

Sampling Methods

$$\langle f \rangle = \int f(x) p(x) dx \quad (1.1)$$

$$\hat{f} = \frac{1}{M} \sum_{m=1}^M f(x^{(m)}) \quad (1.2)$$

$x^{(m)} \sim p(x)$

$$\Rightarrow \langle \hat{f} \rangle = \langle f \rangle$$

$$\text{var}[\hat{f}] = \frac{\sigma^2}{M} \quad \sigma^2 = \langle (f - \langle f \rangle)^2 \rangle \quad \text{var}[\hat{f}] \downarrow \text{as } M \rightarrow \infty.$$

- Issues:- independent samples  $\{x^{(m)}\}_{m=1}^M$  ?  
 $f(x_1)$  large,  $p(x_1)$  small }  $\rightarrow$  interaction of magnitude of  $f$   
 $f(x_0)$  small,  $p(x_0)$  large . } with probability weighting.

Ancestral Sampling in GM: GM joint: -  $p(x) = \prod_{i=1}^d p(x_i | x_{1:(i)})$

- sample following GM flow

- no observed variables.

Logic Sampling in GM

- some nodes instant. at observed values
- sample  $x_i$  and observed agree  $\rightarrow$  continue sampling.
- else discard.

1.0.1 - Sampling and EM: - direct impl. of Bayesian network.  
- frequentist  $\rightarrow$  2-step ML est. with mvtct. posterior.

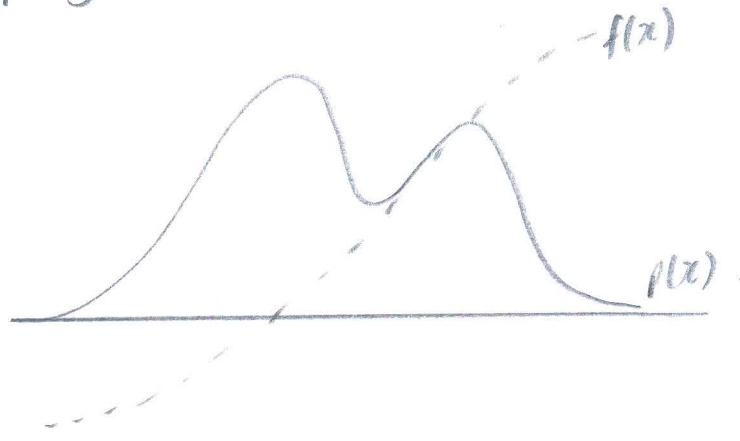
$x_H$  - hidden  
 $x_V$  - visible

$\theta$  - param.

M-step  
optimises :-

ELL

$$\begin{aligned} Q(\theta, \theta_{\text{old}}) &= \mathbb{E}_{p(x_H | x_V, \theta_{\text{old}})} [\ln p(x_H, x_V | \theta)] \\ &= \int p(x_H | x_V, \theta_{\text{old}}) \ln p(x_H, x_V | \theta) d x_H \end{aligned} \quad (1.5)$$



approx 1.5 via finite sum over samples  $\{x_H^{(i)}\}$  drawn from current est.  
for posterior  $p(\theta_H | x_V, \theta_{old})$  :-

$$a(\theta, \theta_{old}) \approx \frac{1}{M} \sum_{m=1}^M \ln p(x_H^{(m)}, x_V | \theta) \quad (1.6)$$

- then optimise  $Q$  in usual way in M-step

(\*) that was Monte Carlo EM

(\*) augment to find mode of posterior over  $\theta$  (MAP estimate) by placing prior  $p(\theta)$  and adding  $\ln p(\theta)$  to  $Q(\theta, \theta_{old})$  in M-step.

(\*) stochastic EM.

(\*) full Bayesian treatment (JP/drawing algo) :-

- sample from post over  $\theta$

- draw from joint post  $p(\theta, x_H | x_V) \rightarrow$  computationally difficult.

- Imputation-post.

I-step: sample from  $p(x_H | x_V) \rightarrow$  cannot directly.

$$\text{- note: } p(x_H | x_V) = \int p(x_H | \theta, x_V) p(\theta | x_V) d\theta$$

- for  $m = 1, \dots, M$  :-

i) sample  $\theta^{(m)}$  from curr est. for  $p(\theta | x_V)$ .

ii) use this to sample  $x_H^{(m)}$  from  $p(x_H | \theta^{(m)}, x_V)$ .

$$\text{P-step: note } p(\theta | x_V) = \int p(\theta | x_H, x_V) p(x_H | x_V) dx_H$$

- use  $\{x_H^{(m)}\}_{m=1}^M$  from I-step to compute revised est. of post. over  $\theta$

given by

$$\text{- } p(\theta | x_V) = \frac{1}{M} \sum_{m=1}^M p(\theta | x_H^{(m)}, x_V)$$

## 1.1. Basic Sampling Algos

- sample from dist'n defined over univariate  $x \in \mathbb{R}$
- Pseudo-random  $\rightarrow$  deterministic, passes tests for randomness
- pseudo-random over  $Unif(0,1)$
- this is basis for other algos.

### 1.1.1. Standard dist'n's

- PRNG - successive app. of trans  $\theta(\cdot)$
- sequence:  $(x^{(1)}, x^{(2)}, \dots, x^{(n)})$  via recurrence  $x^{(n)} = \theta(x^{(n-1)})$
- seed/initial value:  $x^{(1)}$ .
- choice of  $\theta(\cdot)$  - finite rep of digit computers.
- integer representation: -
  - $x \in \{0, \dots, N\}$   $\rightarrow$   $N$ -largest poss. integer m comp.
- interval  $[0,1)$  generated by: -  $\frac{x}{N+1}$

- period?: - smallest integer  $T$ :  
 $x^{(n+T)} = x^{(n)} \quad \forall n$ . or:  $\theta(\cdot)^T$  is identity op.
- largest value which simple gen: -  $x^{(n)} = \theta(x^{(n-1)}) \Rightarrow (N+1) \text{?}$

where initial:  $\theta(x) = (ax+b) \bmod (N+1)$   $\xrightarrow{(a,b) \text{- select with care}}$

- non-uniform - assume source of uniformly dist'n random nos.
- PRNG:  $x \sim U(0,1)$ .

- transform  $x$  via  $f(\cdot)$ :  $y = f(x)$ .  $\xrightarrow{\text{Q:}}$  - is this similar to transforms of r.v.s?  
 $p(y) = p(x) \left| \frac{dx}{dy} \right|$   $p(x) = 1 \text{? } \int p(x) dx = 1 \text{?}$

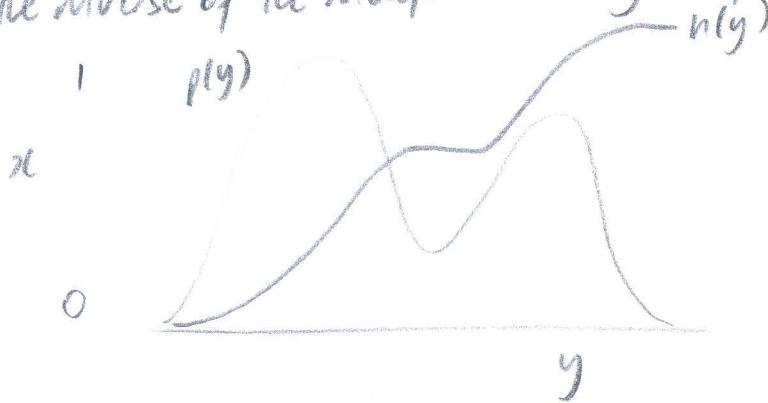
- select  $f(\cdot)$  such that resulting values of  $y$  have desired dist'n  $p(y)$ . 0/51
- missing a step.

-  $x = h(y) \equiv \int_{-\infty}^y p(y') dy'$  (indefinite integral of  $p(y')$ )

*inuse  
transform. tighten this*

$$\Rightarrow y = h^{-1}(x)$$

④ transform uniformly distributed random nos using a function  $f(\cdot)$  that is the inverse of the cumulative integral of desired distri.



e.g. exponential

$$p(y) = \lambda \exp(-\lambda y) \quad (1.12)$$

$$h(y) = 1 - \exp(-\lambda y)$$

- use  $y = h^{-1}(x) = -\lambda^{-1} \ln(1-x)$  when  $y \sim \text{Exp}$ .

e.g. Cauchy

$$p(y) = \frac{1}{\pi} \frac{1}{1+y^2} \quad (1.13)$$

$$y = h^{-1}(x) = \tan(\pi x)$$

multiple variables

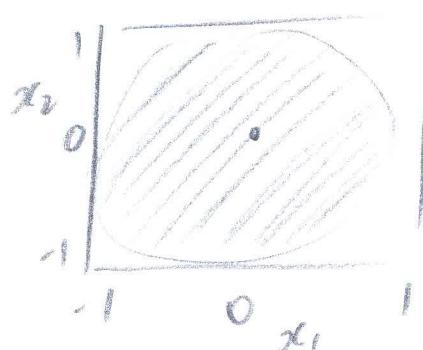
$$p(y_1, \dots, y_d) = p(x_1, \dots, x_d) \left| \frac{\partial(y_1, \dots, y_d)}{\partial(x_1, \dots, x_d)} \right| \quad \textcircled{D} X$$

Box Muller:

- generate uniformly distri  $(x_1, x_2) \in (-1, 1)$

- discard pair unless  $x_1^2 + x_2^2 \leq 1$

- uniform distri



$$p(x_1, x_2) = \frac{1}{\pi} \quad \textcircled{D}$$

• for each pair  $(x_1, x_2)$ , evaluate:-

$$y_1 = x_1 \left( \frac{-2\ln x_1}{n} \right)^{1/2} \quad y_2 = x_2 \left( \frac{-2\ln x_2}{n} \right)^{1/2} \quad (1.15)/(1.16)$$

$$\rho^2 = \tilde{x}_1 + \tilde{x}_2$$

Joint distn:  $f(y_1, y_2) = f(x_1, x_2) \left| \frac{\partial(y_1, y_2)}{\partial(x_1, x_2)} \right|$  (1.17) ?

$$= \left[ \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y_1^2}{2}\right) \right] \left[ \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y_2^2}{2}\right) \right] \quad (1.18)$$

•  $y_1, y_2$  map; zero-mean, unit variance.

• finally use transformation:  $y \rightarrow \delta y + \mu$  to get normal. with  $\mu, \sigma^2$ .

• vector-valued multi-variate Normal, mean  $\mu$ , cov  $\Sigma$ .

• eigenvalue/vector decmp. of  $\Sigma$ :

$$\Sigma u_i = \lambda_i u_i \quad (1.19)$$

• If  $x_i$  are multivariate, normally distn, then

$$y = \mu + \sum_i \lambda_i^{1/2} x_i u_i \text{ has required distn}$$

• use Cholesky decmp of  $\Sigma = W^T m$  practice.

• rejection + importance samp!  $\rightarrow$  multivariate distns.

### 1.12. rejection sampling

• sample from  $p(x) = \frac{1}{Z} \hat{p}(x)$   
is difficult  
(analytic. intract.)

• introduce a proposal distn  $q(x)$ . and constant  $K$ :

$$K q(x) \geq \hat{p}(x) \quad \forall x.$$

$\hat{p}(x)$  - readily evaluated  
 $Z$  - unknown normalising constant.

$Rg(x)$  - comparison function

• generate 2 random nos.

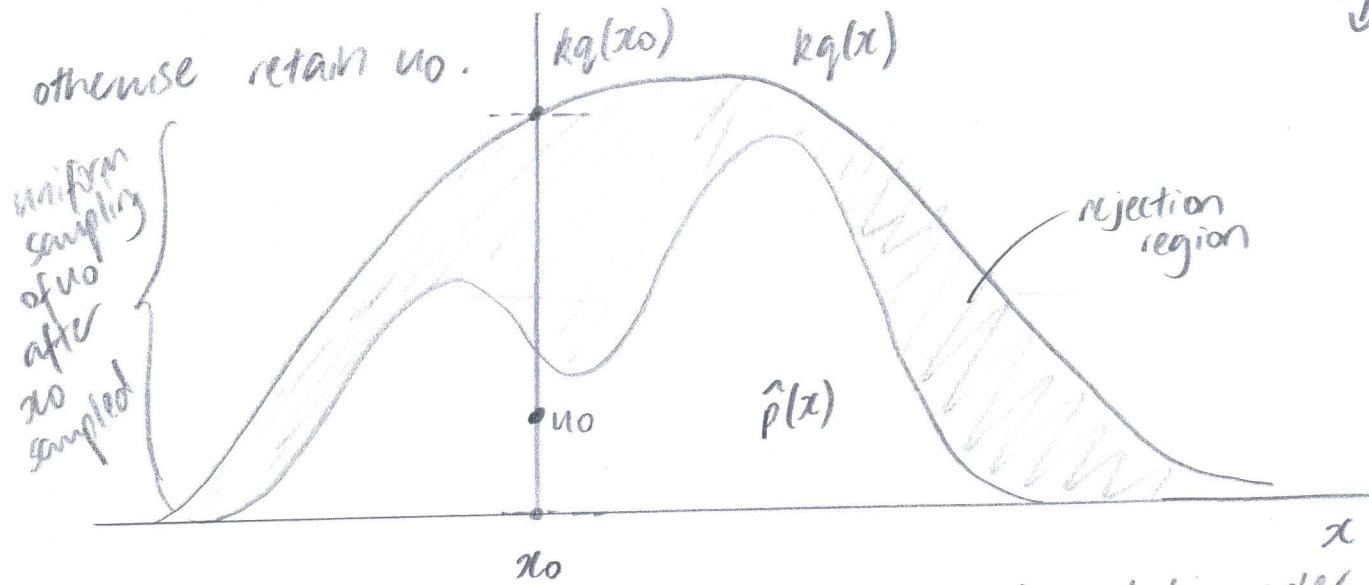
$$x_0 \sim g(x)$$

$$u_0 \sim \text{unif}(0, Rg(x_0)) \quad (\text{no } x_0)$$

•  $(x_0, u_0)$  have uniform distri under curve of  $Rg(x)$

• If  $u_0 > \hat{p}(x_0)$  then reject sample

otherwise retain  $u_0$ .



• Remaining pairs in unshaded region  $\rightarrow$  have uniform distri under the curve of  $\hat{p}(x)$

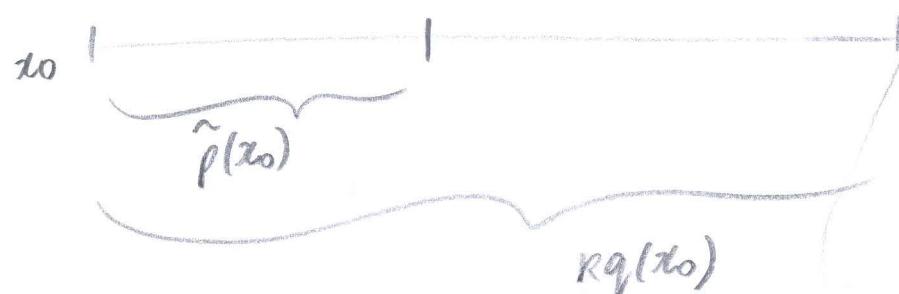
• corresponding  $x$  values distri according to  $p(x)$  ②

(\*) correctness of rejection sampling (from  $p(x)$ )

• original values  $x_0$  generated from  $g(x)$  distri.

• these samples are accepted with probability  $\frac{\hat{p}(x_0)}{Rg(x_0)}$

↳ 0 accept  $\hat{p}(x_0)$  reject.  $Rg(x_0)$



(\*) resulting distn of  $x$  after norm:-

$$\frac{[\hat{p}(x)/Rq(x)]q(x)}{\int [\hat{p}(x)/Rq(x)]q(x)dx} = \frac{R\hat{p}(x)}{\int \hat{p}(x)dx} = \frac{\hat{p}(x)}{\int \hat{p}(x)dx} = p(x).$$

(\*\*) Prob. sample accepted:-  $\text{prob}(\text{accept})$

$$\text{prob}(\text{accept}) = \int [\hat{p}(x)/Rq(x)]q(x)dx = \frac{1}{R} \int \hat{p}(x)dx$$

①  
② ③/53

- fraction of points rejected depends on ratio of area under unnormalised  $(\hat{p}(x))$  compared to area under curve  $Rq(x)$ .

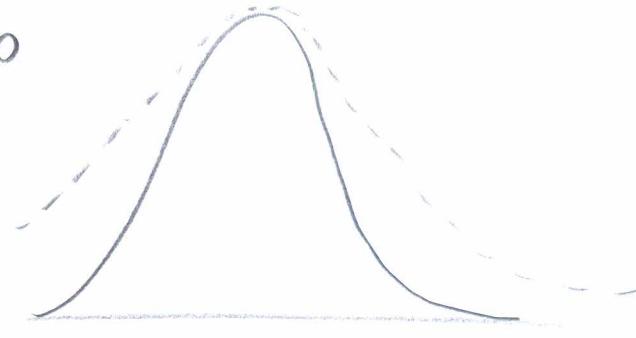
- difficulty seeing this intuitively

④ constant  $R$  must be small as possible subject to constraint that  $Rq(x)$  is nowhere less than  $\hat{p}(x)$

example:

- sample from gamma:-

$$\text{Gam}(x|a) = \frac{x^{a-1} \exp(-ax)}{\Gamma(a)} \quad a > 0$$



- use cauchy as  $q(x)$ .

- generalise cauchy to ensure that nowhere is it smaller than gamma distri.

- transform I.V.Z:-

$$x = b \tan Z + c$$

$$\text{yielding: } q(x) = \frac{R}{1 + \frac{(x-c)^2}{b^2}}$$

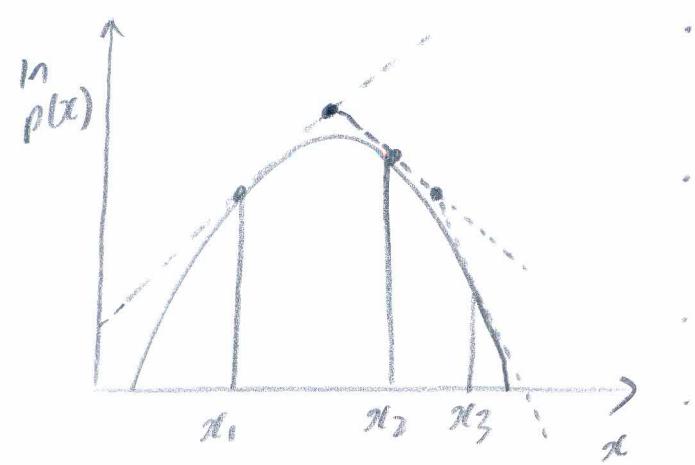
⑤ ⑥/54  
. min rejection rate:-  
- set  $c = a - 1$   $b^2 = 2a - 1$   
- set  $R$  as small as possible while ensuring  $Rq(x) \geq \hat{p}(x) \forall x$

### 1.1.3. Adaptive Rejection Sampling

univariate

log concave  $p(x)$

construction of envelope fn  $q(x)$  on the fly using grid points, tangent lines.



- function  $\ln p(x)$  evaluated at set of grid points
- use target line intersections to construct envelope functions.
- draw samples from envelope distri.

• envelope distri  $q(x)$  - piecewise exponent.:-

$$q(x) = R_i \lambda_i \exp\{-\lambda_i(x - x_{i-1})^2\} \quad x \in (x_{i-1}, x_i]$$

0/55

show this

• usual rejection criterion.

1. Acceptance  $\rightarrow$  draw from desired distri.

2. Rejection  $\rightarrow$  incorporate into grid, compute target, refine envelope.

• more grid points, envelope is better approx of  $p(x)$ ,  $p(\text{reject}) \downarrow$ .

• adaptive rejection Metropolis sampling.

other points on rejection sampling

④ comparison function be close to required distri  $\rightarrow$  minimise rejection rate

(interests of an efficient algorithm)

• dimensionality issues (illustration)

• sample from zero-mean MVG.

$$x \sim N(0, \sigma_p^2 I) \quad x \in \mathbb{R}^d \quad (\text{isotropic})$$

• do this by using an MVG proposal :-

$$Q(x) = N(0, \sigma_Q^2 I)$$

• for  $R Q(x) \geq p(x) \forall x$ ; require  $\sigma_Q^2 > \sigma_p^2$  (see diagram).

• d-dimensions, optimal value of R is given by

$$R = \left(\frac{\sigma_Q}{\sigma_p}\right)^d$$

0/56

show this

Acceptance rate:

- ratio of volumes under  $p(x)$  and  $Rq(x)$ .  
 ○  $\hat{p}(x)$ ?

Since both dists  
normalised, this will be  $\frac{1}{R}$ .

Acceptance rate =  $\frac{1}{R} = \left(\frac{\delta q}{\delta p}\right)^d$

$\delta Q^2 > \delta P^2 \Rightarrow \delta Q > \delta P \Rightarrow \frac{\delta Q}{\delta P} > 1$

As  $d \rightarrow \infty$ , acceptance rate  $\rightarrow 0$  (exponential)

Numerically:  $\delta Q = (1.1 \delta P) \times \delta a = 1.1 \delta P$ !  
 required to und. avg.

Q. Can't get these numbers/back of envelope calc's to work.

Mackay: what value of  $R$  is required if dimensionality  $d=1000$ ?

- density of  $q(x)$  at origin is  $\frac{1}{(2\pi\delta_Q^2)^{d/2}}$

- for  $Rq(x) \geq p(x)$  requires  $\Rightarrow$

$$\frac{R}{(2\pi\delta_Q^2)^{d/2}} \geq \frac{1}{(2\pi\delta_P^2)^{d/2}} \Rightarrow R \geq \frac{(2\pi\delta_Q^2)^{d/2}}{(2\pi\delta_P^2)^{d/2}} = \frac{(2\pi)^{d/2} (\delta_Q^2)^{d/2}}{(2\pi)^{d/2} (\delta_P^2)^{d/2}}$$

$$\Rightarrow R \geq \left(\frac{\delta Q}{\delta P}\right)^d \Rightarrow R \geq \exp\left(d \ln \frac{\delta Q}{\delta P}\right) \quad (\text{exp. lower bound on } R).$$

- set  $d=1000$ ,  $\left(\frac{\delta Q}{\delta P}\right) \approx 1.01 \Rightarrow R \geq \exp(10) \Rightarrow R \approx 22000$

- Acceptance rate for this value of  $R$ ?

- Ratio of volume under  $p(x)$  to volume under  $Rq(x)$ .

- As  $p$  and  $q$  both appropriately normalised i.e.  $\int p(x) dx = 1$  and  $\int q(x) dx = 1$

Then acceptance rate is  $\frac{1}{R}$ .



(0/56)

- Assess why  
you couldn't  
reproduce  
below argument  
- Arithmetic error!

∴ not suitable in high-dimensional cases.

- other issues → multi-modal, sharply peaked, difficult to find good proposal
- exponential decrease in acceptance is a feature + comparison.
- BUT subroutine for more sophis. algos

#### 1.1.4. Importance Sampling

- recall there are two problems in sampling:-

- i) Evaluating expectations under a distri  $p(x)$ .
- ii) Sampling from  $p(x)$

- Monte Carlo sampling solves (i) but not (ii). - assume diff: - draw samples from  $p(x)$  directly with supplement Mackay

- note: we discuss:-

$$\langle f \rangle = \int f(x) p(x) dx \xrightarrow{\text{approx by}} \langle f \rangle \approx \sum_{m=1}^M p(x^{(m)}) f(x^{(m)})$$

- using a discretisation of  $X \rightarrow \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$  uniformly.

- issue: no. of terms grows exponentially with dimensionality of  $x$ . ②  
- contradicts claim about accuracy of estimator and why?

- this is an example of uniform sampling

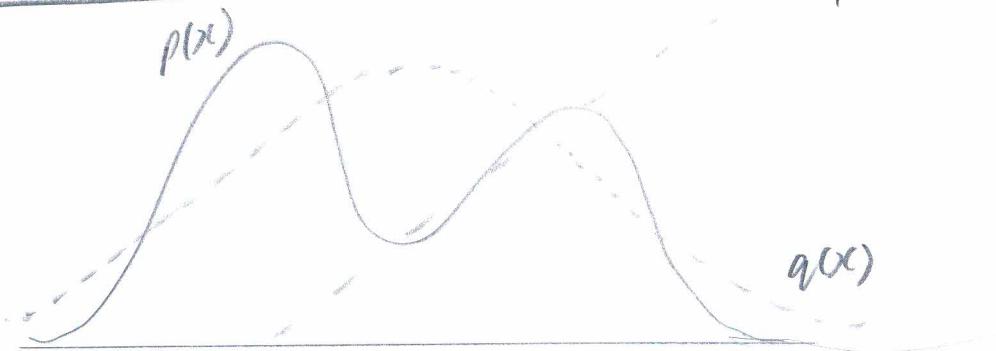
- Jordan - to solve typos and clarity.

③ inefficiency of uniform sampling in high-dim space:  
probability dists of atleast one much mass confined  
in small regions of ~~support~~  $x$ -space.

- in high-dims, only a very small proportion of samples contrib. to sum.

- ideally, select sample points to fall in regions where  $p(x)f(x)$  is large.

④ use a proposal distri  $q(x)$  to draw samples (similar to rejection-accept.)



(\*) draw samples  $\{x^{(m)}\}$  from simple  $q(x)$  and weight correc. terms in summation with  $p(x^{(m)})/q(x^{(m)})$ .

$$\begin{aligned}\langle f \rangle &= \int f(x) p(x) dx \\ &= \int f(x) \frac{p(x)}{q(x)} q(x) dx \quad - \text{IS 'trick' similar to V.I.} \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{p(x^{(m)})}{q(x^{(m)})} f(x^{(m)}) \quad x^{(m)} \sim q(x) \quad (1.28)\end{aligned}$$

Q: computing  $\left\langle \frac{p(x)}{q(x)} f(x) \right\rangle_q$ ?

- importance weights:-  $w_m = \frac{p(x^{(m)})}{q(x^{(m)})}$  - correct bias introduced by sampling from wrong distribution.

- all samples retained.

normalised importance sampling:

-  $p(x)$  can only be evaluated up to a normalisation constant.

$$p(x) = \frac{\hat{p}(x)}{Z_p} \quad \hat{p}(x) - \text{easily eval.}$$

$Z_p$  - unknown norm constant.

- use similar proposal distri  $q(x) = \frac{\hat{q}(x)}{Z_q}$

$$\langle f \rangle = \int f(x) p(x) dx$$

$$= \int f(x) \frac{p(x)}{q(x)} q(x) dx = \int f(x) \frac{\hat{p}(x)/Z_p}{\hat{q}(x)/Z_q} q(x) dx$$

$$= \frac{Z_q}{Z_p} \int f(x) \frac{\hat{p}(x)}{\hat{q}(x)} q(x) dx$$

Hence approx:-

$$\langle f \rangle \approx \frac{2q}{2p} \frac{1}{M} \sum_{m=1}^M r_m f(x^{(m)})$$

$$r_m = \frac{\hat{p}(x^{(m)})}{\hat{q}(x^{(m)})}$$

should be  
 $\hat{p}(x)$  (typo).

Evaluation of  $\frac{\partial p}{\partial q}$  :-

$$\frac{\partial p}{\partial q} = \frac{1}{2q} \int \hat{p}(x) dx = \int \frac{1}{2q} \hat{p}(x) dx = \int p(x) \frac{q(x)}{\hat{q}(x)} dx$$

$$\approx \frac{1}{L} \sum_{m=1}^L r_m$$

② typo?

L=M as same sample set used

Hence:  $\langle f \rangle \approx \frac{\sum_m r_m f(x^{(m)})}{\sum_m r_m}$

③: think there is a typo and  $p(x)$  replaced with  $\hat{p}(x)$ .

$$\frac{\partial p}{\partial q} = \frac{1}{2q} \cdot 2p = \frac{1}{2q} \int \hat{p}(x) dx = \frac{q(x)}{\hat{q}(x)} \int \hat{p}(x) dx = \int \frac{\hat{p}(x)}{\hat{q}(x)} q(x) dx$$

$$\Rightarrow \frac{\partial p}{\partial q} = \left\langle \frac{\hat{p}(x)}{\hat{q}(x)} \right\rangle_{q(x)} = \langle r_m \rangle_q$$

use finite/empirical approx of  $\frac{\partial p}{\partial q}$  :-

$$\frac{\partial p}{\partial q} \approx \frac{1}{L} \sum_{m=1}^L r_m \Rightarrow \frac{\partial p}{\partial q} = \frac{1}{\frac{1}{L} \sum_{m=1}^L r_m}$$

$$\text{As } \langle f \rangle \approx \frac{2q}{2p} \frac{1}{M} \sum_{m=1}^M r_m f(x^{(m)})$$

$$= \frac{\frac{1}{M} \sum_{m=1}^M r_m f(x^{(m)})}{\frac{1}{L} \sum_{m=1}^L r_m} = \frac{\sum_{m=1}^M r_m f(x^{(m)})}{\sum_{m=1}^L r_m}$$

6/58

- a little confused  
 (need to review again when fresher)

- reconciling with 10-708 U3.

- may just be difference mass.  
 i.e. proposal / inform constat.  
 as required.

### caveats

- dependent on how well sampling distri  $q(x)$  matches desired distri  $p(x)$ .
- $p(x)f(x)$  strongly varying, pmass concentrated in small regions of  $X$
- set of importance weights  $\{r_m\}$  drawn by few weights having large values.
- effective sample size  $\ll$  apparent samp size  $M$ .
- even worse  $\rightarrow$  none of samples fall where  $p(x)f(x)$  large.
- Apparent variances of  $r_m$  and  $r_m f(x^{(m)})$  may be small even if est. of exp very wrong.
- sampling distri  $q(x)$  should not be zero or small where  $p(x)$  sig.

(\*) Sampling distri  $q(\cdot)$  continuous densities

- use heavy-tailed  $q(\cdot)$

(\*) provides a mechanism for sampling from BN.

(\*) And also likelihood-weighted sampling

1.1.5. Sampling - Importance resampling (weighted resampling)

- rejection sampling  $\rightarrow$  determine  $R$  in comparison.
- issue:- impractical for many  $p(x), q(x)$  to determine  $R$  which is sufficiently large to guarantee bound on impractically small accept.
- SJR/weighted resampling  $\rightarrow$  use sampling distri, avoid determining  $R$ .

2 stage scheme:

- draw  $N$  samples  $x^{(1)}, \dots, x^{(N)}$  from  $q(x)$ .
- construct weights  $w^{(1)}, \dots, w^{(N)}$ :

$$w^{(n)} = \frac{\hat{p}(x^{(n)}) / q(x^{(n)})}{\sum_{n=1}^N \hat{p}(x^{(n)}) / q(x^{(n)})}$$

1st stage

2nd stage

- 2<sup>nd</sup> set of  $M$  samples drawn from discrete distri  $(x^{(1)}, \dots, x^{(N)})$  with probabilities given by  $(w^{(1)}, \dots, w^{(N)})$ .

$M \gg N$

characterising these  $M$  samples

- the  $M$  samples are only approximately distributed according to  $p(x)$ .
- that is distri becomes correct as  $N \rightarrow \infty$  (converge in distri?)

CDF of resampled values:-

$$P_r(x \leq a) = \sum_{n: x_n \leq a} w^{(n)}$$

$$= \frac{\sum_{n=1}^N \mathbb{I}(x_n \leq a) \hat{p}(x_n) / q(x_n)}{\sum_{n=1}^N \hat{p}(x_n) / q(x_n)}.$$

- take limit as  $N \rightarrow \infty$ , assuming regularity conditions:-

- replace summations with integrals weighted according to orig sampling distri  
 $q(x)$ .  $\downarrow \textcircled{?}$

$$P_r(x \leq a) = \frac{\int [\mathbb{I}(x \leq a) \hat{p}(x) / q(x)] q(x) dx}{\int [\hat{p}(x) / q(x)] q(x) dx}$$

6/59

$$= \frac{\int \mathbb{I}(x \leq a) \tilde{p}(x) dx}{\int \tilde{p}(x) dx} \quad \downarrow \textcircled{?}$$

$$= \int \mathbb{I}(x \leq a) p(x) dx \quad \rightarrow \text{CDF of } p(x).$$

$\rightarrow$  Normalisation of  $p(x)$   
not required.

- fix  $N$ , given initial sample set, resampled values only approx. drawn from desired distri  $p(x)$ .

-  $\hat{p}$  approx. improves (à la rejection sampling) as

$$q(x) \rightarrow p(x)$$

-  $q(x) = p(x)$ , initial samples  $(x^{(1)}, \dots, x^{(N)}) \sim p(x)$

$$\text{and weights } w^{(n)} = \frac{1}{N}.$$

- Resampled values then have desired distri.

## computation of moments:

use original samples and weights:-

$$\langle f(x) \rangle = \int f(x) p(x) dx$$

$$= \int f(x) [\hat{p}(x)/q(x)] q(x) dx$$

---

$$\int [\hat{p}(x)/q(x)] q(x) dx$$

$$\approx \sum_{n=1}^N f(x_n) w_n$$

?) review.  
intuitively  
not 'clicking'



## 1.1. Particle filtering

- use weighted resampling ideas to obtain SMC algorithm - particle filter
- particularly for online applications in which we make posterior inferences/updates in light of new observations.
- HMM + KF  $\rightarrow$  use discrete or near Gaussian distri.
- Particle filtering allows for more complex choices of observation/emission  $p(y_t|x_t)$ , and where posteriors are analytically intractable.
- recall time and measurement updates from Jordan (2003).

time update:  $p(x_t|y_0, \dots, y_t) \rightarrow p(x_{t+1}|y_0, \dots, y_t)$

measurement update:  $p(x_{t+1}|y_0, \dots, y_t) \rightarrow p(x_{t+1}|y_0, \dots, y_{t+1})$

given observed values  $y(t) = (y_1, \dots, y_t)$

draw  $M$  samples  $p(x_t|y(t))$  to eval. exp. of  $f(x)$  wrt posterior.

Bayes:  $\langle f(x_t) \rangle = \int f(x_t) p(x_t|y(t)) dx_t$

$$= \int f(x_t) p(x_t|y_t, y_{(t-1)}) dx_t$$

$$= \int f(x_t) \cdot \frac{p(y_t, x_t|y_{(t-1)})}{p(y_t|y_{(t-1)})} dx_t$$

$$\begin{aligned}
 &= \int f(x_t) \cdot \frac{p(y_t|x_t) p(x_t|y_{t-1})}{\int p(y_t|x_t) p(x_t|y_{t-1}) dx_t} dx_t \\
 &\quad \text{- constant + } x_t \text{ integrated out.} \\
 &= \frac{\int f(x_t) p(y_t|x_t) p(x_t|y_{t-1}) dx_t}{\int p(y_t|x_t) p(x_t|y_{t-1}) dx_t} \\
 &\approx \sum_{m=1}^M \left( \frac{p(y_t|x_t^{(m)})}{\sum_{m=1}^M p(y_t|x_t^{(m)})} \right) f(x_t^{(m)}) \quad \text{where } \{x_t^{(m)}\}_{m=1}^M \sim p(x_t|y_{t-1})
 \end{aligned}$$

• use of C.I.  $p(y_t|x_t, y_{t-1}) = p(y_t|x_t)$   $\circlearrowleft$  - use now this features (check prob rules)

• defining:  $w_t^{(m)} = \frac{p(y_t|x_t^{(m)})}{\sum_{m=1}^M p(y_t|x_t^{(m)})}$  yields

$$\Rightarrow \langle f(x_t) \rangle = \sum_{m=1}^M w_t^{(m)} f(x_t^{(m)}) \quad \text{set of samples}$$

$\circlearrowleft$  (\*) posterior distn  $p(x_t|y_t)$  is represented by  $\{x_t^{(m)}\}_{m=1}^M$  and weights  $\{w_t^{(m)}\}_{m=1}^M$

• weights:  $0 \leq w_t^{(m)} \leq 1$   $\sum_m w_t^{(m)} = 1$   $\circlearrowleft$  (\*) typo: -  $p(x_t|y_t)$  !

• suppose:

- we have a set of:-

$$\begin{cases}
 1) \text{Samples } \{x_t^{(m)}\}_{m=1}^M \sim p(x_t|y_{t-1}) \\
 2) \text{weights } \{w_t^{(m)}\}_{m=1}^M = \frac{p(y_t|x_t^{(m)})}{\sum_{m=1}^M p(y_t|x_t^{(m)})} 
 \end{cases}
 \quad \begin{array}{l}
 \text{obtained} \\
 \text{at time step } t, \\
 \text{represents post. distn} \\
 p(x_t|y_t)
 \end{array}$$

$\circlearrowleft$  use weighted resampling rep. of  $p(x_t|y_{t-1})$ .  
to iteratively infer  $p(x_t|y_{t-1})$ .

• then we observe  $y_{t+1}$  and wish to find weights and samples at time  $(t+1)$  i.e.  $p(x_{t+1}|y_{t+1})$

PTO.  $\rightarrow$

- note have used C.I properties:-  
 $p(x_{t+1}|x_t, y_{t+1}) = p(x_{t+1}|x_t)$  via d-sep.  
 $p(y_{t+1}|x_t, y_{t+1}) = p(y_{t+1}|x_t)$
  - similar to mixture distribution  
    - samples drawn by selecting component  $m$  with prob given by mixing coefficients  $w^{(m)}$  and drawing sample from corresp. component.
  - (W) Jordan (2003) is very brief here! + measured update nuclear  
    - Also using lecture slides 10-708 L13. to supplement.
    - . Summary: - At time step  $t$ , have sample-rep of post  $p(x_t|y_{t+1})$  contained in  $\{w_t^{(m)}\}_{m=1}^M$  and  $\{x_t^{(m)}\}_{m=1}^M$ . (mixture rep.)
      - for next step, draw  $M$  samples from mixture distri. and for each sample use  $y_{t+1}$  to evaluate new weights  $w_{t+1}^{(m)} \propto p(y_{t+1}|x_t^{(m)})$  ?
  - (\*) develop measurement update as nuclear!
  - (W) made a decision to re-present particle filtering, with a view to providing full clarity (effectively merging Jordan (2003) and 10-708 L13).
  - sequential weighted resampling for particle filtering consists of:-
  - 1) starting point - sample based representation of the posterior  $p(x_t|y_{t+1})$  using weights  $\{w_t^{(m)}\}_{m=1}^M$  and samples  $\{x_t^{(m)}\}_{m=1}^M$
  - 2) time-update - computation of  $p(x_{t+1}|y_{t+1})$  using above sampling representation
  - 3) measurement - computation of  $p(x_{t+1}|y_{t+1})$  using  time-update and also newly evaluated weights.
- ← Jordan covers this but this is very vague
- GOAL: Sample from an intractable posterior.

first sample from  $p(x_{t+1}|y_{(t)})$  via Bayes:

time update

$$p(x_{t+1}|y_{(t)}) = \int p(x_{t+1}|x_t, y_{(t)}) p(x_t|y_{(t)}) dx_t$$

$$= \int p(x_{t+1}|x_t) p(x_t|y_{(t)}) dx_t$$

$$= \int p(x_{t+1}|x_t) \underbrace{p(x_t|y_t, y_{(t-1)})}_{\downarrow} dx_t$$

$$= \int p(x_{t+1}|x_t) \cdot \frac{p(y_t, x_t|y_{(t-1)})}{p(y_t|y_{(t-1)})} dx_t$$

$$= \int p(x_{t+1}|x_t) \cdot \frac{p(y_t|x_t)p(x_t|y_{(t-1)})}{\int p(y_t|x_t)p(x_t|y_{(t-1)}) dx_t} dx_t$$

no  $x_t$  terms

$$= \frac{\int p(x_{t+1}|x_t) p(y_t|x_t) p(x_t|y_{(t-1)}) dx_t}{\int p(y_t|x_t) p(x_t|y_{(t-1)}) dx_t}$$

using similar  
principles  
to earlier.

- more

② like mixture  
model!

$$\approx \sum_{m=1}^M \left( \frac{p(y_t|x_t^{(m)})}{\sum_{n=1}^M p(y_t|x_t^{(n)})} \right) p(x_{t+1}|x_t^{(m)})$$

$$= \sum_{m=1}^M w_t^{(m)} p(x_{t+1}|x_t^{(m)})$$

that is; we use representation of posterior  $p(x_t|y_{(t)})$  given by weights and samples:-

1) samples  $\{x_t^{(m)}\}_{m=1}^M \sim p(x_t|y_{(t-1)})$

2) weights  $\{w_t^{(m)}\}_{m=1}^M$

## 1. Starting point

- At time step  $t$ , we denote a sampling-based representation of the posterior  $p(x_t | y_{t-1})$ ; and assume this is available to us.
- Recall this is because  $p(x_t | y_{t-1})$  is analytically intractable, due to it being likely that there is no conjugacy we can exploit (e.g. in K.F.).
- The posterior  $p(x_t | y_{t-1})$  has following expression:-

$$p(x_t | y_{t-1}) = p(x_t | y_t, y_{t-1}) = \frac{p(x_t | y_{t-1}) p(y_t | x_t)}{\int p(x_t | y_{t-1}) p(y_t | x_t) dx_t} \quad \left( = \frac{p(y_t, x_t | y_{t-1})}{p(y_t | y_{t-1})} \right)$$

- We represent  $p(x_t | y_{t-1})$  using:-

$$\left\{ \begin{array}{l} x_t^{(m)} \sim p(x_t | y_{t-1}), \quad w_t^{(m)} = \frac{p(y_t | x_t^{(m)})}{\sum_{m=1}^M p(y_t | x_t^{(m)})} \\ x_t^{(m)} \text{ drawn from } ? \end{array} \right\}_{m=1}^M$$

weights  
computed  
using observation  
distn  $p(y_t | x_t^{(m)})$

## 2. Time update

- Note that:-

$$p(x_{t+1} | y_{t-1}) = \int p(x_{t+1} | x_t) p(x_t | y_{t-1}) dx_t$$

- which is represented by:-

$$p(x_{t+1} | y_{t-1}) \approx \sum_{m=1}^M w_t^{(m)} p(x_{t+1} | x_t^{(m)})$$

(akin to mixture  
model)

- We can draw samples of  $x_{t+1}$  from this, using  $x_t^{(m)}$  and  $w_t^{(m)}$ !

→ see earlier  
for details  
on how C.1. assumptions  
are used.  
+ derivations

## 3. Measurement update

(starting point)

- Note that measurement update is the posterior  $p(x_t | y_t)$  propagated forward by one time-step and using new evidence  $y_{t+1}$  i.e:-

$$p(x_{t+1} | y_{t+1}) = p(x_{t+1} | y_{t+1}, y_t) = \frac{p(x_{t+1} | y_t) p(y_{t+1} | x_{t+1})}{\int p(x_{t+1} | y_t) p(y_{t+1} | x_{t+1}) dx_{t+1}}$$

To get the weighted, resampled version of  $p(x_{t+1}|y_{t+1})$ , we can represent it using:-

$$\left\{ \begin{array}{l} x_{t+1}^{(m)} \sim p(x_{t+1}|y_t) \\ w_{t+1}^{(m)} = \frac{p(y_{t+1}|x_{t+1}^{(m)})}{\sum_{m=1}^M p(y_{t+1}|x_{t+1}^{(m)})} \end{array} \right\}_{m=1}^M$$

- We draw samples  $x_{t+1}^{(m)}$  using the mixture model sample representation of the time update  $p(x_{t+1}|y_t)$ .
- Then compute new weights  $w_{t+1}^{(m)}$  using  $x_{t+1}^{(m)}$ , new evidence  $y_{t+1}$ , and our observation distri  $p(y|x)$ . (plug values in)
- (\*) This includes particle filtering  $\rightarrow$  now see diagram.

Corrected diagram:

