

2 Importance Sampling

Can we do better than the simple Monte Carlo estimator of

$$\theta = E[g(X)] = \int g(x) f(x) dx \approx \frac{1}{m} \sum_{i=1}^m g(X_i)$$

where the variables X_1, \dots, X_m are randomly sampled from f ?

Yes!!

Goal: estimate integrals with lower variance than the simplest Monte Carlo approach.

↳ more efficient estimation.

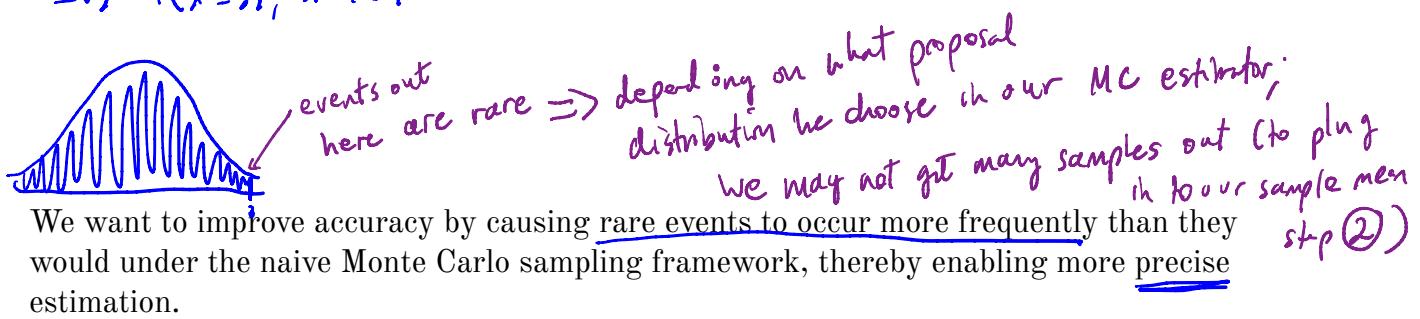
To accomplish this, we will use *importance sampling*.

2.1 The Problem (why would our simple MC estimators have high variance to begin with?)

If we are sampling an event that doesn't occur frequently, then the naive Monte Carlo estimator will have high variance.

Example 2.1 Monte Carlo integration for the standard Normal cdf. Consider estimating $\Phi(-3)$ or $\Phi(3)$. (HW6)

$$\Phi(z) = P(X \leq z), \quad X \sim N(0, 1).$$



We want to improve accuracy by causing rare events to occur more frequently than they would under the naive Monte Carlo sampling framework, thereby enabling more precise estimation.

For very rare events, extremely large reductions in the variance of the MC estimator are possible.

2.2 Algorithm

Consider a density function $f(x)$ with support \mathcal{X} . Consider the expectation of $g(X)$,

$$\underset{\text{parameter of interest}}{\theta} \rightarrow \theta = E[g(X)] = \int_{\mathcal{X}} g(x) f(x) dx.$$

Let $\phi(x)$ be a density where $\phi(x) > 0$ for all $x \in \mathcal{X}$. Then the above statement can be rewritten as

$$\begin{aligned} &\uparrow \text{support of } \phi \text{ includes the support of } f \\ \{x : \phi(x) > 0\} &\supseteq \{x : f(x) > 0\} = \mathcal{X} \end{aligned}$$

$$\begin{aligned} \theta &= E[g(X)] = \int_{\mathcal{X}} g(x) \frac{f(x)}{\phi(x)} \phi(x) dx \\ &= E\left[g(Y) \cdot \frac{f(Y)}{\phi(Y)}\right] \text{ for } Y \sim \phi \end{aligned}$$

ϕ is called the importance sampling function (similar to an envelope in accept-reject),

An estimator of θ is given by the *importance sampling algorithm*:

ϕ MUST be a density!
(integrate to 1 $\Leftrightarrow \int \phi(x) dx = 1$).

→ 1. Sample X_1, \dots, X_m from ϕ

2. Compute

$$\hat{\theta} = \frac{1}{m} \sum_{i=1}^m g(X_i) \frac{f(X_i)}{\phi(X_i)}$$

For this strategy to be convenient, it must be

① easy to sample from ϕ

② easy to evaluate f (and ϕ) even if it is not easy to sample from f .

$X = \text{result of rolling one fair six-sided die.}$

Example 2.2 Suppose you have a fair six-sided die. We want to estimate the probability that a single die roll will yield a 1. $P(X=1)$.

We could

① Roll a die m times

② A point estimate of $P(X=1)$ would be the proportion of ones in the sample.

The variance of this estimator is $\frac{5}{36m}$ if the die is fair.

$$X = \{1, \dots, 6\} \text{ and } f(x) = \begin{cases} \frac{1}{6} & x \in \{1, \dots, 6\} \\ 0 & \text{o.w.} \end{cases}$$

$$\text{Define } Y = \begin{cases} 1 & \text{if } X=1 \\ 0 & \text{o.w.} \end{cases} \Rightarrow Y \sim \text{Bernoulli}\left(\frac{1}{6}\right).$$

$$EY = \sum_{i=1}^6 \mathbb{I}[X=i] \cdot \frac{1}{6} = \frac{1}{6}$$

Expected # of 1's in m rolls

$$E\left(\sum_{i=1}^m Y_i\right) = \sum_{i=1}^m EY = \frac{m}{6}$$

$$\text{Var } Y = p(1-p) = \frac{1}{6} \cdot \frac{5}{6} = \frac{5}{36}$$

Proportion of 1's in the sample:

$$E\left[\frac{\sum Y_i}{m}\right] = \frac{1}{m} \cdot \frac{m}{6} = \frac{1}{6} \quad \text{unbiased estimator}$$

$$\text{Var}\left[\frac{\sum Y_i}{m}\right] = \frac{1}{m^2} \sum \text{Var } Y_i = \frac{1}{m} \text{Var } Y_i = \frac{1}{m} \cdot \frac{5}{36} \quad \begin{array}{l} \nearrow \text{variance} \\ \text{of our} \\ \text{estimator} \end{array}$$

We could consider the "coefficient of variation"

$$\text{So, } CV\left[\frac{\sum Y_i}{m}\right] = \frac{\sqrt{5/(36m)}}{\sqrt{1/6}}$$

$$CV[X] = \frac{\sqrt{\text{Var}(X)}}{E(X)} \quad \begin{array}{l} \swarrow \text{relative} \\ \text{measure of} \\ \text{variation} \\ (\text{chemistry,} \\ \text{physics}). \end{array}$$

Say we want a CV of 5%. Want $\frac{\sqrt{5/(36m)}}{\sqrt{1/6}} = .05 \dots$ solve for $m \dots m = 2000$ rolls!

goal: Can we do better? i.e. use less rolls to get the same $CV = .05$ with a different estimator.

To reduce the # rolls, we could consider biasing the die by replacing the faces bearing 2 and 3 with 1's.

This will increase the probability of rolling a 1 to 0.5, but now we are not sampling from target distribution (a fair die).

We can correct this biasing by

- weighting each roll of 1 by $\frac{1}{3}$

- Let $Y_i = \begin{cases} \frac{1}{3} & \text{if } X=1 \\ 0 & \text{o.w.} \end{cases}$

Now:

$$P(X=1) = \frac{1}{2}$$

$$P(X=2) = P(X=3) = 0$$

$$P(X=4) = P(X=5) = P(X=6) = \frac{1}{6}$$

$$\text{estimator: } \frac{1}{m} \sum_{i=1}^m Y_i$$

$$E\left[\frac{1}{m} \sum_{i=1}^m Y_i\right] = \frac{1}{m} \sum EY_i \stackrel{iid}{=} EY_i = \frac{1}{3} \cdot \frac{1}{2} + 0 \left[0 + 0 + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \right] = \frac{1}{6} \quad \text{unbiased!}$$

But the variance:

$$EY^2 = \left(\frac{1}{3}\right)^2 \cdot \frac{1}{2} = \frac{1}{18}$$

$$\text{Var}\left[\frac{1}{m} \sum_{i=1}^m Y_i\right] \stackrel{\text{ind}}{=} \frac{1}{m^2} \sum \text{Var}Y_i \stackrel{iid}{=} \frac{1}{m} \text{Var}Y_i = \frac{1}{m} \left[\frac{1}{18} - \left(\frac{1}{6}\right)^2 \right] = \frac{1}{36m}$$

So to achieve CV of 5% we would only need

$$\frac{\sqrt{1/(36m)}}{Y_6} = .05$$

: solve for 5

$$\Rightarrow m = 400 \text{ rolls.}$$

- This die rolling example is successful because an importance function (rolling die w/ 3 sides=1) is used to over-sample a portion of the state space that receives lower probability under the target and then importance weighting to correct the bias.

2.3 Choosing ϕ

① support of ϕ must include support of f .

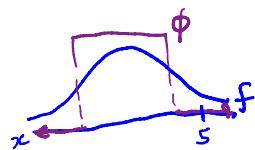
In order for the estimators to avoid excessive variability, it is important that $f(x)/\phi(x)$ is bounded and that ϕ has heavier tails than f .

target/importance
importance weights

If this requirement is not met, some importance weight $\frac{f(x)}{\phi(x)}$ will be huge.

① Example 2.3

If we ignore requirement $\phi(x) > 0$ when $f(x) > 0$

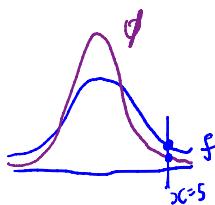


Then $\frac{f(s)}{\phi(s)} = \frac{f(s)}{0}$ unbounded! → breaks requirement ①.

And, we can't draw $x=5$ from ϕ .

② Example 2.4

If we select ϕ with lighter tails than f .



$\frac{f(s)}{\phi(s)}$ will be large if $\phi(s)$ is small.

Thus $x=5$ draw has large weight \Rightarrow approx will be poor.

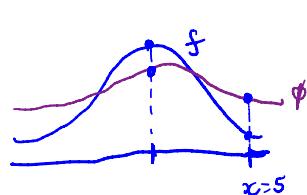
A rare draw from ϕ with much higher density under f than under ϕ will receive a huge weight and inflate the variance of the estimate.

opposite of our goal.

Strategy - choose prefuntion ϕ so that $\frac{f(x)}{\phi(x)}$ is large only when $g(x)$ is small.

Example 2.5

If we select and an appropriate ϕ ,



$\frac{f(0)}{\phi(0)}$ will be large.

$\frac{f(5)}{\phi(5)}$ will be small

We care about something out in the tails \Rightarrow this makes sense to try.

The importance sampling estimator can be shown to converge to θ under the SLLN so long as the support of ϕ includes all of the support of f .

↓
parameter
of interest

2.4 Compare to Previous Monte Carlo Approach

Common goal – estimate an integral $\int h(x)dx$.

Step 1 Do some derivations.

- a. Find an appropriate f and g to rewrite your integral as an expected value.

$$\theta = \int h(x)dx = \int g(x)f(x)dx = E[g(x)] \text{ wrt. } X \sim f.$$

- b. For importance sampling only,

Find an appropriate ϕ to rewrite θ as an expectation with respect to ϕ .

$$\theta = \int g(x) \frac{f(x)}{\phi(x)} \phi(x) dx = E\left[g(x) \frac{f(x)}{\phi(x)}\right], X \sim \phi$$

Step 2 Write pseudo-code (a plan) to define estimator and set-up the algorithm.

Naïve

- For Monte Carlo integration

1. Sample $X_1, \dots, X_m \sim f$

$$2. \hat{\theta} = \frac{1}{m} \sum_{i=1}^m g(X_i)$$

- For importance sampling

1. Sample $X_1, \dots, X_m \sim \phi$

$$2. \hat{\theta} = \frac{1}{m} \sum_{i=1}^m g(X_i) \frac{f(X_i)}{\phi(X_i)}$$

importance weights.

Step 3 Program it.

may also want to approximate $\text{Var}(\hat{\theta})$.

note
 $\phi(x) > 0$ when
 $f(x) > 0$
 Required.

2.5 Extended Example

In this example, we will estimate $\theta = \int_0^1 \frac{e^{-x}}{1+x^2} dx$ using MC integration and importance sampling with two different importance sampling distributions, ϕ . ① ② a) b)

STEP 1 Derive things.

a) Select $X \sim \text{Exp}(1)$. so $f(x) = \begin{cases} e^{-x} & x \geq 0 \\ 0 & \text{o.w.} \end{cases}$

$$\Rightarrow \text{find } g(x). \text{ s.t. } \theta = E[g(X)] = \int_0^\infty g(x) f(x) dx. \Rightarrow g(x) = \frac{1}{1+x^2} \mathbb{I}[x \leq 1]$$

$$\theta = \int_0^1 \frac{e^{-x}}{1+x^2} dx = \int_0^\infty \frac{e^{-x}}{1+x^2} \mathbb{I}(x \leq 1) dx = E\left[\frac{1}{1+x^2} \mathbb{I}[X \leq 1]\right], X \sim \text{Exp}(1).$$

Option ① MC integration w/ $X \sim \text{Exp}(1)$ (w/ 1b step).

Option ② Importance sampling w/

a) $\phi \sim \text{Unif}(0,1)$ $\phi_a(x) = \begin{cases} 1 & 0 \leq x \leq 1 \\ 0 & \text{o.w.} \end{cases}$

b) $\phi \sim \text{Exp}(1)$ rescaled to have support $0 \leq x \leq 1$,

$$\Rightarrow \phi_b(x) = \begin{cases} \frac{e^{-x}}{1+e^{-1}} & 0 \leq x \leq 1 \\ 0 & \text{o.w.} \end{cases}$$

you can check this is a valid pdf by making sure it integrates to 1 ($\phi_b(x) \geq 0$)

$$\int_0^1 \frac{e^{-x}}{1+e^{-1}} dx = 1 - e^{-1}$$

Step 1b Option 2a) $\phi_a(x) = 1 \quad 0 \leq x \leq 1$

$$\begin{aligned} \theta &= E_f[g(X)] = \int_0^\infty g(x) \frac{f(x)}{\phi_a(x)} \phi_a(x) dx \\ &= \int_0^\infty \frac{1}{1+x^2} \mathbb{I}(x \leq 1) \frac{e^{-x}}{1} \cdot 1 \mathbb{I}(0 \leq x \leq 1) dx. \\ &= E\left[\frac{1}{1+x^2} \mathbb{I}(X \leq 1) e^{-X}\right], X \sim \phi_a = \text{Uniform}(0,1). \end{aligned}$$

Option 2b) $\phi_b(x) = \frac{e^{-x}}{1+e^{-1}} \quad 0 \leq x \leq 1$.

$$\begin{aligned} \theta &= E_f[g(X)] = \int_0^\infty g(x) \frac{f(x)}{\phi_b(x)} \phi_b(x) dx = \int_0^\infty \frac{1}{1+x^2} \mathbb{I}(x \leq 1) \frac{-x}{e^{-x}/(1+e^{-1})} \cdot \frac{-x}{\frac{e^{-x}}{1+e^{-1}}} \mathbb{I}(0 \leq x \leq 1) dx \\ &= E\left[\frac{1}{1+x^2} \mathbb{I}(X \leq 1)\right], X \sim \phi_b \end{aligned}$$

STEP 2 Make a plan.

Option 1:

1. Sample X_1, \dots, X_m from $\text{Exp}(1)$.

$$2. \hat{\theta} = \frac{1}{m} \sum_{i=1}^m \left[\frac{1}{1+x_i^2} \cdot \mathbb{I}[X_i \leq 1] \right]$$

Option 2a:

1. Sample X_1, \dots, X_m from $\text{Unif}[0,1]$.

$$2. \hat{\theta} = \frac{1}{m} \sum_{i=1}^m \frac{1}{1+x_i^2} e^{-x_i} \mathbb{I}[X_i \leq 1]$$

This will always hold.

Option 2b:

1. Sample X_1, \dots, X_m from $\phi_b(x) = \frac{e^{-x}}{1+e^{-x}}$, $x \in [0,1]$.

Note this ϕ_b is not a named distribution. \Rightarrow We need to sample from it.

Inverse transform!

$$F_{\phi_b}(x) = \int_0^x \frac{e^{-y}}{1+e^{-y}} dy = -\left. \frac{e^{-y}}{1+e^{-y}} \right|_0^x = \begin{cases} 0 & x < 0 \\ \frac{1-e^{-x}}{1+e^{-x}} & \text{for } x \in [0,1] \\ 1 & x > 1 \end{cases}$$

↳ Find cdf \rightarrow take inverse.

$$u = F_{\phi_b}(x) = \frac{1-e^{-x}}{1+e^{-x}}$$

$$u(1+e^{-1}) = 1-e^{-x}$$

$$e^{-x} = 1-u(1+e^{-1})$$

$$-x = \log(1-u(1+e^{-1}))$$

$$\hat{F}(u) = x = -\log(1-u(1+e^{-1}))$$

$$= -\log(1-u(1-e^{-1}))$$

Algorithm for option 2b

i. Sample $x_1, \dots, x_m \stackrel{iid}{\sim} \phi_b$ using inverse transform method.

a) Sample $U_1, \dots, U_m \stackrel{iid}{\sim} \text{Unif}(0,1)$.

b) Set $x_i = -\log(1 - U_i / (1 + e^{-1}))$ for $i=1, \dots, m$.

$$2. \hat{\theta} = \frac{1}{m} \sum_{i=1}^m \frac{1 + e^{-1}}{1 + x_i^2} \mathbb{I}(x_i \leq 1)$$

this will always hold by definition.

which will be the best? We can compare $h(x) = f(x)g(x)$ to $f(x)$, $\phi_a(x)$, and $\phi_b(x)$.

can look at $\frac{f(x)g(x)}{f(x)}$, $\frac{f(x)g(x)}{\phi_a(x)}$, and $\frac{f(x)g(x)}{\phi_b(x)}$

want to pick the one that is closest to constant.

This will give us the lowest variance.