

Monte Carlo Methods

Useful refs:

Monte Carlo Stat. Methods : Robert & Casella

MCMC : Stochastic Simulation for Bayesian

Inference : Gamerman & Lopes

Ross: Ch. 11

No reqd. text: lecture notes + occasional handouts

Computing assignments: must be typed up.

(easiest in 'R')

Monte Carlo

Basic idea: learn about prob. models by simulating them. Useful for both probability & stat. inference.

Formally: Use pseudo-random (simulated) values from a prob. distr. f to estimate expectations w.r.t. f . Also useful for integration in general.

Suppose we want to find μ where

$$\mu = E_f(g(x)) = \int g(x) f(x) dx$$

(or $\sum_i g(x_i) f(x_i)$ if f is pmf)

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f$ using a computer

Define $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n g(X_i)$

$\therefore \hat{\mu}_n$ is the Monte Carlo estimate of μ .

Strong Law of Large Numbers:

If $E_f |g(x)| < \infty$ then $\hat{\mu}_n \rightarrow \mu$ a.s.

a.s. ('almost sure') convergence: $P(\lim_{n \rightarrow \infty} \hat{\mu}_n = \mu) = 1$

If $E_f |g(x)| < \infty$ and $\sigma^2 = \text{Var}_f(g(x)) < \infty$

then by the Central Limit Thm

$$\sqrt{n} (\hat{\mu}_n - \mu) \rightarrow N(0, \sigma^2) \text{ in distribution}$$

Convergence in distr. (weak convergence, convergence in law):

$X_n \rightarrow X$ in distr. if $F_{X_n}(x) \rightarrow F_X(x)$ as $n \rightarrow \infty$ for all x st. $F_X(x)$ is conts.

$F_X(x) = P(X \leq x)$, the cdf.

Note: Informally, ^{distr. of} $\hat{\mu}_n$ approaches $N(\mu, \sigma^2/n)$
so Monte Carlo error of $\hat{\mu}_n$ decreases as $1/\sqrt{n}$.

Above results hold for multidimensional spaces
(not just 1-D r.v.s).

Accuracy of M.C. estimate is independent
of dimensionality of space sampled.

Thinking about iid Monte Carlo = thinking about
basic statistics.

Midterm 2011

Mean 17.7

Median 18.5

V. roughly : $[20, 25]$: excellent
 $[16, 20)$: fine-good
 ≤ 15 : possible cause for concern

Midterm prob.:

- 10 1. standard M.C. $\leq \frac{(a) \text{ set-up} + (b) - (f) \text{ standard shift (no pts off if error due to (a))}}{(b) - (f)}$
- 4 2. detailed balance / standard M.C. (done in class)
- 2 3. (a) basic P.P.
- 2 (b) cond. P.P. $X, Y, Z \text{ iid Unif}(0, 1)$ } more thought
- 3 (c) memoryless / conditioning for exp. } ~~tricky~~, but like hw
- 2 4. (a) iterated expec.
- 2 (b) using indep. incr. } tricky

Toy example

(Geyer's notes)

$X \sim N(0,1)$, $Y \sim N(0,1)$, independent

What is $P(Y < X^2)$?

$$\text{Let } \mu = P(Y < X^2) = \int \underbrace{\Phi(x^2)}_{P(Y < x^2) \text{ for } N(0,1)} \underbrace{\phi(x)}_{\text{pdf for } N(0,1)} dx$$

Monte Carlo approach: simulate $(X_1, Y_1), \dots, (X_n, Y_n)$
of pairs of $N(0,1)$ r.v.'s

$$\hat{\mu}_n = \frac{\sum_{i=1}^n \mathbb{I}(Y_i < X_i^2)}{n} \quad (\text{sample prop.})$$

This is easy enough that numerical integration also works, but often numerical integration is not an option.

3. If $\mu = E_f(g(x)) < \infty$, we have.

$\hat{\mu}_n \rightarrow \mu$ a.s. (almost surely) by the Strong Law of Large Numbers (S.L.L.N.) That is,

$$P(\lim_{n \rightarrow \infty} \hat{\mu}_n = \mu) = 1. \quad (1.1)$$

Furthermore, if $\sigma^2 = \text{Var}_f(g(x)) < \infty$, we can establish a convergence rate for this estimate, that is we can establish how quickly $\hat{\mu}_n$ converges to μ from the **Central Limit Theorem**:

$$\sqrt{n}(\hat{\mu}_n - \mu) \rightarrow N(0, \sigma^2) \text{ in distribution} \quad (1.2)$$

Example 1: Suppose we want to calculate $\Pr(-1 < X < 0)$ when X is a Normal(0,1) random variable. We could easily do this by Monte Carlo:

- Generate $X_1, \dots, X_n \sim N(0, 1)$.
- Compute the estimate

$$\hat{\mu}_n = \frac{\sum_{i=1}^n 1(-1 < X_i < 0)}{n},$$

which is simply the proportion of times $X_i \in (-1, 0)$ for sampled values X_1, \dots, X_n .

For large enough n , $\hat{\mu}_n$ will be very close to $\Pr(-1 < X < 0)$. Of course, Monte Carlo is not really needed for this toy problem since statistical software can easily calculate such probabilities.

Example 2: Suppose we want to conduct a simple hypothesis test to see if the correlation ρ between two random variables is significant. Assume that the two random variables X and Y come from a bivariate normal distribution, and we observe 30 data points $(X_1, Y_1), \dots, (X_{30}, Y_{30})$, and want to conduct a hypothesis test based on these data. The sample correlation, $\hat{\rho}$, is the test statistic. For these data, $\hat{\rho} = 0.3$. To find the associated p-value, we need to find the probability $P(\hat{\rho} > 0.3)$ under the null hypothesis that there is no correlation ($\rho = 0$), and the alternative that $\rho > 0$. To calculate this

start \rightarrow
here

probability, we would need to know the sampling distribution of $\hat{\rho}$ under the null hypothesis that $\rho = 0$. This null distribution is not easy to calculate, but it is given in Anderson (2003). The sample correlation coefficient $\hat{\rho}$, for a sample of size N from a bivariate Normal with mean $\boldsymbol{\mu}$ and correlation ρ , depends only on ρ and N (not on $\boldsymbol{\mu}$ or the marginal variances), and is given by:

$$f(\gamma) = \frac{2^{N-3}(1-\rho^2)^{0.5(N-1)}(1-\gamma^2)^{-0.5(N-4)}}{\pi\Gamma(N-2)} \sum_{\alpha=0}^{\infty} \Gamma\left(\frac{N+\alpha-1}{2}\right)^2 \frac{(2\rho\gamma)^\alpha}{\alpha!},$$

where $\gamma \in (-1, 1)$. Now, for $\rho = 0$, the sampling distribution simplifies to:

$$f(\gamma) = \frac{2^{N-3}(1-\gamma^2)^{-0.5(N-4)}}{\pi\Gamma(N-2)} \Gamma\left(\frac{N-1}{2}\right)^2.$$

Finding the above distribution is a non-trivial and time consuming problem, and even though we can assume that we did not have to work to find the distribution (since the theory has already been worked out), finding the p-value for this hypothesis test would still involve integrating the above density over the interval $(0.3, \infty)$. However, if we simply use the fact that the distribution of $\hat{\rho}$ only depends on ρ and the sample size N , a Monte Carlo solution to this problem is very simple:

- To draw a single sample of $\hat{\rho}$ from the null distribution, generate a sample $X_1, \dots, X_N \sim N(0, 1)$ and a sample $Y_1, \dots, Y_N \sim N(0, 1)$. (In R you would use the command `xs=rnorm(30,0,1)`, for instance). Find the sample correlation r_1 based on the pairs $(X_1, Y_1), \dots, (X_N, Y_N)$. Repeat this process m times to obtain m sample correlations r_1, \dots, r_m , generated from the null distribution.
- Compute the estimate

$$\hat{\mu}_n = \frac{\sum_{i=1}^m 1(r_k > 0.3)}{m}.$$

For large enough m , this is an accurate estimate of the desired p-value.

Note that the only theory necessary was recognizing that the sampling distribution of $\hat{\rho}$ depends only on ρ and N , which is easy to prove by the invariance of the distribution of $\hat{\rho}$ to affine transforms of a bivariate normal random variable (see Anderson (2003) for details). It was not necessary to derive the complicated formula for the sampling distribution above, nor was it necessary to compute the integral; the p-value is easily estimated through a simple Monte Carlo procedure. Of course, since the exact distribution is available here, and the p-value only involves a 1-dimensional integral, it is possible to do this by using numerical integration procedures. In more complicated situations, Monte Carlo will often be the only solution.

1.3 Monte Carlo standard errors

Informally, (1.2) states that the distribution of $\hat{\mu}_n$ approaches a Normal distribution, $N(\mu, \sigma^2/n)$. Hence, to obtain confidence intervals and error estimates for the estimator $\hat{\mu}_n$, we need to estimate σ^2 . Monte Carlo standard error is an assessment of the error of our Monte Carlo estimator. For the simple independent Monte Carlo scenario above, we can easily estimate σ^2 by the sample variance, $\hat{\sigma}^2$. Then, the estimate of Monte Carlo standard error is simply $\hat{\sigma}/\sqrt{n}$. Since $\hat{\sigma}^2$ is a consistent estimator of σ^2 , we can use Slutsky's Theorem and (1.2) to obtain asymptotic 95% confidence intervals for Monte Carlo estimates in the usual way: $\hat{\mu}_n \pm 1.96\hat{\sigma}/\sqrt{n}$ ((1.2) still holds when σ^2 is replaced by a consistent estimator of σ^2). We note the following:

1. The independence requirement for the samples is unnecessarily restrictive as we will see later on when we discuss Markov chain Monte Carlo methods. The S.L.L.N. will typically hold under similar conditions for the dependent case but the C.L.T. may not hold and estimating σ^2 is typically very difficult.
2. f may be multivariate, i.e., the random variable for which f is the probability distribution may be multidimensional.

If we can simulate draws from f , easy to estimate expectations w.r.t. f .

Generally hard to do for multivariate random quantities.

Few exceptions: multivariate normal, Wishart.

How do we simulate samples from an arbitrary distr. f ? General strategy:

- 1) Build a method to generate $U_1, U_2, \dots \stackrel{\text{iid}}{\sim} \text{Unif}(0,1)$.
- 2) Generate $X \sim f$ using U_1, U_2, \dots .

1) Uniform pseudo-random^{number} generator: starting from initial value u_0 , called the 'seed', produce a deterministic sequence of values in $[0,1]$, u_1, u_2, \dots by $u_k = d(u_{k-1}, u_{k-2}, \dots)$ for some function d . This sequence imitates a sequence of iid uniform r.v.s.

end 200

This sequence is not random but is an acceptable pseudo-random sequence if it passes a set of tests s.t. the hypothesis $u_1, \dots, u_n \stackrel{iid}{\sim} \text{Unit}(0,1)$ is reasonable e.g. Marsaglia's 'Die-Hard' tests.

Typical uniform random # generator: of form:

Initial value: X_0 (seed)

$$X_{n+1} = \underbrace{(a X_n + c) \text{ modulo } m}_{\text{remainder from dividing } aX_n + c \text{ by } m} \quad a, c, m: \mathbb{Z}^+$$

so $X_n \in \{0, \dots, m-1\}$

X_n/m approximation to $\text{Unit}(0,1)$.

Lots of generators, increasingly sophisticated:
Knuth's TAOC, Marsaglia's Super-Duper,
Mersenne-Twister etc. see `help(Random.seed)`.

Properties of unif. r.v. generators:

- 1) Repetitive: after enough draws, pattern repeats itself. # draws before repetition = period.
- 2) Not independent.
- 3) Not quite uniform continuous (actually discrete).

Lots of v. good uniform r.v. generators in existing stat. software — we will generally not worry abt. it, (except if using parallel computing. Why: Start M.C. on different machines — results not as indep. as we would like.)

E.g. R uses Mersenne-Twister w/ period $2^{19937} - 1$.

Note: advantage of deterministic behavior: reproducibility.

In R: Random Seed: can look at seed.

2) How do we generate $X \sim f$ using uniform r.v.'s?

For important r.v.'s special techniques exist.

See Ross's book for normal, gamma, χ^2 , beta, et.c.

More general approach: Inverse CDF / inverse transformation method

Let $U \sim \text{Unit}(0,1)$. If F is a cdf
 $X = F^{-1}(U)$ has cdf F . $F^{-1}(u) = \{x : F(x) = u\}$

Pf: $\text{Prob}(F^{-1}(U) \leq x) = \text{Prob}(U \leq F(x))$
 $= F(x) \quad \therefore U \sim \text{Unit}(0,1)$.

Very limited use: only if univariate r.v. and F has explicit inverse. E.g. Cauchy, Exponential, Weibull, Logistic.
but not normal, beta, and many gammas.

Drawing $X \sim f$ is generally hard, especially in high dimensions.

Further complication: Often normalizing constants are unknown.

Rejection (Accept-Reject) Sampling

Goal: Simulate $X \sim f(x)$ when we have $h(x)$ s.t. $h(x)/c_1 = f(x)$, i.e.

only know $f(x)$ up to a ~~const~~ normalizing constant c_1 (we don't know c_1).

Suppose we have a simpler pdf $g(x)$ s.t. it is easy to draw $Y \sim g(x)$.

and we know $\sup_x \frac{h(x)}{g(x)} < K$ for some $K < \infty$

and we know K ; with $h(x)/c_2 = g(x)$

Rejection sampler: for $i=1, \dots, N$:

Simulate $Y_1, \dots, Y_N \stackrel{\text{iid}}{\sim} g$

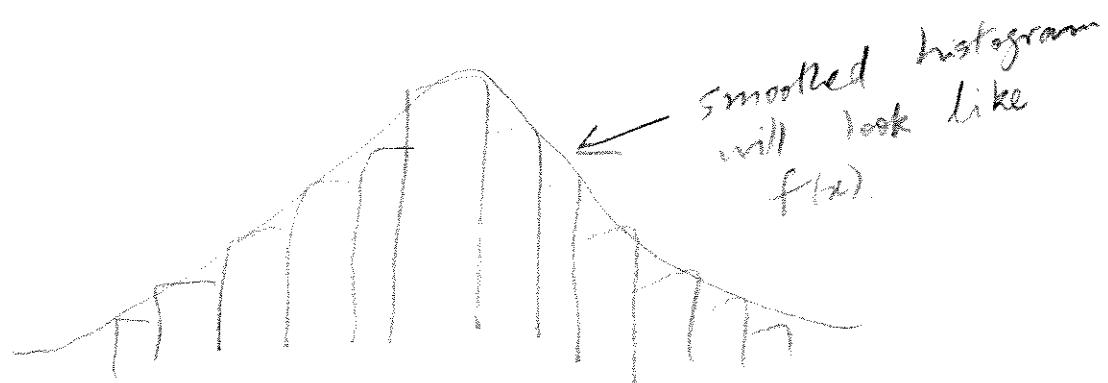
1. Draw $Y_i \sim g$
2. Draw $U \sim \text{Unit}(0, 1)$
3. If $U \leq \frac{h(Y_i)}{K g(Y_i)}$ return Y_i , else do not.

Values returned by above algorithm are draws from $f(x)$.

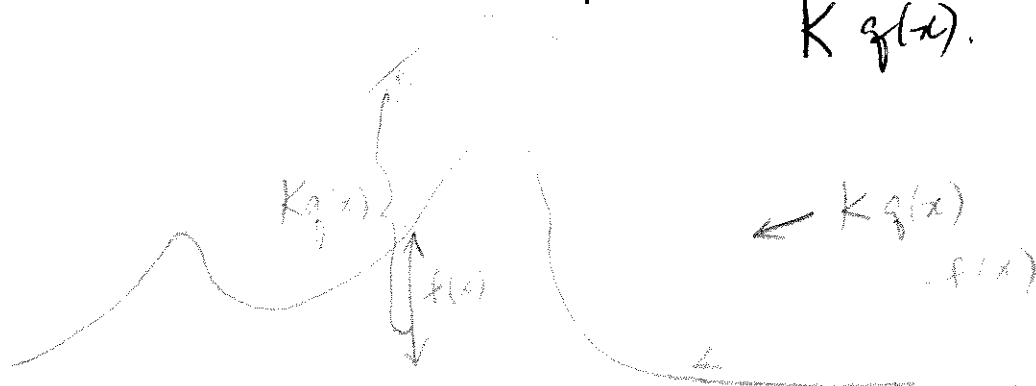
Intuition for rejection sampler

If we could draw directly from $f(x)$.

$X_1, \dots, X_n \stackrel{iid}{\sim} f(x)$. Large n .



Suppose $f(x)$ is more complicated, but we have envelope $Kg(x)$.



If we accepted all draws from g would get shape of outer density

But if we ^{accept} reject samples w/ prob $\frac{f(x)}{Kg(x)}$,

can 'recover' inner density.

Proof for rejection sampler ($1-D_n$ case). ^{controls}

Let $X \sim q(x)$ and $U \sim \text{Unit}(0,1)$

Accept X if $U \leq \frac{h(x)}{K r(x)}$ where $\frac{r(x)}{c_2} = q(x)$,
 $\frac{h(x)}{c_1} = f(x)$

and $\sup_x \frac{h(x)}{r(x)} \leq K < \infty$

Claim: X accepted by this alg. has distr. $f(x)$

Pf: $P(X \leq x | X \text{ accepted})$

$$= P(X \leq x | U \leq \frac{h(x)}{K r(x)})$$

$$= P(X \leq x, U \leq \frac{h(x)}{K r(x)}) / P(U \leq \frac{h(x)}{K r(x)})$$

iterated expectations =

$$= \frac{E_q \{ P(X \leq x, U \leq \frac{h(x)}{K r(x)} | X) \}}{E_q \{ P(U \leq \frac{h(x)}{K r(x)} | X) \}} = \frac{E_q \{ P(X \leq x | X) P(U \leq \frac{h(x)}{K r(x)} | X) \}}{E_q \{ \frac{h(x)}{K r(x)} \}}$$

(condl. indep.)

$$= \frac{E_q \{ \underbrace{I(X \leq x)}_{\text{either } X \leq x \text{ or } \text{reject (given } x)} \frac{h(x)}{K r(x)} \}}{E_q \{ \frac{h(x)}{K r(x)} \}} = \frac{E_q \{ I(X \leq x) \frac{c_1 f(x)}{K c_2 q(x)} \}}{E_q \{ \frac{c_1 f(x)}{K c_2 q(x)} \}}$$

uses fact that $\frac{h}{r} < \infty$, same for $\frac{f}{q} < \infty$

$$= \frac{\int_{-\infty}^{\infty} \frac{f(x)}{q(x)} q(x) dx}{\int_{-\infty}^{\infty} \frac{f(x)}{q(x)} q(x) dx}$$

$$(\because \frac{h}{r} < \infty, \frac{f}{q} < \infty) \Rightarrow \frac{I(X \leq x)}{F(x)} \quad F(x), \text{ cdf corresponding to } f(x). \text{ [norm. constant not needed]}$$

Result holds much more generally (multivariate, discrete etc.)

Computing tip #1 likelihoods
v. useful when evaluating ratio of pdf's etc.
v. useful to take logs whenever possible

~~For eg. $\frac{g(x)f(x)}{g(x)} = \exp \{ \log g(x) + \log f(x) - \log g(x) \}$~~

Fewer problems w/ round-off error, ~~instability of~~
numerical instability.

For eg.
Rej. sampling: Instead of ~~comparing~~ evaluating:

Is $U \leq \frac{f(x)}{K g(x)}$?

Try: Is $\log(U) \leq \log f(x) - \log g(x) - \log K$?

Computing tip #2: use apply, sapply (eg see
help and Rlinks).

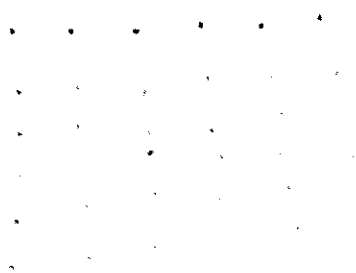
Tip #3: printing out intermediate values,
debugging: use `cat("this var =", varname, "\n")`
OR `print(varname)`.

E.g. Where normalizing constant is unknown

E.g.1. Binary Markov Random Field ^{Look up: Ising Model from stat mechanics}

N individuals 'pixels' arranged in a lattice. Each individual is in one of 2 state: $+1$ or -1 (eg. black or white)

Ω = space of all configurations, each config
= (s_1, \dots, s_N) , where $s_i = +1$ or -1 .



Suppose model: $P(S = (s_1, \dots, s_N)) = \frac{1}{Z} \exp\{-J U(s)\}$
where $J > 0$, fixed.

And $U(s) = \sum_{i \sim j} I(s_i \neq s_j)$ $i \sim j$ if i is a nbr of j

Z is a normalizing constant.

$Z = \sum_{s \in \Omega} \exp\{-J U(s)\} \Rightarrow$ have to sum over 2^N possible configurations. $N=10 \Rightarrow 1.27 \times 10^3$

$U(s) = \#$ of unlike nbr pairs in configuration

This model favors smooth images, especially for large J .

We would want a simulation strategy that does not require normalizing constant to be known.
MC 9

E.g. 2. Simple Bayesian model.

$L(Y|\theta)$: likelihood of data Y given parameters θ .

$p(\theta)$: prior distr. on θ .

Inference based on posterior, $\pi(\theta|Y)$

$$\pi(\theta|Y) \propto \frac{L(Y|\theta) p(\theta)}{\int L(Y|\theta) p(\theta) d\theta} \leftarrow \text{fn. of } \theta$$

$$\propto \frac{L(Y|\theta) p(\theta)}{Z} \leftarrow \begin{array}{l} \text{constant} \\ \text{w.r.t. } \theta \end{array}$$

$Z \leftarrow$ unknown and potentially complicated especially if lots of parameters.

Very rarely is $\pi(\theta|Y)$ tractable / available in closed form

Recap. For rejection sampler we need:

1) q s.t. $\sup_x \frac{h(x)}{q(x)} < \infty$ $\left(\frac{h(x)}{c_1} = f(x)\right)$

q should have heavier tails than h .

2) Value of K s.t. $\left(\frac{h(x)}{c_2} = g(x)\right)$
 $\sup_x \frac{h(x)}{h(x)} < K, K < \infty.$

~~For~~ For efficient sampler, try to find smallest
 K satisfying above req.

3) q must match h well or else
alg. will be v. inefficient.

4) Must be easy to simulate from q , i.e.
need to ~~draw~~ draw values from q quickly for alg.
to be efficient. end 6/27/07

Note: ① Do not need to know normalizing
constant of $f(x)$. In fact, do not need
normalizing constant of $g(x)$. Why is this useful?

~~If g is multivariate normal, evaluating~~

② Can show acceptance probability = $\frac{c}{K}$