

## Lecture 21 - scribbles

Tuesday, November 21, 2017

14:17

today: • finish Gibbs sampling  
• variational methods

(ctd for Gibbs sampling)

GS is MH with time varying proposal

suppose we pick  $i$  at time  $t$

then proposal is  $q_t(x' | x^{(t-1)}) = p(x'_i | x_{-i}^{(t-1)}) \underbrace{\delta(x'_{-i}, x_{-i}^{(t-1)})}_{\text{force rest to be constant}}$

acceptance ratio:

$$a(x' | x^{(t-1)}) = \frac{\overbrace{q_t(x^{(t-1)} | x')}^{\cancel{p(x'_i | x'_{-i})}} \overbrace{p(x')}^{\cancel{p(x'_{-i} | x_{-i}^{(t-1)})}}}{\overbrace{q_t(x' | x^{(t-1)})}^{\cancel{p(x'_i | x'_{-i})}} \overbrace{p(x^{(t-1)})}^{\cancel{p(x_{-i}^{(t-1)} | x_{-i}^{(t-1)})}}} \rightarrow \frac{\cancel{p(x'_i | x'_{-i})} p(x')}{\cancel{p(x'_i | x'_{-i})} \cancel{p(x_{-i}^{(t-1)} | x_{-i}^{(t-1)})}} \rightarrow \frac{p(x'_i) p(x_{-i}^{(t-1)})}{p(x_{-i}^{(t-1)}) p(x'_i)}$$

$= 1$   $\nearrow$  always accepts?

convergence of G-S:

Let  $A$  be <sup>Markov</sup> transition kernel of one full cycle of G-S. (i.e.  $n$  steps)

$\rightarrow$  homogeneous M.C.

$A$  is irreducible and aperiodic because  $A_{ii} > 0 \forall i$

(since can get to any state with  $n$  steps)

(supposing  $p(x_i | x_{-i}) > 0$   
 $\forall x_i, x_{-i}$ )

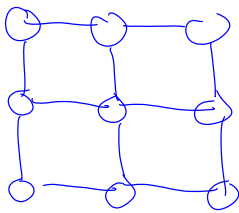
$$\Rightarrow A^t \pi_0 \xrightarrow{t \rightarrow \infty} p$$

⊗ also works for random scan (pick  $i \sim \text{Unif}(1:n)$  at each step)

example: G-S. for Ising model:

Ising model  $x_i \in \{0, 1\}$

UGM:



$$p(x) = \frac{1}{Z(\eta)} \exp\left(\sum_i \eta_i x_i + \sum_{i,j \in E} \eta_{ij} x_i x_j\right)$$

for Gibbs sampling,

want to compute  $p(x_i | x_{-i}) \propto p(x_i, \overbrace{x_{-i}}^{\text{fixed}})$

$$= \exp\left(\eta_i x_i + \sum_{j \in N(i)} \eta_{ij} x_i x_j + \text{rest}\right)$$

$\Rightarrow$  renormalize to get conditional:

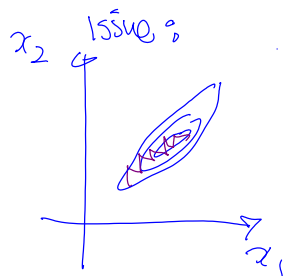
$$p(x_i=1 | x_{-i}) = \frac{\exp(\eta_i + \sum_{j \in N(i)} \eta_{ij} x_j)}{1 + \exp(\dots)}$$

$$= \sigma\left(\eta_i + \sum_{j \in N(i)} \eta_{ij} x_j\right)$$

↑  
sigmoid  $\sigma(z) \triangleq \frac{1}{1 + \exp(-z)}$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Variants:



a) block Gibbs sampling

$$p(z_A | z_{-A})$$

↑  
set (block)

$$U A_i = V$$

b) Fast-Blockwelized Gibbs sampling: marginalize out some variables first

eg. say can compute  $p(z, y)$  from  $p(x, y, z)$

(used in LDA (latent Dirichlet allocation))

overall

⊗ need to be able to sample from  $p(x_i | x_{-i})$  to run G.-S.

(otherwise, need M.H.)

Gibbs is easy when

discrete graph. model

cts. DGM with "conjugate distributions"

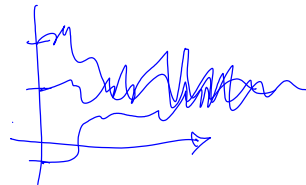
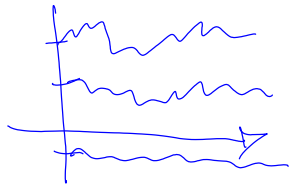
\* <sup>good</sup> proposal is an art.

look at software: STAN (C++) (Andrew Gelman)

(BUGS)

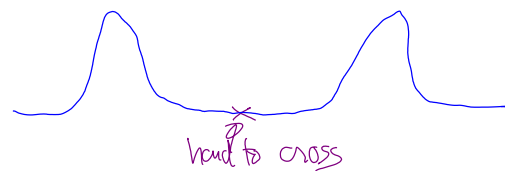
## diagnostic of mixing:

monitor mixing by running independent chains



"sticky chain"  
→ slow mixing

usually slow mixing comes to difficulty to move between modes



→ annealing methods help this

proposal looks like  $\frac{1}{Z} \exp(-\frac{E(x)}{T_k})$

example: "Annealed importance sampling"

physics analogy

energy

$\downarrow$   
 $\frac{1}{T_k}$

temperature

→ high temperature  
⇒ more exploration

## Variational methods

general idea: say we want to approximate  $\theta^*$

then, express it as soln to optimization problem

$$\Theta^* = \underset{\Theta \in \Theta}{\operatorname{argmin}} f(\Theta) \quad ] \text{ OPT}$$

idea: approximate  $\Theta^*$  by approximating OPT

linear algebra example: want to compute soln to  $Ax=b$  i.e.  $x=A^{-1}b$

$$\min_x \|Ax-b\|^2$$

variational EM:

recall the EM trick      latent variable model  $p(x, \overset{\text{unobserved}}{z})$

$$\log p(z|\epsilon) \geq \mathbb{E}_q[\log \frac{p(x, z|\epsilon)}{q(z)}] \triangleq f(q, \epsilon)$$

$$\underbrace{\log p(z|\epsilon)}_{f(q, \epsilon) + \underbrace{KL(q||p(z|x, \epsilon))}_{\text{KL}(q||p(z|x, \epsilon))}}$$

$$\log p(z|\epsilon) - f(q, \epsilon) = KL(q(\cdot) || p(\cdot|x, \epsilon))$$

$$\text{E-step: } \underset{q \in \text{all distributions}}{\operatorname{argmax}} f(q, \epsilon^{(t)}) \Leftrightarrow \underset{q}{\operatorname{argmin}} KL(q || p(z|x, \epsilon^{(t)}))$$

"a" variational approximation for the E-step:

$$\text{do } q_{\text{approx}}^{(t+1)} = \underset{q \in \underbrace{\Omega_{\text{simple}}}_{\text{source of approximation}}}{\operatorname{argmin}} KL(q(\cdot) || p(\cdot|x))$$

(still get lower bound on  $\log p(z|\epsilon^{(t)})$  but no more monotonically guarantees)

$\rightarrow$  to approximate  $p(z|x, \epsilon^{(t)})$

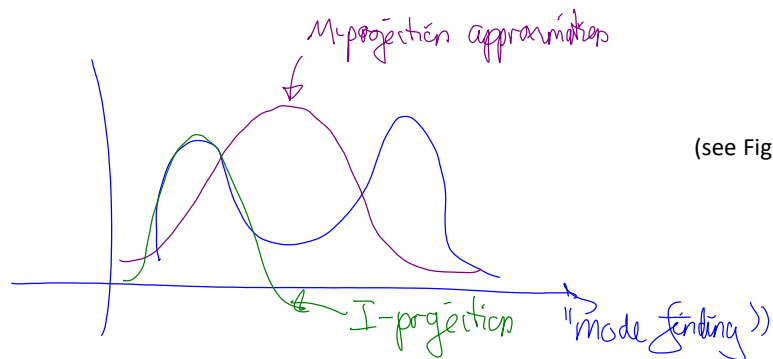
approximate MZP :  $\arg \max_{\theta \in \Theta} \mathbb{E}_{q_{\text{approx}}^{\theta}} [\log p(x, z | \theta)]$

more generally, using  $\arg \min_{q \in \mathcal{Q}} KL(q || p)$  is a variational approach to approximate  $p$

note :  $\uparrow$   
I-projection ; if  $q$  is simple, can compute  $\mathbb{E}_q [\log \frac{q}{p}]$

alternative :  $\arg \min_{q \in \mathcal{Q}} KL(p || q)$  M-projection

"motivation for" EP algorithm  $\rightarrow$  moment matching  
expectation propagation



(see Figure 10.2 in Bishop)

Mean-field approximation:

(section 10.1 in Bishop)

let's suppose that  $p(z)$  is in exponential family

$$z_1, \dots, z_p \quad p(z) = \exp(n^T T(z) - A(n))$$

mean field approximation :  $Q_{MF} = \{q(z) = \prod_i q_i(z_i)\}$

set of fully factorized distributions

$$\begin{aligned} KL(q \| p) &= \mathbb{E}_q \left[ \log \frac{q(z)}{p(z)} \right] \\ &= -n^T \mathbb{E}_q [T(z)] + A(n) + \underbrace{\sum_z q(z) \log q(z)}_{\sum_i \left( \sum_{z_i} q(z_i) \log q(z_i) \right)} \end{aligned}$$

coordinate descent on  $q_i$ 's :

fix  $q_j$  for  $j \neq i$

minimize with respect to  $q_i$   $KL(q_i q_{-i} \| p)$

$$= -\mathbb{E}_{q_i} \left[ \underbrace{n^T \mathbb{E}_{q_{-i}} [T(z)]}_{\triangleq f_i^0(z_i)} \right] + \text{cst.} + \sum_{z_i} q_i(z_i) \log q_i(z_i)$$

add Lagrange multiplier for  $\sum_{z_i} q_i(z_i) = 1 \quad \rightsquigarrow \quad + \lambda (1 - \sum_{z_i} q_i(z_i))$

$$\frac{\partial}{\partial q_i(z_i)} = 0 \Rightarrow -f_i(z_i) + \log q_i(z_i) + 1 - \lambda = 0$$

$$q_i^*(z_i) \propto \exp(f_i(z_i))$$

general mean field update when target  $p$  is in exp. family

$$q_i^{(t+1)}(z_i) \propto \exp(n^T \mathbb{E}_{q_{-i}^{(t)}} T(z))$$

Ising model example:

$$T(z) = \begin{matrix} (z_i)_{i \in V} & z_i \in \{0,1\} \\ (z_i z_j)_{\{i,j\} \in E} \end{matrix}$$

$$\mathbb{E} q_{Ti}(z_j) = q_j(z_j=1) \triangleq \mu_j$$

$$\mathbb{E} q_{Ti}[z_i z_j] = z_i \mu_j$$

$$\begin{aligned} \eta^T \mathbb{E} q_{Ti}^{(t)} T(z) &= \eta_i z_i + \sum_{j \neq i} \eta_j \overbrace{\mathbb{E} q_{Ti}^{(t)}[z_j]}^{\mu_j^{(t)}} \\ &\quad + \sum_{j \in N(i)} \eta_{ij} \underbrace{\mathbb{E} q_{Ti}^{(t)}[z_i z_j]}_{z_i \mu_j^{(t)}} + \text{rest (no } z_i) \end{aligned}$$

result  $q_i^{(t+1)}(z_i) \propto \exp(\eta_i z_i + z_i \sum_{j \in N(i)} \mu_j^{(t)})$

$$\boxed{\mu_i^{(t+1)} = \sigma\left(\eta_i + \sum_{j \in N(i)} \eta_{ij} \mu_j^{(t)}\right)}$$

MF update for  $q_i(z_i)$   
[with parameter  $\mu_i$ ]

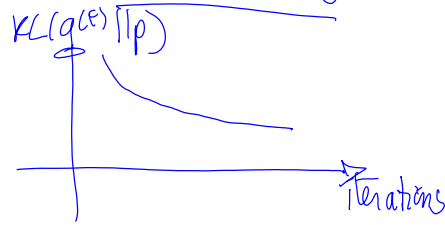
Compare with G-S update where  $z_i^{(t+1)} = 1$  with prob.  $\sigma\left(\eta_i + \sum_{j \in N(i)} \eta_{ij} z_j^{(t)}\right)$

⊗ here min  $KL(q||p)$   
 $q \in Q$

$KL(\cdot||p)$  is a convex function of  $q$   
but  $Q_{MF}$  is a non-convex constraint set



$\Rightarrow$  can get stuck in local minima  
but can monitor progress by



pros & cons of variational methods vs. Sampling

⊕ optimization based  
 $\Rightarrow$  often faster  
& easier to debug

⊖ noisy  $\Rightarrow$  hard to debug  
mixing problem for chains

⊖ biased estimate

$$\mathbb{E}_{q(x)}[f(z)] \neq \mathbb{E}_p[f(z)]$$

⊕ unbiased estimate

$$\mathbb{E}[\mathbb{E}_{q(x)}[f(z)]] = \mathbb{E}_p[f(z)]$$

with respect to  
random sample