

## Sampled representations:

crucial, useful fact

$$x_i \sim q(x)$$

then

$$\frac{1}{N} \sum f(x_i) \rightarrow \int f g dx$$

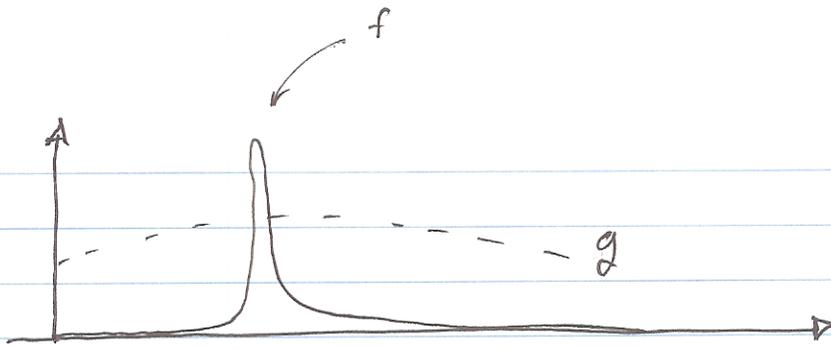
This gives a method to represent Prob densities by samples.

→ it's a representation, because we can compute expectations

~~Notice this~~

But we may not get good est's:

(2)



$x_i \sim g$  will most likely give a high variance estimate of  $E_g[f] = \int fg dx$

Importance sampling:

assume we know  $h$ , which is "similar" to  $fg$  (big when  $fg$  is big).

1)  $x_i \sim h$

2)  $\frac{1}{N} \sum f(x_i) \cdot \frac{g(x_i)}{h(x_i)} \rightarrow \int fg dx$

will be a better estimate  
(lower variance)

This suggests a representation

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

Representing a posterior from a prior

we have

$$\begin{array}{ccc}
 P(x|\theta) & & P(\theta) \\
 \uparrow \text{likelihood} & & \uparrow \text{prior}
 \end{array}$$

- assume we have  $\theta_i \sim p(\theta)$

- want to represent

$$\begin{aligned}
 p(\theta|x) &= \frac{p(x|\theta) p(\theta)}{p(x)} \\
 &= \frac{p(x|\theta) p(\theta)}{\int p(x|\theta) p(\theta) d\theta}
 \end{aligned}$$

Now

$$\text{if } w_i = p(x/\theta_i)$$

then

$$\frac{1}{N} \sum_i w_i \rightarrow \int p(x/\theta) p(\theta) d\theta$$

and

$$\frac{1}{N} \sum_i f(\theta_i) w_i \rightarrow \int f(\theta) p(x/\theta) p(\theta) d\theta$$

So

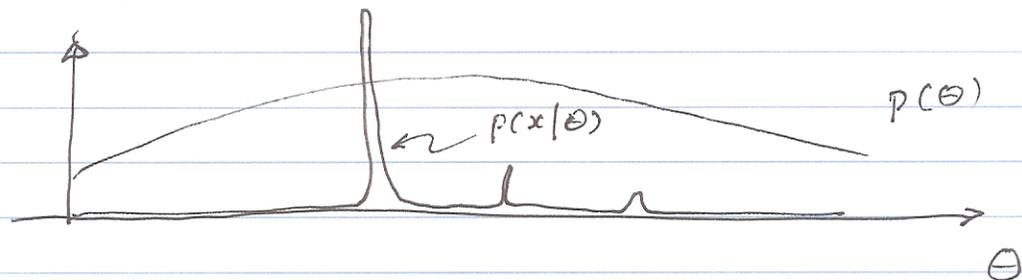
$$\frac{\sum_i f(\theta_i) w_i}{\sum_i w_i} \rightarrow \int f(\theta) p(\theta|x) d\theta$$

This gives one amazingly simple recipe  
for representing a posterior!

---

But:

it's often not very good.



In particle filtering, the effect that results is sometimes called "sample impoverishment"

- Many  $w_i$  are v. small, few are big
- ⇒ very few  $\theta_i$  contribute to the estimate of the integral in any significant way.

→ test:

look at variance of  $w_i$

## How do we get samples?

- Uniform, Normal, Some <sup>(few)</sup> others:

there are standard algs

- Finite mixtures of above:

- Draw a sample from mixture weights

- then from component

## Rejection sampling:

- wish to draw a sample from  $p(x)$

- know how to draw from  $q(x)$

- $p(x) \leq c q(x)$

Algorithm:

- repeat until accepted
- draw  $x_i \sim q$
- draw  $y \sim U[0, 1]$
- if  $y \leq \frac{p(x_i)}{c q(x_i)}$  accept
- else go again.

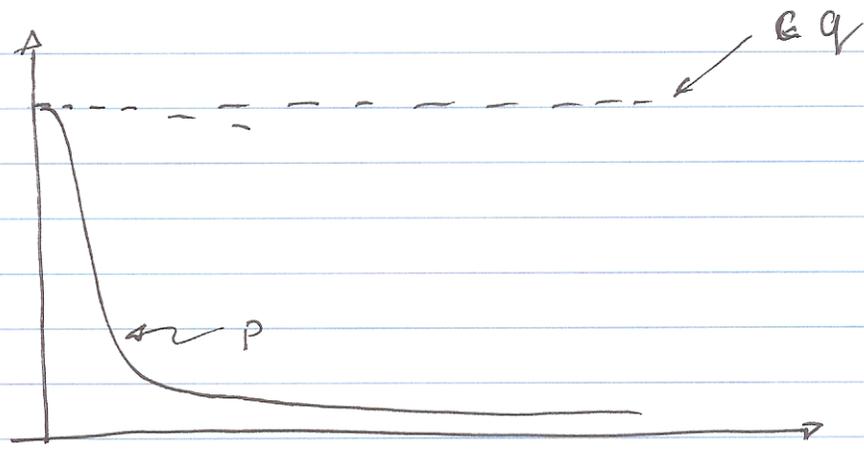
1) Why does this work?

$$\begin{aligned}
 P\{\text{sample at } x_i\} &= P\{\text{generate sample at } x_i \text{ and accept}\} \\
 &= q(x_i) \cdot \frac{p(x_i)}{c q(x_i)}
 \end{aligned}$$

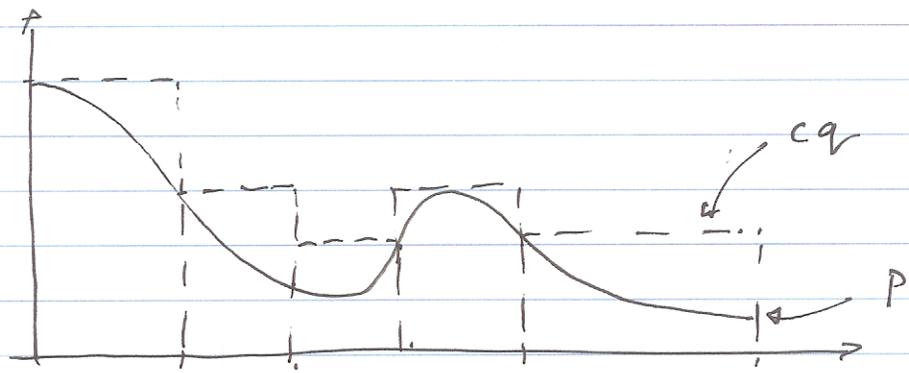
So  $P\{\text{sample at } x_i\} \propto p(x_i)$

2) How efficient is it?

- Depends on relationship between  $p, q$ .



many rejects in this case



here we have multiple intervals and draw from c q as a mixture

⑨

- Drawing a sample from a Discrete distribution



- draw  $x \sim U[0, 1]$  to point location

Or:

build a binary tree,  
to rejection sampling

Or:

build a binary tree,  
use a ~~binary~~ biased sample.

# Gibbs sampling :

- assume we want to draw a sample from a multivar distrib.

$$p(u, v)$$

- assume  $p(u|v)$  is easy  
 $p(v|u)$  is easy.

- Procedure: start with  $(u_0, v_0)$

- draw  $u_i \sim p(u|v_{i-1})$

- $v_i \sim p(v|u_i)$

- do this many times

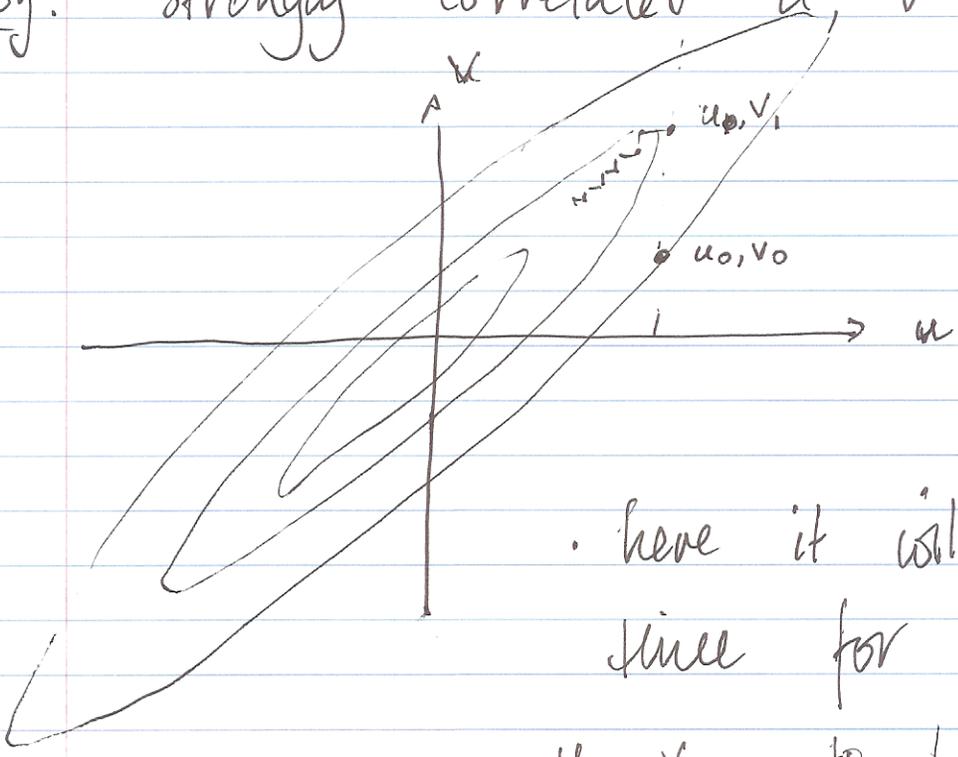
$$(u_n, v_n) \sim p(u, v)$$

(if we iterate often enough).

eg  $p(u, v) \propto \exp\left[-\frac{1}{2}\left(\frac{u^2}{\sigma_1^2} + \frac{v^2}{\sigma_2^2}\right)\right]$

- each stage is a draw from a univariate gaussian
- not much real advantage, because we could draw  $(u, v) \sim N(0, Id)$  then scale.

eg: strongly correlated  $u, v$



• here it will take a long time for  $u_n, v_n$  to have high prob.

# Markov chain

sequence of random vars,

$X_i \dots$  such that

$$\begin{aligned}
 P(X_{i+1} = x \mid X_i = \dots \mid X_0 = \dots) \\
 = P(X_{i+1} = x \mid X_i = \dots)
 \end{aligned}$$

And this transition is stationary.

Discrete case:

Write  $P(X_{i+1} = x_u \mid X_i = x_v) = P_{uv}$

Some matrix ↑

Now, start with

$$P(X_0 = x_w) = \pi_w$$

$$\begin{aligned}
 P(X_1 = x_s) &= \sum_{x_w} P(X_1 = x_s \mid X_0 = x_w) P(X_0 = x_w) \\
 &= \sum_w P_{sw} \cdot \pi_w = P \pi
 \end{aligned}$$

Similarly:

$$\begin{aligned}
P(X_i = x_s) &= \sum_{w_{i-1}, \dots} P_{s w_{i-1}} P_{w_{i-1} w_{i-2}} \dots P_{w_i w_0} \pi_{w_0} \\
&= P^k \pi
\end{aligned}$$

• Now  $P$  has at least 1 unit eigenvalue  
 (because  $P^T$  does,  $P^T \mathbf{1} = \mathbf{1}$ )

• Assume only 1; all others  $< 1$

• then  $P^k v \rightarrow \pi_s$  such  
 that  $P \pi_s = \pi_s$

Stationary distribution

This yields following procedure to sample from discrete dist  $\Pi_S$

- Build MC with  $\Pi_S$  as stationary dist
- "run the chain"

eg to sample from

$$\frac{1}{2}$$

$$\frac{1}{2}$$

chain:

$$\begin{bmatrix} 1-\epsilon & \epsilon \\ \epsilon & 1-\epsilon \end{bmatrix}$$

- Something seems like a problem here  $\rightarrow$
- No guarantee that  $x_i \parallel x_{i-1}$ !

Q: How do we build such a chain?

- Design requirement

$$\sum_{x_w} P(X_{i+1} = x_u | X_i = x_w) \pi(X_i = x_w) = \pi(X_{i+1} = x_u)$$

or

$$\sum P_{uw} \pi_w = \pi_u$$

This is a weak condition, rather hard to achieve by choice of  $P$

- Detailed balance

(stronger condition)

$$P_{uw} \pi_w = P_{wu} \pi_u$$

For each  $u, w$

Notice this implies above. →

Now we can build an MC in a continuous domain with stationary

Dist  $\pi(x)$

1) obtain Proposal process  $P(x \rightarrow y)$

(this is  $P(X_{i+1} = y | X_i = x)$ )

- $P(x \rightarrow y)$  should ~~have~~ satisfy some support conditions.

• Good enough if

$$\pi(y) > 0 \Rightarrow P(x \rightarrow y) > 0$$

$$P(x \rightarrow y) > 0 \Rightarrow P(y \rightarrow x) > 0$$

2) Let we achieve detailed balance

by a form of rejection sampling

Algorithm

given  $x_i$

• draw  $x_{i+1}^P \sim P(x_i \rightarrow y)$

• compute

$$\alpha = \min\left(1, \frac{P(x_i \rightarrow x_{i+1}^P) \pi(x_{i+1}^P)}{P(x_i \rightarrow x_i^P) \pi(x_i)}\right)$$

$$x_{i+1} = \begin{cases} x_{i+1}^P & \text{with prob } \alpha \\ x_i & \text{" " " } 1-\alpha. \end{cases}$$

Thm: for large enough  $N$ ,

$$x_N \sim \pi$$

whatever  $x_0$

# Amazingly powerful algorithm

+ • can draw samples from unnormalized dists

+ • easy; general

? + • can sometimes prove that chains are "fast mixing" (i.e. forget start point quickly)

? ~~+~~ • Samples are correlated

- • Correlation may be hard to est.

- • Can be hard to tell if sampler is "sticky"

- • "Sticky" samplers look better.