

2/4/13

Review

- * Learning in BNs

- * ML estimation from complete data

examples $\{(x_1^{(t)}, x_2^{(t)}), \dots, x_n^{(t)}\}_{t=1}^T$

$$P_{ML}(x_i = x | p_{ai} = r) = \text{count}(x_i = x, p_{ai} = r)$$

$$\text{count}(p_{ai} = r)$$

$$= \sum_{t=1}^T I(x^{(t)}, x) I(p_{ai}^{(t)}, r)$$

I: indicator functions

test for equality

t	x_1	x_2	\dots	x_n
1	1	2	5	3
2	2	?	0	3
3	?	1	1	8
4	?	?	5	5
T	2	6	3	9

? - hidden values

- * ML estimation from incomplete data

Examples $t = 1, 2, \dots, T$

Hidden nodes $H^{(t)}$

Visible nodes $V^{(t)}$

choose CPTs to maximize log-likelihood :

$$\ell = \sum_{t=1}^T \log P(V = V^{(t)})$$

$$= \sum_t \log \sum_h P(V = V^{(t)}, H = h)$$

$$= \sum_t \log \sum_h \prod_{i=1}^n P(x_i = x_i | p_{ai} = r) \Big|_{V = V^{(t)}, H = h}$$

} How to
maximize?

- * EM algorithm

Iterative procedure to maximize $\sum_t \log P(V = V^{(t)})$

for incomplete data in terms of CPTs of BN.

- * Pseudo code

1) Initialize all CPTs with (possibly) random values.

2) Do until log-likelihood steps increasing :

(A) E-step : compute posterior probabilities

$$P(x_i = x_i, p_{ai} = r | V = V^{(t)}) \text{ for } t = 1, \dots, T$$

(B) M-step : update CPTs (non-inference algorithm)

$$P(x_i = x_i | p_{ai} = r) \leftarrow \frac{\sum_t P(x_i = x_i, p_{ai} = r | V^{(t)})}{\sum_t P(p_{ai} = r | V^{(t)})}$$

Intuitions:

- expected statistics under $P(H|V^{(t)})$ are filling in "missing values" of incomplete data.
- expected counts are substituting for observed counts in complete data case.

Iterate E&M steps until convergence. Why iterate?

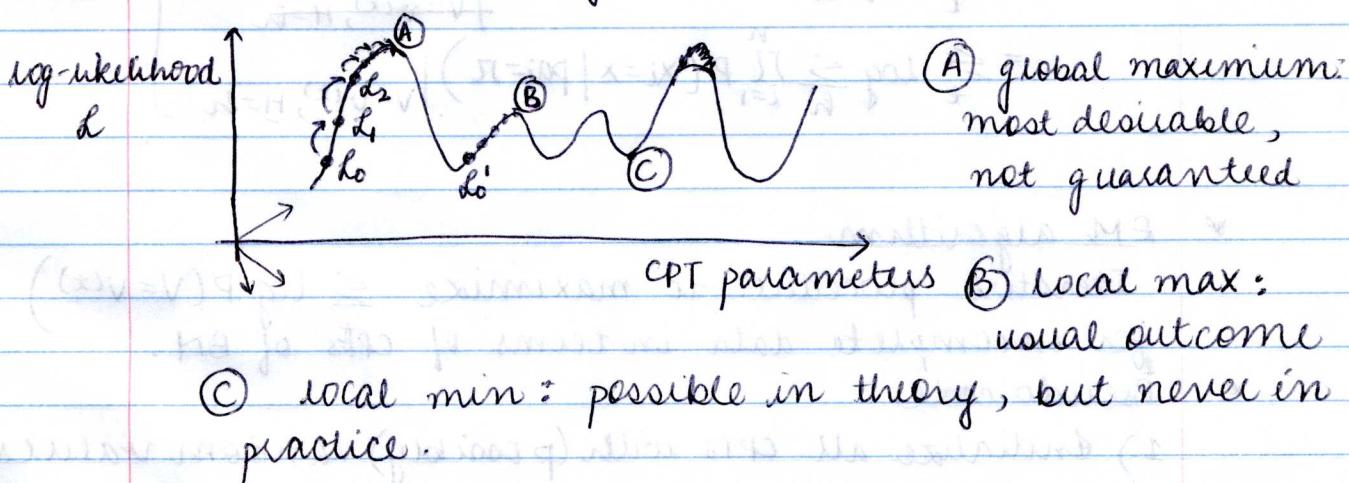
- RHS depends on current CPTs.

Key properties

- monotonic convergence
Each iteration of EM improves log-likelihood
- $L = \sum_t \log P(V^{(t)})$

If L_k is log-likelihood at k^{th} iteration, then
 $L_k \geq L_{k-1}$.

converges to stationary point of log-likelihood
where gradient vanishes



- Why is EM popular?
- No tuning parameters: no step sizes, learning rates, no backtracking,

Example $(A) \rightarrow (B) \rightarrow (C)$ A and C are observed (visible)
B is hidden.

* Posterior probability

$$P(B=b | A=a, C=c) = P(C=c | B=b, A=a) P(B=b | A=a)$$

$$\leq_{b'} P(C=c | B=b', A=a) P(B=b' | A=a)$$

Bayes rule

$= P(C=c | B=b) P(B=b | A=a)$ conditional

$\sum_{b'} P(C=c | B=b') P(B=b' | A=a)$ independence

Shorthand: $P(b/a, c) = P(B=b | A=a, C=c)$

* Incomplete data set $\{(a_t, c_t)\}_{t=1}^T$

t	A	B	C
1	a ₁	?	c ₁
2	a ₂	?	c ₂
T	a _T	?	c _T

log-likelihood:

$$L = \sum_t \log P(A=a_t, C=c_t)$$

$$= \sum_t \log P(A=a_t, C=c_t)$$

$$= \sum_t \log \sum_b P(A=a_t, B=b, C=c_t) \text{ marginalization}$$

$$= \sum_t \log \left\{ \sum_b P(a_t) P(b|a_t) P(c_t|b) \right\} \text{ product rule and CI}$$

General EM algorithm:

$$P(x_i=x | p_{ai}=\pi) \leftarrow \frac{\sum_t P(x_i=x, p_{ai}=\pi | V^{(t)})}{\sum_t P(p_{ai}=\pi | V^{(t)})}$$

Now apply to this example...

$$M\text{-step: } P(B=b | A=a) \leftarrow \sum_t P(A=a, B=b | A=a_t, C=c_t) / \sum_t P(A=a | A=a_t, C=c_t)$$

$$\text{Simplify: } P(b/a) \leftarrow \sum_t I(a, a_t) P(b/a_t, c_t)$$

We showed how to compute these from CPTs using Bayes rule.

$$P(C=c | B=b) \leftarrow \sum_t P(B=b, C=c | A=a_t, C=c_t) / \sum_t P(B=b | A=a_t, C=c_t)$$

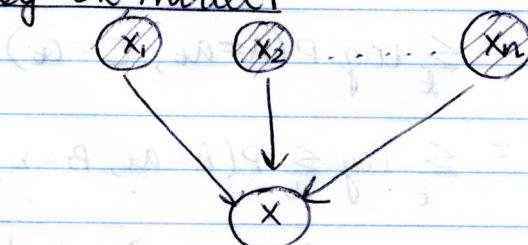
$$\text{Simplify: } P(c/b) \leftarrow \sum_t I(c, c_t) P(b/a_t, c_t) / \sum_t P(b/a_t, c_t)$$

computed in terms of CPTs using Bayes rule.

Aside:

$$\begin{aligned} & P(B=b, C=c | A=a_t, C=c_t) \\ &= P(C=c | A=a_t, C=c_t) P(B=b | A=a, C=c_t, C=c) \\ &= I(c, c_t) P(B=b | A=a_t, C=c_t) \end{aligned}$$

Noisy-OR model



disease $x_i \in \{0, 1\}$
symptom $y \in \{0, 1\}$

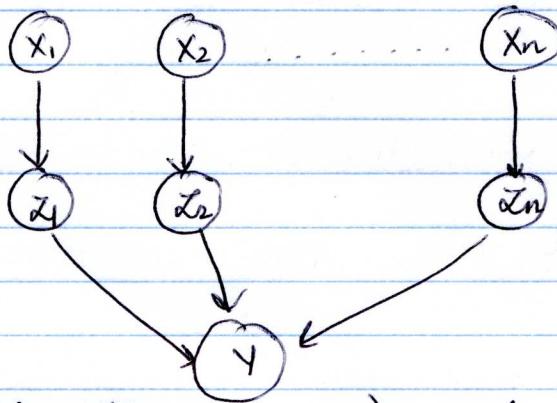
$$P(y=1 | x_1, x_2, \dots, x_n) = 1 - \prod_{i=1}^n (1-p_i)^{x_i} \text{ where } p_i \in [0, 1]$$

* for complete data $\{(x_t, y_t)\}_{t=1}^T$ how to estimate $p_i \in [0, 1]$?

Note: NOISY-OR is a "parametric" model of CPT.

No simple, "closed-form" MLE estimate for $p_i \in [0, 1]$

* Alternative formulation



$$P(Y=1 | z_1, z_2, \dots, z_n) = \text{OR}(z_1, z_2, \dots, z_n) \text{ deterministic}$$

$P(z_i=1 | x_i=0) = 0$ copy x_i if $x_i=0$

$P(z_i=1 | x_i=1) = p_i$ flip x_i with prob $1-p$
if $x_i=1$

Are these models equivalent?