

Integration

We often need to compute:

① $E_f g(x)$ — expectations / moments

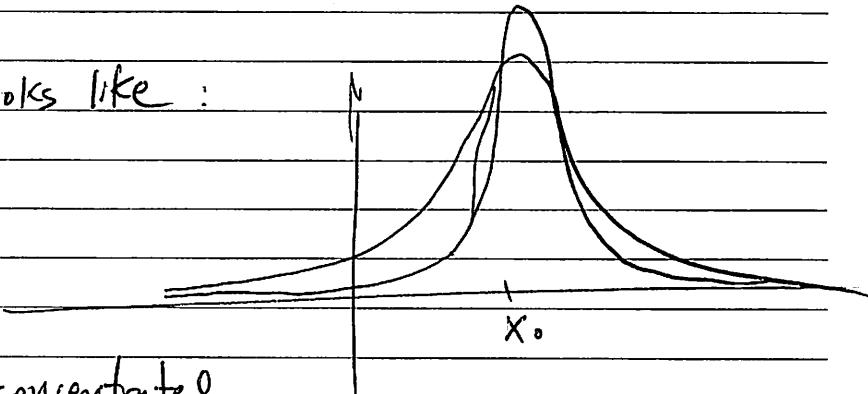
② $\int p(x|\theta) dx$ — normalizing constants for densities.

\Rightarrow In general this can be very hard, but there are many special cases.

Suppose we want to calculate:

$$\int g(x) dx \quad \text{and } g \in L^2$$

Suppose g looks like:



So g is highly concentrated about x_0 .

Then we could say (?)

$$\int g(x) dx \doteq g(x_0) \varepsilon \quad ?$$

$$\begin{aligned} x^{1/2} \\ \int g(x) = \\ x^a \end{aligned}$$

Suppose we want

$$\int_a^b g(x) dx = \int_a^b e^{h(x)} dx \quad [h(x) = \log g(x)]$$

Suppose
h achieves
its max at x_0

$$= \int_a^b \exp\left(h(x_0) + h'(x_0)(x-x_0) + \frac{h''(x_0)(x-x_0)^2}{2}\right) dx$$

$$= \int_a^b \exp\left(h(x_0) + \frac{h''(x_0)(x-x_0)^2}{2}\right) dx$$

$$= \exp(h(x_0)) \int_a^b \exp\left(\frac{h''(x_0)}{2}(x-x_0)^2\right) dx$$

$$= \cancel{\exp(h(x_0))} \exp\left(-\frac{1}{2} \frac{(x-x_0)^2}{h''(x_0)}\right)$$

$$= \exp(h(x_0)) \int_a^b \exp\left(-\frac{1}{2} \frac{(x-x_0)^2}{h''(x_0)}\right) dx$$

$N(x_0, -h''(x_0)^{-1})$ density

$$= \underbrace{\exp(h(x_0))}_{= g(x_0)} \left[\frac{\sqrt{2\pi}}{-h''(x_0)^{1/2}} \right] \left[\Phi(b|x_0, -h''(x_0)^{-1}) - \Phi(a|x_0, -h''(x_0)^{-1}) \right]$$

$$= 1 \text{ if } b = \infty \\ a = -\infty$$

$$= \cancel{g(x_0)}$$

$$h'' \leftarrow \left(\frac{g'}{g}\right)' \\ \frac{gg'' - (g')^2}{g^2}$$

$$g \\ \frac{g''}{g} - \frac{g'^2}{g^2}$$

$$\mathbb{E}(\theta) = \cancel{\int_{-\infty}^{\infty} \theta f(y|\theta) \pi(\theta)}$$

$$\mathbb{E}\theta = \int_{-\infty}^{\infty} \theta p(\theta|y) d\theta$$

$$= \cancel{\int_{-\infty}^{\infty} \theta f(y|\theta) \pi(\theta) d\theta}$$

$$\int \cancel{\theta f(y|\theta) \pi(\theta)} d\theta$$

$$= \int_{-\infty}^{\infty} \theta \exp(\log f(y|\theta) \pi(\theta))$$

$$\underbrace{\int \exp(\log f(y|\theta) \pi(\theta))}_{h(\theta|y)} d\theta$$

Need $\hat{\theta}$

$$= \int_{-\infty}^{\hat{\theta}} \theta \exp(h(\hat{\theta}|y) + h''(\hat{\theta}|y)(\theta - \hat{\theta})^2/2) d\theta$$

$$\int \exp(h(\hat{\theta}|y) + h''(\hat{\theta}|y)(\theta - \hat{\theta})^2/2) d\theta$$

$$= \int_{-\infty}^{\hat{\theta}} \theta \exp(h''(\hat{\theta}|y)(\theta - \hat{\theta})^2) d\theta$$

$$\int \exp\left(\frac{h''(\hat{\theta})}{2}(\theta - \hat{\theta})^2\right) d\theta$$

$$= \int_{-\infty}^{\hat{\theta}} \theta \frac{\sqrt{2\pi}}{|h''(\hat{\theta})|^{1/2}} \mathcal{C}(\theta|\hat{\theta}, -h''(\hat{\theta})^{-1}) d\theta$$

$$\int \frac{\sqrt{2\pi}}{|h''(\hat{\theta})|^{1/2}} \mathcal{C}(\theta|\hat{\theta}, -h''(\hat{\theta})^{-1}) d\theta$$

$$= \hat{\theta}$$

Laplace approx
posterior mean
is the
posterior mode.

See Lange Chap 4.6 for Laplace Approx.

Monte Carlo

Suppose we want to compute $\mathbb{E}_f[h(x)]$ for $h: \mathbb{R}^K \rightarrow \mathbb{R}$

$$\mathbb{E}_f[h(x)] = \int h(x) f(x) dx$$

If we can simulate $x_1, \dots, x_n \stackrel{iid}{\sim} f$, then by the LLN,

$$\frac{1}{n} \sum_{i=1}^n h(x_i) \xrightarrow{\text{sum}} \mathbb{E}_f[h(x)] \quad \text{integral}$$

Furthermore, $\text{Var}[\frac{1}{n} \sum_{i=1}^n h(x_i)] = \frac{1}{n^2} \sum \text{Var}(h(x_i))$

~~$$\text{Var}[\frac{1}{n} \sum_{i=1}^n h(x_i)] = \sqrt{\frac{1}{n^2} \sum \text{Var}(h(x_i))}$$~~

Notice that the variance doesn't depend on dimension of x . Thus, $\text{Var} \sim \frac{1}{n}$. This is very important, both good and bad.

Applications

① Monte Carlo / Simulation studies. We have a method that estimates a parameter $h(x) = \hat{\beta}$ and we want to explore performance of $\hat{\beta}$.

② $\mathbb{E}_f[h(x)]$ might be a posterior mean (f is post. dist.)

③ Might want $\int h(x) dx = \int \frac{h(x)}{f(x)} f(x) dx$.

$$\approx \frac{1}{n} \sum \frac{h(x_i)}{f(x_i)}, x_1, \dots, x_n \stackrel{iid}{\sim} f$$

We need to be able to simulate numbers, i.e. from some dist. f

Random Number Generation

Most popular (and simplest) are linear congruential generators.

Let X_0 be some starting value, called the "seed".
Then generate a sequence (for $n=0, 1, 2, \dots$)

$$X_{n+1} = (aX_n + b) \bmod m$$

a = multiplier

b = increment

m = modulus

For uniform RN, just let $U_{n+1} = \frac{X_{n+1}}{m}$

Ideally, X_1, X_2, \dots will hit every number from 0 to $m-1$ before repeating

of steps until repeat is "period".
A "maximal period generator" has period m.

Setting a, b, and m is a very tricky business.

For example, this is bad

$$X_{n+1} = (2X_n + 0) \bmod 2^{32}$$

One good set is

$$a = 106, b = 1283, m = 6075$$

There are quadratic and cubic PRNGs, but they require more work and not much better, i.e.

$$X_{n+1} = (a X_n^2 + b X_n + c) \bmod m$$

LCG are not useful for things like cryptography.
Any polynomial generator can be broken.

RNG are critical for stream ciphers which

~~breakable~~
Use linear feedback shift registers (LFSR)
"break + better"

All PRNGs produce deterministic sequences that "look" random

One can check "randomness" of sequence with a test for uniformity

- Kolmogorov-Smirnov test
- Chi-square test

Marsaglia has "die hard" tests

Also, many PRNGs generate good 1-D sequences but do not look random in > 1 dimension

It may be stuck on a hyperplane.

Uniforms are good but we need values from some density f . First, assume you can simulate $\text{Unif}(0, 1)$

④ Integral transform / rescaling

exponential(1)

$$\cancel{f(x) = e^{-x}}$$

$$f(x) = e^{-x}$$

$$F(x) = \cancel{\frac{1}{\lambda}e^{-\lambda x}} = e^{-x}$$

$$F(x) = 1 - e^{-x}$$

$$F(y) = \cancel{1 - e^{-y}}$$

$$\log(Y) = \cancel{-\lambda^{-1} \ln U}$$

$$F^{-1}(u) = \cancel{-\lambda^{-1} \ln u}$$

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

$$① U \sim \text{Unif}(0, 1)$$

$$② X \sim F^{-1}(U), \text{ where } f(x) = \int_0^x f(t) dt \text{ cdf.}$$

⑤ Transformation

$$① Z_1 = \sqrt{-2 \log U_1} \cos(2\pi U_2) \quad \begin{matrix} \text{i.i.d} \\ \sim N(0, 1) \end{matrix}$$

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

$$② \text{Beta}(\alpha, \beta) \quad \frac{X}{X+Y} \sim \text{Beta}(\alpha, \beta)$$

$$\text{where } X \sim \text{gamma}(\alpha, 1)$$

$$Y \sim \text{gamma}(\beta, 1)$$

$$③ M \sim \text{gamma}(\alpha, \beta)$$

$$Y | M \sim \text{Poisson}(M)$$

$$Y \sim \text{negative binomial}$$

Most standard dist. are implemented in R.

$$p(y) = \underbrace{\sum_{m=0}^{\infty} p(y|m) p(m)}_{\text{Poisson gamma}}$$

Multivariate Normal

Want $X \sim N(\mu, \Sigma)$

Let $\Sigma = L^T L$ (cholesky decomp)

① Simulate $Z \sim N(0, I)$

② $X = \mu + WZ \sim N(\mu, \Sigma)$

Sieve sampling

Rejection sampling

A way of generating samples/obs. from f by thinning out obs/samples from a candidate density g . ("Random thinning")

Suppose f is our target density and we can evaluate it.

let g be our candidate density that we can simulate from.

let X_f be the support of f and X_g be support of g and assume $X_f \subset X_g$

Assume $C = \sup_{x \in X_f} \frac{f(x)}{g(x)} < \infty$

and that we can calculate C

Rejection Sampling

① Simulate $u \sim \text{uniform}(0,1)$

② Simulate $X \sim g(x)$

③ If $u \leq \frac{f(x)}{cg(x)}$, "accept" X

otherwise, ~~go back to 1~~ and X

go back to ①.

let $X_1, X_2, \dots \sim \underset{\text{iid}}{g}$

$X_1 \ X_2 \ X_3 \ X_4 \ X_5 \ X_6 \ X_7 \ \dots$

0	0	1	0	1	0	0	...
reject	reject	accept	accept				
geometric(p)							

iid coin flips

sample size of 2.

of flips until acceptance is geometric w/
success probability $1/c$.

$$\text{Pf: } \text{P}(X \text{ accepted}) = \text{P}\left(u \leq \frac{f(x)}{cg(x)}\right)$$

$$= \int \text{P}\left(u \leq \frac{f(x)}{cg(x)} \mid X=x\right) g(x) dx$$

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

$$= 1/c$$

The dist. of accepted values is f .

$$P(X \leq t | X \text{ accepted})$$

$$= \frac{P(X \leq t, X \text{ accepted})}{P(X \text{ accepted})}$$

$$= E_g [1\{X \leq t\} 1\{X \text{ accepted}\}]$$

$$= c E_g [E_u [1\{X \leq t\} 1\{u \leq \frac{f(x)}{cg(x)}\} | X]]$$

$$= c E_g \left(1\{X \leq t\} \frac{f(x)}{c g(x)} \right)$$

$$= c \int 1\{X \leq t\} \frac{f(x)}{c g(x)} g(x) dx$$

$$= \int_{-\infty}^t f(x) dx = F(t)$$

Notes:

- ① Only need to know f or g up to a constant of proportionality
- ② Any number $c' > c$ will work, but will be less efficient.
- ③ Operations can (and should) be performed on log scale

- ② $\text{Gaussian } u \sim \text{unif}(0,1), X \sim g$
- ③ Accept X if $u \leq f(x)/g(x)$
- ④ Update $C^* = \max_C \left\{ C, \frac{f(x)}{g(x)} \right\}$
- ⑤ Go to step 2, set $C = C^*$

$$X_{n+1} = k_n - \alpha h(k_n)$$

$$h(x_n), h(x_{n+1})$$

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

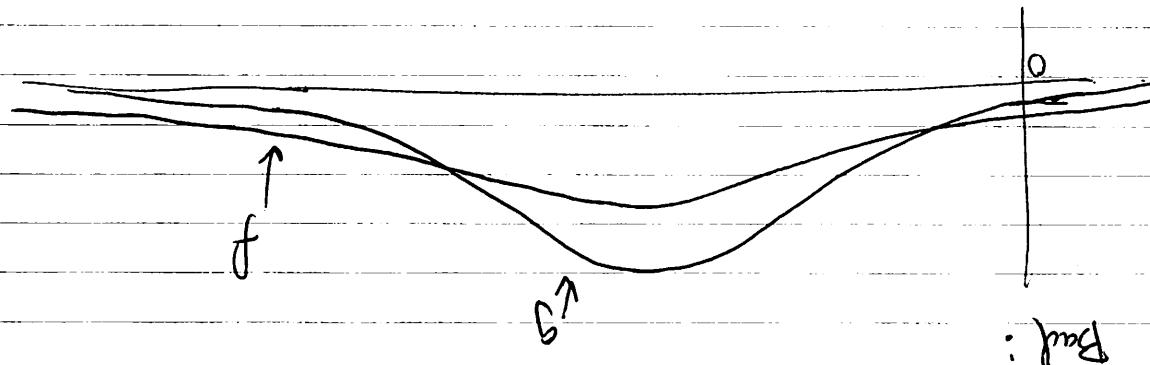
① Guess $C > 1$

Empirical step rejection sample (affo, 02)

Can we estimate C ? Yes!

Why if we count calls $C = \sup_{x \in \mathbb{R}} f(x)/g(x) \downarrow \infty$

As $x \rightarrow \infty, g(x) \uparrow 0$ faster than $f(x) \uparrow 0$.



Paul:

⑤ Whether $C = \infty$ depends (usually) on the shape of the curve, which must be heavier than the tails.

④ The higher the dimension of f, g the more difficult rejection sampling will be.

ESUP

Extra assumption is needed for ~~ESUP~~.

$$C = \frac{f(x_c)}{g(x_c)} \text{ for some } x_c \in X_f$$

& (sup is achievable)

Thus is satisfied if g has heavier tails than f .

~~Caffo proved that~~
~~If f is discrete~~

let X_1, X_2, X_3, \dots iid g

$X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5 \quad X_6 \quad \dots$

True C

Y_i	0	0	1	0	1	0
	accept			accept		

 \hat{C}

\hat{Y}_i	1	0	1	0	0	0
	accept		accept	accept	reject	accept
	error		error			

Caffo showed that

① ~~P~~ $P(Y_i \neq \hat{Y}_i \text{ i.o.}) = 0$ if f is discrete

~~Def:~~ By assumption 3, $\exists x_c \in X_f$ s.t. $C = \frac{f(x_c)}{g(x_c)}$.

Let $\gamma = \min_i \{X_i = x_c\}$ where $X_i \sim g$.

then $\gamma \sim \text{geometric}(g(x_c))$. Once $\hat{C} = C$, algorithms are the same.

$$P(Y_i \neq \hat{Y}_i) \leq P(\gamma \geq i) = (1 - g(x_c))^{i-1}$$

$$\Rightarrow \sum_{i=1}^{\infty} P(Y_i \neq \hat{Y}_i) < \infty$$

Coupling
lemma \Rightarrow

② For continuous f , it's trickier. In general,

$$P(Y_i \neq \tilde{Y}_i) = O(i^{-1})$$

But if $\log(f/g)$ is smooth (twice differentiable) and unimodal (w/mode at x_c), then

$$P(Y_i \neq \tilde{Y}_i) = O(i^{-2})$$

Suggests a "burn-in", ~~or~~ discard first several vectors.

More on burn-in with MCMC

Often we want to simulate from

$$p(\theta|y) \propto \underbrace{f(y|\theta)}_{\text{likelihood}} \underbrace{\pi(\theta)}_{\text{prior}}$$

\Rightarrow use π as the candidate dist.

$$C = \sup_{\theta} \frac{p(\theta|y)}{\pi(\theta)} = \sup_{\theta} f(y|\theta) = \underset{\substack{\text{not made} \\ \text{MLE}}}{f(y|\hat{\theta})}$$

Rejection sampling is

① Simulate $U \sim \text{uniform}(0,1)$

② Simulate $\theta \sim \pi$

③ Accept θ if $U \leq \frac{p(\theta|y)}{C\pi(\theta)} = \frac{f(y|\theta)}{f(y|\hat{\theta})}$

Ex: $Y_1, \dots, Y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

Prior $\pi(\mu) = N(0, B) \xrightarrow{\text{BIG}}$

$$\pi(\sigma^2) = \text{IG}(\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} z^{-\alpha-1} e^{-\beta/z} = f(z)$$

$$\text{MLE: } \hat{\mu} = \bar{Y}, \quad \hat{\sigma}^2 = S^2 = \frac{1}{n} \sum (Y_i - \bar{Y})^2$$

$$\begin{aligned} C = f(y | \hat{\mu}, \hat{\sigma}^2) &= \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{1}{2\sigma^2} \sum (Y_i - \hat{\mu})^2} \\ &= \left(\frac{1}{(2\pi)^{\frac{n}{2}}} \left(\frac{\sum (Y_i - \bar{Y})^2}{n} \right)^{-n/2} e^{-\frac{n}{2}} \right) \end{aligned}$$

① Simulate $U \sim \text{Uniform}(0, 1)$

② Simulate $\mu \sim N(0, B), \frac{1}{\sigma^2} \sim \text{gamma}(\alpha, \beta)$

③ Accept (μ, σ^2) if

$$U \leq \frac{f(y | \mu, \sigma^2)}{f(y | \hat{\mu}, \hat{\sigma}^2)} = \left(\frac{\sum (Y_i - \bar{Y})^2}{n \sigma^2} \right)^{n/2} \exp \left(-\frac{1}{2\sigma^2} \sum (Y_i - \mu)^2 + \frac{n}{2} \right)$$

Given a sample from $p(\theta | y)$ we can compute confidence/posterior intervals, mean, variance, etc.

What does ESUP do? If searches for MLE by random search.

It is MCMC,
but not Gibbs.
exactly

$$\mathbb{E}[X] \text{ iff } \mathbb{E}[X|Y]$$

Data Augmentation (Tanner + Wong '87)

We have some data y_1, \dots, y_n , and we propose a model

$$y \sim p(y|\theta)$$

$$\theta \sim \pi(\theta)$$

So we have a likelihood $p(y|\theta)$ and a prior $\pi(\theta)$.

We want the posterior $p(\theta|y)$.

$$p(\theta|y) = p(y|\theta) \pi(\theta)$$

BUT because of some missing data z , it is difficult to evaluate $p(\theta|y)$.

Let (y, z) be the "complete data" and we want to know

$$p(\theta|y) = \int p(\theta|y, z) p(z|y) dz \quad (1)$$

"complete data predictive density
posterior" for z .

↳ easy to calculate

Notice also that we have

$$p(z|y) = \int p(z|y, \theta) p(\theta|y) d\theta \quad (2)$$

Substituting, ~~we get~~ (2) into (1), we get

$$\begin{aligned} p(\theta|y) &= \int p(\theta|y, z) \left[\int p(z|y, \theta') p(\theta'|y) d\theta' \right] dz \\ &= \int [p(\theta|y, z) p(z|y, \theta)] dz p(\theta|y) d\theta' \end{aligned}$$

System of
Integral
equations

DA algorithm

Roughly: At iteration i

① Sample $z_1, \dots, z_m \sim p(z|y)$

② Estimate $p_{i+1}(\theta|y) = \frac{1}{m} \sum_{j=1}^m p(\theta|y, z_j)$

Repeat until $\|p_{i+1}(\theta|y) - p_i(\theta|y)\| < \varepsilon$.
complete data posterior

① Pick some initial $p_0(\theta|y)$. At step i

②a Generate $\theta \sim p_i(\theta|y)$

②b Generate $z \sim p(z|\theta, \theta)$

Repeat ②a - ②b m times to get z_1, \dots, z_m .

③ Let $p_{i+1}(\theta|y) = \frac{1}{m} \sum p(\theta|z_j, z_j)$

mixture of conditonal densities.

Monte Carlo estimate of

$$\int p(\theta|y, z) p(z|y) dz$$

DA not so straightforward when there are more than 2 "missing" components.

Use Gibbs Sampling for more complicated problems

$$\text{Let } K(\theta, \theta') = \int p(\theta|y, z) p(z|y, \theta') d\theta'$$

$$p(\theta|y) = \int \underset{(1)}{K(\theta, \theta')} \underset{(2)}{p(\theta'|y)} d\theta'$$

(fixed point system) If $p(\theta'|y)$ is the
 ~~$\pi(\theta)$ and $\pi(\theta')$~~ the posterior, then it
 will map to itself.

Let T be a functional which maps a function $g(\theta)$ to

$$Tg(\theta) = \int K(\theta, \theta') g(\theta') d\theta'$$

If we take some initial value $\pi_0(\theta|y)$, what happens when we iterate ~~$\pi_{t+1}(\theta|y) = T\pi_t(\theta|y)$~~ , i.e.

~~$\pi_{t+1}(\theta|y) = T\pi_t(\theta|y)$~~

$$= \int K(\theta, \theta') \pi_0(\theta'|y) d\theta'$$

Does the sequence $\{\pi_i(\theta|y)\}$ converge to anything?

~~(1) $\pi_i(\theta|y)$ the true posterior is the~~

(1) $\{\pi_i(\theta|y)\} \rightarrow p(\theta|y)$ monotonically (always getting closer)
 (true posterior)

(2) $p(\theta|y)$ is the unique solution to system of integral equations

(3) $\pi_i(\theta|y) \rightarrow p(\theta|y)$ linearly.

Ex. 3-stage hierarchical model

$$\cancel{y|\alpha} \quad y|\alpha \sim N(\alpha, 1)$$

$$\alpha|\theta \sim N(\theta, 1)$$

$$\theta \sim \pi(\theta) = N(0, 1)$$

$$p(y, \alpha, \theta) = \frac{p(y|\alpha, \theta)}{N(\alpha, 1)} \frac{p(\alpha|\theta)}{N(\theta, 1)} \pi(\theta)$$

~~p($\alpha|y, \theta$)~~

$$p(\alpha|y, \theta) \propto p(y|\alpha, \theta) p(\alpha|\theta)$$

$$= N\left(y + \frac{1}{2}(\theta - y), \frac{1}{2}\right)$$

$$p(\theta|y, \alpha) \propto p(y|\alpha, \theta) p(\alpha|\theta) \pi(\theta)$$

$$N(\quad)$$

① Find some $p_i(\theta|y)$ [Normal?]

②a Sample $\theta \sim p_i(\theta|y)$

②b Sample $\alpha \sim p(\alpha|y, \theta)$

Repeat m times to get $\alpha_1, \dots, \alpha_m$

③ Let $P_{\text{fit}}(\theta|y) = \frac{1}{m} \sum p_i(\theta|y, \alpha)$

~~...~~

EM

Likelihood

Target

Method Handle z
toNeed $(z|y)$

Std. errors

Rate of convergence linear

Monotonicity

DA

Bayesian

Posterior dist.

sample (conditional)

complete E

sample from

No, by default

Yes, have entire posterior
(and prior)

linear

✓

Both methods "impute" missing data via the specification of the complete data model.

Once $p(y, z)$ is specified, everything else is determined.

$$\hat{M}_n = \sum_{i=1}^n \frac{f(x_i)}{g(x_i)} h(x_i) = \frac{\sum w(x_i) h(x_i)}{\sum w(x_i)}$$

With rejection sampling, we can sample from f given a candidate density g .

What if we want to estimate $E_f[h(x)]$ for some $h: \mathbb{R}^K \rightarrow \mathbb{R}$?

Obvious way: Sample $X_1, \dots, X_n \sim f$ and use
 $E_f[h(x)] \approx \frac{1}{n} \sum_{i=1}^n h(X_i) = \hat{M}_n$

$$\sqrt{n}(\hat{M}_n - M_n) \xrightarrow{\text{D}} N(0, \sigma^2)$$

But sampling from f is hard, so we use RS with cand. dens. g .

In order to obtain samples of size n , we need on avg. $C \times n$ samples from g , where

$$C = \sup \frac{f}{g} \quad (\text{assumed } < \infty)$$

We reject $(C-1) \times n$ of the samples from g !
 on avg.

Those samples belong in the domain of f ,
 but they may be over/under-represented

e.g. if g has heavier tails, there will be too many extreme values — RS rejects those.

But maybe we can down/up weight values to get ~~the~~ the right answer.

Note that

$$\mathbb{E}_f[h(x)] = \mathbb{E}_g\left[\frac{f(x)}{g(x)} h(x)\right].$$

If $x_1, \dots, x_n \sim g$ then

$$\mathbb{E}_f[h(x)] \approx \frac{1}{n} \sum_{i=1}^n \frac{f(x_i)}{g(x_i)} h(x_i)$$

$$= \frac{1}{n} \sum w_i h(x_i) = \tilde{M}_h$$



importance weights

\Rightarrow assumes we can compute f and g exactly.

Notice that if $f=g$, our estimator is just

$$\frac{1}{n} \sum_{i=1}^n h(x_i).$$

~~So~~ If the density of X_i w.r.t g is larger than its density w.r.t f , then $h(x_i)$ is down weighted in the sum (vice versa)

N

Comparison

Rejection Sampling : Sample directly from f , take averages

Importance Sampling : Sample from g , reweight by f/g

For computing expectations IS is much more efficient because there is no rejection

~~What if we only know~~

RS: Let $C = \sup f/g$. Then

sample $X_1, \dots, X_n \sim g$ and $U_1, \dots, U_n \sim \text{unif}(0,1)$.

$$\hat{M}_h = \frac{\sum \mathbb{1}\{U_i \leq \frac{f(x_i)}{g(x_i)}\} h(x_i)}{\sum \mathbb{1}\{U_i \leq \frac{f(x_i)}{g(x_i)}\}}$$

IS: Sample $X_1, \dots, X_n \sim g$

$$\tilde{M}_h = \frac{1}{n} \sum \frac{f(x_i)}{g(x_i)} h(x_i)$$

① IS estimator (\tilde{M}_h) is "smoother" than \hat{M}_h and should have lower variance

② RS requires $\sup f/g < \infty$. But IS does not need C.

③ IS computations cannot be done on log scale

What if we only know $f^* \neq c_1 f$ and $g^* \neq c_2 g$. Then use

$$\hat{M}_h^* = \frac{\frac{1}{n} \sum_{i=1}^n \frac{f^*(x_i)}{g^*(x_i)} h(x_i)}{\frac{1}{n} \sum_{i=1}^n \frac{f^*(x_i)}{g^*(x_i)}} \quad (\text{ratio estimator})$$

Since f^* and g^* are in numerator and denominator, constants drop out.

$\hat{M}_h^* \rightarrow M_h$ by Slutsky theorem

Ex. Importance Sampling for Bayesian sensitivity analysis.

We have data y with a likelihood $L(y|\theta)$ and a prior for θ $\pi(\theta|\psi_0)$ where ψ_0 is

a known hyperparameter. The posterior for θ is

$$p(\theta|y, \psi_0) \propto L(y|\theta) \pi(\theta|\psi_0).$$

Suppose we want to obtain much energy obtain

$\theta_1, \dots, \theta_n$, a sample from $p(\theta|y, \psi_0)$, and we can compute posterior mean $E[\theta] = \frac{1}{n} \sum \theta_i$. What

if we want to calculate different values of $E(\theta_i|y, \psi_0)$? We do not need to resample $\theta_1, \dots, \theta_n$, just

reweight them. ~~let's plot~~ Let ψ be

the new hyperparameter, $f = p(\theta|y, \psi)$, $g = p(\theta|y, \psi_0)$.

$$E[\theta|y, \psi] \approx \frac{\sum \theta_i \frac{f(\theta_i)}{g(\theta_i)}}{\sum \frac{f(\theta_i)}{g(\theta_i)}} = \frac{\sum \theta_i \frac{p(\theta_i|y, \psi)}{p(\theta_i|y, \psi_0)}}{\sum \frac{p(\theta_i|y, \psi)}{p(\theta_i|y, \psi_0)}} \quad \begin{matrix} \text{have a sample} \\ \text{from this} \end{matrix}$$

$$= \frac{\sum \theta_i \frac{L(y|\theta) \pi(\theta|\psi)}{L(y|\theta) \pi(\theta|\psi_0)}}{\sum \frac{L(y|\theta) \pi(\theta|\psi)}{L(y|\theta) \pi(\theta|\psi_0)}} = \frac{\sum \theta_i \frac{\pi(\theta|\psi)}{\pi(\theta|\psi_0)}}{\sum \frac{\pi(\theta|\psi)}{\pi(\theta|\psi_0)}}$$

$$= \frac{\sum \theta_i \frac{\pi(\theta|\psi)}{\pi(\theta|\psi_0)}}{\sum \frac{\pi(\theta|\psi)}{\pi(\theta|\psi_0)}} = \frac{\sum \frac{\pi(\theta|\psi)}{\pi(\theta|\psi_0)}}{\sum \frac{\pi(\theta|\psi)}{\pi(\theta|\psi_0)}}$$

$$\frac{E[\theta|\psi_0]}{p(\theta|y, \psi_0)} = \int \theta p(\theta|y, \psi_0) d\theta$$

Ex. Calculating Marginal likelihoods.

Suppose we have $f(y|u)$, the dist. of y given some random effect u and $h(u|\theta)$ the dist. of random effects for parameter θ .

If

$$y_{ij} \sim N(\mu + u_i, \sigma^2)$$

$$u_i \sim N(0, \theta)$$

We want to maximize

$$L(\theta) = \int f(y|u) h(u|\theta) du = E_h [f(y|u)]$$

↳ Integrate out random effects.

Suppose we simulate u_1, u_2, \dots, u_n from a candidate dist $\underbrace{h(u|\theta_0)}_{\text{wrong dist}}$. Then

$$\hat{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \frac{h(u_i|\theta)}{h(u_i|\theta_0)} f(y|u_i)$$

Could use a more general candidate dist $g(u|\theta_0)$.

$$\hat{L}(\theta) = \frac{1}{n} \sum \frac{h(u_i|\theta)}{g(u_i|\theta_0)} f(y|u_i)$$

$$g(u|\theta) = f(y|u) h(u|\theta)$$

Marginal likelihood

Ideal carbamate is proportional to
 $f(y|u) h(u|\theta)$ which is
 $p(u|y, \theta)$. Since we want to maximize \hat{L} ,
the best carbamate is $p(u|y, \hat{\theta})$ where
 $\hat{\theta}$ maximizes \hat{L} . But that solves our
problem. So try (Geyer 1990)

- ① Sample $\{u_i\}$ from $p(u|y, \theta_0)$
- ② Max \hat{L} to get θ_1
- ③ Set $\theta_0 = \theta_1$. Go to ①

MCEM

MCEM ("Ascent-based") (Caffo et al 2008)

Given γ, z , the E-step computes

$$Q_t = \mathbb{E}_h \left[\log g(\gamma, z | \theta) \mid \gamma, \theta_0 \right]$$

$h(z | \gamma, \theta_0)$ → simulate from vsj rejection sampling

$$\tilde{Q}_t \approx \frac{1}{M_t} \sum_{i=1}^{M_t} \log g(\gamma, z_i | \theta) \xrightarrow{\text{LN}} Q_t$$

as $M_t \rightarrow \infty$

How do we know that θ^* which max \tilde{Q} has the ascent property, $Q(\theta^* | \theta_0^*) \geq Q(\theta_0^* | \theta_0^*)$?

$$\Delta Q = Q(\theta^* | \theta_0^*) - Q(\theta_0^* | \theta_0^*)$$

$$\approx \Delta \tilde{Q} = \tilde{Q}(\theta^* | \theta_0^*) - \tilde{Q}(\theta^* | \theta_0^*)$$

$$\sqrt{M_t} (\Delta \tilde{Q} - \Delta Q) \sim N(0, \sigma^2) \quad [\text{Caffo shows}]$$

$\hat{\sigma}^2$ depends on IS or RS.

Update the sample size

$$\Delta \tilde{Q}(\hat{\theta}_{t+1}^{(test)}, \theta_t^*) \sim N\left(\Delta Q(\theta_{t+1}^*/\theta_t^*) \frac{\hat{\tau}}{M_{t+1}}\right)$$

$$M_{t+1} = \frac{\hat{\tau}^2 (z_\alpha + z_\beta)^2}{\Delta \tilde{Q}(\theta_t^*/\theta_{t-1}^*)^2}$$

For IS we do not need $\sup f/g < \infty$ but what is required?

Recall Cramer's Theorem (Delta Method)

Let y_1, \dots, y_n be s.t. $\sqrt{n}(\bar{y} - \mu) \rightarrow N(0, \Sigma)$

where $\mu = E[\bar{y}]$. Let $\psi: \mathbb{R}^K \rightarrow \mathbb{R}$ be a differentiable map, then

$$\sqrt{n}(\psi(\bar{y}) - \psi(\mu)) \rightarrow N(0, \psi'(\mu)^T \Sigma \psi'(\mu))$$

Importance Sampling

Let $y_i = \begin{pmatrix} h(x_i) w_i \\ w_i \end{pmatrix}$ where $x_i \sim g$
 $w_i = \frac{f(x_i)}{g(x_i)}$

$$\text{So } E_g y_i = \begin{pmatrix} E h(x_i) \\ 1 \end{pmatrix} = \begin{pmatrix} \mu_h \\ 1 \end{pmatrix}. \quad \text{---}$$

Note that we can estimate $\text{Var}(y_i)$ consistently with

$$\hat{\Sigma} = \begin{bmatrix} \text{Sample Var}\{h(x_i)w_i\} & \text{Sample Cov}\{h(x_i)w_i, w_i\} \\ \text{Sample Cov}\{h(x_i)w_i, w_i\} & \text{Sample Var}\{w_i\} \end{bmatrix}$$

Let $V(a, b) = \frac{a}{b}$. Then $\hat{\mu}_h = V(\sum h(x_i)w_i, \sum w_i)$.

$$\text{Also } V'(a, b) = \left(\frac{1}{b}, -\frac{a}{b^2} \right)$$

Variance estimates for IS

Cramers' Theorem: Suppose y_1, y_n are

$$\text{st } \sqrt{n}(\bar{y} - \mu) \rightarrow N(0, \Sigma)$$

where $\mu = \mathbb{E} \bar{Y}$. Let $g: \mathbb{R}^K \rightarrow \mathbb{R}$ be a differentiable map. Then

$$\sqrt{n}(g(\bar{y}) - g(\mu)) \rightarrow N(0, g'(\mu)^T \Sigma g(\mu))$$

Let f is target density

g is proposal density

~~$\alpha \neq \frac{f}{g}$~~

$$\mu_h = \mathbb{E}_f h(x)$$

Let $y_i = \begin{pmatrix} h(x_i) w_i \\ w_i \end{pmatrix}$ where $x_i \sim g$
 and $w_i = \frac{f(x_i)}{g(x_i)}$

~~$\mathbb{E}_g y_i$~~

$$\mathbb{E}_g y_i = \begin{pmatrix} \mathbb{E}_g [h(x_i) \frac{f(x_i)}{g(x_i)}] \\ \mathbb{E}_g [\frac{f(x_i)}{g(x_i)}] \end{pmatrix} = \begin{pmatrix} \mathbb{E}_f h(x) \\ 1 \end{pmatrix}$$

$$\text{Let } m\begin{pmatrix} a \\ b \end{pmatrix} = \frac{a}{b} \quad m'\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \frac{1}{b} \\ -\frac{a}{b^2} \end{pmatrix}$$

$$m(\bar{Y}) = m\left(\frac{\frac{1}{n}\sum h(x_i)w_i}{\frac{1}{n}\sum w_i}\right) = \frac{\frac{1}{n}\sum h(x_i)w_i}{\frac{1}{n}\sum w_i}$$

$$\text{Note } E_g Y_i = \begin{pmatrix} E_g\left[h(x_i) \frac{P(x_i)}{g(x_i)}\right] \\ E_g\left[\frac{P(x_i)}{g(x_i)}\right] \end{pmatrix} = \begin{pmatrix} E_f[h(\cancel{x})] \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} \mu_h \\ 1 \end{pmatrix}$$

$$E\bar{Y} = \begin{pmatrix} \mu_h \\ 1 \end{pmatrix}$$

$$m(E\bar{Y}) = \mu_h$$

We can estimate $\text{Var}(Y_i)$ consistently w/M

$$\hat{\Sigma} = \begin{pmatrix} \hat{\text{Var}}\{h(x_i)w_i\} & \hat{\text{Cov}}\{h(x_i)w_i, w_i\} \\ \hat{\text{Cov}}\{h(x_i)w_i, w_i\} & \hat{\text{Var}}\{\cancel{h(x_i)}w_i\} \end{pmatrix}$$

Cramer's Thm says:

$$\sqrt{n} \left(\frac{m(\bar{Y}) - \mu_h}{m'(\bar{Y})^\top \hat{\Sigma} m'(\bar{Y})} \right) \rightarrow N(0, 1)$$

Cramer's Theorem says:

$$\sqrt{n} \left(\frac{V(\bar{Y}) - E[M_h]}{V'(\bar{Y})^T \hat{\Sigma} V'(\bar{Y})} \right) \rightarrow N(0, 1)$$

And

$$V'(\bar{Y})^T \hat{\Sigma} V'(\bar{Y}) =$$

$$n \left(\frac{\sum h(x_i) w_i}{\sum w_i} \right)^2 \left(\frac{\sum h(x_i)^2 w_i^2}{(\sum h(x_i) w_i)^2} - 2 \frac{\sum h(x_i) w_i^2}{(\sum h(x_i) w_i)(\sum w_i)} + \frac{\sum w_i^2}{(\sum w_i)^2} \right)$$

\Rightarrow We need

$$E_g[h(x)^2 w^2] = E\left[\left(h(x) \frac{f(x)}{g(x)}\right)^2\right] < \infty$$

$$E[w^2] = E\left(\frac{f(x)}{g(x)}\right)^2 < \infty$$

$$E[h(x) w^2] = E\left[h(x) \left(\frac{f(x)}{g(x)}\right)^2\right] < \infty$$

All true if $\frac{f(x)}{g(x)}$ is bounded, which is required for RS.

$$\text{Close to } \hat{\theta} \sum_{i=0}^{\infty} \alpha^i = \frac{1}{1-\alpha} \quad Q(\theta | \theta_t) = 0$$

$$\hat{\theta}_{t+1} = g(\theta_t) = \arg \max_g Q(\theta | \theta_t)$$

$$g(\theta_t) \doteq g(\theta_{t-1}) + J(\theta_{t-1})(\theta_t - \theta_{t-1})$$

$$\theta_{t+1} \doteq \theta_t + J(\theta_{t-1})(\theta_t - \theta_{t-1})$$

$$d_{t+1} \doteq J(\theta_{t-1}) d_t$$

$$d_{t+2} \doteq J(\theta_t) d_{t+1} = \underbrace{J(\theta_t) J(\theta_{t-1})}_{J(\hat{\theta})} d_t$$

$$d_{t+j+1} \doteq J_t(\hat{\theta}) d_{j+1}$$

$$\hat{\theta} = \theta_j + \sum_{l=1}^{\infty} d_{l+j} = \theta_j + (\theta_{j+1} - \theta_j) + (\theta_{j+2} - \theta_{j+1}) + \dots$$

$$= \theta_j + \left(\sum_{l=0}^{\infty} J_l(\hat{\theta}) \right) d_{j+1} \quad \xrightarrow{\text{Eigenvalues of } J(\hat{\theta}) \text{ are all } < 1}$$

$$\theta_{j+1}^* = \theta_j + (1 - J(\hat{\theta}))^{-1} d_{j+1} \quad (\text{Aitken acceleration})$$

$$(I - J(\hat{\theta}))^{-1} = I_x I_y^{-1}$$

↳ $-Q''(\hat{\theta}/\hat{\theta})$

~~Part C~~

$$I_y = I_x - \left(E \left[s(x/\hat{\theta})^+ s(x/\hat{\theta})' \right] - s(y/\hat{\theta}) s(y/\hat{\theta})' \right)$$

Non-monotone, numerical instability, diff. f?

SQUAREM

Super linear
Convergence?

① $\theta_1 = g(\theta_0)$

② $\theta_2 = g(\theta_1)$

③ $r = \theta_1 - \theta_0$

④ $v = \theta_2 - \theta_1 - r$

⑤ Compute $\alpha = -\frac{\|r\|}{\|v\|}$

⑥ $\theta' = \theta_0 - 2\alpha r + \alpha^2 v$

⑦ $\theta_0 = g(\theta')$

⑧ Converge?

Go to ①