

# Lecture 22 - scribbles

Friday, November 24, 2017  
13:30

today: - finish variational methods  
• ML for graphs models

## Structured mean field:

context: suppose target  $p$  is UGM with graph  $G=(V,E)$

(see [https://metacademy.org/graphs/concepts/structured\\_mean\\_field](https://metacademy.org/graphs/concepts/structured_mean_field) for pointers)

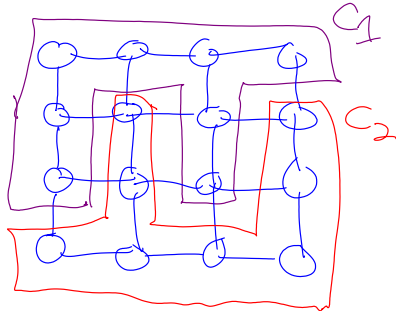
general idea: use  $q(z) = \prod_{j=1}^k q_j(z_{C_j})$

where  $C_1, \dots, C_k$  partition  $V$

and  $q_j$ 's are tractable distributions  
(for example tree UGM)

to do this, you choose  $C_j$

st.  $G$  restricted to  $C_j$  is a low treewidth UGM



⊗ suppose  $q(z) = \prod_j q_j(z_{C_j})$

where  $q_j$  is arbitrary distribution over  $z_{C_j}$

graph  $G$   
restricted  
on nodes  
in  $C_j$

it's easy to show that

$q_j^*(z_{C_j})$  is UGM with structure  $G|C_j$

i.e. when do  $\min_{q \in \mathcal{Q}} KL(q||p)$

⊗ we can rederive mean field ~~opt~~ as before using coordinate descent on  $KL(q||p)$

to get  $q_j^{(t+1)}(z_j) \propto \exp(n^T \mathbb{E}_{q_{7C_j}^{(t)}}(T(z)))$

this will be VGM  $G \setminus C_j$

Ising model example:

$$n^T T(z) = \sum_i n_i z_i + \sum_{\{i,j\} \in E} n_{ij} z_i z_j \quad z_i \in \{-1, 1\}$$

$$\mathbb{E}_{q_{7C_k}^{(t)}} [n^T T(z)] = \sum_{i \in C_k} \left[ n_i z_i + z_i \left( \sum_{j \in N(i) \setminus C_k} n_{ij} \underbrace{\mathbb{E}_{q_{7C_k}^{(t)}}(z_j)}_{\text{effect of rest of graph}} \right) \right]$$

$$+ \sum_{\substack{\{i,j\} \in E \\ i,j \in C_k}} n_{ij} z_i z_j + \text{rest}$$

does not depend on  $z_i$  for  $i \in C_k$

convex optimization viewpoint on variational inference

Wainwright & Jordan FTML 2008 [paper](#)

$p$  in exponential family

$$KL(q||p) = \mathbb{E}_q \left[ \log \frac{q(x)}{p(x)} \right] = -n^T \underbrace{\mathbb{E}_q [T(x)]}_{\triangleq \mu(q)} + A(n) - H(q) = -(-\mathbb{E}_q [\log q(x)])$$

$$KL(q||p) \geq 0$$

(recall that  $\nabla_n A(n) = \mu(n)$ )

$$\Rightarrow A(n) \geq n^T \mu(q) + H(q) \text{ for any } q$$

$$A(n) = \sup_{q \in \text{all dist.}} n^T \mu(q) + H(q) \stackrel{\text{Wainwright showed}}{=} \sup_{\mu \in \mathcal{M}} n^T \mu + H(n(\mu))$$

$\mu = (\dots)$

since  $\mu(q) \triangleq \mathbb{E}_q[T(x)]$

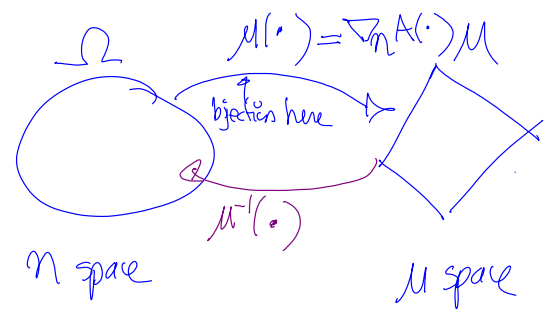
$$M = \{ \mu : \exists \text{ a dist. } q \text{ s.t. } \mathbb{E}_q[T(x)] = \mu \} \quad (\text{convex set})$$

$\triangleq$  Marginal polytope

(e.g. for Ising model:  $T(z) = \begin{pmatrix} (z_i)_{i \in V} \\ (z_i z_j)_{i,j \in E} \end{pmatrix}$ )

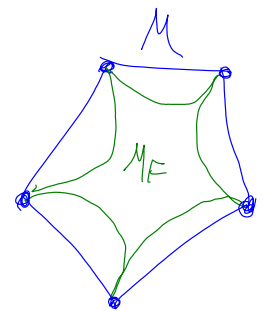
for a minimal exponential family, Wainwright & al. showed

that there is a one-to-one correspondence between  $\mu \in \text{int}(M)$  and  $\eta \in \text{int}(\Omega)$



$$A(\eta) = \sup_{\mu \in M} \eta^T \mu + H(\eta | \mu)$$

intractable for Ising model



idea: approximate  $\mu$  &  $H$  to get tractable approximation

mean field approximate is to restrict  $M$  to  $\mathbb{E}_q[T(x)]$  where  $q$  is fully factorized

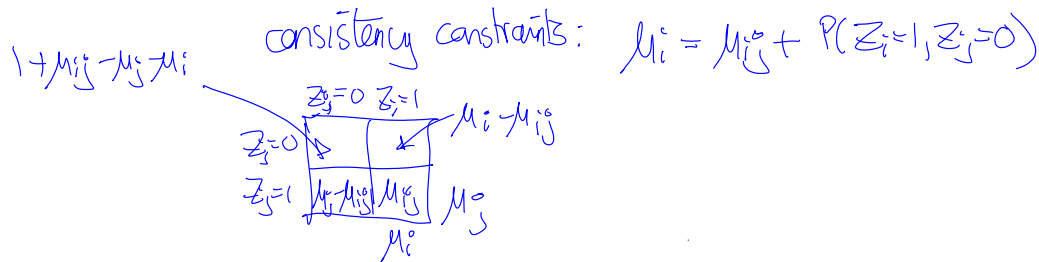
recall that for Ising model

here for mean field  $\approx$

$$(M_i)_{i \in V}$$

$$M_{ij} \triangleq P(z_i=1, z_j=1) \quad i, j \in E$$

$$\{ M_{ij} = M_i \circ M_j \}$$



⊗ outer bound approximation:

$$\sup_{\mu \in L} nT\mu + H_{\text{Berne}}(\mu)$$

local consistency polytope  $M \subseteq L$

$$\text{which says e.g. that } \sum_k P(z_i, z_j=k) = P(z_i)$$

→ approximation of entropy as if it was a tree

↓  
recall that for a tree,  
we had  $P(z) = \prod_i P(z_i) \prod_{i,j \in E} \frac{P(z_i, z_j)}{P(z_i)P(z_j)}$

other approaches: TRW approach → gives a convex upper bound on  $A(n)$  optims.

aside:  $f^*(\mu) \triangleq \sup_n nT\mu - f(n)$  Fenchel-conjugate

$$\text{when } f \text{ is convex, } f^{**}(n) = f(n)$$

$$\text{it turns out that } A^*(\mu) = \begin{cases} H(n\mu) & \text{for } \mu \in \text{int}(M) \\ +\infty & \text{for } \mu \notin M \end{cases}$$

$$\bullet \text{ Duro } n(\text{data}, n) = \leq nT\langle z_i \rangle - nA(n)$$

$$\bullet \log p(\text{data}; n) = \sum_i n^T J(x_i) - nA(n)$$

Estimation of parameters for PGM:

DGM: parametric family  $\mathcal{P}_{\Theta} = \left\{ p_{\Theta}(x) = \prod_i p(x_i | x_{\pi_i}, \Theta_i) \right\}$

independent parametrization

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

$\Theta \in \Theta = \Theta_1 \times \dots \times \Theta_M$

ie. no tying of parameters

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

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

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

example: for discrete R.V.  $\Rightarrow \hat{\Theta}_j^{\text{ML}} = \frac{\#(x_j = k, x_{\pi_j} = \text{something})}{\#(x_{\pi_j} = \text{something})}$

$\otimes$  If have latent variable (i.e. unobserved variables)

$\Rightarrow$  use EM

(1)  $\rightarrow$  M  $\rightarrow$

# UGM

example for exponential family:

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

$$\exp(\eta_c^T T_c(x)) = \psi_c(x)$$

gradient ascent on log-likelihood:

$$\frac{1}{n} \sum_{i=1}^n \log p(x^{(i)}|m) = \sum_c \eta_c^T \left( \frac{1}{n} \sum_{i=1}^n T_c(x^{(i)}) \right) - A(\eta)$$

$$\nabla_{\eta_c} ( \quad ) = \hat{\mu}_c - \mu_c(m)$$

↳  $\mathbb{E}_{p(x|m)} [T_c(x)]$

to compute this, need inference

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

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

perhaps use approximate inference

- ↳ variational
- ↳ sampling