

# Review

- Learning in BNs
- ML estimation from complete data

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

$$P_{ML}(X_i = x | p_{q_i} = \pi) = \frac{\text{count}(X_i = x, p_{q_i} = \pi)}{\text{count}(p_{q_i} = \pi)}$$

$$= \frac{\sum_{t=1}^T I(x^{(t)}, x) I(p_{q_i}^{(t)}, \pi)}{\sum_{t=1}^T I(p_{q_i}^{(t)}, \pi)}$$

$I$ : indicator function  
tests for equality

- ML estimation from incomplete data

examples  $t=1, 2, \dots, T$   
Hidden nodes  $H^{(t)}$   
Visible nodes  $V^{(t)}$

$t$	$X_1$	$X_2$	...	$X_n$
1	1	2	5	3
2	2	?	0	3
3	?	1	1	8
⋮	?	?	5	5
$T$	2	6	3	9

? : hidden values

Choose CPTs to maximize log-likelihood

$$\mathcal{L} = \sum_{t=1}^T \log P(V = v^{(t)})$$

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

$$= \sum_t \log \sum_h \prod_{i=1}^n P(X_i = x | p_{q_i} = \pi) \Big|_{V = v^{(t)}, H = h}$$

How to maximize?

## • Expectation-Maximization (EM) algorithm

Iterative procedure to maximize  $\sum_t \log P(V=v^{(t)})$   
for incomplete data in terms of CPTs of BN.

### • Algorithm:

- (1) Initialize all CPTs with (possibly) random values
- (2) Do until log-likelihood stops increasing:

E-step: compute posterior probabilities (inference)

$$P(X_i=x, pa_i=\pi \mid V=v^{(t)}) \quad \text{for } t=1, 2, \dots, T$$

M-step: update CPTs

$$P(X_i=x \mid pa_i=\pi) \leftarrow \frac{\sum_t P(X_i=x, pa_i=\pi \mid V=v^{(t)})}{\sum_t P(pa_i=\pi \mid V=v^{(t)})}$$

### Intuition

- expected statistics under  $P(H \mid V^{(t)})$  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.

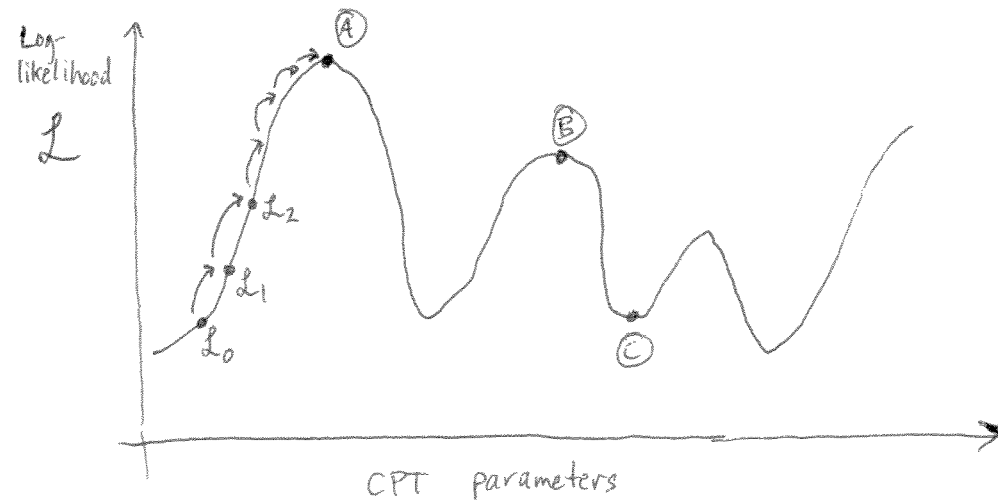
### Key properties

- monotonic convergence

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

If  $\mathcal{L}_k$  is log-likelihood at  $k^{\text{th}}$  iteration, then  $\mathcal{L}_k \geq \mathcal{L}_{k-1}$ .

• converges to stationary point of log-likelihood where gradient vanishes



- (A) global max: most desirable, not guaranteed
- (B) local max: usual outcome
- (C) local min: possible in theory, but never in practice

• no tuning parameters: no step sizes, learning rates, backtracking

Example



A, C are observed (visible), B is hidden

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

Log-likelihood:

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

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

$$= \sum_t \log \left\{ \sum_b P(a_t) P(b|a_t) P(c_t|b) \right\}$$

t	A	B	C
1	a <sub>1</sub>	?	c <sub>1</sub>
2	a <sub>2</sub>	?	c <sub>2</sub>
⋮	⋮	⋮	⋮
T	a <sub>T</sub>	?	c <sub>T</sub>

marginalization

product rule, CI, shorthand

General EM algorithm:

E-step

M-step  $P(X_i=x | p_{\pi}=\pi) \leftarrow \frac{\sum_t P(X_i=x, p_{\pi}=\pi | V^{(t)})}{\sum_t P(p_{\pi}=\pi | V^{(t)})}$

Now apply to this example...

• First: posterior probability  $P(H|V)$

Shorthand:  $P(B=b|A=a_t, C=c_t) \rightarrow P(b|a_t, c_t)$

$$P(b|a_t, c_t) = \frac{P(c_t|a_t, b)P(b|a_t)}{P(c_t|a_t)} \quad \text{Bayes rule}$$

$$= \frac{P(c_t|b)P(b|a_t)}{\sum_{b'} P(b', c_t|a_t)} \quad \begin{array}{l} \text{CI} \\ + \\ \text{marginalization} \end{array}$$

$$= \frac{P(c_t|b)P(b|a_t)}{\sum_{b'} P(c_t|b')P(b'|a_t)} \quad \begin{array}{l} \text{product rule +} \\ \text{CI} \end{array}$$

• Update CPTs  $P(A)$ ,  $P(B|A)$ ,  $P(C|B)$

$$P(A=a) = \frac{\sum_t \text{count}(A=a)}{T}$$

M-step:

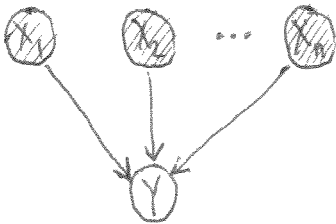
$$P(B=b|A=a) \leftarrow \frac{\sum_t P(A=a, B=b | A=a_t, C=c_t)}{\sum_t P(A=a | A=a_t, C=c_t)}$$

$$\text{Simplify: } P(b|a) \leftarrow \frac{\sum_t I(a, a_t) P(b|a_t, c_t)}{\sum_t I(a, a_t)} \quad \leftarrow \text{posterior computed above, in terms of CPTs}$$

$$P(C=c|B=b) \leftarrow \frac{\sum_{t=1}^I P(B=b, C=c | A=a_t, C=c_t)}{\sum_{t=1}^I P(B=b | A=a_t, C=c_t)}$$

Simplify:  $P(c|b) \leftarrow \frac{\sum_t I(c, c_t) P(b|a_t, c_t)}{\sum_t P(b|a_t, c_t)} \rightarrow$  posterior computed above

### Noisy-OR model



e.g. disease  $X_i \in \{0, 1\}$   
 symptom  $Y \in \{0, 1\}$

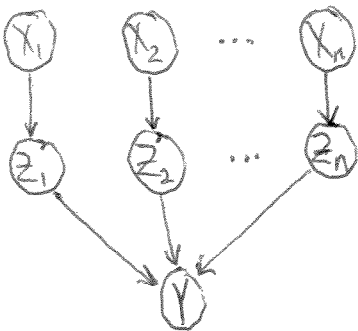
$$P(Y=1 | X_1, \dots, X_n) = 1 - \prod_{i=1}^n (1 - p_i)^{X_i} \quad \text{where } p_i \in [0, 1]$$

- From complete data  $\{(\vec{x}_t, y_t)\}_{t=1}^T$ , how do we estimate  $p_i \in [0, 1]$ ?

Noisy-OR is a "parametric" model of CPT ( $n$  values parameterize  $2^n$  CPT entries)

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

- Alternative formulation (motivation: use EM to estimate  $p_i$ )



$$P(Y | Z_1, \dots, Z_n) = \text{OR}(Z_1, \dots, Z_n)$$

logical-OR (deterministic)

$$P(Z_i=1 | X_i=1) = p_i$$

$$P(Z_i=1 | X_i=0) = 0$$

Equivalently: 
$$P(Z_i=0 | X_i) = (1 - p_i)^{X_i} = \begin{cases} 1 - p_i, & X_i=1 \\ 1, & X_i=0 \end{cases}$$

What is  $P(Y=1|\vec{X})$  in this new model?

$$P(Y=1|\vec{X}) = \sum_{\vec{Z} \in \{0,1\}^n} P(Y=1, \vec{Z}|\vec{X})$$

marginalization

$$= \sum_{\vec{Z}} P(Y=1|\vec{Z}, \vec{X}) P(\vec{Z}|\vec{X})$$

product rule

$$= \sum_{\vec{Z}} P(Y=1|\vec{Z}) P(\vec{Z}|\vec{X})$$

conditional independence

$$= \sum_{\vec{Z} \neq \vec{0}} P(\vec{Z}|\vec{X})$$

because  $Y = \text{OR}(\vec{Z})$

$$= 1 - P(\vec{Z} = \vec{0}|\vec{X})$$

from normalization

$$= 1 - \prod_{i=1}^n P(Z_i=0|X_i)$$

conditional independence

$$= 1 - \prod_{i=1}^n (1-p_i)^{X_i}$$

Same as original Noisy-OR BN!

• Posterior probability

1 if  $Y=1$   
0 if  $Y=0$

$p_i$  if  $X_i=1$   
0 if  $X_i=0$

$$P(Z_i=1|\vec{X}, Y) = \frac{P(Y|\vec{X}, Z_i=1) P(Z_i=1|\vec{X})}{P(Y|\vec{X})}$$

$P(Y|\vec{X})$

$\leftarrow 1 - \prod_{i=1}^n (1-p_i)^{X_i}$  if  $Y=1$

$$= \frac{Y p_i X_i}{1 - \prod_{i=1}^n (1-p_i)^{X_i}}$$