

7.12

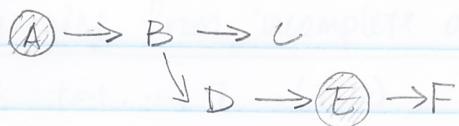
Review - Inference

- * Tractable in polytrees — singly connected networks ; no loops
- * Key tools:

Bayes Rule
Marginalization
Product Rule
Conditional Independence

} (and conditionalized version of these)

Example



How to compute $P(B|E, A)$?

$$P(B|E, A) = \frac{P(E|B, A) P(B|A)}{P(E|A)}$$

Bayes rule: use it to rewrite RHS in terms of probabilities where nodes are conditioned on ancestors

$$P(E|B, A) = \sum_d P(E, D=d | B, A)$$

use marginalization to introduce parents b/c CPTs always express $P(\text{child} | \text{parents})$

$$P(E, D=d | B) = P(D=d | B) P(E | D=d, B)$$

conditional product rule gets you even closer to CPTs, independence

b/c CPTs predict one node at time.

Review - Learning from complete data

* examples $t=1, 2, \dots, T$

variables X_1, X_2, \dots, X_n

data set $\{(X_1^t, X_2^t, \dots, X_n^t)\}_{t=1}^T$ (IID)

* choose CPTs to maximize log-likelihood $L = \sum_t \log P(X_1^t, X_2^t, \dots, X_n^t)$

* ML estimates

$$P_{\text{ML}}(X_i = x | \text{pa}_i = \pi) = \frac{\text{count}(X_i = x, \text{pa}_i = \pi)}{\text{count}(\text{pa}_i = \pi)} = \frac{\sum_{t=1}^T I(X_i^t = x) I(\text{pa}_i^t = \pi)}{\sum_{t=1}^T I(\text{pa}_i^t = \pi)}$$

Review - Learning from incomplete data

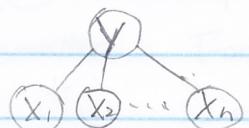
* examples $t=1, \dots, T$ (IID)

hidden nodes $H^{(t)}$

visible nodes $V^{(t)}$

Why useful?

Ex: document clustering



$Y \in \{1, 2, \dots, K\}$ document topic $H^{(t)} = \{Y\}$

$X \in \{0, 1\}$ Does i^{th} word in dictionary

appear in document? $V^{(t)} = \{x_1, x_2, \dots, x_n\}$

Question: can we learn such a model from a data set of unlabeled documents?

Data:

t	x_1	x_2	\dots	x_n	Y
1	0	1	\dots	1	?
2	0	0	\dots	1	?
3	0	1	\dots	0	?
\vdots					
T	1	0	\dots	0	?

* Choose CPTs to maximize log-likelihood

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

$$= \sum_{t=1}^T \log \sum_h P(V^{(t)}, H=h)$$

$$= \sum_{t=1}^T \log \underbrace{\sum_h \prod_{i=1}^n P(X_i=x | \text{pa}_i=\pi)}_{\text{Much more difficult to compute and optimize}} \quad \begin{array}{l} \text{evaluate there} \\ V^{(t)}, H=h \end{array}$$

No "closed-form" solution for CPTs.

Alternative = seek an iterative solution.

Iterative solution is known as EM algorithm:

* EM algorithm:

i) Initialization

Choose values for elements of CPTs in BN

(should be valid probabilities, could be poor guesses)

2) Iterative steps — Repeat until convergence.

(E) E-step — Inference

expectation For each example $t=1, 2, \dots, T$ compute posterior distribution $P(H=h | V^{(t)})$

In particular, for each node, compute

$$P(X_i=x, \text{pa}_i=\pi | V=V^{(t)})$$

(M) M-step = update CPTs to increase log-likelihood.

$$\text{maximization } P(X_i=x | \text{pa}_i=\pi) \leftarrow \frac{\sum_{t=1}^T P(X_i=x, \text{pa}_i=\pi | V^{(t)})}{\sum_{x_i} \left[\sum_{t=1}^T P(X_i=x, \text{pa}_i=\pi | V^{(t)}) \right]}$$

Simplify =

① nodes w/ parents:

$$P(X_i=x | \text{pa}_i=\pi) \leftarrow \frac{\sum_{t=1}^T P(X_i=x, \text{pa}_i=\pi | V^{(t)})}{\sum_{\pi} P(\text{pa}_i=\pi | V^{(t)})}$$

② root node (no parents)

$$P(X_i=x) \leftarrow \frac{1}{T} \sum_{t=1}^T P(X_i=x | V^{(t)})$$

Note:

- RHS of these update depends on current values of CPTs
- Inference is key subroutine of learning.

Intuition:

- Use posterior distribution $P(H | V^{(t)})$ to "fill in" missing values of hidden nodes
- expected statistics of posterior distribution $P(H | V^{(t)})$ substitute for observed counts that we had in complete data case
- Compare to ML estimate in complete data case

$$P_{ML}(X_i = x | \text{par}_i = \pi) = \frac{\text{count}(X_i = x, \text{par}_i = \pi)}{\text{count}(\text{par}_i = \pi)} \\ = \frac{\sum_{t=1}^T I(X_i^{(t)}, x) I(\text{par}_i^{(t)}, \pi)}{\sum_{t=1}^T I(\text{par}_i^{(t)}, \pi)}$$

} special case of EM
when all data is observed.

Key properties of Expectation Maximization (EM) algorithm

1) Monotonic convergence

Each iteration of EM improves log-likelihood

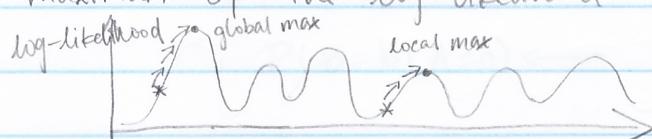
$$\mathcal{L} = \sum_{t=1}^T \log P(V^{(t)})$$

If \mathcal{L}_k is log-likelihood at k^{th} iteration, then

$$\mathcal{L}_k \geq \mathcal{L}_{k-1} \quad \swarrow \quad \text{Really helpful debugging diagnostic}$$

2) Converges in general to local (but not global)

maximum of the log-likelihood $\mathcal{L} = \sum_t \log P(V^{(t)})$



Final output of EM depends on model initialization

3) No tuning parameters, no learning rates, no back-tracking

Example



A, C are observed (visible) nodes

B is hidden node

- First, consider inference in this BN.

$$\begin{aligned} \text{Posterior: } P(B|A, c) &= \frac{P(c|B, A) P(B|A)}{P(c|A)} \quad \text{Bayes Rule / conditionality} \\ &= \frac{P(c|B) P(B|A)}{\sum_b P(c|B=b) P(B=b|A)} \quad \text{Independence / normalization} \end{aligned}$$

Bottom line: We can compute $P(b|a, c)$ in terms of model's CPTs.

* How do we learn from incomplete data $\{(a_t, c_t)\}_{t=1}^T$

- EM algorithm

$$\begin{aligned} \text{Log-likelihood } L &= \sum_t \log P(A=a_t, C=c_t) \\ &= \sum_t \log \sum_b P(A=a_t, B=b, C=c_t) \\ &= \sum_t \log \sum_b P(a_t) P(b|a_t) P(c_t|a_t, b) \end{aligned}$$

Reminder: general step of EM algorithm

$$P(X_i=x | p_{ai}=\pi) \leftarrow \frac{\sum_{t=1}^T P(X_i=x, p_{ai}=\pi | V^{(t)})}{\sum_{t=1}^T P(p_{ai}=\pi | V^{(t)})}$$

Now specialize to our example $V^{(t)} = (a_t, c_t)$

- CPT at node B:

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

Simplify:

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

↓ We've computed already
in terms of CPTs

↑ Read off
from data

- CPT at node C

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

Simplify:

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

Example

Noisy-OR Model

$x_1 \ x_2 \ \dots \ x_n$ diseases $x_i \in \{0, 1\}$

$\downarrow \downarrow \downarrow \downarrow$ symptom $y \in \{0, 1\}$

$$P(Y=1 | x_1, x_2, \dots, x_n) = 1 - \prod_{i=1}^n (1-p_i)^{x_i}$$

with $p_i \in [0, 1]$

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

Note: Noisy-OR is a "parametric" model of CPT; there is no simple, closed-form ML estimate for $p_i \in [0, 1]$ that maximizes.

$$\sum_t \log P(Y=y^{(t)} | x_1=x_1^{(t)}, \dots, x_n=x_n^{(t)})$$

* Alternative formulation =

$x_1 \ x_2 \ \dots \ x_n \quad x_i \in \{0, 1\}$ noisy copy of parent

$\downarrow \downarrow \downarrow \downarrow$

$$z_1 \ z_2 \ \dots \ z_n \quad z_i \in \{0, 1\}$$

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

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

$$P(z_i=1 | x_i=1) = p_i$$

$\downarrow \downarrow \downarrow \downarrow$

$y \in \{0, 1\}$

$$Y = OR(z_1, z_2, \dots, z_n)$$

$$P(Y=1 | z_1, \dots, z_n) = \begin{cases} 1 & \text{if } \sum_{i=1}^n z_i > 0 \\ 0 & \text{if } z_i = 0 \text{ for all } i. \end{cases}$$

$$P(z_i=0 | x_i=0) = 1, \quad P(z_i=1 | x_i=0) = 0$$

$$P(z_i=1 | x_i=1) = p_i$$

Equivalently:

$$P(z_i=0 | x_i) = (1-p_i)^{x_i} = \begin{cases} 1-p_i & \text{if } x_i=1 \\ 1 & \text{if } x_i=0 \end{cases}$$

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

$$\begin{aligned} P(Y=1 | \vec{x}) &= \sum_{\vec{z} \in \{0,1\}^n} P(Y=1, \vec{z} | \vec{x}) \quad \text{marginalization} \\ &= \sum_{\vec{z}} P(Y=1 | \vec{z}, \vec{x}) P(\vec{z} | \vec{x}) \quad \text{product rule} \\ &= \sum_{\substack{\vec{z} \in \{0,1\}^n \\ \vec{z} \neq (0,0,0,\dots,0)}} P(Y=1 | \vec{z}) P(\vec{z} | \vec{x}) \quad \begin{matrix} \text{cond.} \\ \text{independence} \end{matrix} \\ &= \sum_{\vec{z} \neq (0,0,0,\dots,0)} P(\vec{z} | \vec{x}) + P(Y=1 | \vec{z}=\vec{0}) P(\vec{z}=\vec{0} | \vec{x}) \\ &\quad \text{logical-OR, } Y = \text{OR}(\vec{z}) \\ &= 1 - P(\vec{z}=\vec{0} | \vec{x}) \quad \text{normalization} \\ &= 1 - \prod_{i=1}^n P(z_i=0 | x_i) \quad \text{product rule} \\ &= 1 - \prod_{i=1}^n (1-p_i)^{x_i} \end{aligned}$$

key insight:

alternative BN has same expression for $P(Y=1 | \vec{x})$!!

→ estimate p_i parameters for noisy-OR CPT

using EM in "expanded" network

First, let's do inference in this model:

$$P(z_i=1 | \vec{x}, y) = \frac{P(Y=1 | z_i=1, \vec{x}) P(z_i=1 | \vec{x})}{P(Y=1 | \vec{x})} \quad \begin{matrix} \{0 \text{ if } x_i=0 \\ p_i \text{ if } x_i=1\} \end{matrix}$$

only care if $y=1$,
b/c otherwise numerator
vanish.

$$= \frac{y \cdot p_i \cdot x_i}{1 - \prod_{i=1}^n (1-p_i)^{x_i}}$$

true for $x_i \in \{0,1\}$, $y \in \{0,1\}$

* (conditional) log-likelihood of data $\{(\vec{x}_t, y_t)\}_{t=1}^T$

$$L = \sum_t \log P(Y_t | \vec{x}_t)$$

$$\begin{aligned} &= \sum_{t=1}^T [y_t \log P(Y=1 | \vec{x}_t) + (1-y_t) \log P(Y=0 | \vec{x}_t)] \\ &= \sum_{t=1}^T [y_t \log \left(1 - \prod_{i=1}^n (1-p_i)^{x_i^{(t)}} \right) + (1-y_t) \log \prod_{i=1}^n (1-p_i)^{x_i^{(t)}}] \end{aligned}$$

$$\mathcal{L}(p_1, p_2, \dots, p_n) = \sum_{t=1}^T y_t \log \left(1 - \prod_{i=1}^n (1-p_i)^{x_{it}} \right) + \sum_{t=1}^T \sum_{i=1}^n (1-y_t) x_{it} \log (1-p_i)$$

y_t, x_{it} are data observed

p_i is parameter — how to choose?

Note:

Complicated nonlinear expression of p_i .

How to optimize? EM to the rescue

Shorthand: let $T = \text{total } \# \text{ of examples}$

$$\text{let } T_i = \sum_{t=1}^T x_{it} = \text{count}(x_i = 1)$$

EM update rule for $p_i = P(\vec{z}_i=1 | \vec{x}_i=1)$

$$P(\vec{z}_i=1 | \vec{x}_i=1) \leftarrow \frac{\sum_{t=1}^T P(z_i=1, x_i=1 | \vec{x}_t, y_t)}{\sum_{t=1}^T P(x_i=1 | \vec{x}_t, y_t)}$$

Simplify =

$$P(\vec{z}_i=1 | \vec{x}_i=1) \leftarrow$$

$$\frac{\sum_{t=1}^T I(x_{it}, 1) P(\vec{z}_i=1 | \vec{x}_t, y_t)}{\sum_{t=1}^T I(x_{it}, 1)} \quad \begin{array}{l} \text{have already from} \\ \text{Bayes rule} \end{array}$$

EM update
for Noisy OR

$$p_i \leftarrow \frac{p_i}{T_i} \sum_{t=1}^T \frac{y_t x_{it}}{1 - \prod_{i=1}^n (1-p_i)^{x_{it}}}$$

$$\frac{y_t \cdot p_i \cdot x_{it}}{1 - \prod_{i=1}^n (1-p_i)^{x_{it}}}$$

This update rule, applied in parallel to all $\{p_i\}_{i=1}^n$, will monotonically increase $\mathcal{L} = \sum_t \log P(y_t | \vec{x}_t)$

Recall from HW4.

Punigram(w)

Pbigram($w'|w$)

$$P_{\text{mixture}}(w'|w) = (1-\lambda) P_{\text{unigram}}(w') + \lambda P_{\text{bigram}}(w'|w) \text{ where } \lambda \in [0, 1]$$

More generally:

$$P_{\text{mixture}}(w'|w) = (1-\lambda(w)) P_{\text{unigram}}(w') + \lambda(w) P_{\text{bigram}}(w'|w) \text{ where } \lambda(w) \in [0, 1]$$

Next Lecture: This is also a hidden variable model!

We can use EM to estimate λ .