Importance Sampling

Used for a number of purposes:

- Variance reduction

- Allows for difficult distributions to be sampled from.

- Sensitivity analysis

- Reusing samples to reduce computational burden.

Idea is to sample from a different distribution that picks points in "important" regions of the sample space.

Want

$$E\left[f\left(X\right)\right] = \int f\left(x\right)g\left(x\right)dx$$

Instead of sampling from density (or probability mass function) $g\left(x\right)$, sample from a distribution with density (or pmf) $h\left(x\right)$.

Since we are sampling from the "wrong" distribution we have to make adjustments in our estimator.

$$E_g\left[f(X)\right] = \int f(x)g(x)\,dx$$

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

$$= E_h\left[f(X)\frac{g(X)}{h(X)}\right]$$

This suggests the following estimation scheme

1) Sample $x_1,\ldots,x_n$ from $h(x)$.

2) Calculate weights

$$w_i = \frac{g(x_i)}{h(x_i)}$$

3) Use estimator

$$\hat{\mu}_{f,IS} = \frac{1}{n}\sum_{i=1}^{n}w_i f(x_i)$$

So instead of a regular average, this estimator is a weighted average.

So points that occur more often under $h(x)$ than $g(x)$ get downweighted and those that occur less often get upweighted.

Notice that $\hat{\mu}_{f,IS}$ is an unbiased estimate of $E_g\left[f(X)\right]$ regardless of which proposal distribution $h(x)$ as long as $h(x)$ has the same support as $g(x)$, i.e.

$$g(x) > 0 \text{ implies that } h(x) > 0$$

Note that $h(x) > 0$ can be allowed to occur when $g(x) = 0$, though doing this tends to be inefficient (but there are times you want to do this).

Since $\hat{\mu}_{f,IS}$ is unbiased, the main idea is to pick a distribution $h(x)$ that reduces the variance.

$$\text{Var}_h\left(\frac{f(X)g(X)}{h(X)}\right) = E_h\left[\left(\frac{f(X)g(X)}{h(X)}\right)^2\right] - \mu_f^2$$

To do this, we want $h(x)$ to look like $f(x)g(x)$, i.e. make

$$\frac{f(x)g(x)}{h(x)}$$

look like a constant.

The optimal $h(x)$ satisfies

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

Note that this usually can't be determined, due to the normalizing constant.

However this does give us a motivation for picking $h(x)$.

Example: Monte Carlo Evaluation of a Likelihood Ratio (Genetics Example)

Assume that you have a missing data model where $X = (X_{obs}, X_{mis})$. Then the observed data likelihood ratio satistifies

$$l(\theta_1, \theta_0) = \frac{L(\theta_1)}{L(\theta_0)} = \frac{p_{\theta_1}(X_{obs})}{p_{\theta_0}(X_{obs})}$$

$$= E_{\theta_0} \left[ \frac{p_{\theta_1}(X_{obs}, X_{mis})}{p_{\theta_0}(X_{obs}, X_{mis})} \Big| X_{obs} \right]$$

This can be estimated by sampling $z_1, \ldots, z_n$ from $p(X_{mis} | X_{obs})$ calculating

1)  $f(z_i) = \dfrac{p_{\theta_1}(X_{obs}, z_i)}{p_{\theta_0}(X_{obs}, z_i)}$

2)  $\hat{l}(\theta_1, \theta_0) = \dfrac{1}{n} \sum_{i=1}^{n} f(z_i)$

Suppose that you are interested in getting $l(\theta_2, \theta_0)$, based on this Monte Carlo estimate.

This can be done with the importance sampling estimate

$$\hat{l}(\theta_2, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} f(z_i) \frac{p_{\theta_2}(X_{obs}, z_i)}{p_{\theta_1}(X_{obs}, z_i)}$$

This can be shown to be an unbiased estimator of $l(\theta_2, \theta_0)$.

Genetics example:

Observed Data Model

$$(Y_1, Y_2, Y_3, Y_4) \sim \text{Multi}\left(197, \left(\frac{\lambda}{4}, \frac{1-\lambda}{4}, \frac{1-\lambda}{4}, \frac{2+\lambda}{4}\right)\right)$$

$$g(Y|\lambda) = \left(\frac{\lambda}{4}\right)^{Y_1} \left(\frac{1-\lambda}{4}\right)^{Y_2+Y_3} \left(\frac{2+\lambda}{4}\right)^{Y_4}$$

Complete Data Model

$$(X_1, X_2, X_3, X_4, X_5)$$

$$\sim \text{Multi}\left(197, \left(\frac{\lambda}{4}, \frac{1-\lambda}{4}, \frac{1-\lambda}{4}, \frac{\lambda}{4}, \frac{1}{2}\right)\right)$$

$$g(X|\lambda) = \left(\frac{\lambda}{4}\right)^{X_1+X_4} \left(\frac{1-\lambda}{4}\right)^{X_2+X_3} \left(\frac{1}{2}\right)^{X_5}$$

As seen before $X_4 | Y_4 \sim \text{Bin}\left(Y_4, \frac{\lambda}{2+\lambda}\right)$

The complete data likelihood ratio satisfies

$$\frac{g(Y,X|\lambda_1)}{g(Y,X|\lambda_0)} = \left(\frac{\lambda_1}{\lambda_0}\right)^{Y_1+X_4} \left(\frac{1-\lambda_1}{1-\lambda_0}\right)^{Y_2+Y_3}$$

Note that this implies the importance sample weight satisfies

$$w_i = c(Y_1, Y_2, Y_3, \theta_2, \theta_1) \left(\frac{\theta_2}{\theta_1}\right)^{z_i}$$

In this case $\hat{l}(\lambda_2, \lambda_0)$ has the form

$$\hat{l}(\lambda_2, \lambda_0) = \frac{c(Y_1, Y_2, Y_3, \lambda_2, \lambda_1)}{n} \sum_{i=1}^{n} f(z_i) \left(\frac{\lambda_2}{\lambda_1}\right)^{z_i}$$

As in many problems, the desired sampling distribution doesn't need to be known exactly, but only up to the normalizing constant (i.e. $l(x) = cg(x)$).

Importance sampling still works fine in this case.

1) Sample $x_1, \ldots, x_n$ from $h(x)$.

2) Calculate weights

$$w_i = \frac{l(x_i)}{h(x_i)}$$

3) Use estimator

$$\hat{\mu}_{f,IS} = \frac{1}{W} \sum_{i=1}^{n} w_i f(x_i)$$

$$= \frac{1}{n\bar{w}} \sum_{i=1}^{n} w_i f(x_i)$$

where $W = \sum w_i$.

Properties of this estimator

$$E_h \left[ w(X) f(X) \right] = \int \frac{l(x)}{h(x)} f(x) h(x) dx$$

$$= \int f(x) c g(x) dx$$

$$= c E_g \left[ f(X) \right]$$

$$E_h \left[ w \right] = \int \frac{l(x)}{h(x)} h(x) dx$$

$$= \int c g(x) dx = c$$

As this is a ratio estimator, it is no longer unbiased, but it is consistent.

In addition, when $c$ is known, this estimator is often preferred to the unbiased one discussed last time as it often has a smaller mean square error (to be discussed later).

Efficiency of importance sampling

Effective sample size

$$\text{ESS}(n) = \frac{n}{1 + \text{var}_h \left( w(X) \right)}$$

Since the weights are usually only known up to the normalizing constant, $\text{var}_h\big(w(X)\big)$ needs to be estimated by the coefficient of variation of the unnormalized weights

$$\text{cv}^2(w) = \frac{\sum_{i=1}^{n}(w_i - \bar{w})^2}{(n-1)\,\bar{w}^2}$$

Assume for what follows that $c = 1$ (and this is known).

We have two estimators

$$\tilde{\mu} = \frac{1}{n}\sum_{i=1}^{n} w_i f(x_i) \qquad\qquad \text{(unbiased)}$$

$$\hat{\mu} = \frac{1}{n\bar{w}}\sum_{i=1}^{n} w_i f(x_i) \equiv \frac{\tilde{\mu}}{\bar{w}} \qquad \text{(ratio)}$$

Let $Z = w(X)f(X)$. Then by the delta method

$$E_h\big[\hat{\mu}\big] \approx E_h\left[\bar{Z}\Big(1 - (\bar{w}-1) + (\bar{w}-1)^2\Big)\right]$$

$$\approx \mu - \frac{\text{cov}_h(w, Z)}{n} + \frac{\mu\,\text{Var}_h(W)}{n}$$

In addition,

$$\text{Var}_h\left(\hat{\mu}\right)$$

$$\approx \frac{1}{n}\left(\mu^2 \text{Var}_h\left(w\right) + \text{Var}_h\left(Z\right) - 2\mu \text{Cov}_h\left(w, Z\right)\right)$$

For the unbiased estimator,

$$E_h\left[\tilde{\mu}\right] = \mu \text{ and } \text{Var}_h\left(\tilde{\mu}\right) = \text{Var}_h\left(Z\right)/n$$

Thus

$$\text{MSE}\left(\tilde{\mu}\right) = \frac{\text{Var}_h\left(Z\right)}{n}$$

and

$$\text{MSE}\left(\hat{\mu}\right) = \left(E_h\left[\hat{\mu}\right] - \mu\right)^2 + \text{Var}_h\left(\hat{\mu}\right)$$

$$= \frac{\text{MSE}_h\left(\tilde{\mu}\right)}{n}$$

$$+ \frac{1}{n}\left(\mu^2 \text{Var}_h\left(w\right) - 2\mu \text{Cov}_h\left(w, Z\right)\right) + O\left(n^{-2}\right)$$

Thus the ratio estimate is to be preferred when

$$2\mu \text{Cov}_h\left(w, Z\right) > \mu^2 \text{Var}_h\left(w\right)$$

(assuming $\mu > 0$). That is when $w(X)$ and $w(X)f(X)$ are highly correlated.

In addition, the formula for $\text{Var}_h(\hat{\mu})$ implies (whole bunch of steps omitted)

$$\text{Var}_h(\tilde{\mu}) \approx \text{Var}_g(f(X))\big(1 + \text{Var}_h(w(X))\big)\big/n$$

Rearranging this gives

$$\frac{\text{Var}_g(f(X))}{\text{Var}_h(w(X)f(X))} \approx \frac{1}{1 + \text{Var}_h(w(X))}$$

One way of thinking of this statement is that $n$ samples from the proposal distribution $h(x)$ is worth $n\big/\big(1 + \text{Var}_h(w(X))\big)$ samples drawn from the target distribution $g(x)$.

The nice thing with this rule of thumb is that it doesn't depend on the function being integrated.

Thus the effective sample size is a useful measure of the efficiency of the method when different functions $f$ are investigated with a single sample.

One consequence of this is that we want to keep the coefficient of variation of the weights well behaved.

One way of doing this is to have the proposal distribution $h(x)$ be heavier tailed than the target distribution $g(x)$. This will help minimize

$$E_h\left[\left(\frac{g(X)}{h(X)}\right)^2\right] = \int\left(\frac{g(x)}{h(x)}\right)^2 h(x)\, dx$$

$$= E_g\left[\frac{g(X)}{h(X)}\right]$$

which will keep $\mathrm{Var}_h\left(g(X)/h(X)\right)$ small.

For example, use a $t$ distribution with moderate degrees of freedom instead of a normal.

Marginalization in importance sampling

Let $g(X_1, X_2)$ and $h(X_1, X_2)$ be two probability densities where the support of $g$ is a subset of the support of $h$ (e.g. $g(X_1, X_2) > 0$ implies $h(X_1, X_2) > 0$).

Then

$$\mathrm{Var}_h\left(\frac{g(X_1, X_2)}{h(X_1, X_2)}\right) \geq \mathrm{Var}_h\left(\frac{g(X_1)}{h(X_1)}\right)$$

where $g(X_1)$ and $h(X_1)$ are the respective marginal distributions.

$$E_h\left[\frac{g(X_1, X_2)}{h(X_1, X_2)}\Big| X_1\right]$$

$$= \int \frac{g(X_1, X_2)}{h(X_1, X_2)} h(X_2 | X_1) dX_2$$

$$= \int \frac{g(X_1, X_2)}{h(X_1) h(X_2 | X_1)} h(X_2 | X_1) dX_2$$

$$= \frac{g(X_1)}{h(X_1)} \int g(X_2 | X_1) dX_2 = \frac{g(X_1)}{h(X_1)}$$

Thus

$$\text{Var}_h\left(\frac{g(X_1,X_2)}{h(X_1,X_2)}\right) \geq \text{Var}_h\left(E\left[\frac{g(X_1,X_2)}{h(X_1,X_2)}\Big| X_1\right]\right)$$

$$= \text{Var}_h\left(\frac{g(X_1)}{h(X_1)}\right)$$

In addition

$$\text{Var}_h\left(\frac{g(X_1,X_2)}{h(X_1,X_2)}\right) - \text{Var}_h\left(\frac{g(X_1)}{h(X_1)}\right)$$

$$= E\left[\text{Var}_h\left(\frac{g(X_1,X_2)}{h(X_1,X_2)}\Big| X_1\right)\right]$$

The implication of this result is that where possible, minimize the number of variables you sample as this will increase $\text{ESS}(n)$.

However the computational burden must also be considered.

For example, if

$$\mathrm{Var}_h \left( \frac{g(X_1, X_2)}{h(X_1, X_2)} \right) = 2\,\mathrm{Var}_h \left( \frac{g(X_1)}{h(X_1)} \right)$$

but the computational time involved in sampling $h(X_1)$ is 4 times the time involved sampling $h(X_1, X_2)$, sampling over the second space is to be preferred.
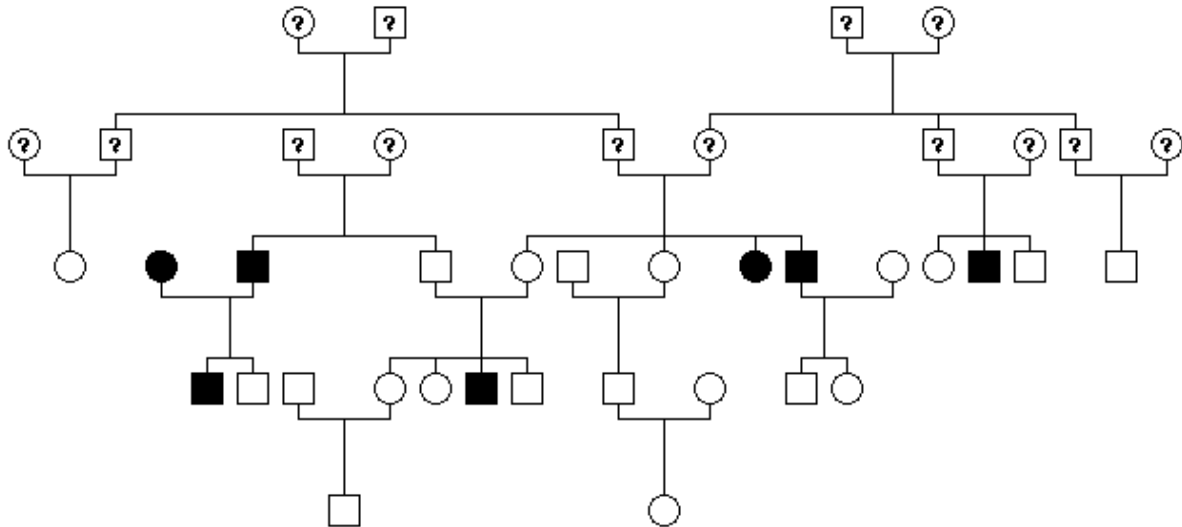
Similarly, Rao-Blackwellization can be with importance sampling, giving an estimate of the form

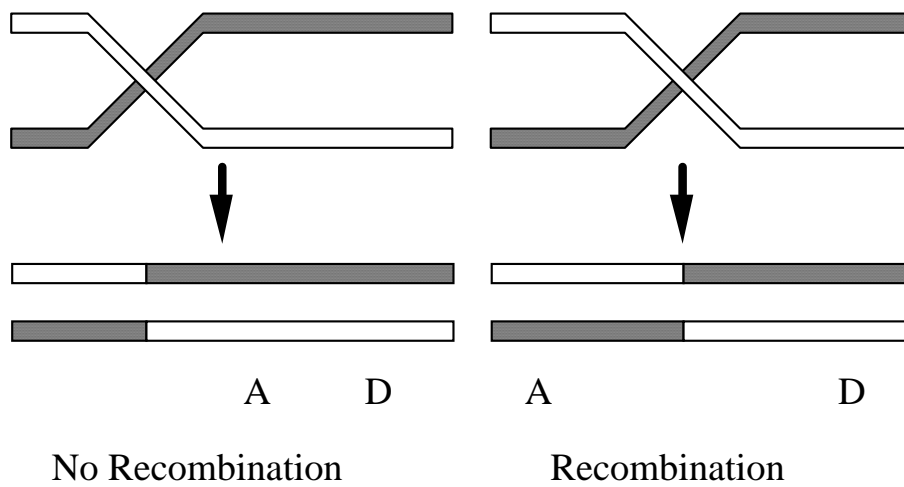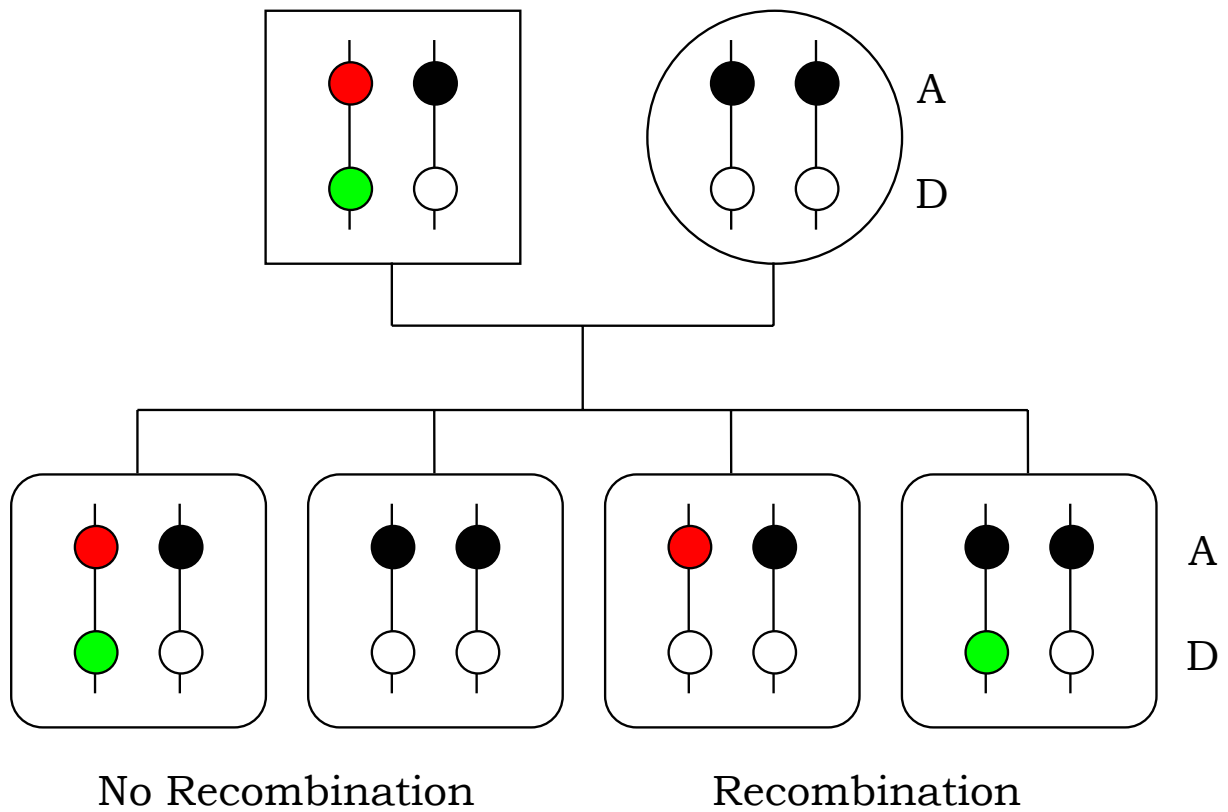$$\breve{\mu} = \frac{1}{W} \sum_{i=1}^{n} w_i E_g \left[ f(X_1, X_2) \big| x_{1,i} \right]$$

While this is a consistent estimate, it may not have a smaller variance than the non-Rao-Blackwellized estimate. However in many situations it should.

The following is an example where it did help.

Estimating recombination fractions with pedigree data



- 41 members

- $n = 27$ (nonfounders), $f = 14$ (founders)

- 8 markers from chromosome 19

- #alleles ranges from 6 to 8

- 14 members in top 2 generation have no marker data

- Want to use pedigree to estimate the distances between the 8 markers

No Recombination          Recombination



No Recombination          Recombination

Can use the recombination fraction (P[recombination]) as a measure of distance.

Data Structure

- $\mathbf{x}_l$ = data on locus $l$

  o $\mathbf{x}_M = (\mathbf{x}_1, ..., \mathbf{x}_m)$ $m$ markers

  o $\mathbf{x}_D$ = disease / trait data

  o $\mathbf{x} = (\mathbf{x}_1, ..., \mathbf{x}_m, \mathbf{x}_D)$

- $\mathbf{y}_l$ = haplotype for locus $l$

  o $\mathbf{y} = (\mathbf{y}_1, ..., \mathbf{y}_m, \mathbf{y}_D)$

  o Allele information with parental source

- $\mathbf{z}_l$ = inheritance vector for locus $l$

  o $\mathbf{z} = (\mathbf{z}_1, ..., \mathbf{z}_m, \mathbf{z}_D)$

  o Indicators of whether allele inherited came from the grandmother or grandfather

Assume that $\mathbf{z}_l$ is a vector of length $2n$. Then an estimate of $\theta_j$, the recombination fraction between markers $j$ and $j + 1$ is

$$\theta_j = \frac{1}{2n} \sum_{i=1}^{2n} I\left(z_{j,i} \neq z_{j+1,i}\right)$$

assuming that $\mathbf{z}$ is known.

Unfortunately, for most data sets, $\mathbf{z}$ can't be determined with certainty. However it is possible to do simulation and estimate the recombination fractions using Monte Carlo EM (MCEM)

As part of this procedure is it necessary to calculate $E\left[\mathbf{z}_l \middle| \mathbf{x}, \boldsymbol{\theta}\right]$ via Monte Carlo. There are two approaches to doing this.

1)  Sample $\mathbf{y}^{(i)}$ from $p_{\boldsymbol{\theta}}\left(\mathbf{y} \middle| \mathbf{x}\right)$ by importance sampling

   Sample $\mathbf{z}^{(i)}$ from $p_{\boldsymbol{\theta}}\left(\mathbf{z} \middle| \mathbf{x}, \mathbf{y}^{(i)}\right)$

   Estimate $E\left[\mathbf{z}_l \middle| \mathbf{x}, \boldsymbol{\theta}\right]$ by

$$\frac{1}{W}\sum_{i=1}^{m} w_i \mathbf{z}_l^{(i)}$$

2)  Sample $\mathbf{y}^{(i)}$ from $p_{\boldsymbol{\theta}}\left(\mathbf{y} \middle| \mathbf{x}\right)$ by importance sampling

   Estimate $E\left[\mathbf{z}_l \middle| \mathbf{x}, \boldsymbol{\theta}\right]$ by

$$\frac{1}{W}\sum_{i=1}^{m} w_i E\left[\mathbf{z}_l \middle| \mathbf{y}^{(i)}, \boldsymbol{\theta}\right]$$

Both approaches estimate exactly the same quantity, but the 2nd approach is much more efficient.

First the second approach has a smaller variance.

In addition, calculating the expected value is actually faster than sampling the **z**'s. (usually doesn't work out this way).