

Lecture 21 - Gibbs sampling; variational methods

Tuesday, November 22, 2022 2:58 PM

today :

- Gibbs sampling
- variational methods

Gibbs sampling alg.

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

example of applications:

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

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

(UGM)

Cyclic Gibbs sampling alg.: nodes $i=1, \dots, n$

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

for $t=1, \dots,$

- pick $i = (t \bmod n) + 1 \rightarrow 1:n \setminus \{i\}$
- sample $x_i^{(t)} \sim p(x_i = \cdot | X_{\setminus i} = x_{\setminus i}^{(t-1)})$
true conditional on $X_{\setminus i}$ as proposal
- set $x_j^{(t)} = x_j^{(t-1)}$ for $j \neq i$

end

★ Gibbs sampling is M.H. with a time varying proposal
suppose we pick i at time t

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

M.H. acceptance ratio:

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

for $x_{\setminus i}$ to stay constant

always accept?

$$\rightarrow p(x_i | \cancel{x_{\neq i}^{(t-1)}}) \delta(x_{\neq i}^{(t)}, x_{\neq i}^{(t-1)})$$

always accepts?

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 & aperiodic because $A_{ii} > 0$

$\Rightarrow p(x_i | x_{\neq i}) > 0$ $\forall x_i \neq x_{\neq i}$ $\begin{cases} A_{ij} > 0 \\ i \neq j \end{cases}$

\Rightarrow since can get to any state with n "flips"

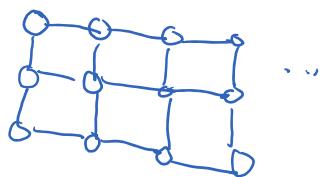
\Rightarrow by detailed balance, p is stationary dist. of chain ...

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

example : G.S. for Ising model

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

UGM :



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

[mininal exp. family representation]

for G.S.,

$$\text{want to compute } p(x_i | x_{\neq i}) = p(x_i | x_{\text{neighbors}(i)}) \stackrel{\text{by cond. indep}}{=} \propto \exp(m_i x_i + \sum_{j \in N(i)} m_{ij} x_i x_j + \text{rest})$$

\Rightarrow renormalize to get conditional :

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

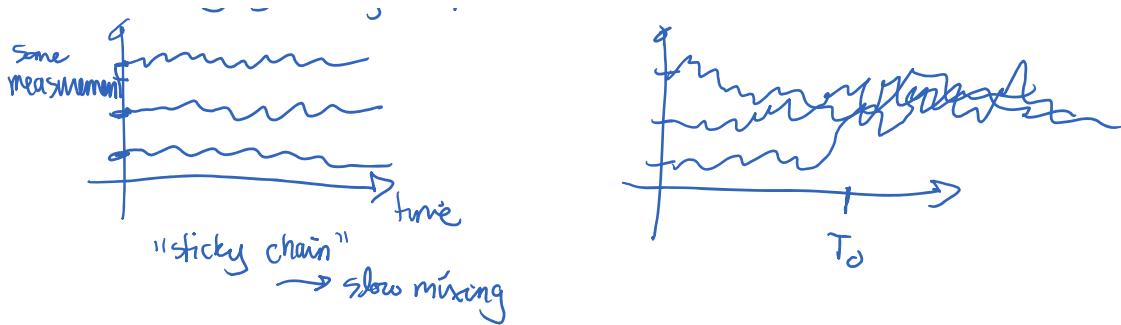
$$p(x_i^{(t)} = 1 | x_{\neq i}^{(t-1)}) = \sigma(m_i + \sum_{j \in N(i)} m_{ij} x_j^{(t-1)})$$

diagnostic of mixing

monitor mixing by running independent chains

some measurement

from a blank slate



16h02

Variational methods

general idea: say we want to approximate θ^*

then, express it as solution to opt. problem

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{arg\,min}} f(\theta) \quad] \text{OPT}$$

Idea: approximate θ^* via an approximation to OPT

linear algebra example:

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

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

Variational EM (motivation for objective)

recall EM trick

latent variable $p(x, z | \theta)$
unobserved \downarrow

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

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

EM step: $\underset{\substack{q \in \text{all distributions} \\ \text{on } z}}{\operatorname{arg\,max}} \mathcal{L}(q, \theta^{(t)}) \Leftrightarrow \underset{q}{\operatorname{arg\,min}} KL(q(\cdot) || p(\cdot | x, \theta^{(t)})$

⇒ a variational approx. for E step:

$$\text{do } Q_{\text{approx}}^{(t+1)} = \underset{\substack{q \in \text{simple} \\ \Phi}}{\operatorname{arg\,min}} KL(q || p(\cdot | x, \theta^{(t)})$$

sample of approx. → to approximate $p(z|x, \theta^{(t)})$

approximate M step:

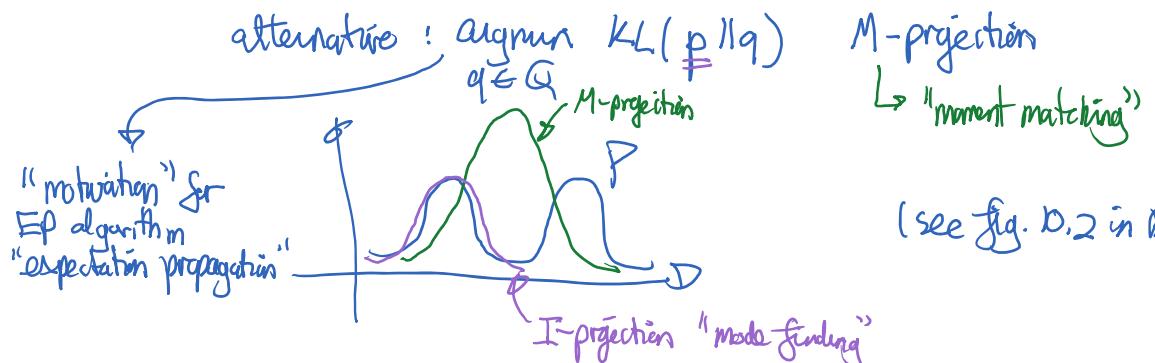
$$\theta^{(t+1)} = \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E}_{q_{\text{approx}}^{(t)}} [\log p(z|x, \theta)]$$

$$= \mathbb{E}_{q_{\text{approx}}^{(t)}} [\log p(z|x, \theta)]$$

still a lower bound on $\log p(x|\theta)$
 but lose monotonicity guarantee
 on $\log p(z|q^{(t)})$ 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, $\mathbb{E}_q [\log p/q]$
 can compute



Mean-field approximation (section 10.1 in Bishop)

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

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

mean field approximation, $Q_{MF} = \{q(z) = \prod_i q_i(z_i)\}$
 set of fully factorized dist.

$$KL(q||p) = \mathbb{E}_q [\log \frac{q}{p}]$$

$$= -\eta^T \mathbb{E}_q [T(z)] + A(\eta) + \sum_z q(z) \log q(z)$$

$$\sum_i \underbrace{\sum_z q_{zi}(z_i) q_i(z_i) \log q_i(z_i)}_{1 \cdot \sum_i} = \sum_i \sum_z q_i(z_i) \log q_i(z_i)$$

$$\sum_j (\prod_i q_{zi}(z_i)) \sum_z \log q_{ji}(z_j)$$

Coordinate descent on q_i 's

$$\text{fix } q_j \text{ for } j \neq i; \min_{\text{w.r.t. to } q_i} KL(q_i \| p) = -\mathbb{E}_{q_i} [m^T \mathbb{E}_{q_{\bar{i}}}[t(z)]] + \text{const.} + \sum_{z_i} q_i(z_i) \log q_i(z_i)$$

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

$$\partial_{\lambda q_i(z_i)} = 0 \rightarrow -f_i(z_i) + \log q_i(z_i) + 1 - \lambda = 0 \Rightarrow q_i(z_i) \propto \exp(f_i(z_i))$$

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

$$q_i^{(t+1)}(z_i) \propto \exp(m^T \mathbb{E}_{q_{\bar{i}}^{(t)}}[t(z)])$$

I sing model:

$$T(z) \Rightarrow (z_i)_{i \in V} \quad z_i \in \{0, 1\}$$

$$(z_i z_j)_{i, j \in E}$$

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

$$\mathbb{E}_{q_{\bar{i}}^{(t)}}(z_i z_j) = z_i \mu_j^{(t)}$$

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

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

$$M_i^{(t+1)} = \sigma(m_i + \sum_{j \in N(i)} m_{ij} \mu_j^{(t)})$$

parameter for $q_j^{(t)}$
MF update for $q_i(z_i)$
[with parameter M_i]

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