# Week 8: Monte Carlo (MC) MATH-517 Statistical Computation and Visualization

Linda Mhalla

2024-11-08

## Introduction

 $MC\equiv$  repeated random sampling to mimic outcome of random process and produce numerical results such as

- generating draws from complicated distributions and/or domains
- integration
  - calculation of moments or confidence intervals
  - high-dimensional densities in Bayesian settings
- optimization
  - mode evaluation

**Basic idea**: If we can sample from a process or mimic its outcomes, we can learn a lot about it by doing statistics on the simulated samples (as opposed to analyzing the process itself)

MC methods  $\equiv$  simulation-based statistical techniques/inference

Linda Mhalla

## Introduction

Gambling experiments have random outcomes - hence "Monte Carlo"



Method initially developed by Stanislaw Ulam and John von Neumann for the Manhattan Project (to estimate integrals)

# Example

 $(X_1,Y_1)^\top,\ldots,(X_N,Y_N)^\top$  a sample from the standardized bivariate Gaussian distribution

$$\mathcal{N}\left(\begin{pmatrix}0\\0\end{pmatrix},\begin{pmatrix}1&\rho\\\rho&1\end{pmatrix}\right)$$

 $\bullet$  We want to test  $H_0: \rho = \rho_0$  against  $H_1: \rho \neq \rho_0$ 

• Statistic: 
$$\hat{
ho}_N = N^{-1} \sum_{n=1}^N X_n Y_n$$

Since the data generation process is fully determined under  $H_0$ , we can simulate data to approximate the sampling distribution and thus also the p-value

To test a hypothesis, we only need to simulate data under  $H_0$ 

But how to draw samples from a specific distribution?

Linda Mhalla

Week 8: Monte Carlo (MC)

## Example: Estimation of $\pi$

```
\frac{\text{area of the circle}}{\text{area of the square}} = \frac{\pi r^2}{2r \times 2r} = \frac{\pi}{4}
          \Pr(\text{point inside the circle}) = -
piVal <- function(nPoints = 300){</pre>
   x <- runif(nPoints,-1,1)</pre>
   y <- runif(nPoints,-1,1)</pre>
   result <- ifelse(x<sup>2</sup>+y<sup>2</sup> <= 1, TRUE, FALSE)</pre>
   4*sum(result)/nPoints
}
Ν
  <- 1000
pi est <- replicate(N, piVal())</pre>
mean(pi est) #to get uncertainty
```

```
[1] 3.141707
```

## Section 1

# Random Number Generation (RNG)

# RNG

True randomness is hard to come by. Historically:

- dice, cards, coins
- physical processes
- census data, tables, etc.

Practical reasons not to use "truly" random numbers: debugging and reproducibility

John von Neumann: pseudo-RNG

- $\bullet\,$  approximates the desired dist. for  $N\to\infty\,$
- cannot be predicted
- pass a set of independence tests
- repeatability ( $\Rightarrow$  reproducibility)
- long cycle (before it starts repeating) and fast sampling

Uniformity and independence tests needed to assess quality of pseudo-RNG

Week 8: Monte Carlo (MC)





# Cornerstone: Generating from $\mathcal{U}[0,1]$

Assume now we can generate numbers from the  $\mathcal{U}[0,1]$  distribution

• e.g., the linear congruential method (LCG)

$$X_n=(aX_{n-1}+c)\,\mathrm{mod}\,m,\quad n=1,2,\ldots,$$

where a, c, m, and  $X_0$  are cleverly chosen to fulfill the pseudo-RNG requirements, i.e., maximize period, speed, and "randomness"

- X<sub>0</sub> is the seed
- produces integers between 0 and m-1

• 
$$U_n = X_n/m \stackrel{.}{\sim} U(0,1)$$

Bad example:  $m=2^{31}$ ,  $a=2^{16}+3$ , and c=0  $\Rightarrow$  IBM's RANDU

**Example:** a = 13, c = 0, and m = 64

# LCG

#### Choice of parameters

- $\bullet \ \, {\rm For} \ m=2^b {\rm , \ and} \ c\neq 0$ 
  - longest possible period  $P = m = 2^b$  is achieved if c is relative prime to m and a = 1 + 4k, where k is an integer

• For 
$$2, m = 2^b$$
, and  $c = 0$ 

- longest possible period  $P=m/4=2^{b-2}$  is achieved if the seed  $X_0$  is odd and a=3+8k or a=5+8k, for  $k=0,1,\ldots$
- For m a prime and c = 0
  - longest possible period P=m-1 is achieved if the multiplier a has property that smallest integer k such that  $a^k-1$  is divisible by m is  $k={\rm m}-1$

Now, better and much more complicated algorithms are available

- every piece of software has its favorite pseudo-RNG
- outside the scope of the course (see, e.g., shift-register generators or Wichmann-Hill generator)

Linda Mhalla

#### Transforms

**Lemma.** (Inverse Transform.) Let  $U \sim \mathcal{U}(0,1)$  and F be a distribution function and  $F^{-1}$  the quantile (or generalized inverse) function. Then  $X = F^{-1}(U) \sim F$ .

**Proof**: Simply  $P(X \le x) = P\{F(X) \le F(x)\} = P\{U \le F(x)\} = F(x).$ 



The inverse transform method is general, but not almighty:

- distribution/quantile functions can be complicated/unknown - e.g.,  $\mathcal{N}(0,1)$ 

Often, simpler relationships can be used: diagram

 $\bullet$  still, there is no arrow there between  $\mathcal{U}(0,1)$  and  $\mathcal{N}(0,1),$  generating  $\mathcal{N}(0,1)$  is actually a bit tricky...

## Transforms

Lemma. (Box-Muler transform.) Let  $U_1, U_2 \sim \mathcal{U}(0,1)$  be independent. Then

$$Z_1 = \sqrt{-2\log(U_1)}\cos(2\pi U_2) \quad \& \quad Z_2 = \sqrt{-2\log(U_1)}\sin(2\pi U_2)$$

are two independent standard Gaussian random variables

**Explanation:** points from the  $\mathcal{N}(0, I_2)$ , when written in polar coordinates, have an angle  $\theta \sim U[0, 2\pi)$  which is independent of the radius R

Again, software uses its favorite relationships

- e.g., R has tabulated F and  $F^{-1}$  for  $\mathcal{N}(0,1)$  to a high precision and actually uses the inverse transform, because evaluating trigonometric functions is rather expensive (slow)
- ?rnorm  $\Rightarrow$  rnorm, pnorm, qnorm, dnorm help

12/33

# **Rejection Sampling**

 ${\bf Setup}:$  we know how to simulate from a  $proposal\ g,$  we want to simulate from a  $target\ f$ 

- let  $supp(f)\subset supp(g),$  i.e.,  $f(x)>0\Rightarrow g(x)>0$
- $\bullet \,$  let there be c>1 such that  $\forall x:\,f(x)\leq c\,g(x),$  i.e.,  $\sup_{x}\frac{f(x)}{g(x)}=c<\infty$

**Algorithm**: (to draw a single sample X from f)

#### Example:

- $\bullet \ \mathcal{U}(0,1) \ \mathrm{proposal}$
- $\mathcal{B}(2,5)$  target
- $c \approx 2.5$



2024-11-08

Week 8: Monte Carlo (MC)

# **Rejection Sampling**

Does the algorithm really sample from f?

$$\begin{split} P(X \le x) &= P\left\{Y \le x \left| U \le \frac{1}{c} \frac{f(Y)}{g(Y)} \right\} = \frac{P\{Y \le x \land U \le t(Y)\}}{P\{U \le t(Y)\}} \\ &= \frac{\int_{-\infty}^{x} \int_{0}^{t(y)} du \, g(y) dy}{\int_{-\infty}^{+\infty} \int_{0}^{t(y)} du \, g(y) dy} = \frac{\int_{-\infty}^{x} t(y) g(y) dy}{\int_{-\infty}^{+\infty} \int_{0}^{t(y)} du \, g(y) dy} = \frac{\int_{-\infty}^{x} t(y) g(y) dy}{\int_{-\infty}^{+\infty} \frac{1}{c} f(y) dy} \\ &= \frac{\frac{1}{c} F(x)}{\frac{1}{c}} = F(x) \end{split}$$

The rejection sampling algorithm above is again quite general, but it needs

- $\bullet$  a good proposal g
  - as close as possible to the target density (small values of c), leading to
  - high acceptance probability  $P\{U \leq t(Y)\} = 1/c$
- $\bullet\,$  fast evaluation of f and g

# Example: $\mathcal{N}(0,1)$ again

**Goal**: Simulate data from the standard Gaussian target using the double exponential (Laplace) proposal, i.e.,

$$f(x)=\frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}\quad \&\quad g(x)=\frac{\alpha}{2}e^{-\alpha|x|}, \text{ where } \alpha>0, \quad x\in\mathbb{R}$$



Another way of obtaining  $\mathcal{N}(0,1)$  from  $\mathcal{U}(0,1)\text{:}$ 

$$\begin{split} \mathcal{U}(0,1) &\longrightarrow Exp(1) \\ Exp(1) &\longrightarrow DbExp(1) \\ DbExp(1) &\longrightarrow N(0,1) \end{split}$$

Note:  $\alpha = 1$  minimizes the value of  $c = \sup f(x)/g(x)$  and hence maximizes the acceptance probability

Linda Mhalla

# Section 2

#### Numerical Integration

#### Deterministic Quadrature Approaches

**Goal**: approximate  $J = \int_a^b f(x) dx$  via composite quadrature rules *Quadrature method*: evaluate the function on a grid

$$S_K=\{(b-a)/K\}\sum_{k=1}^K f(t_k)$$

- $\bullet$  if f is nice (smooth),  $S_K \to J$  for  $K \to \infty$
- ullet corresponds to integrating a local constant interpolation of f
- local linear interpolation (*trapezoidal rule*) or local quadratic (*Simpson's rule*) are also well known



#### Naive Monte Carlo

We consider the more general integral

$$J = \int_{\mathcal{X}} m(x) f(x) dx = \mathbb{E}_f \big\{ m(X) \big\} \quad \text{for} \quad X \sim f$$

 $\Rightarrow$  generate  $X_1,\ldots,X_N \overset{i.i.d.}{\sim} f$  and approximate J by

$$\widehat{\bar{J}}_N = N^{-1} \sum_{n=1}^N m(X_n)$$

• Provided  $\int |m(x)| f(x) dx < \infty$ , unbiased and consistent due to SLLN • monitoring convergence via CLT-based confidence intervals:

$$\sqrt{N}\frac{\widehat{J}_N-J}{\sigma(m)}\stackrel{.}{\sim}\mathcal{N}(0,1),\quad\text{where}\quad\sigma^2(m)=\int m^2(x)f(x)dx-J^2$$

(Naive Monte Carlo estimate  $v_N^2$  of  $\sigma^2(m)$  is used in practice)

 $\Rightarrow$  MC integration methods give a probabilistic error bound versus a deterministic one with numerical quadrature rules

Linda Mhalla

Week 8: Monte Carlo (MC)

#### Naive Monte Carlo

$$\lim_{N \to \infty} \Pr\left(-3\frac{\sigma(m)}{\sqrt{N}} \le \frac{1}{N} \sum_{n=1}^{N} m\left(X_n\right) - J \le 3\frac{\sigma(m)}{\sqrt{N}}\right) \approx 99.73\%$$

 $\Rightarrow$  to attain an error of  $\epsilon>0$  with 99.73%, we need

$$N = \frac{9}{\epsilon^2} \sigma^2(m) = O(1/\epsilon^2)$$

- This is independent of the dimension (depends on the smoothness of *m* though)!
- Deterministic numerical integration methods suffer the curse of dimensionality:  $N=O(1/\epsilon^d),$  roughly
- Monte Carlo estimate for the integral converges slowly to the true value (slower than quadrature methods for smooth functions in low dimensions)
- Beware of rare events: if f has heavy tails,  $\hat{\bar{J}}_N$  can be a bad estimate and we need huge N to get small  $v_N \Rightarrow$  techniques for variance reduction

# Importance Sampling

We often can not simulate directly from f and we require sophisticated approaches. Rewrite % f(x) = 0

$$J := \int_{\mathcal{X}} m(x) f(x) dx = \int_{\mathcal{X}} m(x) \frac{f(x)}{g(x)} g(x) dx = \int_{\mathcal{X}} m(x) w(x) g(x) dx$$

with g a density whose support contains that of f and  $w(x) \geq 0$  the importance weighting function

Thus, by sampling  $X_1,\ldots,X_N \stackrel{\mathrm{i.i.d.}}{\sim} g,$  we can approximate J by

$$\widehat{J}_N:=N^{-1}\sum_{n=1}^N m(X_n)w(X_n)$$

Idea:

- $\bullet\,$  Use a simpler proposal distribution g from which we can generate
- $\bullet\,$  Candidates generated from g fall within the domain of f
- Reweight the observations generated from it when taking the mean

Linda Mhalla

## Importance Sampling: Intuitive Explanation

 ${\bf Key}:$  integrating f amounts to integrating f/g under sampling from g

$$J=\int_{\mathcal{X}}m(x)f(x)dx=\mathbb{E}_g\{m(X)w(X)\}$$

- when f is flat (all regions are equally important), use either the naive MC (with uniform sample) or deterministic approaches that need only small samples
- when f is not flat, using a "good" g allows us to encode which regions are important ⇒ "importance sampling" (vs rejection sampling)

Of course, it is not always easy to find a "good" g which

- ullet has a similar shape than f and
- from which we can easily sample

As we will see, when  $\mathcal{X}=\mathbb{R},$  it is important to match the decay of the tails between the target and reference measures

Linda Mhalla

Week 8: Monte Carlo (MC)

# Importance Sampling: Properties

• unbiased and the variance is given by

$$\begin{split} \operatorname{var}(\widehat{J}_N) &= \frac{1}{N} \bigg\{ \int_{\mathcal{X}} m^2(x) \frac{\widehat{f}(x)}{g(x)} f(x) dx - J^2 \bigg\} \\ &= N^{-1} \int_{\mathcal{X}} \bigg\{ m(x) \frac{f(x)}{g(x)} - J \bigg\} m(x) f(x) dx \end{split}$$

which is small if  $g(x)\approx m(x)f(x)J^{-1}$  or  $g(x)\propto m(x)f(x)$ 

 $\Rightarrow$  good choices of g can yield huge improvements in efficiency, e.g.,

- g with similar shape than mf
- $\bullet\,$  maximum principle ( g(x) and m(x)f(x) take their maximum at the same value)
- $\bullet \,\,g$  from the same distribution family as f

For instance, reduction in variance is

$$\operatorname{var}(\widehat{\bar{J}}_N) - \operatorname{var}(\widehat{J}_N) = \frac{1}{N} \int_{\mathcal{X}} m^2(x) \bigg\{ 1 - \frac{f(x)}{g(x)} \bigg\} f(x) dx$$

Linda Mhalla

- $\bullet\,$  Most difficult aspect to importance sampling is in choosing a good sampling density g
- Need to be very careful as it is possible to choose g according to some heuristics, but that results in a variance increase
- It is possible to have an importance sampling estimator with infinite variance
  - e.g., if  $g \ {\rm puts}$  too little weight relative to  $f \ {\rm on}$  the tails of the distribution

Task: Approximately calculate P(2 < X < 6) for the target distribution  $X \sim f$  using a reference g

Gaussian target, Exponential reference - densities (left), their ratio (middle), importance sampling error (right)



Approximately calculate P(2 < X < 6) for the target distribution  $X \sim f$  using a reference g

Cauchy target, Exponential reference - densities (left), their ratio (middle), importance sampling error (right)



# Examples

Approximately calculate P(2 < X < 6) for the target distribution  $X \sim f$  using a reference g

Cauchy target, Gaussian reference - densities (left), their ratio (middle), importance sampling error (right)



- $\bullet$  the tails of Cauchy and Gaussian distributions are too different  $\Rightarrow$  importance sampling performs poorly
- if we can simulate from Gaussian, we can simulate directly from Cauchy:  $Z_1, Z_2 \sim \mathcal{N}(0, 1)$  independent  $\Rightarrow Z_1/Z_2 \sim Cauchy(0, 1)$

# Variance Reduction

Accuracy of MC integration is assessed by the estimator's efficiency/variance (assuming efforts of simulation are similar)

There are ways to tweak the sampling scheme in order to reduce the variance

- importance sampling (we have seen above)
- antithetic variables (to follow)
- stratified sampling (to follow)
- quasi-random sampling and control variates (see the supplementary notes)
- many other techniques: latin hypercube sampling, ratio estimator, etc

**Remark**: When comparing several different estimators via simulations, the same simulated datasets should be used for all the estimators

#### Variance Reduction: Antithetic Variables

Idea: Introduce negative correlation between pairs of replications and rely on

$$\mathrm{var}(f_1+f_2) = \mathrm{var}(f_1) + \mathrm{var}(f_2) + 2\mathrm{cov}(f_1,f_2)$$

Given two i.i.d. samples  $X_1,\ldots,X_N\sim f$  and  $Y_1,\ldots,Y_N\sim f$  , consider the estimator

$$\tilde{J}_N = \frac{1}{2N} \sum_{n=1}^N \{m(X_n) + m(Y_n)\} = \frac{1}{2} (\hat{\bar{J}}_N^X + \hat{\bar{J}}_N^Y)$$

Then,

$$\mathrm{var}(\tilde{J}_N) = \frac{1}{2} \mathrm{var}(\hat{\bar{J}}_N) \{1 + \mathrm{corr}(\hat{\bar{\mathbf{J}}}_N^{\mathrm{X}}, \hat{\bar{\mathbf{J}}}_N^{\mathrm{Y}})\}$$

 $\Rightarrow \tilde{J}_N$  is more efficient than the naive MC (with sample of size 2N) if  $m(X_n)$  and  $m(Y_n)$  are negatively correlated

*Basic result:* if g(u) is monotonic on 0 < u < 1, then

$$\operatorname{corr}\{g(U),g(1-U)\}<0$$

Thus  $F^{-1}(U) \mbox{ and } F^{-1}(1-U)$  are negatively correlated with distribution F

# Variance Reduction: Stratified Sampling

- Break sampling space into strata and sample appropriate number of observations in each
- Compute the naive MC estimator in each stratum and sum over all strata

We want to estimate  $J = \mathbb{E}\{m(X)\}$ . Let W be another r.v. s.t.

- $p_i = \Pr(W \in \Delta_i)$ , with  $\Delta_i \subset \mathbb{R}$ , is easily computed
- we know how to generate  $X \mid W \in \Delta_i$

Assuming  $\Delta_1, \dots, \Delta_m$  is a partition of  $\mathbb{R}$  and denoting  $J_j = \mathbb{E}\{m(X) \mid W \in \Delta_j\}$  and  $\sigma_j^2 = \operatorname{var}\{m(X) \mid W \in \Delta_j\}$ , we have

$$J = \mathbb{E}[\mathbb{E}\{m(X) \mid W\}] = \sum_{j=1}^m p_j J_j, \quad \text{and} \ \widehat{J}_{str,N} = \sum_{j=1}^m p_j \widehat{\bar{J}}_{j,N_j}$$

# Variance Reduction: Stratified Sampling

• Method relies on conditional variance:

$$\operatorname{var}\{m(X)\} = \mathbb{E}[\operatorname{var}\{m(X)|I\}] + \operatorname{var}[\mathbb{E}\{m(X) \mid I\}]$$

Thus,

$$\mathrm{var}\{m(X)\} \geq \mathbb{E}[\mathrm{var}\{m(X)|I\}] = \sum_{j=1}^m p_j \sigma_j^2$$

For instance, assuming proportional allocation, i.e.,  $N_j = N p_j$ , then  $\operatorname{var}(\widehat{J}_{str,N}) = \sum_{j=1}^p p_j^2 \frac{\sigma_j^2}{N_j} = \frac{\sum_{j=1}^m p_j \sigma_j^2}{N} \leq \operatorname{var}(\widehat{\bar{J}}_N)$ 

- $\bullet$  Variance reduction substantial if I accounts for a large fraction of the variance of m(X)
- Variance reduction depends on the allocation in each stratum. Thus, we can minimize  $var(\widehat{J}_{str,N})$  subject to  $\sum_{j=1}^m N_j = N$

Donald Knuth (1997, 3rd ed.) The Art of Computer Programming, vol. 2 Robert & Casella (2010) Introducing Monte Carlo methods with R

# Feedback for the mini-project

- Good points
  - original datasets
  - going the extra mile to dig deeper in the data like creating new variables
  - nice introductions and good referencing
- Points to improve
  - be careful when interpreting results (do not jump on conclusions), e.g., correlation coefficients, latent confounders
  - remember to log-transform when needed ...
  - captions missing!!!
  - code appearing in the text
  - bad sectioning of the text (or no sectioning!)
  - bad citations or missing references

Go to Assignment 6 for details.