid \lambda, y) f(\lambda \mid y) .$$

Then $f(\theta \mid \mu, \lambda, y)$ is the product of univariate normal densities $\theta_i \mid \mu, \lambda, y \sim \text{N}((\lambda y_i + a\mu)/(a + \lambda), a\lambda/(a + \lambda))$. Now $f(\mu \mid \lambda, y)$ is also a normal density $\text{N}(\bar{y}, (a + \lambda)/K)$. Finally,

$$f(\lambda \mid y) \propto \frac{1}{\lambda^{b+1}(a + \lambda)^{(K-1)/2}} \exp \left\{ -\frac{c}{\lambda} - \frac{1}{2(a + \lambda)} \sum_{i=1}^K (y_i - \bar{y})^2 \right\} .$$

Thus λ is a linchpin variable. It is left as an exercise to show that one can implement an efficient accept-reject algorithm using an $\text{IG}(b, c)$ proposal distribution to sample from the $\lambda \mid y$ marginal.

Linchpin samplers will typically not be useful when the dimension of the linchpin variable is too large. The following example illustrates this.

Example 2.9. Consider a version of the so-called Bayesian lasso. Let X be a known $m \times p$ design matrix and assume $\lambda > 0$ is known. Assume $a, b > 0$ are known and

$$\begin{aligned} Y \mid \beta, \gamma &\sim \text{N}_m(X\beta, \gamma^{-1}I_m) \\ \nu(\beta \mid \gamma) &= \left(\frac{\lambda\gamma}{4} \right)^p \exp \left\{ -\frac{\lambda\gamma}{2} \|\beta\|_1 \right\} \\ \gamma &\sim \text{Gamma}(a, b) . \end{aligned}$$

This hierarchy gives rise to a posterior density which has conditional

$$\gamma \mid \beta, y \sim \text{Gamma} \left(p + a + \frac{n}{2}, b + \frac{\lambda \|\beta\|_1 + \|y - X\beta\|_2^2}{2} \right)$$

and marginal

$$f(\beta \mid y) \propto \left(1 + \frac{\lambda \|\beta\|_1 + \|y - X\beta\|_2^2}{2b} \right)^{-(a+p+n/2)}.$$

In principle, one can construct an accept-reject sampler for sampling from the marginal of $\beta \mid y$ when X is full rank. However, the method is so inefficient as to be useless. Moreover, in situations where the lasso may be useful X is often not of full rank or p is large so that this linchpin variable sampler is not useful.

2.2 Estimation Theory for Monte Carlo

The estimation theory of Monte Carlo largely coincides with classical large sample statistics. As such this will not be a comprehensive review, instead it is focused on a few key ideas.

2.2.1 Monte Carlo estimation

Suppose there is a Monte Carlo sample $X_1, \dots, X_m \stackrel{iid}{\sim} F$ and the goal is to evaluate an expectation with respect to F

$$\mu_h := E_F[h(X)] = \int_{\mathcal{X}} h(x) F(dx).$$

If $E_F|h(X)| < \infty$, then the strong law of large numbers (SLLN) obtains so that, with probability 1, as $m \rightarrow \infty$,

$$\mu_m := \frac{1}{m} \sum_{i=0}^{m-1} h(X_i) \rightarrow \mu_h. \quad (3)$$

Thus the sample mean is an asymptotically valid estimator of μ_h .

Estimation of quantiles follows much the same pattern as estimation of expectations. Set $V = h(x)$ and let F_V denote the distribution of V . If $0 < q < 1$, the q th quantile is

$$\xi_q = F_V^{-1}(q) = \inf\{v : F_V(v) \geq q\}.$$

If F_V is absolutely continuous and has continuous density function f_V satisfying $0 < f_V(\xi_q) < \infty$, then ξ_q is the unique solution of $F_V(y-) \leq q \leq F_V(y)$.

If X_1, \dots, X_m are iid F , set $Y_i = h(X_i)$ for $i = 1, \dots, m$. Let $Y_{m(j)}$ be the j th order statistic of Monte Carlo sample. Then an estimator of ξ_q is

$$\xi_{m,q} = Y_{m(j)} \quad \text{where} \quad j-1 < mq \leq j \quad (4)$$

and, with probability 1,

$$\xi_{n,q} \rightarrow \xi_q \quad \text{as} \quad m \rightarrow \infty. \quad (5)$$

2.2.2 Monte Carlo Error

An obvious question is when should the simulation terminate? That is, when is θ_m a “good” estimate of θ ?

In estimating θ with θ_m , even for large values of the Monte Carlo sample size, m , there will be an unknown Monte Carlo error, $\theta_m - \theta$. Thus θ_m will be more valuable if there is a measure of its accuracy and it is reported. This can be accomplished with an interval estimator through the approximate sampling distribution of the Monte Carlo error. The interval estimator will allow assessment of the Monte Carlo error in the sense that it can describe the confidence in the number of significant figures reported. For example, suppose $\theta_{1000} = 123.45$. Then there are 5 significant figures reported in the estimate. But suppose that $100(1-\alpha)\%$ interval estimate is $[122.01, 124.91]$ so at the chosen confidence level only two of the five reported significant figures are trusted since values such as $\theta = 122$ or $\theta = 125$ would be plausible upon rounding. If this is not a sufficient level of precision, then this would prompt increasing the Monte Carlo sample size.

Consider the estimators μ_m and $\xi_{q,m}$; more general cases will be considered later. If μ_m estimates an expectation μ_h , then a central limit theorem (CLT) holds under classical conditions: If $E_F[h^2(X)] < \infty$, then, as $m \rightarrow \infty$,

$$\sqrt{n}(\mu_m - \mu_h) \xrightarrow{d} N(0, \text{var}_F[h(X)]) . \quad (6)$$

Moreover, $\text{var}_F[h(X)]$ can be consistently estimated with the sample variance

$$S_m^2 = \frac{1}{m-1} \sum_{i=0}^{m-1} [h(X_i) - \mu_m]^2$$

and consequently it is easy to calculate a *Monte Carlo standard error* (MCSE), s_m/\sqrt{m} . An MCSE can be used to produce a $100(1 - \alpha)\%$ confidence interval for the unknown value μ_h in the usual way: if $t_{m-1, \alpha/2}$ denotes a quantile from a Student's t distribution with $n - 1$ degrees of freedom, then

$$\mu_m \pm t_{n-1, \alpha/2} \frac{s_n}{\sqrt{m}}.$$

The width of the interval then conveys an idea of the accuracy of the Monte Carlo approximation—that is, how many significant figures can be trusted.

Now consider estimating the quantile ξ_q with the sample quantile $\xi_{m,q}$. Under the standing assumptions on F , as $m \rightarrow \infty$,

$$\sqrt{m}(\xi_{q,m} - \xi_q) \xrightarrow{d} N(0, q(1 - q)/f(\xi_q)^2).$$

Hence constructing a confidence interval for ξ_q will require estimation of $f(\xi_q)$. If f can be evaluated, then $f(\xi_{m,q})$ would suffice. Otherwise one would need to estimate the density at the point $\xi_{n,q}$ to obtain $\hat{f}(\xi_{m,q})$. A $100(1 - \alpha)\%$ confidence interval for the unknown value ξ_q in the usual way: if $z_{\alpha/2}$ denotes a quantile from a standard normal distribution, then

$$\xi_{q,m} \pm z_{\alpha/2} \frac{q(1 - q)}{\hat{f}(\xi_{m,q})}.$$

2.2.3 Generalized Monte Carlo Error

The typical Monte Carlo experiment is multivariate in two directions. That is, the sample X_1, \dots, X_m consists of d -dimensional random vectors while the goal is to estimate several, say p , features of F .

For definiteness suppose $h : \mathbf{X} \rightarrow \mathbb{R}^p$ for some $p \geq 1$ so that μ_h is a p -dimensional vector. Of course, if h is the identity mapping, then $p = d$, but in general the relative size of d and p is context dependent. Let $Y_j = h(X_j)$. As long as the Monte Carlo sample size satisfies $m > p$, describing the

variability in the sample is based on the sample covariance matrix

$$S_m = \frac{1}{m-1} \sum_{j=1}^m (Y_j - \mu_m)(Y_j - \mu_m)^T .$$

Then the *mean-centered ellipsoid of concentration* is defined by

$$E_k = \{y : (y - \mu_m)^T S_m^{-1} (y - \mu_m) \leq k^2\} \quad (7)$$

and consists of the points which are k units from the sample mean in squared Mahalanobis distance. The ellipsoid will have axes oriented along the eigenvectors of S_m having lengths proportional to the corresponding eigenvalues.

Notice that E_k will allow investigation of the bivariate marginals insofar as μ_m and S_m allow. For example, if the sample is approximately normally distributed and $k^2 = \chi_2^2(\alpha)$, then the ellipse will contain approximately $\alpha\%$ of the sample points. Consider the following two examples whose results are plotted in Figure 1.

Example 2.10. Consider a bivariate normal distribution, specifically $N_2(\mu, \Sigma)$ with $\mu = (1, 1)^T$ and

$$\Sigma = \begin{pmatrix} 1 & -.5 \\ -.5 & 2 \end{pmatrix} .$$

Simulation from this target distribution is easily accomplished with standard software.

Example 2.11. Suppose $Y_1, \dots, Y_m \sim N(\mu, \theta)$ and suppose the prior on (μ, θ) is proportional to $1/\sqrt{\theta}$. The posterior density $q(\mu, \theta|y)$ is proper if $m \geq 3$. Then $\mu \mid \theta, y \sim N(\bar{y}, \theta/n)$ and $\theta \mid y \sim \text{IG}(n/2, ns^2/2)$ so that θ is a linchpin variable and hence simulation from the posterior is straightforward.

While ellipsoids of concentration are especially useful in two dimensions, in general it is helpful to have univariate summaries of the variability in the Y_j . Two such summaries are total sample variance and generalized variance. The total sample variance is defined as $\text{trace}(S_m)$, that is, the sum of the diagonal elements of S_n . The generalized sample variance is $\det(S_n)$, which is proportional to the hypervolume of the p -dimensional ellipsoid of concentration, E_k , since

$$\text{Volume}(E_k) = \frac{2\pi k^p \Gamma(p/2)}{p} \det(S_m)^{1/2} .$$

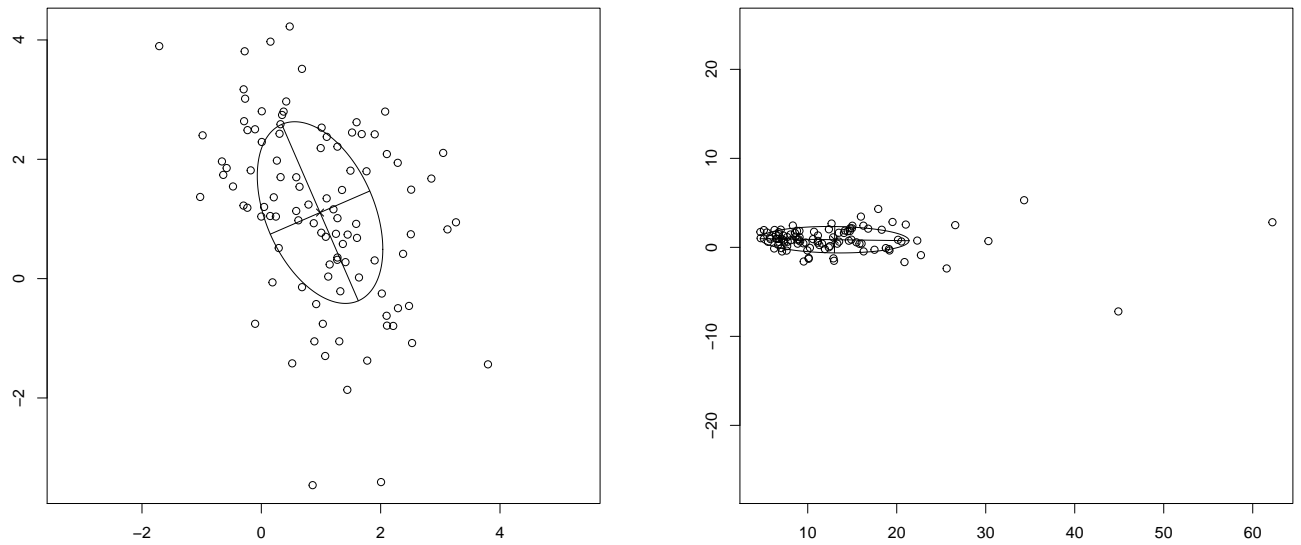


Figure 1: Ellipse of concentration with major and minor axes for two examples where the ellipse is formed by taking $k = \chi_2^2(.4)$. Left: A plot of 100 samples from the bivariate normal model in Example 2.10. Notice that 38 observations are in the ellipse. Right: A plot of 100 samples from the posterior in Example 2.11. Notice that 68 observations are in the ellipse.

Despite the appealing geometrical interpretation of generalized variance neither it nor the total sample variance tell the entire story. An examination of the ordered sample eigenvalues, say $\hat{\lambda}_1, \dots, \hat{\lambda}_p$, of S_m is advised since the total sample variance can be expressed as

$$\text{trace}(S_m) = \sum_{i=1}^p \hat{\lambda}_i$$

while generalized variance may be expressed as

$$\det(S_m) = \prod_{i=1}^p \hat{\lambda}_i.$$

Thus the total sample variance may be large due to one large eigenvalue, while generalized variance may be small due to one tiny eigenvalue.

2.2.4 Confidence Regions

Using univariate confidence intervals would require adjustment for multiplicity, perhaps a Bonferroni correction. Adjusting individual confidence intervals can work well when p is small, but will be overly conservative if p is even somewhat large. Indeed smaller regions are possible even when $p = 2$. The Monte Carlo estimation process can be better understood by using an explicit multivariate approach.

Assume the Monte Carlo sample size is larger than the number of estimands, that is, $m > p$. Using a Cramér-Wold device argument yields a multivariate CLT as in (6), that is, as $m \rightarrow \infty$,

$$\sqrt{n}(\mu_m - \mu_h) \xrightarrow{d} N_p(0, V) \quad (8)$$

where $V = \text{var}_F[h(X)]$ is a $p \times p$ positive definite matrix which can be estimated with the sample covariance

$$S_m = \frac{1}{m-1} \sum_{j=1}^m (h(X_j) - \mu_m)(h(X_j) - \mu_m)^T.$$

Then a confidence region for μ_h is defined by the ellipsoid

$$m(\mu_m - \mu_h)^T S_m^{-1} (\mu_m - \mu_h) \leq k^2$$

and k is chosen to ensure a desired coverage probability. More specifically, if $k^2 = \chi_{p,1-\alpha}^2$, then the ellipsoid will have approximate coverage probability $1 - \alpha$. If $\lambda_1, \dots, \lambda_p$ are the eigenvalues of S_m ,

then the ellipsoidal region is oriented along axes determined by the corresponding eigenvectors and whose lengths are proportional to $\sqrt{\lambda_i}$.

While ellipsoidal regions have minimal volume, they do not provide convenient marginal interpretations. It is also desirable to handle the simultaneous estimation of both means and quantiles. This will be addressed later.

References

- [1] Angus, J. E. (1994). The probability integral transform and related results. *SIAM Review*, 36:652–654.
- [2] Devroye, L. (1986). *Non-uniform Random Variate Generation*. Springer-Verlag Inc.
- [3] Fishman, G. S. (1996). *Monte Carlo: Concepts, Algorithms, and Applications*. Springer, New York.
- [4] Geyer, C. J. (1994). On the convergence of Monte Carlo maximum likelihood calculations. *Journal of the Royal Statistical Society, Series B*, 56:261–274.
- [5] Jones, G. L., Haran, M., Caffo, B. S., and Neath, R. (2006). Fixed-width output analysis for Markov chain Monte Carlo. *Journal of the American Statistical Association*, 101:1537–1547.
- [6] Robert, C. P. and Casella, G. (2004). *Monte Carlo Statistical Methods*. Springer, New York, second edition.