

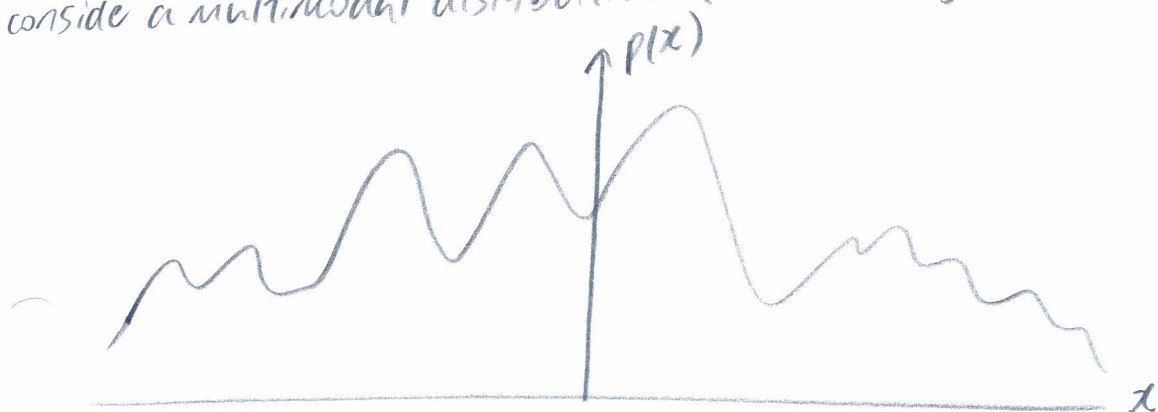
U3 - Approx Inference: Monte Carlo and sequential Monte Carlo

(*) How to represent a joint/marg.

$$\mathbb{E}_P[f(x)] = \int f(x)p(x)dx$$

• sample based representation:

• consider a multimodal distribution (and density): -



(?) How to integrate?

- It may not be feasible (analytically)? to integrate this.

- Sol: Draw samples from distri

- computationally: will use a uniform PRNG to draw a sample $u(0,1)$ and transform to get Gaussian.

(*) Draw N samples $x^{(i)} \sim P$:-

$i=1, \dots, N$

$$\text{compute } \mathbb{E}_P[f(x)] = \frac{\sum_{i=1}^N f(x^{(i)})}{N}$$

(*) original distri P expressed as a set of samples

} Monte Carlo methods

(*) Monte Carlo methods

(*) replace integration with summation (sample-based averages)

(*) Asymptotically exact \rightarrow more samples drawn ($N \rightarrow \infty$) arbitrarily close to expectation we wish to approx.

Challenges:

- 1) How to draw samples from distri?
- How to do this efficiently?
from an arbitrary
- 2) How to make better use of samples?
- asymptote - how many samples can we draw ...

(*) Naive sampling

$$P(x_A, \dots, x_M) = P(x_B)P(x_A|x_B) \dots$$

$$\underline{X}^{(k)} = \begin{bmatrix} x_A \\ \vdots \\ x_M \end{bmatrix} \sim P$$

- How to draw samples?

- traverse graph

(*) EX walks through realisation of sample in BN
- sample i.v.s. from B.N.

(*) Rows of table correspond to a 'sample' from B.N.

(*) Computation of queries: check samples where condition is met.

(*) Rare events \rightarrow conditioning on rare events?

(*) Shows the limitations of sampling

(*) Monte Carlo methods (rejection sampling)

- attempt to get around above deficiencies

$$\pi(x) = \frac{\pi'(x)}{Z}$$

$\pi(x)$ - difficult to sample from

$\pi'(x)$ - easy to evaluate
(e.g. unnormalised clique pot. product)

- sample from simpler distri $Q(x)$

⑧: $x^* \sim Q(x)$ (draw samples from Q)
 - Accept x^* with probability $\frac{\pi'(x^*)}{K Q(x^*)}$ K -arbitrary constant

(*) Ratio of unnormalised part of given (evaluable dist.) and a constant multiplied by proposal distribution.

(*) with these two steps:-

$$x^* \sim \pi(x)$$

ex: How to prove this?

$$P(x^*) = \frac{Q(x^*) \frac{\pi'(x^*)}{K Q(x^*)}}{\int Q(x^*) \frac{\pi'(x^*)}{K Q(x^*)} dx^*} = \frac{K \pi'(x^*)}{K \int \pi'(x^*) dx^*} = \pi(x)$$

(*) A solution to sampling from a distri $\pi(x)$ which is difficult to normalise (using a proposal distrib)

Pitfall: ①: Geometric intuition / explanation of too many rejections ②

(*) Rejection sampling

- Conceptual experiment to illustrate rejection sampling
 - Small differences between proposal distri and distri of interest caused by scaling constant can cause a large difference in rejection rate

$$K = \left(\frac{\sigma_q}{\sigma_p} \right)^d \approx \frac{1}{20,000}$$

(*) The actual vol. of 2 dists are vastly different in high-dim space \Rightarrow reject samples most of time even though dists are close

(*) Solution: Don't use one Q -distri (proposal distri) to cover the target distri, but use a piece-wise envelope

- (*) Adaptive rejection sampling is an example of this
- So rejection sampling doesn't work well in a high-dim space

(*) Unnormalised importance sampling

- Draw, as before, samples $x^{(i)} \sim Q(x)$ from a proposal distri. N times
- Instead of accepting/rejecting as earlier, compute weight based on new sample

$$w^{(i)} = \frac{p(x^{(i)})}{Q(x^{(i)})}$$

- Store $\{x^{(i)}, w^{(i)}\}_{i=1}^N$ as representation of true distri.

samples

(*) Why is P suddenly available/tractable?

- Assume in simplest case that P can be eval (up to norm. constant)

(*) Assess proof, picture.

- Proposal distri introduced as a dummy in proof
- (*) Unnormalised as weights directly comp. from likelihood ratio (do not sum to one).

Q: What if I can only evaluate true distri up to a constant? (2)

e.g. no exact value of norm. constant
(define to be arbitrary α)

(*) Normalised importance sampling

- only can evaluate $p'(x) = \alpha P(x)$ (e.g. MRF)
- Get around normalisation constant:-

(*) take expectation of the ratio

$$r(x) = \frac{p'(x)}{Q(x)} \Rightarrow \langle r(x) \rangle_Q = \int \frac{p'(x)}{Q(x)} Q(x) dx = \int p'(x) dx = \alpha$$

- expectation of ratio is normalisation constant

$$\alpha = \frac{r^{(1)}}{\sum_{i=1}^N r^{(i)}}$$

- Now draw $x^{(i)} \sim Q(x)$ $i=1, \dots, N$ (N samples)

- Compute $r^{(i)}$ for each i^{th} sample.

- Compute α via $\alpha = \sum_{i=1}^N r^{(i)}$

(A3) - confused here
- review derivations

Proof:

(*) Compute expectations of function using samples drawn as above

(*) Instructor presentation lost me. (*) don't need to know norm constant or part function of target

(*) normalised vs unnormalised importance sampling

(A4) - review at home (instructor does not emphasise).

(*) Efficiency of likelihood weighting

- So far we have used likelihood weighting \rightarrow may not get 'true' answer in peculiar/path scenarios

- Consider practical implementation of MCMC. When do we stop sampling?
Typically stop when stability/convergence observed as empir. strategy for controlling iteration/alg duration.

- In a technically equivalent way; can view this as terminating when variance \downarrow of current solution.

- makes sense. (i.e. emergent behaviour through variance)

(*) So variance in a sense can be viewed as measure of closeness to the answer.
(of current algo output)

(*) In importance sampling; we can have a pathological case: -

(A5): Review; watching + record intuition

(*) Example of poor proposal distri (note we do not know it's poor);
does not envelop true distri.

(*) Hoping weights $w^{(i)}$ will correct this.

Solution: use heavy tailed $Q \rightarrow$ (but many samples wasted)

- weighted resampling from $\{x^{(i)}, w^{(i)}\}_{i=1}^N$

- we resample from $\{w^{(i)}\}_{i=1}^N$ N' times where $N' \gg N$.

(*) 'Amplify' distn using resampled weighted resampling (A6: review weighted resampling)

(*) weighted resampling

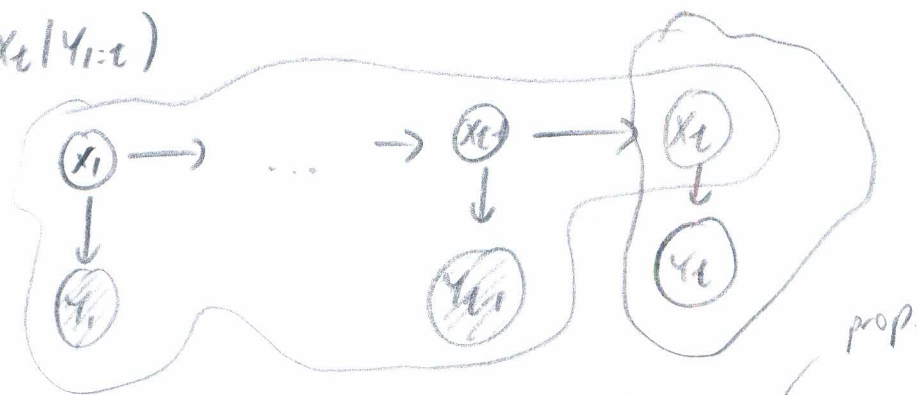
(*) particle filtering

- use resampling idea to elegant, efficient inference.
- make a fast sampling based algorithm for SS/MS.
- KF demanding to implement in high-dim space; or non-Gaussian (Taylor approx)

(*) Sketch of particle filtering

- use KF to establish recursive procedure.

$P(x_t | Y_{1:t})$



- Break down $P(x_t | Y_{1:t}) \rightarrow 2$ parts.

$$P(x_t | Y_{1:t}) = P(x_t | y_t, Y_{1:t-1}) = \frac{P(x_t | Y_{1:t-1}) P(y_t | x_t)}{\int P(x_t | Y_{1:t-1}) P(y_t | x_t) dx_t}$$

'weights'

- inspect $P(x_t | Y_{1:t-1})$ and use this as a proposal distrib. from which we sample

- represent $P(x_t | Y_{1:t})$ using above decomp: -

$$\left\{ \begin{aligned} \{x_t^{(i)}\} &\sim P(x_t | Y_{1:t-1}) ; w_t^{(i)} = \frac{P(y_t | x_t^{(i)})}{\sum_{i=1}^N P(y_t | x_t^{(i)})} \end{aligned} \right\}$$

(Hence $P(x_t | Y_{1:t})$ is expressed as weighted resamples from proposal distri.

(*) Goal: - use weighted resample representation of $P(x_t | Y_{1:t})$ to iteratively infer $P(x_{t+1} | Y_{1:t+1})$

(*) sequential weighted resampler

time update: keep evidence; propagate forward with no new measurement

$$p(x_{t+1} | Y_{1:t}) = \int \underbrace{p(x_{t+1} | x_t)}_{(i)} \underbrace{p(x_t | Y_{1:t})}_{(ii)} dx_t$$

→ given by weighted res. rep.

• i) given via model

(*) replace integration with sampling:-

$$p(x_{t+1} | Y_{1:t}) = \sum_{i=1}^N w_t^{(i)} p(x_{t+1} | x_t^{(i)})$$

- weighted sum of transition prob. cond. on samples of previous states.

(*) similar to mixture model (sampling from).

(*) use new evidence:-

measurement update:- $p(x_{t+1} | Y_{1:t+1}) = \frac{\overbrace{p(x_{t+1} | Y_{1:t})}^{(*)} \overbrace{p(Y_{t+1} | x_{t+1})}^{\text{new prop.}}}{\underbrace{\int p(x_{t+1} | Y_{1:t}) p(Y_{t+1} | x_{t+1}) dx_{t+1}}_{\text{reweights}}}$

(*) note (*) is computed from time update

- replace by weighted resampling version:-

⇒ $p(x_{t+1} | Y_{1:t+1})$ has representation:-

$$\left\{ x_{t+1}^{(i)} \sim p(x_{t+1} | Y_{1:t}), w_{t+1}^{(i)} = \frac{p(Y_{t+1} | x_{t+1}^{(i)})}{\sum_{i=1}^N p(Y_{t+1} | x_{t+1}^{(i)})} \right\} \text{reweighted}$$

(*) particle filtering graph (illustrating resampling)

- size of balls correspond to weights

(*) - review

(*) Particle filtering
gives online results for predictive distri

(*) PF for SSMs (switching)

- Allows drawing from more complex distri e.g. switching SSM
- multiple hidden distri.

(*) Rao-Blackwellised sampling

(*) Review slides (not emphasized)

(*) Statistical properties of unnormalized/normalised importance sampling.

- unnormalised importance sampling is unbiased

$$\text{i.e. } \mathbb{E}_Q[f(x)w(x)] = \mathbb{E}_P[f(x)]$$

- normalised importance sampling is biased for finite samples e.g. $M=1$

$$\mathbb{E}_Q\left[\frac{f(x')r(x')}{\sum r(x')}\right] \neq \mathbb{E}_P[f(x)]$$