

Lecture 1.

Chp 1.1, 1.2. Liu's Book.

Note Title

4/2/2006

Goal: Teach Monte Carlo methods with some related material on optimization.

Text Book: "Monte Carlo Strategies in Scientific Computing". J.S. Liu.

Grading: 4 Homework Assignments + Final.

Motivation: Many problems in Statistics can be formulated as probabilistic inference, or as optimization.

$$\text{Find } \hat{\underline{x}} = \underset{\underline{x}}{\text{ARG-MIN}} E(\underline{x}) \quad \text{MIN } E(\underline{x}) = E(\hat{\underline{x}})$$
$$\hat{\underline{x}} = \underset{\underline{x}}{\text{ARG-MAX}} P(\underline{x} | \underline{d})$$

Or as evaluating an integral

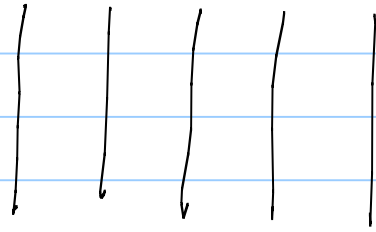
$$I = \int_D g(\underline{x}) d\underline{x}$$

Monte Carlo gives a way to perform this
History: Developed after the 2nd world war. The Manhattan project. Computers available.

Page 2 Count de Buffon (1707-1788) use sampling to estimate π .
Early Sampling Experiment.

Needle length $l < D$

Drop needle at
random.



parallel lines

spacing D between lines

probability that needle will intersect a
line is $\frac{2l}{\pi D}$. (check $\frac{2l}{\pi D} < 1$, because $l < D$).

Let p_N be the proportion of "intersects"
in N samples (i.e. drop the needle N
times, count no. times it intersects line - say M ,
set $p_N = M/N$).

$$\lim_{N \rightarrow \infty} p_N = \frac{2l}{\pi D}$$

$$\text{Hence } \pi = \lim_{N \rightarrow \infty} \frac{2l}{p_N D} \quad //$$

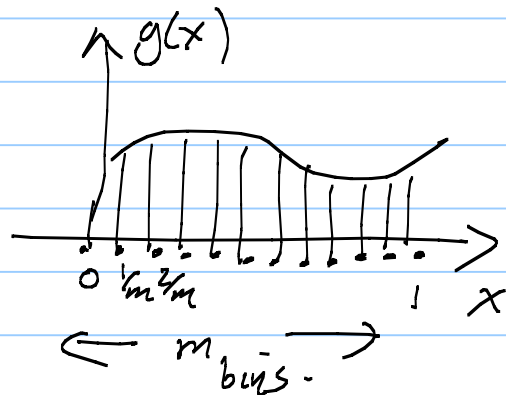
Note: Sampling has been around for 1,000's of years.
For example, Banks/Kings count number of coins in a pile by selecting
a few coins at random, weighing them; then weigh the pile to estimate total
number of coins.

Page 3

Integration

$$I = \int_0^1 g(x) dx$$

$$D = [0, 1]$$



Riemannian Approximation

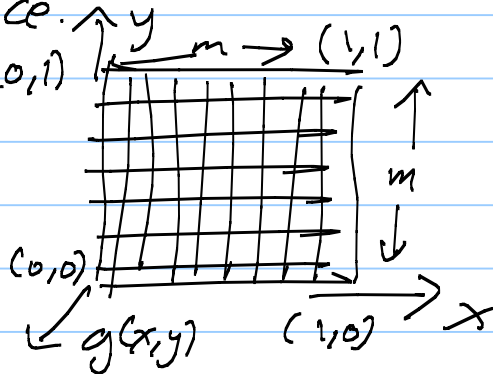
Approximate I by $\tilde{I}_m = \frac{1}{m} \{g(1/m) + g(2/m) + \dots + g(1)\}$

The typical error is $O(m^{-1})$.
As $m \rightarrow \infty$, $\tilde{I}_m \rightarrow I$. Not so bad.

But no. of bins increases exponentially with the dimensionality of the space.

In 2-dimensions
 $D = [0, 1]^2$

$$\tilde{I}_m = \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m g\left(\frac{i}{m}, \frac{j}{m}\right)$$



$\tilde{I}_m \rightarrow I$, as $m \rightarrow \infty$

But requires evaluating m^2 points to get error $O(m^{-1})$

In n -dimensions, to get $O(m^{-1})$ error requires evaluating m^n points \rightarrow too many!

Page 4.

Monte Carlo approximation

$$I = \int_{\Omega} \pi(x) g(x) dx, \quad \underline{x} \text{ in } n\text{-dimensional space}$$

Note: If we let $\pi(x)$ be the uniform distribution

$$\pi(\underline{x}) = 1/|\Omega|, \quad \underline{x} \in \Omega \quad \text{so that } \int_{\Omega} \pi(\underline{x}) d\underline{x} = 1$$

then $I = \frac{1}{|\Omega|} \int_{\Omega} g(\underline{x}) d\underline{x}$. (E.g. if $\Omega = [0,1]^2$, $\pi(x,y) = 1$ for $0 \leq x \leq 1, 0 \leq y \leq 1$.)

Use Monte Carlo (MC) to draw m independent & identically distributed (i.i.d) samples from $\pi(x)$:

$$\underline{x}^{(1)}, \dots, \underline{x}^{(m)}$$

Approximate I by $\hat{I}_m = \frac{1}{m} \{g(\underline{x}^{(1)}) + \dots + g(\underline{x}^{(m)})\}$

By the law of large numbers

$$\lim_{m \rightarrow \infty} \hat{I}_m = I, \quad \text{with probability 1.}$$

($\forall \epsilon > 0, \lim_{m \rightarrow \infty} P(|\hat{I}_m - I| > \epsilon) = 0$)

By central limit theorem: $\hat{I}_m = I + \epsilon/\sqrt{m}$, $\epsilon \sim N(0, \sigma^2)$

$$\sqrt{m}(\hat{I}_m - I) \rightarrow N(0, \sigma^2) \quad \text{— error is } O(m)^{-1/2}$$

where $N(0, \sigma^2)$ is a zero-mean Gaussian with variance σ^2 :

$$\sigma^2 = \int_{\Omega} d\underline{x} \pi(\underline{x}) (g(\underline{x}) - \bar{g})^2, \quad \text{with } \bar{g} = \int_{\Omega} d\underline{x} \pi(\underline{x}) g(\underline{x})$$

Note: error is $O(m^{-1/2})$ with only m computations and independent of the dimensions of the space.

Page 5.

The Monte Carlo approach gives better error rates with far fewer samples!

But

(1.) We need to be able to draw samples from $\pi(x)$ — this is not easy (the main purpose of the course is to show how to sample from probability distributions.)

(2.) The variance σ^2 may be very large. Recall Error $\sim \sigma/\sqrt{m}$.

There are ways to reduce the variance
E.G. Importance Sampling. Suppose $\pi(x)$ is uniform.
Then get i.i.d samples $x^{(1)}, \dots, x^{(m)}$ from distribution $\pi_1(x)$ that puts more probability on important parts of D .

Estimate
$$\hat{I}_m = \frac{1}{m} \sum_{j=1}^m \frac{g(x^{(j)})}{\pi_1(x^{(j)})} \quad \left| \quad \mathbb{E}_{\pi_1} \left(\frac{g(x)}{\pi_1(x)} \right) = \int_D g(x) dx$$

As before
$$\sqrt{m} (\hat{I}_m - I) \rightarrow N(0, \bar{\sigma}_{\pi_1}^2)$$
 ← constant

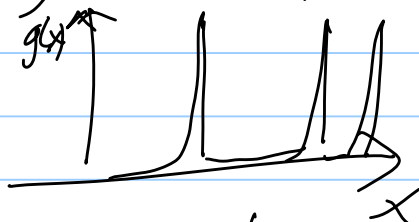
with
$$\bar{\sigma}_{\pi_1}^2 = \text{var}_{\pi_1} \left(\frac{g(x)}{\pi_1(x)} \right)$$
. If $\pi_1(x) = K g(x)$
then $\bar{\sigma}_{\pi_1}^2 = 0$, ideal!

(6) Best distribution to sample from is $g(x)$ (if $\pi(x)$ is uniform) but this may not be possible.

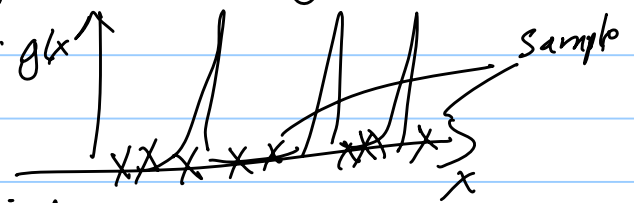
|| Some distributions are easy to sample from but others are very difficult. ||

Notice that if $\pi_1(x)$ and $g(x)$ are very different, then the variance $\bar{g}_\pi^2 = \text{var}_{\pi_1} \left(\frac{g(x)}{\pi_1(x)} \right)$ may be huge — so sampling might be even worse than using the Riemann approximation.

For example, in 1-D suppose $g(x)$ is very 'spiky' and you sample from a uniform distribution $\pi_1(x) = 1/D$



There is very low probability that the samples from $\pi_1(x)$ will lie on the spikes — so your estimate of $\int g(x) dx$ will often be bad.



Moral \rightarrow you need to pick a sampling distribution $\pi_1(x)$ which is close to $g(x)$.

So the "miracle" of MC — the ability to estimate an integral with error independent of dimension — occurs only if you know a lot about the function $g(x)$ that you want to integrate. (described in more detail in other lectures)