

In regression and other settings the primary inference objectives can often be expressed in the form

$$E(f) = \int f(z) p(z) dz \quad (1)$$

For instance in regression we might have

$$z = [\beta, \sigma^2], \quad f(z) = p(y_{\text{new}} | x_{\text{new}}, z)$$

with

$p(z) = p(z | X, Y)$  the posterior distribution over the model parameters, meaning that

$$E(f) = \int p(y_{\text{new}} | x_{\text{new}}, z) p(z | X, Y) dz$$

is the posterior predictive distribution of  $y_{\text{new}}$  given  $x_{\text{new}}$  and the post. distribution  $p(z | X, Y)$

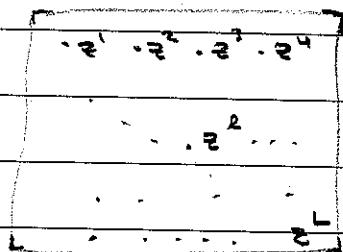
Today: approaches for generating "samples" from  $p(z)$  so that (1) can be approximated as

$$E(f) \approx \frac{1}{L} \sum f(z^{(l)})$$

with  $z^{(l)} \sim p(z)$

First intuitive (but bad) approach:

Grid the domain of  $p(z)$

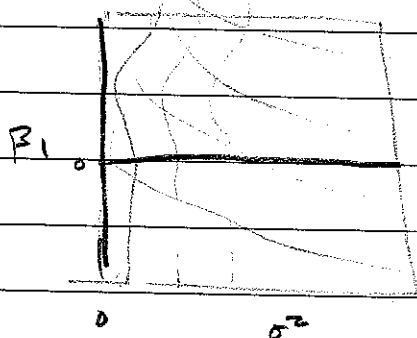


consider  $z \in \mathbb{R}^2$

and compute

$$E(f) \approx \sum_{\ell=1}^L p(z^{(\ell)}) f(z^{(\ell)})$$

ie. Bayesian regression through origin  
with prior like:



$$p(z^0 | \lambda, \mu_{\beta}, \sigma^2) = \lambda \exp(-\lambda \|z^0\|^2) \times N(\beta_1 | \mu_{\beta}, \sigma_{\beta}^2)$$

But gridding requires the number of points to grow exponentially in the dimension of  $z$

Better: choose points  $z^{(\ell)}$  only where  $p(z)$  large\* ... but how?

\* ideally where  $f(z)p(z)$  large!

One way: Importance sampling

Assume that we have  $q(z)$  which is easy to draw samples from

$$\text{i.e. } q(z) = N(z | \mu_q, \Sigma_q)$$

We can rewrite

$$E(f) = \int f(z) p(z) dz$$

$$= \int f(z) \frac{p(z)}{q(z)} q(z) dz$$

$$\approx \sum_{k=1}^L \frac{p(z^k)}{q(z^k)} f(z^k)$$

$$\text{where } z^k \sim q \text{ and } w_k = \frac{p(z^k)}{q(z^k)}$$

are known as importance weights

Remember:

1) We're doing this because  $p(z)$  is difficult to sample from:

When does this happen?

- When parametric form is not easily transformed or doesn't have an easy analytic description.

$$\text{eg. } p(z | x, y) = \frac{p(z | x, y) p(z)}{\int p(z | x, y) p(z) dz}$$

What if we could get away with knowing  $p(z)$  up to a normalizing constant only? i.e.

$$p(z|x, y) \propto p(z|x, y) p(z)$$

then we could, for instance, do posterior inference without needing to be able to compute the posterior normalizing constant.

To be most general let's consider

$$p(z) = \frac{\tilde{p}(z)}{Z_p} \quad \text{and} \quad q(z) = \frac{\tilde{q}(z)}{Z_q}$$

↑ is proposal  
only known up to a  
normalizing constant too

Remember that we can sample from  $q(z)$  even if we only know it up to a normalizing constant, i.e.  $q(z) = \frac{\tilde{q}(z)}{Z_q}$ .  
sampling from  $\tilde{q}(z)$  is the same as sampling from  $q(z)$

In this situation we can write

$$\begin{aligned} E(f) &= \int f(z) p(z) dz \\ &= \int f(z) \frac{\tilde{p}(z)}{Z_p} \frac{\tilde{q}(z)}{Z_q} q(z) dz \\ &= \frac{Z_p}{Z_q} \int f(z) \frac{\tilde{p}(z)}{\tilde{q}(z)} q(z) dz \end{aligned}$$

$$\approx \frac{z_a}{z_p} \frac{1}{L} \sum_{l=1}^L \tilde{r}_l f(z^l)$$

where

$$\tilde{r}_l = \frac{\tilde{p}(z^l)}{\tilde{q}(z^l)} \quad \text{and} \quad z^l \sim q$$

But now what about  $\frac{z_a}{z_p}$ ?

$$\frac{z_p}{z_a} = \frac{1}{z_a} \int \tilde{p}(z) dz$$

$$z_a = \frac{\tilde{q}(z)}{\tilde{q}(z)} = \text{is constant so}$$

$$\begin{aligned} \frac{z_p}{z_a} &= \frac{\int \tilde{p}(z) dz}{z_a} = \int \frac{\tilde{p}(z)}{\tilde{q}(z)} q(z) dz \\ &\approx \frac{1}{L} \sum_{l=1}^L \tilde{r}_l \end{aligned}$$

So  $E(f) \approx \sum w_l f(z^l)$  where  $z^l \sim q$

$$\text{and } w_l = \frac{\tilde{r}_l}{\sum_m \tilde{r}_m} = \frac{\tilde{p}(z^l)/\tilde{q}(z^l)}{\sum_m \tilde{p}(z^m)/\tilde{q}(z^m)}$$

So now, we can do posterior inference numerically without computing any integrals of the Bayesian normalization kind explicitly.

What's wrong?

if  $\tilde{q}$  doesn't align well (tough to know in advance), most weights will be small and the approximation will be dominated by a few (or worse, no samples) -- no diagnostics to help

A solution: get better samples from  $p(z)$  or  $\tilde{p}(z)$ . One way (not without it's own set of problems) MCMC

Simplest recipe: Metropolis algorithm

Ingredients:  $\tilde{p}(z)$ , distribution to sample from  
 $q(z|z^T)$  "proposal" distribution with  
 $q(z|z^T) = q(z^T|z)$   
i.e. symmetric

Algorithm—

Set  $\tau = 0$ , pick  $z^{\tau=0}$  arbitrarily in domain  
Repeat Forever {  
Propose  $z^*$  from  $q(z^*|z)$   
"Accept"  $z^*$  w.p.  
 $A(z^*, z^T) = \min\left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^T)}\right)$   
Set  $\tau = \tau + 1$   
If  $z^*$  "accepted" set  
 $z^\tau = z^*$   
otherwise set  
 $z^\tau = z^{\tau-1}$

Note: if ratio is  $> 1$   
sample is always  
accepted

Claim  $\{z^\tau\}_{\tau=0}^T$  are samples  
from  $\tilde{p}$  and correspondingly  
from  $p$  as well.

Justification of Metropolis sampling

Realize:  $z^0, \dots, z^T, \dots, z^T$   $T \rightarrow \infty$  forms  
a Markov chain with a given transition  
function

$$z^{t+1} = \begin{cases} z^* \text{ w.p. } A(z^*, z^t) \\ z^t \text{ w.p. } 1 - A(z^*, z^t) \end{cases}$$

Remember Markov chains have "equilibrium"  
distributions.

Markov chain (1<sup>st</sup> order): A first order Markov  
chain is defined to be a series of R.V.s

$$z^1, \dots, z^u \text{ st.}$$

$$p(z^{u+1} | z^1, \dots, z^u) = p(z^{u+1} | z^u)$$

We can specify a Markov chain by  
giving  $p(z^0)$  and the transition probabilities

$$T_u(z^u | z^{u+1}) = p(z^{u+1} | z^u)$$

A Markov chain is homogenous if the  
transition probabilities are the same for  
all  $u$ .

The marginal probability of a particular variable  $z^{m+1}$  can be expressed in terms of the marginal probability of the previous variable in the chain by

$$p(z^{m+1}) = \sum_{z^m} p(z^{m+1} | z^m) p(z^m)$$

A dist.  $p^*(z)$  is invariant or stationary w.r.t. the Markov chain if each step in the chain leaves it unchanged

For a homogeneous Markov chain with transition probs  $T(z', z)$ , the dist  $p^*(z)$  is invariant if

$$p^*(z) = \sum_{z'} T(z', z) p^*(z')$$

Intuition

$$z = T z \quad \begin{array}{c} z' \\ \begin{bmatrix} .5 & .25 & .25 \\ .25 & .5 & .25 \\ .25 & .25 & .5 \end{bmatrix} \end{array} \quad \begin{array}{c} z' \\ \begin{bmatrix} .3 \\ .3 \\ .3 \end{bmatrix} \end{array}$$

$z = T z$  is eigenvector of  $T$  with eigenvalue 1



A sufficient (but not necessary) condition for ensuring that  $p^*(z)$  is invariant under the Markov chain is to choose transition probabilities to satisfy "detailed balance", defined by

$$p^*(z) T(z, z') = p^*(z') T(z', z)$$

To show that a transition probability that satisfies detailed balance w.r.t. to a particular dist will leave that dist. invariant note that

$$\begin{aligned} \sum_{z'} p^*(z') T(z', z) &= \sum_{z'} p^*(z) T(z, z') \quad \text{1st line from} \\ &= p^*(z) \sum_{z'} T(z, z') \\ &= p^*(z) \sum_{z'} p(z'|z) \\ &= p^*(z) \end{aligned}$$

starting from here we get

Remember the goal: construct a Markov chain s.t. "running" the Markov chain yields samples from a  $p^*(z)$  of our choice.

This means picking  $T(z', z)$  s.t.  $p^*(z)$  is the distribution we're interested in drawing samples from. (at a minimum)

## Metropolis Algorithm -

$$A(z^*, z^T) = \min\left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^T)}\right)$$

We have to show that  $p(z)$  is the invariant distribution of the Markov chain defined by  $T(z^{T+1}, z^T)$

$$= \begin{cases} z^*, & z^* \sim q(z^* | z^T) \text{ w.p. } \min\left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^T)}\right) \\ z^T & \text{w.p. } 1 - \min\left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^T)}\right) \end{cases}$$

We can do this by demonstrating that the Markov chain defined by the Metropolis algorithm satisfies detailed balance, i.e.

$$\begin{aligned} p(z) T(z, z') &= p(z') T(z', z) \\ p(z) q(z' | z) A(z' | z) &= \min(p(z) q(z' | z), \frac{p(z') q(z | z') p(z)}{p(z)}) \\ &= \min(p(z) q(z' | z), q(z' | z) p(z')) \\ &= \min(p(z') q(z' | z), p(z) q(z' | z)) \\ \text{but } q(z' | z) &= q(z | z') \text{ by assumption and} \\ &= p(z') q(z | z') \min\left(1, \frac{p(z)}{p(z')}\right) \\ &= p(z') q(z | z') A(z | z') = T(z', z) \end{aligned}$$

which shows that  $p(z)$  is an invariant dist. of the Metropolis Alg.

$$T(z, z') \quad P^*(z)$$

$$z' \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad \begin{bmatrix} p_1 \\ p_2 \end{bmatrix}$$

$$T(z', z) \quad P^*(z')$$

$$= z \begin{bmatrix} a & c \\ b & d \end{bmatrix} \quad \begin{bmatrix} p_1 \\ p_2 \end{bmatrix}$$

$$2p_1 a + p_2 b = p_1 a + p_2 c$$

$$p_1 c + p_2 d = p_1 b + p_2 d$$

$$a + c = 1$$

$$a + b = 1$$

$$b + d = 1$$

$$c + d = 1$$

$$b = c, \quad c = b \quad /$$

$$T(z, z') \quad P(z)$$

$$z' \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

$$= z \begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

$$\text{s.t.} \quad \begin{aligned} a+b+c &= 1 \\ d+e+f &= 1 \\ g+h+i &= 1 \end{aligned}$$

$$\begin{aligned} \cancel{a}p_1 + bp_2 + cp_3 &= \cancel{a}p_1 + dp_2 + gp_3 \\ d\cancel{p_1} + \cancel{e}p_2 + fp_3 &= b\cancel{p_1} + e\cancel{p_2} + hp_3 \\ g\cancel{p_1} + h\cancel{p_2} + i\cancel{p_3} &= c\cancel{p_1} + f\cancel{p_2} + i\cancel{p_3} \end{aligned}$$

$$\begin{aligned} a-d &= 0 \\ a+d+g &= 2 \\ b+e+h &= 2 \\ c+f+i &= 2 \end{aligned}$$

$$bp_2 + cp_3 + dp_1 = b\cancel{p_1} + dp_2 + gp_3$$

$$b(p_2 - p_1) = dp_2 + gp_3 + cp_3 + d\cancel{p_1}$$

$$b = \frac{dp_2 + gp_3 + cp_3 + d\cancel{p_1}}{p_2 - p_1}$$

$$T_P = T'_P$$

if  $T$  orthogonal  
and symmetric then

$$T^{-1} = T' \quad \text{and} \quad T^{-1}T_P = T' T_P$$

$$P = P$$

We also require that for  $n \rightarrow \infty$ , the dist  $P(\mathbf{z}^n)$  converges to the desired dist  $p^*(\mathbf{z})$  irrespective of choice of  $p(\mathbf{z}^{(0)})$

This prop. is called ergodicity and the invariant dist of an ergodic Markov chain is known as its equilibrium dist. An ergodic Markov chain can have only one equilibrium dist.

It can be shown that: a homogeneous Markov chain will be ergodic subject to weak restrictions on the invariant dist. and the transition probs.

← Fill

The Metropolis Algorithm is

- 1) homogeneous
- 2) ergodic
- 3) and has  $p^*$  as its invariant dist.

It is a general <sup>Markov chain</sup> transition rule that is parameterized by the distribution of interest  $p(\mathbf{z})$  and results in an ergodic Markov chain whose equilibrium dist is  $p(\mathbf{z})$ .

Uses: sample from posterior dist for use in Monte Carlo integrals like posterior predictive and/or posterior inference.

Metropolis Hastings :  $q$  not symmetric

$$p(z) q(z'|z) A(z', z) = p(z) q(z'|z) \min\left(1, \frac{p(z') p(z|z')}{p(z) q(z'|z)}\right)$$

$$= \min(p(z) q(z|z), p(z') p(z|z'))$$

$$= \min(p(z|z') p(z'), p(z) q(z'|z))$$

$$= p(z') q(z|z') \min\left(1, \frac{p(z) q(z'|z)}{p(z') q(z|z')}\right)$$

$$= p(z') q(z|z') A(z|z') \quad \checkmark$$