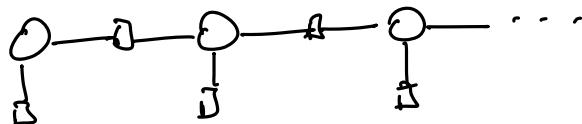


LECTURE 6 - INFERENCE (CONT.)

1. Belief Propagation

last time : BP on Ising chains

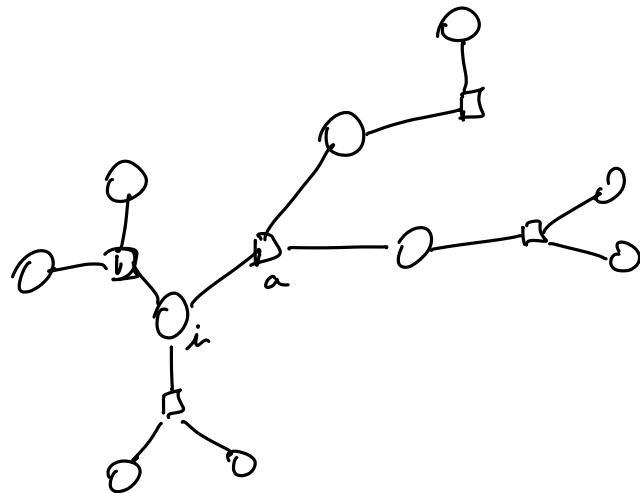


→ forward ($v_{\rightarrow j}$) and backward ($v_{j \leftarrow}$) messages computed using dynamic programming.

BP on trees

$$G = (V, A, E)$$

$$E \subseteq V \times A$$



$$p(x) = \frac{1}{Z} \prod_{a \in A} \phi_a(x_{N(a)})$$

→ At iteration t , two messages for each edge (i, a)

- $v_{i \rightarrow a}^{(t)}(x_i)$
- $\hat{v}_{a \rightarrow i}^{(t)}(x_i)$

→ BP / sum-product equations:

$$(*) \quad v_{j \rightarrow a}^{(t+1)}(x_j) \propto \prod_{b \in N(j) \setminus \{a\}} \hat{v}_{b \rightarrow j}^{(t)}(x_j)$$

$$(**) \quad \hat{v}_{a \rightarrow j}^{(t)}(x_j) \propto \sum_{x_{N(a) \setminus \{j\}}} \phi_a(x_{N(a)}) \prod_{b \in N(a) \setminus \{j\}} v_{b \rightarrow a}^{(t)}(x_b)$$

→ Estimate marginals:

$$v_i^{(t)}(x_i) \propto \prod_{a \in N(i)} \hat{v}_{a \rightarrow i}^{(t-1)}(x_i)$$

→ Algorithm:

Initialize $v_{i \rightarrow a}^{(0)}(\cdot)$ as i.i.d. with some dist. p_0

For $t=0, \dots, t_{\max}$

- For each $(j, a) \in \bar{E}$,

compute $\hat{v}_{a \rightarrow j}^{(t)}$ using $(**)$

- For each $(j, a) \in \bar{E}$

compute $v_{a \rightarrow j}^{(t+1)}$ using $*$

- If no message changed:

return current messages.

return "not converged"

Theorem (BP is exact on trees)

Let $t^* = \text{maximum distance between any two variable nodes}$
(diameter)

1. BP updates converge after at most t^* iterations,
irrespective of initial condition.

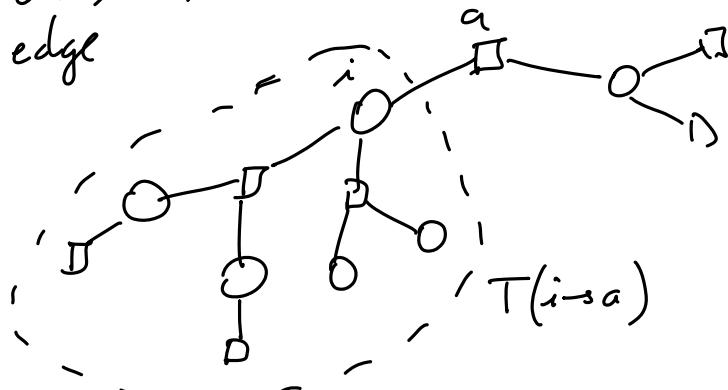
For any $(i,a) \in E$, $t > t^*$, $v_{i \rightarrow a}^{(t)} = v_{i \rightarrow a}^*$ (fixed-point
 $\hat{v}_{a \rightarrow i}^{(t)} = \hat{v}_{a \rightarrow i}^*$ messages)

2. The fixed-point messages provide exact marginals.

i.e. $p(x_i) = v_i^* (x_i)$ for all $i \in V$.

Proof (sketch):

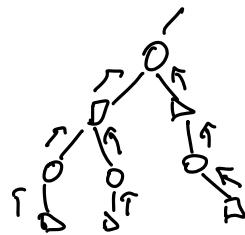
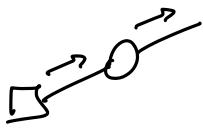
| For an edge $(i,a) \in E$, let $T(i \rightarrow a)$ be the subtree
| rooted at this edge



| Denote by $t^*(i \rightarrow a)$ the depth of the tree $T(i \rightarrow a)$.

| We can show by induction on $t^*(i \rightarrow a)$ that
| the message $v_{i \rightarrow a}^{(t)}$ converges to the marginal of x_i
| in the subgraph $T(i \rightarrow a)$ after $t^*(i \rightarrow a)$ steps.

(exercise = similar to DP recursion in the chain)



The marginals are then given by.

$$p(x_i) \propto \sum_{x_{V \setminus \{i\}}} \prod_a \phi_a(x_{N(a)})$$



$$\propto \prod_{b \in N(i)} \sum_{x_{N(b) \setminus \{i\}}} \phi_b(x_{N(b)}) \prod_{j \in N(b) \setminus \{i\}} v_{j \rightarrow b}^*(x_j)$$

$$= v_i^*(x_i)$$

□

Remarks:

- conditioning can be done using the same BP equations by adding factors $\mathbb{I}\{x_i = x_i^*\}$

$$p(x | x_B = x_B^*) \propto \prod_a \phi_a(x_{N(a)}) \cdot \mathbb{I}\{x_B = x_B^*\}$$

$$\stackrel{\text{O}_i}{\downarrow} \mathbb{I}\{x_i = x_i^*\}$$

- complexity depends - linearly on depth of tree
- exponentially on size of factors

■ BP on general graphs - ("Loopy BP")

- Note that one can run the BP algorithm on any graph
- In general, BP does not converge to the correct marginals.
- Yet, it often works well in practice ("Loopy BP")
- It can be shown that the fixed points of the BP equations correspond to the stationary points of a quantity known as the "Bethe Free Energy", which is often a good approximation of the Free Energy $\log Z$

■ Max-product and Min-sum

- Sometimes we care about optimization instead of marginals

Ex.: MAP inference

$$\max_{\mathbf{x}} p(\mathbf{z} | \mathbf{x} = \mathbf{x}_0)$$

$\stackrel{\mathbf{z}}$
(image reconstruction, sequence decoding)

- Similar algorithms to BP can be derived for such optimization problems.

→ Max-product : $(\Sigma, \tau) \mapsto (\max, \tau)$

$$v_{j \rightarrow a}^{(t+1)}(x_j) \propto \prod_{b \in N(j) \setminus \{a\}} \hat{v}_{b \rightarrow j}^{(t)}(x_j)$$

$$\hat{v}_{a \rightarrow j}^{(t)}(x_j) \propto \max_{x_{N(a) \setminus \{j\}}} \phi_a(x_{N(a)}) \prod_{b \in N(a) \setminus \{j\}} v_{b \rightarrow a}^{(t)}(x_b)$$

max-marginals : $v_i^{(t)}(x_i) \propto \prod_{a \in \Lambda(i)} \hat{v}_{a \rightarrow i}^{(t-1)}(x_i)$

→ Min-sum : $(\Sigma, \tau) \mapsto (\min, \Sigma)$

$$\max_{\Sigma} \prod_a e^{-E_a(x_{n(a)})} \Leftrightarrow \min_{\Sigma} \sum_a E_a(x_{n(a)})$$

Similar updates, defined up to an additive constant
(replaces normalization)

[2] Learning in latent-variable models

Latent-variable models: ex: $\begin{matrix} \textcircled{z} \\ \downarrow \\ \textcircled{x} \end{matrix}$ (unobs.)
(observed)

Goal: estimate parameters θ $p(x, z; \theta)$ when z is not observed

Maximum likelihood estimator :

$$\max_{\theta} \sum_i \log p(x_i; \theta)$$

$$= \max_{\theta} \sum_i \log \int p(x_i, z; \theta) dz$$

depending on discrete or continuous z .

→ this problem can be highly non-convex!

Q: How can we solve this?

■ K-means

→ clustering with discrete assignments

- $x_i \in \mathbb{R}^P, i=1, \dots, m$: data
- $\mu_h \in \mathbb{R}^P, h=1, \dots, K$: cluster centers (parameters)
- $z_i \in \{1, \dots, K\}$: cluster assignments

→ Objective :

$$\min_{\mu, z} \left\{ L(\mu, z) := \sum_{i=1}^m \|x_i - \mu_{z_i}\|^2 \right\}$$

Algorithm

Repeat until convergence

$$\rightarrow z_i = \arg \min_h \|x_i - \mu_h\|^2$$

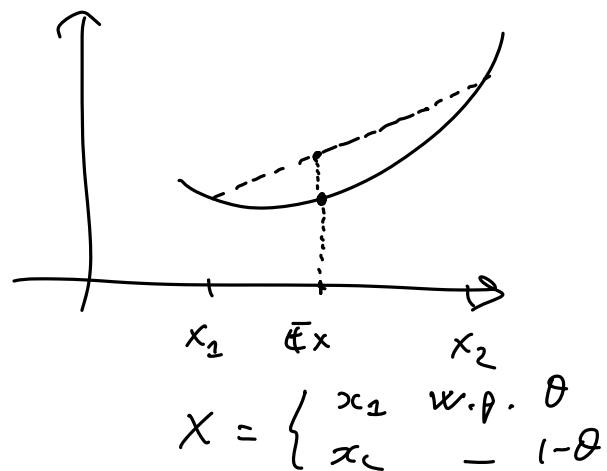
$$\rightarrow \mu_h = \frac{1}{|\{i : z_i = h\}|} \sum_{i : z_i = h} x_i$$

→ alternates between minimization of $L(\mu, z)$ w.r.t.
 z and μ -

■ EM algorithm

- Recall Jensen's inequality :

→ if φ is convex, then $\varphi(\bar{x}) \leq \mathbb{E} \varphi(x)$



→ we have equality ($\varphi(\bar{x}) = \mathbb{E} \varphi(x)$)
if and only if $X = x_0$ almost surely (constant)

- “Variational” lower bound :

$$\log p(x; \theta) = \log \sum_z p(x, z; \theta)$$

$$= \log \sum_z q(z) \frac{p(x, z; \theta)}{q(z)} \quad (\text{for some } q(z) \geq 0 \text{ and } \sum_z q(z) = 1)$$

$$= \log \mathbb{E}_q \left[\frac{p(x, z; \theta)}{q(z)} \right]$$

$$\stackrel{(*)}{\geq} \mathbb{E}_q \left[\log \frac{p(x, z; \theta)}{q(z)} \right]$$

(Jensen's, since \log is concave)

We thus have:

$$\log p(x; \theta) \geq \mathbb{E}_q [\log p(x, z; \theta)] - \mathbb{E}_q [\log q(z)]$$

$=: L(\theta, q)$

→ Note that we have equality in (*) when $q(z) \propto p(z|x; \theta)$
i.e. $q(z) = p(z|x; \theta)$

→ E-M algorithm (Expectation-Maximization)

For $t > 0$:

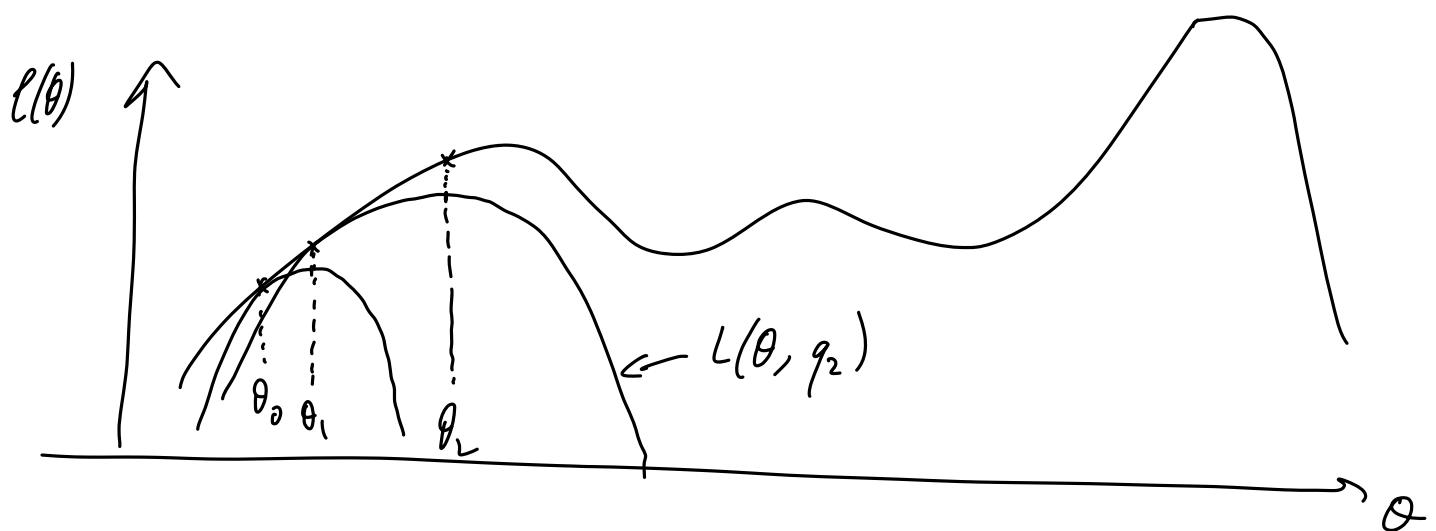
$$q_{t+1} \in \arg \max_q L(\theta_t, q) \quad \left[= p(z|x; \theta_t) \right]$$

(E-step)

$$\theta_{t+1} \in \arg \max_{\theta} L(\theta, q_{t+1}) \quad \left[= \arg \max_{\theta} \mathbb{E}_{q_{t+1}} [\log p(x, z; \theta)] \right]$$

(M-step)

"complete"
log-likelihood



Remarks:

- corresponds to alternating maximizations on $L(\theta, q)$
- $\ell(\theta) = \max_q L(\theta, q) = L(\theta, p(\cdot | x; \theta))$

$\ell(\theta_t)$ is non-decreasing $\Rightarrow \theta_t$ converges to stationary point (not necessarily global maximum)

- First example of a variational method (more next week)
- The gap between $\log p(x; \theta)$ and $L(\theta, q)$ corresponds to a Kullback-Leibler divergence (relative entropy)

$$\Delta_{KL}(p, q) = \bar{E}_p \left[\log \frac{p(x)}{q(x)} \right]$$

$$\ell(\theta) - L(\theta, q) = \log p(x; \theta) - \bar{E}_q \left[\log \frac{p(x; \theta) p(z|x; \theta)}{q(z)} \right]$$

$$= \bar{E}_q \left[\log \frac{q(z)}{p(z|x; \theta)} \right]$$

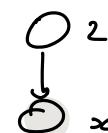
$$= \Delta_{KL}(q, p(\cdot | x; \theta))$$

Example: Gaussian Mixtures

- $z \sim \pi$
- $x|z=k \sim \mathcal{N}(\mu_k, \Sigma_k)$

$$\theta = (\pi, \{\mu_k\}, \{\Sigma_k\})$$

$$\sum_k \pi_k = 1$$



$\rightarrow n$ points x_i , corresponding latents: z_i

$\rightarrow \underline{E\text{-step}}:$

$$\underbrace{p(z_i=k|x_i; \theta)}_{=: \tau_{ik}} \propto \mathcal{N}(x_i|\mu_k, \Sigma_k) \cdot \pi_k$$

$\rightarrow \underline{M\text{-step}}:$

- complete log-likelihood:

$$\sum_i \log p(x_i, z_i; \theta) = \sum_{i=1}^m \sum_{k=1}^K \mathbb{I}\{z_i=k\} \log \pi_k$$

$$+ \sum_{i=1}^m \sum_{k=1}^K \mathbb{I}\{z_i=k\} \log \mathcal{N}(x_i|\mu_k, \Sigma_k)$$

$$\mathbb{E} \left[\sum_i \log p(x_i, z_i; \theta) | x_i; \theta_t \right] = \sum_{i=1}^m \sum_{k=1}^K p(z_i=k | x_i; \theta_t) \overline{\tau_{ik}^t} \log \pi_k$$

$$+ \sum_{i=1}^m \sum_{k=1}^K p(z_i=k | x_i; \theta_t) \log \mathcal{N}(x_i|\mu_k, \Sigma_k)$$

$$\Rightarrow \begin{cases} \pi_k^{t+1} \propto \sum_{i=1}^m \tau_{ik}^t \\ \mu_k^{t+1} = \frac{\sum_i \tau_{ik}^t x_i}{\sum_i \tau_{ik}^t} \\ \Sigma_k^{t+1} = \frac{\sum_i \tau_{ik}^t (x_i - \mu_k^{t+1})(x_i - \mu_k^{t+1})^\top}{\sum_i \tau_{ik}^t} \end{cases}$$