

## Lecture 21 - Gibbs sampling; variational methods

Tuesday, November 21, 2023 2:30 PM

- today:
  - Gibbs sampling
  - variational methods

### Gibbs sampling alg.

↳ M.H. with a clever choice of proposal  $q_t(x'|x)$  [clever:  $a(x'|x) \approx 1$ ]

examples of applications

$$\hookrightarrow \text{UGM} : p(x) = \prod_{c \in C} \psi_c(x_c)$$

$$\begin{aligned} &\text{difficult conditional in DGM} & \tilde{p}(x) = p(x_{E^c}, \bar{x}_E) \delta(x_E, \bar{x}_E) \\ & & \text{or } p(x|\bar{x}_E) \end{aligned}$$

(UGM)

cyclic (Gibbs sampling alg): nodes  $i=1, \dots, n$

start at some  $x^{(0)}$

for  $t=1, \dots,$

- pick  $i = (t \bmod n) + 1$
- sample  $x_i^{(t)} \sim p(x_i = \cdot | x_{\setminus i}^{(t-1)} = x_{\setminus i}^{(t-1)})$
- set  $x_j^{(t)} = x_j^{(t-1)}$  for  $j \neq i$

$\{x_i\}_{i=1}^n$

true conditioned on  $x_i$  as a proposal

end

⊗ Gibbs sampling is M.H. with a time varying proposal

suppose we pick  $i$  at time  $t$

for  $x_i$  to stay constant

$$\text{then the proposal is } q_t(x'|x^{(t-1)}) = p(x_i'|x_{\setminus i}^{(t-1)}) \delta(x_{\setminus i}', x_{\setminus i}^{(t-1)})$$

M.H. acceptance ratio:

$$a(x'|x^{(t-1)}) = \frac{q_t(x'|x^{(t-1)}) p(x')}{q_t(x'|x^{(t-1)}) p(x^{(t-1)})} = \frac{p(x_i'|x_{\setminus i}^{(t-1)}) \delta(x_{\setminus i}', x_{\setminus i}^{(t-1)})}{p(x_i^{(t-1)}) p(x_{\setminus i}'|x_{\setminus i}^{(t-1)})} = 1$$

$$q_t(x' | x^{(t-1)}) p(x^{(t-1)}) \xrightarrow{\text{accept}} p(x_{\tau_i}^{(t-1)}) p(x_i^{(t)} | x_{\tau_i}^{(t-1)})$$

$$\hookrightarrow p(x_i^{(t)} | x_{\tau_i}^{(t-1)}) \delta(x_{\tau_i}^{(t-1)}, x_i^{(t)})$$

//  
always accept

### Convergence of GS.:

- Let  $A$  be a Markov transition kernel of one full cycle of Gibbs sampling (i.e.  $n$  steps)

→ homogeneous M.C.

If suppose that  $\underline{p(x) > 0 \forall x}$   $\Rightarrow A$  is irreducible  $\rightarrow$  aperiodic because  $A_{ii} > 0$

$$\Rightarrow p(x_i | x_{\tau_i}) > 0$$

$\forall x_i \in \mathcal{X}_{\tau_i}$

$$A_{k,k} > 0$$

$$A_{k,k}$$

since can get to any state with h "flips"

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

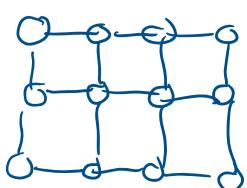
by detailed balance,  $p$  is stationary dist. of chain

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

### Example: GS. for Ising model

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

UGM:



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

[minimal exp. family representation]

for GS.

want to compute  $p(x_i | x_{\tau_i}) \stackrel{\text{by cond. indep.}}{=} p(x_i | x_{\text{neighbors}(i)})$

$$\propto \exp \left( m_i x_i + \sum_{j \in N(i)} m_{ij} x_i x_j + \text{rest} \right)$$

$\Rightarrow$  renormalize to get conditional

$$p(x_i=1 | x_{\tau_i}) \propto \frac{\exp(m_i + \sum_{j \in N(i)} m_{ij} x_j) \exp(\text{rest})}{(1 + \exp(m_i + \sum_{j \in N(i)} m_{ij} x_j)) \exp(\text{rest})}$$

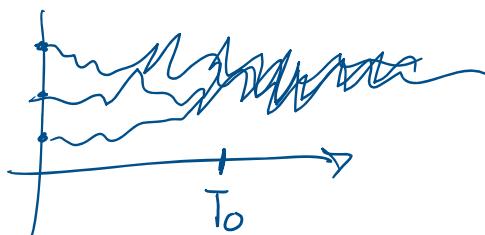
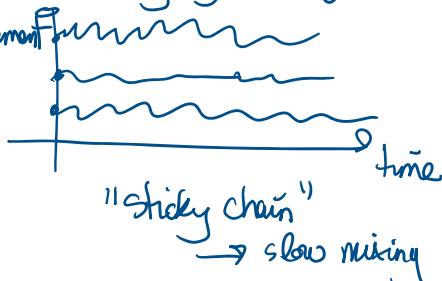
$$\boxed{\text{--- (t) --- (t-1) --- / --- / --- (t+1) ---}}$$

$$P(x_i=1 | x_{-i}) = \sigma\left(\eta_i + \sum_{j \in N(i)} \eta_{ij} x_j\right)$$

## diagnostic of mixing

monitor mixing by running independent chains

## Some metallurgy



15h16

bad at 15h40

## Variational methods

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

Then, express it as a solution to opt. problem

$$\hat{\theta}^* = \arg \min_{\theta \in \Theta} f(\theta) \quad ] \text{OPT}$$

Idea: approximate  $\hat{G}^*$  via an approximation of  $OPT$

## Linear algebra example:

say want sol'n to  $Ax=b$  (i.e.  $x^* = A^{-1}b$ )

$$x^* = \underset{x}{\operatorname{arg\,min}} \|Ax - b\|^2$$

## Variational EM (motivation for objective)

recall EM trick

latent  $p(x, z | G)$

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

$$\log p(x|z) = \log p(x|z) - \log \frac{p(x|z)}{q(z)} = \log p(x|z)$$

$$\log p(x|z) - \mathcal{J}(q, \theta) = KL(q(\cdot) || p(z|x, \theta))$$

E step :  $\underset{\substack{q \in \text{all distributions} \\ z}}{\operatorname{argmax}} \mathcal{J}(q, \theta^{(t)}) \Leftrightarrow \underset{q}{\operatorname{argmin}} KL(q(\cdot) || p(\cdot | z, \theta^{(t)}))$

(\*) a variational approx. for E step :

$$\text{do } q_{\text{approx}}^{(t+1)} = \underset{q \in Q_{\text{simple}}}{\operatorname{argmin}} KL(q || p(\cdot | z, \theta^{(t)}))$$

$\uparrow$  source of approx  $\rightarrow$  to approximate  $p(z|x, \theta^{(t)})$

approximate M step :

$$\theta^{(t+1)} = \underset{\theta \in \Theta}{\operatorname{argmax}} \underbrace{\mathbb{E}_{q_{\text{approx}}^{(t+1)}} [\log p(x, z | \theta)]}_{\mathcal{J}(q_{\text{approx}}^{(t+1)}, \theta)}$$

still a lower bound on  $\log p(x)$

but lose monotonicity guarantee  
in  $\log p(x | \theta^{(t)})$  as fct. of  $t$

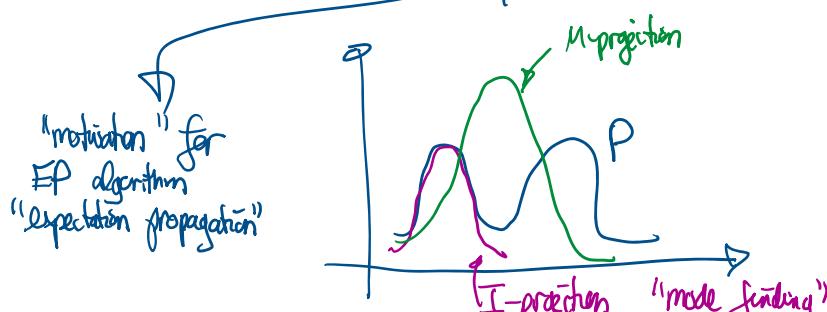
more generally, using  $\underset{q \in Q}{\operatorname{argmin}} KL(q || p)$  is a variational approach to approximate  $P$

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

alternative:  $\underset{q \in Q}{\operatorname{argmin}} KL(p || q)$

M-projection

$\hookrightarrow$  "moment matching"



(see fig. 10.2 in Bishop)

expectation propagation



(see fig. 10.2 in Bishop)

## Mean-field approximation (section 10.1 in Bishop)

Let's suppose that  $p(z)$  is in exp. family

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

$$\text{mean field approximation} \quad Q_{MF} = \left\{ q(z) = \prod_i q_i(z_i) \right\}$$

set of fully factorized dist.

$$\begin{aligned}
 KL(q \parallel p) &= \mathbb{E}_q[\log q/p] \\
 &= -m^T \mathbb{E}_q[T(z)] + A(m) + \sum_z q(z) \log q(z) \\
 &\quad \underbrace{\left( \mathbb{E}_q[T(z)] \right)}_{\sum_j q_j(z_j)} \leq \log q_j(z_j) \\
 &\quad \sum_i \sum_z q_{zi}(z_i) q_i(z_i) \log q_i(z_i) \\
 &\quad \underbrace{\sum_z 1 \cdot \sum_i q_i(z_i)}_{\sum_i q_i(z_i)} \\
 &\quad \sum_i q_i(z_i) \log q_i(z_i)
 \end{aligned}$$

coordinate descent on  $q_i$ 's

$$\text{fix } q_j \text{ for } j \neq i; \quad \min_{\text{w.r.t. } q_i} KL(q_i, q_{zi} \parallel p)$$

$$\begin{aligned}
 &= -\mathbb{E}_{q_i} \left[ n^T \mathbb{E}_{q_{zi}}(T(z)) \right] + \text{const.} + \sum_i q_i(z_i) \log q_i(z_i) \\
 &\triangleq f_i(z_i)
 \end{aligned}$$

(like MaxENT) add Lagrange multiplier for  $\sum_i q_i(z_i) = 1 \rightarrow \lambda (1 - \sum_i q_i(z_i))$

$$\begin{aligned}
 \frac{\partial}{\partial q_i(z_i)} &\Rightarrow \Rightarrow -f_i(z_i) + \log q_i(z_i) + 1 - \lambda = 0 \\
 &\Rightarrow q_i^*(z_i) \propto \exp(f_i(z_i))
 \end{aligned}$$

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

$$q_i^{(t+1)}(z_i) \text{ or } \exp(n^T \mathbb{E}_{q_{zi}}[T(z)])$$

## Ising model

$$T(z) \Rightarrow (z_i)_{\substack{(z_i z_j) \\ \{i,j\} \in E}} \quad z_i \in \{0,1\}$$

$$\mathbb{E} q_{\tau_i}^{(t)}(z_j) = q_j^{(t)}(z_j=1) \triangleq \mu_j^{(t)}$$

$$\mathbb{E} q_{\tau_i}^{(t)}(z_i z_j) = z_i \mu_j^{(t)}$$

$$m^T \mathbb{E} q_{\tau_i}^{(t)} [T(z)] = m_i z_i + \underbrace{\sum_{j \neq i} n_{ij} \mathbb{E} q_{\tau_i}^{(t)} \frac{z_j}{\mu_j^{(t)}}}_{\text{no } z_i \rightarrow \text{constant}} + \underbrace{\sum_{j \in N(i)} n_{ij} \mathbb{E} q_{\tau_i}^{(t)} [z_i z_j]}_{z_i \mu_j^{(t)}} + \text{rest}$$

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

parameter for  $q_j^{(t)}$

$$m_i^{(t+1)} = \sigma(m_i + \sum_{j \in N(i)} n_{ij} z_j^{(t)})$$

MF update for  $q_i^{(t)}$  with parameter  $m_i$

compare with G.S. update

$$\text{where } z_i^{(t)} = 1 \text{ with prob. } \sigma(m_i + \sum_{j \in N(i)} n_{ij} z_j^{(t)})$$