

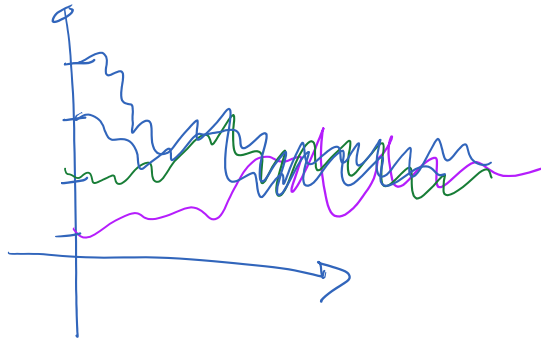
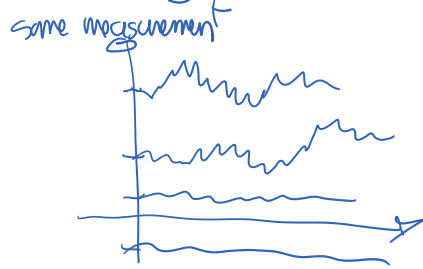
Lecture 22 - scribbles

Friday, November 23, 2018 13:28

today: variational methods
mean field dg.
estimation in PGM

(sampling chd): diagnostic of mixing

monitor mixing by running independent chains



"sticky chain"
→ slow mixing

usually slow mixing comes from difficulty to move between modes



→ annealing methods help this

physics analogy
energy
↓ ↓

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

"temperature"
↑
high temperature
⇒ more exploration

example: "Annealed importance sampling"

14th 47

Variational methods

general idea: say we want to approximate Θ^*

then, express it as solution to optimization problem

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

idea: approximate Θ^* by approximating OPT

Linear algebra example:

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

$$\underset{x}{\operatorname{argmin}} \|Ax-b\|^2$$

Variational EM (motivation)

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

$$\log p(x|\theta) \geq \mathbb{E}_q \left[\log p(x, z|\theta) \right]_{q(z)} \triangleq J(q, \theta)$$

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

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

⊗ a variational approximation for the E-step

$$\text{do } q_{\text{approx}}^{(t+1)} = \underset{q \in \mathcal{Q}_{\text{simple}}}{\text{argmin}} \text{KL}(q \| p(\cdot | x, \epsilon^{(t)}))$$

(still get a lower on $\log p(z|\epsilon^{(t)})$ but no monotonically guarantees for approximate EM)

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

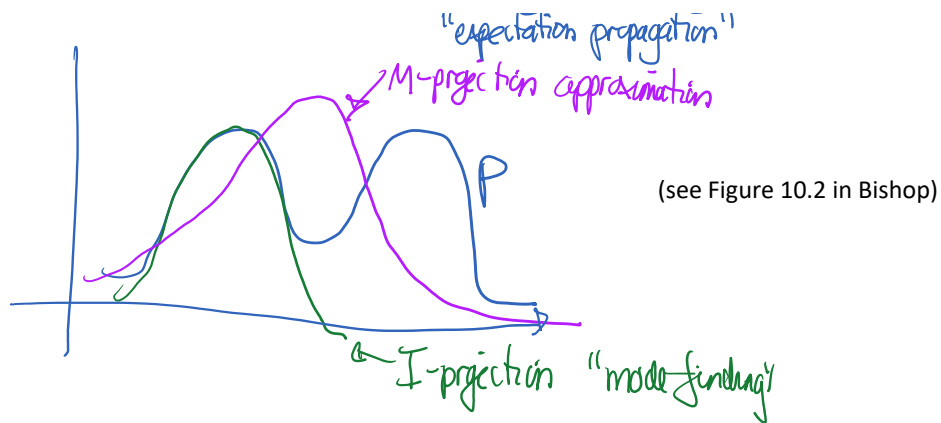
$$\text{approximate M step: } \mathcal{J} = \underset{\theta \in \Theta}{\text{argmax}} \mathbb{E}_{q_{\text{approx}}^{(t)}} [\log p(x, z | \theta)]$$

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

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

alternative: $\underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(p \| q)$ M-projection

"motivation" EP algorithm "expectation propagation" \rightarrow moment matching
 \rightarrow M-projection approximation



Mean field approximation: (section 10.1 in Bishop)

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

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

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

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

coordinate descent on q_i 's:

fix q_j for $j \neq i$

minimize w.r. to q_i $KL(q_i \circ q_{-i} \| p)$

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

add Lagrange multiplier

$$\text{for } \sum_{z_i} q_i(z_i) = 1 \quad \rightarrow \quad + \lambda (1 - \sum_{z_i} q_i(z_i))$$

$$\frac{\partial}{\partial q_i(z_i)} = 0 \quad \Rightarrow$$

$$-f_i(z_i) + \log q_i(z_i) + 1 - \lambda = 0$$

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

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

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

14h40

Ising model example

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

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

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

Model

$$\mathbb{E} q_{7i} L(z_i \in \mathcal{Z}_j) = z_i \mu_j^0$$

$$\begin{aligned} \eta^T \mathbb{E} q_{7i}^{(t)} T(z) &= m_i z_i + \sum_{j \neq i} \eta_j \underbrace{\mathbb{E} q_{7i}^{(t)} [z_j]}_{\mu_j^{(t)}} + \\ &+ \sum_{j \in N(i)} m_{ij} \underbrace{\mathbb{E} q_{7i}^{(t)} [z_i z_j]}_{z_i \mu_j^{(t)}} + \text{rest} \\ &\quad \text{(no } z_i) \end{aligned}$$

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

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

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

Compare with

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

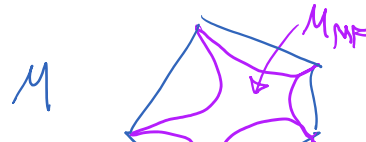
⊗ here $\min_{q \in Q_{MF}} KL(q||p)$

• $KL(\cdot||p)$ is a convex of q

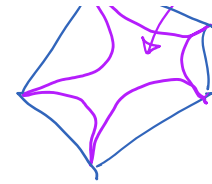
• but Q_{MF} is a non-convex constraint set $\rightarrow \mu_{ij} = \mu_i \mu_j$

\Rightarrow can get stuck in local minima

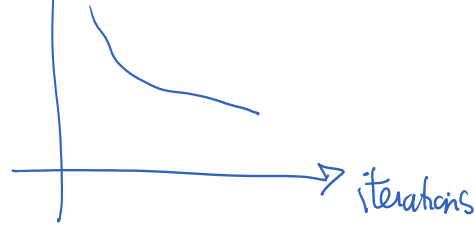
[see lecture 22 last year for "marginal polytope" def. $\underbrace{\text{convex}}_{\text{set of feasible mean parameters}}$]



M



but can monitor progress by $KL(q^{(t)} || p)$



pros & cons of variational methods vs

(+) optimization based
 ⇒ often faster & easier to debug

(-) baised estimate

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

Sampling

(-) noisy ⇒ hard to debug
 mixing problem for chains

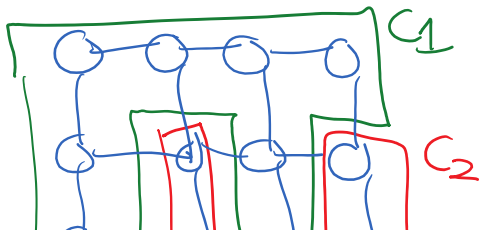
(+) unbiased estimate

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

↳ with respect to random sample

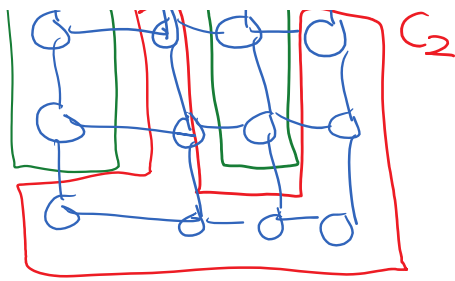
Structured mean field: (see lect. 22 last year)

idea $q(z) = \prod_{j=1}^k q_j(z_{C_j})$



where C_1, \dots, C_k partition of V
 and q_j 's are tractable distributions
 (for example tree UGM)

[lecture 22 Fall 2017 link](#)



(for example tree UGM)

Estimation of parameters for PGM

DGM: parametric family $P_{\Theta} = \sum p_{\Theta}(x) = \prod_i p(x_i | x_{\pi_i}, \Theta_i)$

$\Theta = (\Theta_1, \dots, \Theta_{|V|})$

$\Theta \in \mathcal{H} = \mathcal{H}_1 \times \dots \times \mathcal{H}_{|V|}$

independent parameterization

ie. no tying of parameters

\Rightarrow MLE decouples in $|V|$ independent problems

$\{x^{(i)}\}_{i=1}^n$ $p(\text{data} | \Theta) = \prod_{i=1}^n p(x^{(i)} | \Theta) = \prod_{i=1}^n \prod_{j=1}^{|V|} p(x_j^{(i)} | x_{\pi_j}^{(i)}, \Theta_j)$

$\log [\quad] = \sum_{j=1}^{|V|} \underbrace{\left(\sum_{i=1}^n \log p(x_j^{(i)} | x_{\pi_j}^{(i)}, \Theta_j) \right)}_{f_j(\Theta_j)}$

example: for discrete R.V. $\Rightarrow \Theta_j^{ML} = \text{proportion of observations}$

$\#(x_j = k, x_{\pi_j} = \text{something})$
 $\#(\text{obs} = \text{something})$

(arg = something)

⊗ if have latent variable (i.e. unobserved variable)

⇒ use E.M.

UGM:

example for exp family

$$\rightarrow \exp(\eta_c^T T_c(x_c)) = \psi_c(x_c)$$

$$p(x|n) = \exp\left(\sum_c \eta_c^T T_c(x_c) - A(n)\right)$$

gradient ascent on log-likelihood,

$$\frac{1}{n} \sum_{i=1}^n \log p(x^{(i)}|n) = \sum_c \eta_c^T \left(\underbrace{\frac{1}{n} \sum_{i=1}^n T_c(x_c^{(i)})}_{\hat{\mu}_c} \right) - \frac{1}{n} A(n)$$

$$\nabla \eta_c () = \hat{\mu}_c - \underbrace{\mu_c(n)}_{\mathbb{E}_{p(x|n)} [T_c(x_c)]}$$

to compute this, need inference

e.g. Ising model $T_{ij}(x_i, x_j) = x_i x_j$

$$\mathbb{E}[T_{ij}] = \mu_{ij}^0 = p(x_i=1, x_j=1 | n)$$

perhaps use approximate inference \leftarrow variational

15h26

sampling