# Monte Carlo methods

## Stat 580

# Monte Carlo integration

# Monte Carlo integration

- a numerical approximation for expectation

- often useful for multidimensional problems that require the estimation of $\mu = E\{h(X)\}$, where $X$ is a random vector and $h$ is a function

- simplest approach: approximate $\mu$ by $\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^{m} h(X_i)$ where $X_1, \ldots, X_n$ are iid copies of $X$

- properties:

  - $\hat{\mu}_m$ is consistent: by the Strong Law of Large Numbers, with probability 1, $\hat{\mu}_m \to \mu$ as $m \to \infty$

  - $\hat{\mu}_m$ is unbiased: $E(\hat{\mu}_m) = \mu = E\{h(X)\}$.

  - $\text{Var}(\hat{\mu}_m) = \text{Var}\{h(X)\}/m$ and can be estimated by

$$\widehat{\text{Var}}(\hat{\mu}_m) = \frac{1}{m(m-1)} \sum_{i=1}^{m} \{h(X_i) - \hat{\mu}_m\}^2.$$

# Monte Carlo integration

- we will see methods:

  1. that are applicable when $X_1, \ldots, X_m$ cannot be easily sampled

  2. to reduce $\mathrm{Var}(\hat{\mu}_m)$

- MC integration can be used to evaluate a "usual" integral $I = \int_{\mathcal{X}} H(x)dx$

  1. the idea is to "factorize" $H(x) = f(x)h(x)$ with $f(x)$ as a pdf with support $\supseteq \mathcal{X}$ (we take $h(x) = 0$ if $x \notin \mathcal{X}$.)

  2. approximate $I$ by $\frac{1}{m} \sum_{i=1}^{m} h(X_i)$, where $X_1, \ldots, X_m$ are iid with pdf $f(x)$

# Example

We want to compute

$$\int_{-\infty}^{\infty} \log|x| e^{-\frac{(x+1)^2}{8}} \, dx.$$

Set

$$f(x) = \frac{1}{\sqrt{8\pi}} e^{-\frac{(x+1)^2}{8}}$$

and

$$h(x) = \sqrt{8\pi} \log|x|.$$

Note that $f(x)$ is the pdf for $\mathcal{N}(-1, 4)$, so the integral can be approximated by $\frac{\sqrt{8\pi}}{n} \sum_{i=1}^{m} \log|X_i|$ with $\{X_i\}$ iid $\sim \mathcal{N}(-1, 4)$.

# Example

In a particular missing data problem, suppose $Y \sim p_Y(y|\boldsymbol{\theta})$ but only $X = M(Y)$ is observed, where $M$ is a (known and nonrandom) many-to-fewer mapping.

## Frequentist approach:

To estimate $\boldsymbol{\theta}$ by the MLE, you need the likelihood

$$L(\boldsymbol{\theta}|\boldsymbol{x}) = p_X(\boldsymbol{x}|\boldsymbol{\theta}) = \int_{M(y)=x} p_Y(y|\boldsymbol{\theta})dy.$$

- This integral is often hard to compute analytically

- Using MC integration

  - one designs a suitable function $h$ and a distribution $P$
  - then sample $Y_1, \ldots, Y_n$ from $P$ and get

$$L(\boldsymbol{\theta}|\boldsymbol{x}) \approx \frac{1}{m} \sum_{i=1}^{m} h(Y_i).$$

# Example

## Bayesian approach:

- Suppose $\theta$ has a prior $\pi(\theta)$.

- We are interested in the posterior distribution of $\theta$.

  - If the full data $Y = y$ is observed, by Bayes formula, the posterior is

    $$p_{\theta|Y}(\theta|y) \propto p_Y(y|\theta)\pi(\theta).$$

  - If only $X = x$ is observed, the posterior is

    $$p_{\theta|X}(\theta|x) = \int p_{\theta|Y}(\theta|y)p_{Y|M(Y)}(y|x)dy.$$

    - If $p_{\theta|Y}(\theta|y)$ has a closed form, a MC integration method may sample $Y_1, \dots, Y_m \sim p_{Y|M(Y)}(y|x)$ and approximate

      $$p_{\theta|X}(\theta|x) \approx \frac{1}{m} \sum_{i=1}^{m} p_{\theta|Y}(\theta|Y_i).$$

# Importance sampling

# Importance sampling

- Same setup: want to estimate

$$\mu = \int h(x)f(x)dx, \quad f(x) \text{ is a pdf.}$$

  (the integral is taken over the region where the integrand is positive)

  - but it is difficult to sample from $f$

- Rewrite

$$\mu = \int h(x)\frac{f(x)}{g(x)}g(x)dx,$$

  where $g$ is a pdf such that $g(x) > 0$ whenever $f(x) > 0$.

- Let $X$ have density $g(x)$.

- Then $\mu = E\{h(X)w^*(X)\}$ with $w^*(X) = \frac{f(X)}{g(X)}$.

# Importance sampling

- Consider the following steps:

  1. Generate $X_1, \ldots, X_m$ iid $\sim g(x)$.

  2. Estimate $\mu$ by $\hat{\mu}_g = \frac{1}{m} \sum_{i=1}^{m} h(X_i) w^*(X_i)$.

- $\hat{\mu}_g$ is the importance sampling (IS) estimator of $\mu$ associated with $g$.

- $w^*(X_i)$'s are referred to as importance ratios

- Note that $\hat{\mu}_g$ is a weighted sum of $h(X_i)$.

- If $f = g$, then $\hat{\mu}_g$ is the ordinary MC estimator.

# Example

We want to compute $\mu = E(U^5)$ where $U \sim \text{Unif}(0, 1)$.

- the straightforward MC estimator: $\hat{\mu} = \frac{1}{m} \sum_{i=1}^{m} U_i^5$

  - oversample data $U_i^5$ near the origin and undersample the data near 1

  - $\text{Var}(\hat{\mu}) = 0.0631/m$

- Use IS to put more weights near 1:

  - use $g(x) = 5x^4$ for $0 < x < 1$

  - the IS estimator: $\hat{\mu}_g = \frac{1}{n} \sum_{i=1}^{m} X_i^5 w^*(X_i)$ where $w^*(X_i) = 1/(5X_i^4)$

  - $\text{Var}(\hat{\mu}_g) = 0.00794/m$ (verify!)

  - resulting a variance reduction of 98.74%!

- the IS can be used as a variance reduction technique!

# Properties

1. $\hat{\mu}_g$ is unbiased for $\mu$.

2. $\mathrm{Var}(\hat{\mu}_g) = \frac{1}{m}\mathrm{Var}\{h(X)w^*(X)\}$ .

- To reduce the variance of $\hat{\mu}_g$, $g(x)$ should be in proportion to $h(x)f(x)$ as much as possible.

# Properties

To show this, we need to reduce

$$E\left[\left\{h(X)\frac{f(X)}{g(X)}\right\}^2\right] = \text{Var}(\hat{\mu}_g) + \mu^2.$$

We can use

$$E\left[\left\{h(X)\frac{f(X)}{g(X)}\right\}^2\right] \geq \left[E\left\{h(X)\frac{f(X)}{g(X)}\right\}\right]^2.$$

The equality holds if $h(x)\frac{f(x)}{g(x)}$ is a constant. That is, when $g(x) \propto h(x)f(x)$. (Why is it called importance sampling?)

# When $f$ is only known up to a constant

- That means, $f(x) = cq(x)$ with $c > 0$ unknown, then

$$\mu = \frac{E\{h(X)w^*(X)\}}{E\{w^*(X)\}}$$

with $w^*(X) = \frac{q(X)}{g(X)}$.

- In this case, standardized weights have to be used in IS:

  1. Generate $X_1, \ldots, X_m$ iid from $g(x)$.

  2. Estimate $\mu$ by $\hat{\mu}_g = \frac{1}{m} \sum_{i=1}^{m} h(X_i)w(X_i)$ where $w(X_i) = \frac{w^*(X_i)}{\sum_{i=1}^{m} w^*(X_i)}$.

# Processes

- Let $\{S_1, S_2, \ldots\}$ be a discrete time stochastic process.

- For example, $S_i$ = asset price $S(t_i)$ at time $t_i$.

- In general, to estimate $\mu = E\{h(S_1, \ldots, S_n)\}$, one treats $(S_1, \ldots, S_n)$ as a random vector $S$ with density $f$.

- Then IS estimates $\mu$ by $\hat{\mu}_g = \frac{1}{m} \sum_{i=1}^{m} h(X_i)\frac{f(X_i)}{g(X_i)}$, where $X_1, \ldots, X_m$ are iid with density $g(x)$, and $g(x) > 0$ whenever $f(x) > 0$.

- It is useful to think of $X_i = (X_{i,1}, \ldots, X_{i,n})$ as a sample path of a process (up to the $n$-th index).

- From this point, IS uses a process different from $S$ to estimate $\mu$.

- An important issue is how to design the auxiliary process to take advantage of any special dependence structure of $S_1, S_2, \ldots$.

# Markov processes

- $S_1, \ldots, S_n$ form a Markov process; i.e., we can decompose their joint density as

$$f(x_1, \ldots, x_n) = f_1(x_1) \prod_{j=2}^{n} f_j(x_j | x_{j-1}),$$

  where $f_1(x_1)$ is the density of the initial distribution and $f_j(x_j | x_{j-1})$ is the transition density at step $j$.

- To utilize the Markov property in IS, the auxiliary density $g(x_1, \ldots, x_n)$ may also have a similar structure as $f$:

$$g(x_1, \ldots, x_n) = g_1(x_1) \prod_{j=2}^{n} g_j(x_j | x_{j-1})$$

- Now in IS, the weight for $X = (X_1, \ldots, X_n)$ is

$$w(X) = \frac{f_1(X_1)}{g_1(X_1)} \prod_{j=2}^{n} \frac{f_j(X_j | X_{j-1})}{g_j(X_j | X_{j-1})}.$$

# Markov processes

- The IS in this case is, for $i = 1, \ldots, m$:

  1. Draw $X_1 \sim g_1$.

  2. For $j = 2, \ldots, n$, draw $X_j \sim g_j(x_j | x_{j-1})$.

  3. Set $h_i = h(X_1, \ldots, X_n)$ and

  $$w_i = \frac{f_1(X_1)}{g_1(X_1)} \prod_{j=2}^{n} \frac{f_j(X_j | X_{j-1})}{g_j(X_j | X_{j-1})}.$$

- The IS estimator is

  $$\hat{\mu}_g = \frac{1}{m} \sum_{i=1}^{m} h_i w_i.$$

# Control variates

# Control variates

- We still want to compute $\mu = E\{h(X)\}$.

- Suppose we know the exact value of $\theta = E\{c(Y)\}$, where $c$ is a function of another random variable $Y$.

- The simple MC estimators for $\mu$ and $\theta$ are, respectively,

$$\hat{\mu}_{\mathrm{MC}} = \frac{1}{n} \sum_{i=1}^{n} h(X_i), \qquad \hat{\theta}_{\mathrm{MC}} = \frac{1}{n} \sum_{i=1}^{n} c(Y_i).$$

- Of course, for $\theta$, $\hat{\theta}_{\mathrm{MC}}$ is unnecessary. However, it can be helpful for the estimation of $\mu$.

- How? Suppose $h(X)$ and $c(Y)$ are positively correlated. (If not, this method is not applicable.)

# Control variates

- If we see $\hat{\theta}_{MC} > \theta$, then due to the positive correlation, $\hat{\mu}_{MC}$ is more likely to be $> \mu$.

- Then we can decrease the value of $\hat{\mu}_{MC}$ to obtain a better estimate.

- To be specific, suppose we can sample $(X_1, Y_1), \ldots, (X_n, Y_n)$ iid.

- The control variate estimator for $\mu$ is

$$\hat{\mu}_{CV} = \hat{\mu}_{MC} - b(\hat{\theta}_{MC} - \theta),$$

  where $b$ is a constant.

- $\hat{\mu}_{CV}$ is unbiased and consistent (as MC estimators are unbiased and consistent).

- If $b = 0$, then $\hat{\mu}_{CV} = \hat{\mu}_{MC}$, so the optimal control variate estimator is at least as good as $\hat{\mu}_{MC}$.

- Given a control variate $c(Y)$, need to choose $b$ (choosing $c(Y)$ is harder)

# Choice of $b$

- How should we choose $b$?

  ○ for any given $b$,

  $$\mathrm{Var}(\hat{\mu}_{\mathrm{CV}}) = \mathrm{Var}(\hat{\mu}_{\mathrm{MC}}) + b^2 \mathrm{Var}(\hat{\theta}_{\mathrm{MC}}) - 2b\mathrm{Cov}(\hat{\mu}_{\mathrm{MC}}, \hat{\theta}_{\mathrm{MC}}).$$

  ○ the minimum of $\mathrm{Var}(\hat{\mu}_{\mathrm{CV}})$ happens when $b = b^*$, where

  $$b^* = \frac{\mathrm{Cov}(\hat{\mu}_{\mathrm{MC}}, \hat{\theta}_{\mathrm{MC}})}{\mathrm{Var}(\hat{\theta}_{\mathrm{MC}})} = \frac{\mathrm{Cov}\{h(X), c(Y)\}}{\mathrm{Var}\{c(Y)\}}.$$

  ○ in practice $b^*$ is unknown, but we can estimate it.

  ○ plug $b^*$ into $\mathrm{Var}(\hat{\mu}_{\mathrm{CV}})$ to get $\mathrm{Var}(\hat{\mu}_{\mathrm{CV}}^{\mathrm{opt}})$, and we can show the variance reduction factor is

  $$\frac{\mathrm{Var}(\hat{\mu}_{\mathrm{CV}}^{\mathrm{opt}})}{\mathrm{Var}(\hat{\mu}_{\mathrm{MC}})} = 1 - \rho^2,$$

  where $\rho$ is the correlation coefficient between $h(X)$ and $c(Y)$.

# Estimation of $b^*$

- the optimal $b^*$ can be estimated by

$$\hat{b}_n = \frac{\widehat{\text{Cov}}\,(\hat{\mu}_{\text{MC}}, \hat{\theta}_{\text{MC}})}{\widehat{\text{Var}}\,(\hat{\theta}_{\text{MC}})},$$

where

$$\widehat{\text{Var}}\,(\hat{\theta}_{\text{MC}}) = \frac{1}{n(n-1)} \sum_{i=1}^{n} \left\{ c(Y_i) - \hat{\theta}_{\text{MC}} \right\}^2$$

and

$$\widehat{\text{Cov}}\,(\hat{\mu}_{\text{MC}}, \hat{\theta}_{\text{MC}}) = \frac{1}{n(n-1)} \sum_{i=1}^{n} \{ h(X_i) - \hat{\mu}_{\text{MC}} \} \{ c(Y_i) - \hat{\theta}_{\text{MC}} \}.$$

- $\hat{b}_n$ is the slope of the least-squares regression line for $(h(X_i), c(Y_i))$, $i = 1, \dots, n$

# General idea

- Overall, the general idea for the control variate method is to search $c(Y)$ such that

  1. $E\{c(Y)\}$ is known.

  2. the scatterplot of $(h(X_i), c(Y_i))$ shows strong correlation.

- In practice, $\hat{\mu}_{\text{MC}}$ and $\hat{\theta}_{\text{MC}}$ often depend on the same random variable; i.e., $Y_i = X_i$.

- It is possible to use more than one control variate; i.e.,

$$\hat{\mu}_{\text{CV}} = \hat{\mu}_{\text{MC}} - b_1(\hat{\theta}_{1,\text{MC}} - \theta_1) - b_2(\hat{\theta}_{2,\text{MC}} - \theta_2).$$

# Example

Let $\mu = E(e^U)$ where $U \sim \text{Unif}(0, 1)$. Theoretical study of CV estimator with $b^*$ (when $n = 1$):

- Use $U$ as the control variate

- $E(U) = 1/2, \text{Cov}(e^U, U) = 1 - (e - 1)/2 = 0.14086$ and $\text{Var}(U) = 1/12$

- the CV estimator: $\hat{\mu}_{CV} = e^U - b^*(U - 1/2)$, where $b^* = 12(0.14086)$

- $\text{Var}(\hat{\mu}_{CV}) = 0.0039$ (verify!)

- resulting a variance reduction of 98.4% when compared to $Var(e^U) = 0.242$

# Monte Carlo Testing

# Monte Carlo testing

- We use Monte Carlo testing when

    - we cannot calculate the null distribution of the test statistic

    - but we can simulate from the null hypothesis.

- We "simulate" the null distribution of the test statistic.

# Example

Let $X_1, \ldots, X_n \sim \mathcal{N}(\mu, \sigma^2)$. Suppose we want to test $H_0 : \mu = 3, \sigma^2 = 4$ against the alternative $H_A : \mu \neq 3, \sigma^2 = 4$.

- This is a simple $z$-test. The test statistic is

$$\bar{X} = \frac{X_1 + X_2 + \cdots + X_n}{n}.$$

- To use Monte Carlo testing,

  - we generate $X_1^*, \ldots, X_n^*$ from $\mathcal{N}(3, 4)$, the distribution in the null hypothesis.

  - We then calculate $\bar{X}^* = \frac{\sum X_i^*}{n}$ and repeat this procedure $m$ times (e.g., set $m = 999$) to get $m$ $\bar{X}^*$ values.

  - This approximates the distribution of the test statistic using $m$ simulated values of $\bar{X}$.

  - If $\bar{X}$ is amongst the smallest 2.5% or the largest 2.5% of these $\bar{X}^*$ values, we reject $H_0$.