

INT-APPROX. INF: MC/MC

- recap of Monte Carlo
- generate samples from p.d. $\pi^{(t)}$
- estimate expectations of functions under a p.d. $E[f(x)]$ under $\pi(x)$.

Ex. How does this relate to inference covered in course?

e.g. $p(x|e)$, $\log(x)$, \mathbb{Z}

marginal likelihood partition
query given evidence

- $p(x|e) = \mathbb{P}[\mathcal{I}(x)|e]$
- p. $\log(x), \mathbb{Z}$; recall that the normalisation constants correspond to weights of samples. (①) - using sampling can address inf. tasks

(1) Limitations of Monte Carlo

- related to quality of proposal distri Q .
- tried approach to boost acceptance rates of samples gen by Q .
- cause they (samples) are mostly accepted; they visit different regions of distributions which are multi-modal etc.
- focus on augmenting Q via more aggressive Q ; which updates adaptive.

(2) Markov Chain Monte Carlo

- An example of importance sampling with a bad proposal $Q(x)$
- However the samples are manipulated; it will be hard to 'cover' the distri $\pi(x)$ using $Q(x)$.
- A long tail
- use adaptive Q
- instead of $Q(x')$; use $Q(x'|x)$ where x' - new state being sampled
 x - previous sample.

② AS x changes,

$Q(x'|x)$ can also change (as a function of x')

- see diagram

- draw sample $x^{(1)}$; generate $Q(x^{(2)}|x^{(1)})$ i.e. proposal condit. on $x^{(1)}$
 - (centred on $x^{(1)}$)
 - draw sample $x^{(2)}$; generate $Q(x^{(3)}|x^{(2)})$
 - probability mass shifts to the right progressively; allowing us to better capture the target distri.
- Ex: will study how to operationalise this high-level idea; ensure correctness to make sure it is 'correct'. i.e. after many samples, they are coming from true target distri P .

(ii) Metropolis-Hastings ③ (transitional kernel)

- draw sample x' from $Q(x'|x)$, where x is the previous sample.
- The new sample x' is accepted or rejected with some prob. $A(x'|x)$.
- Acceptance prob.
$$A(x'|x) = \min\left(1, \frac{p(x') Q(x|x')}{p(x) Q(x'|x)}\right)$$
- $A(x'|x)$ ensures that after sufficiently many draws; samples will come from true distri $p(x)$.
- $A(x'|x)$ - ratio of importance weights
 - $\frac{p(x')}{Q(x'|x)}$ - importance weight for x'
 - $\frac{p(x)}{Q(x|x')}$ - importance weight for x
- This ratio can be smaller or greater than 1

(*) If ratio > 1, accept the sample with prob 1.

(*) MH Algorithm

(*) Review pseudocode

- Simple
- If need to remember acceptance rate for new samples
- Illustration of how a proposal distinguishes both modes of multimodal dist.
- (*) In practice; we don't know P
- (*) When we reject a sample; take older example and put in Neg.
re-use proposal

(*) Review this example

- (*)^(*) not very likely under proposal, but we are 'lucky' to sample it.
- However, it gets 'boost' due to high prob under target distri P .
- (*) Therefore it has a better importance weight and is accepted.
- (*) Using this dynamic; we amplify the tend of moving to 'poor modes'
- (*) Under target distri P ; even though the proposal doesn't give a good chance of accepting that point.

(*) (*) Adaptive proposal allows us to capitalise and get a good sample.

- (*) Previous examples: fixed proposal + some kind of augmentation
- (*) Can still sample from bimodal target with unimodal proposal; as the latter 'waves around'!

(*) Theoretical aspects of MCMC

- (*) Burn-in: first 1,000/10,000 samples have to be thrown away.

- Nature of 'burn-in' is a little mysterious (art)

- (*) Why are MH samples guaranteed to be from $f(x)$ - prove

(*) MC and MCMC :-

- Adaptively change Q based on previous example.
- In principle depends on all previous examples; but we assume a 1st order Markov property.

(*) Markov chain:-

- sequence of r.v.s. $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ with M.P.:-
 $P(x^{(n)} = x | x^{(1)}, \dots, x^{(n-1)}) = P(x^{(n)} = x | x^{(n-1)})$

⑥) $P(x^{(n)} = x | x^{(n-1)})$ - transition kernel

- x can be a scalar or vector

(*) Assume that next state only depends on previous state (?)

(*) Homogeneous-time invariant transition matrix.

Hence $T(x'|x)$ is not indexed by time.

(*) MC concepts

Ex: what is effect of conditioning on previous sample ③-review

⑥) $\pi^{(n)}(x') = \sum \pi^{(n)}(x) T(x'|x)$

MCMC: Hope for a protocol generating samples from a target distri.

\Rightarrow Stationary: $\pi(x') = \sum_x \pi(x) T(x'|x) \quad \forall x'$

- distribution after transitioning (using kernel) leads to a distri that is identical to 'starting point'.

- MCMC generates constants on T .

- not every T allows for this stationarity

Ex: when do we have a stationary distri?

- characterise properties (should be familiar)

(*) Irreducibility

④ review intuition on these properties.

(*) Aperiodicity

(*) Ergodicity (irreducible + aperiodic) ⑥

\hookrightarrow Allows stationary distri.

Ex: what transition kernel T satisfies these conditions?

(*) reversible MC

- An MCMC is reversible if there exists a distri $\pi(x)$ such that the detailed balance equation is satisfied

$$\pi(x') T(x|x') = \pi(x) T(x'|x)$$

- probability $x' \rightarrow x$ is same as $x \rightarrow x'$.

(*) reversible MCMCs always have a stationary distri

- simple proof: start with detailed balance; derive stationary.

(*) Why does Metropolis-Hastings algorithm work?

- MH satisfies detailed balance

(*) review the proof (from setting $Q = TA$)

(*) - MH algorithm defines a transition kernel that satisfies detailed bal.

(*) - MH algorithm defines a transition kernel that satisfies detailed bal.

- because of the way we defined our proposal Q and acceptance A ;

MH allows target distri to be the stationary distri

- it has theoretical guarantees

drawbacks

- No guarantees on when convergence occurs; only that it does occur.

- knowing when to start burn-in \rightarrow inform

- rely on heuristics

(*) Gibbs Sampling

- special case of MCMC.

$$x^{(t=1)} = \left(\begin{array}{c} x_1^{(t=1)} \\ x_2^{(t=1)} \\ \vdots \\ x_n^{(t=1)} \end{array} \right) \quad x^{(t=2)} = \left(\begin{array}{c} x_1^{(t=2)} \\ x_2^{(t=2)} \\ \vdots \\ x_n^{(t=2)} \end{array} \right)$$

- simplification of MCMC.

- $Q(x^{(t=2)} | x^{(t=1)})$

- note for $x^{(t=2)}$; we use $Q(\cdot)$ in

such a way that $x_i^{(t=2)}$ is a new value; $x_{i+1}^{(t=2)} = x_{i+1}^{(t=1)} \forall i$

$$\cdot x_1^{(t+1)} \sim p(x_1 | x_2^{(t+1)}, x_3^{(t+1)}, \dots, x_n^{(t+1)})$$

This $p(\cdot)$ is the proposal distri $Q(x^{(2)} | x^{(1)})$

- (*) Sample next example such that all but one dimension/component is same as previous example.
- (*) Next different dimension will be sampled conditional on all other dimensions that have already been sampled.

⑥ Note we can further reduce $p(x_1 | x_2^{(t+1)}, x_3^{(t+1)}, \dots, x_n^{(t+1)})$ in a graphical model.

By using $p(x, MB(x_1))$ where MB is Markov blanket

- ⑦ A specialised way to draw the next sample until you only change one dimension and keep the others; the condition of those.
- Ex: interesting property; this proposal $Q(\cdot)$ leads to an acceptance of 1
- (*) Gibbs sampling \rightarrow ⑥ - review
- (*) distinction between epoch \rightarrow Is it one dim or the whole vector
that is refreshed to what as one 'epoch'?

- varies in literature

(*) Gibbs sampling - example

- previous alarm network

↳ Not many effective samples

- $t=0$ - Random guess/initial start. point

⑦ Review this (mcclurin)

- (*) Exploit C.I. prop of MB ; rather than high-dim proposal of all states
- ⑦ As acceptance rate is 1, we use all the examples?

(*) Topic Models:- Collapsed Gibbs

(*) Gibbs sampling became well-known in ML due to mentioned paper

(Griffiths + Seager)

- DA: used MFA in VI.

- NL paper: integrate out variables; only sample word-topic assignments z

(*) (a) - review plate; algo

- algo: for all variables $\underline{z} = z_1, z_2, \dots, z_n$

- draw $z_i^{(t+1)}$ from $P(z_i | \underline{z}_{-i}, w)$

- well $\underline{z}_i = z_1^{(t+1)}, z_2^{(t+1)}, \dots, z_{i-1}^{(t+1)}, z_{i+1}^{(t+1)}, \dots, z_n^{(t+1)}$

$P(z_i | \underline{z}_{-i}, w)$:

$$P(z_i=j | \underline{z}_{-i}, w) \propto \frac{n_{-i,j}^{(w)} + \beta}{n_{-i,j}^{(t)} + w\beta} \quad \frac{n_{-i,j}^{(d)} + \alpha}{n_{-i}^{(d)} + 1\alpha}$$

'word-topic' 'doc-topic'

- two types of weighted counts (product).

- word-topic \rightarrow no. of words appearing
 within a particular topic over all words
 topic-specific word frequency

- doc-topic \rightarrow no. of topics appearing within a particular
 document over all topics
 document specific topic freq.

(a) - detailed in slides; how nos. can be evaluated from zs.

(*) collapsed gibbs illustration

- shows indexed words of documents

- z_i refers to z_i of document j (z_{ij}).

(*) Heation 1: Given from initialisation

- z_1 ? - 1st word of 1st document \rightarrow topic?

- use proposal dist. of one single word given all others

$$P(z_i=j | \underline{z}_{-i}, w) \propto$$

- (*) Draw $z_{1,1}^{(2)}$ (can draw conditioning on $z_{2,1}^{(1)}, z_{3,1}^{(1)}, z_{4,1}^{(1)} \dots z_{50,5}^{(1)}$) - (*) incorrect indexing here (A10)
- (*) draw $z_{2,1}^{(2)}$ conditioning on $\{z_{3,1}^{(1)}, z_{4,1}^{(1)}, \dots, z_{50,5}^{(1)}\}$ - but only idea

(*) large-scale topic models implement sophisticated indexing programs
(to store and retrieve efficiently)

Ex: extremely simple, liked by practitioners

- allows for additional variables in TM design
- can index word in different ways e.g. introduce sentiment index s

$z_{i,j,s}$ where s indexes sentiment.

(*) Gibbs sampling as special case of MH.

Gibbs proposal: $Q(x_i, x_{-i} | x_i, x_{-i}) = P(x'_i | x_i)$

x_{-i} : all variables except x_i

(A11) Review maths here

Ex: now have 2 algorithms with proof of correctness

(*) Practical aspects of MCMC

- proposal distn quality $Q(x' | x)$ \rightarrow acceptance rate, autocorr
- when to stop burn-in? \rightarrow sample values vs time
log-likel. -||-

(*) Acceptance Rate

(*) issue is you don't know P ; you have to rely on samples blindly to stop from $Q(\cdot)$

(*) slides illustrate some shortcomings & tradeoff

(*) There are tradeoffs between exploration (variance) and acceptance.

(*) Some guidelines are proposed

- see papers for provenance Muller (1993), Roberts, Gelman (1994)

(*) Autocorrelation function

(*) MCMC samples are always correlated (you are using previous examples to generate transition for new).

- autocorrelation function:- $R_x(k) = \frac{\sum_{t=1}^{n-k} (x_t - \bar{x})(x_{t+k} - \bar{x})}{\sum_{t=1}^{n-k} (x_t - \bar{x})^2}$

②

(*) When samples are correlated; there exist formulae to compute effective sample size

(*) E.g. for 2 samples where $\pi^{(1)}$ is twice the magnitude of the 2^{nd} (does this count as 2 'samples'?) - really 1.

- no info other than linear scaling (\approx).

Ex: Compute autocorrelation function to encourage good proposal.
- want low autocorrelation

(*) Monitoring convergence

- track interesting statistics

- e.g. if you want Gaussian parameter estimates
- monitor.

(*) note well-mixed and poorly-mixed chains

(*) In practice; use multiple chains due to poor random guess^{initial}s
(rather like multiple initialisations of EM)

(*) Issue:- high dim LMS (difficult to visualise all r.v. classes at once)

(*) log-likelihood vs Time.

- see profile of log-likelihood over time
- these are main MCMC diagnostic tools

(*) that covers:-

- i) Metropolis-Hastings algo
- ii) Gibbs Sampling

(**) how to choose the random variable of interest to be sampled conditioned on rest (ordering)

- ex: no formal ordering / theory behind it.
- guided by implementation

(#) Supplementary

(#) key struggle of MCMC \rightarrow random walk behavior
(the spirit of Wiener!)

- random walk in MCMC

- symmetric, limited bandwidth proposal $Q(x_{\text{new}} | x_{\text{old}})$

Today? - large bandwidth Gaussian proposal with wide variance
- low acceptance rate. not

Ex: given access to the ground truth target distri $P(\cdot)$ which is known in practice; do we have a preference for the direction?

- we want to explore modalities of target distri
(goes against R.W. behavior)

- All standard vanilla MCMC are r.w. behaviors based on proposal distri

- can we steer the behavior towards mode?

(#) gradients/derivatives is only feature that gives you info on direct. ✓

(#) Hamiltonian Monte Carlo

- Hamiltonian \rightarrow sum of potential, kinetic energy

Ex: interesting corrections (not going to review physics).

\rightarrow see how this yields new sampling method

- introduce more r.v.s into auxiliary distribution

- introduce kinetic energy K

$$(*) \textcircled{D}: \frac{dx_i}{dt} = \frac{\partial H}{\partial p_i} \quad \frac{dp_i}{dt} = -\frac{\partial H}{\partial x_i} \quad (\text{alt. notation})$$

- Radford Neal, Mackay

ex: point is to define change of state as gradient of loss function

(*) average gradient info into transitions.

- position and momentum updates $p(\cdot), q(\cdot)$.
- (*) multiple mechanisms for updates.
 - leapfrog method preferred to avoid over/undershooting.

(B) MCMC from Hamiltonian Dynamics \textcircled{D}

- q, p re variable of interest

Define: $P(q, p) = \frac{1}{Z} \exp(-U(q)/T) \exp(-K(p)/T)$

where $U(q) = -\log[\pi(q)L(q|D)]$ (Bayesian setup)
 $K(p) = \sum_{i=1}^d \frac{p_i^2}{2m_i}$ (kinetic term auxiliary f.v.)

(*) use Hamiltonian dynamics to propose next step

- given q_0 (starting state)

- draw $p_0 \sim N(0, I)$

- use L steps to propose next state

- accept/reject based on change in Hamiltonian

(A13) - review
accepted
reject/reject
based on Ham

$$\textcircled{D}: \min[1, \exp(-U(q^*) + U(q) - K(p^*) + K(p))]$$

(*) broaden generalisation -

combine optimisation + MCMC
gradients

\ sampling

(*) 2D Gaussian

- Hamiltonian MCMC vs R.W. Metropolis

(*) Better coverage of target distri

(*) Better mixing of samples

R.W. - high correlation does not allow us to move away from initial prop.

(*) see empirical results on 2D/100D Gaussian comparing Hamil MCMC, RW.

(*) Langevin dynamics

- Analogous to Gibbs being special case of M-H

- Langevin MCMC is special case of Hamiltonian MCMC

- use more sophisticated updates than leapfrog using 2nd order updates

- Eg. constraints on regions you don't want to enter

- can fully augment optim. technique to make proposal more efficient.

ex: to improve acceptance rate, better mixing ; incorporate optimisation techniques to improve proposals

(*) variational inference

- can be incorporated with MCMC.

- optimisation based method yielding approx of distri based on evidence.

- Render proposal distri a variational distri ; subject it to optimisation

- use M-H for acceptance/rejection of samples.

(*) This is all classical material on PgMS

- representation in OEM view

- Exact inf & learning (with, without latent)

- Approx Inference

(*) next: use PGMS as a foundation for thinking about modern form
of deep learning

