Bayesian Networks III D-separation and EM

Required reading:

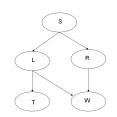
• Bishop online chapter 8, read all of section 8.2

Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

February 23, 2009

Bayesian Networks Definition



Parents	P(W Pa)	P(¬W Pa
L, R	0	1.0
L, ¬R	0	1.0
¬L, R	0.2	0.8
¬L, ¬R	0.9	0.1

A Bayes network represents the joint probability distribution over a collection of random variables

A Bayes network is a directed acyclic graph and a set of CPD's

- Each node denotes a random variable
- Edges denote dependencies
- CPD for each node X_i defines P(X_i / Pa(X_i))
- The joint distribution over all variables is defined as

$$P(X_1 ... X_n) = \prod_i P(X_i | Pa(X_i))$$

Pa(X) = immediate parents of X in the graph

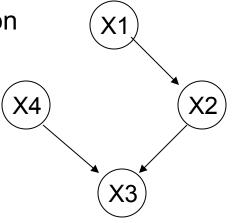
Inference in Bayes Nets

- In general, intractable (NP-complete)
- For certain cases, tractable
 - Assigning probability to fully observed set of variables
 - Or if just one variable unobserved
 - Or for singly connected graphs (ie., no undirected loops)
 - Variable elimination
 - Belief propagation
- For multiply connected graphs
 - Junction tree
 - Loopy belief propagation
- Sometimes use Monte Carlo methods
 - Generate many samples according to the Bayes Net distribution, then count up the results
- Variational methods for tractable approximate solutions

Conditional Independence, Revisited

- We said:
 - Each node is conditionally independent of its non-descendents, given its immediate parents.
- Does this rule give us <u>all</u> of the conditional independence relations implied by the Bayes network?
 - No!
 - E.g., X1 and X4 are conditionally indep given {X2, X3}
 - But X1 and X4 not conditionally indep given X3





Easy Network 1: Head to Tail

prove A cond indep of B given C?

ie., p(a,b|c) = p(a|c) p(b|c)

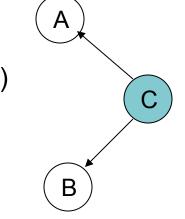
$$P(ablc) = \frac{P(abc)}{P(c)} = \frac{P(a) P(cla) P(blc)}{P(cc)}$$

$$P(ablc) = \frac{P(ablc)}{P(ablc)} = \frac{P(ablc)}{P(blc)}$$

let's use p(a,b) as shorthand for p(A=a, B=b)

Easy Network 2: Tail to Tail

prove A cond indep of B given C? ie., p(a,b|c) = p(a|c) p(b|c)



Easy Network 3: Head to Head

prove A cond indep of B given C? ie., p(a,b|c) = p(a|c) p(b|c)

but
$$p(ab) = p(a)p(b)$$

 $P(a,b) = \begin{cases} p(a,b,c_i) \\ = \\ p(a)p(b) \end{cases} p(c_i|ab)$
 $p(ab) = p(a)p(b) \begin{cases} p(c_i|ab) \\ p(c_i|ab) \end{cases} = 1$

let's use p(a,b) as shorthand for p(A=a, B=b)

Easy Network 3: Head to Head

prove A cond indep of B given C? NO!

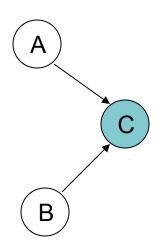
Summary:

- p(a,b)=p(a)p(b)
- p(a,b|c) NotEqual p(a|c)p(b|c)

Explaining away.

e.g.,

- A=earthquake
- B=breakIn
- C=motionAlarm



X and Y are conditionally independent given Z, if and only if X and Y are D-separated by Z.

[Bishop, 8.2.2]

Suppose we have three sets of random variables: X, Y and Z

X and Y are **D-separated** by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is **blocked**

A path from variable A to variable B is **blocked** if it includes a node such that either

- arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z
- 2. the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

X and Y are <u>**D-separated**</u> by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is **blocked**

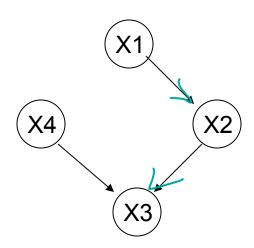
A path from variable A to variable B is **blocked** if it includes a node such that either

- 1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z
- the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

X1 indep of X3 given X2? Yes

X3 indep of X1 given X2? Yes

X4 indep of X1 given X2? Y=5



X and Y are **D-separated** by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is **blocked** by Z

A path from variable A to variable B is **blocked** by Z if it includes a node such that either

1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z

2. the arrows meet head-to-head at the node, and neither the

X4 indep of X1 given X3? No $P(xy \times 11 \times 3 \times 2) = P(xy \times 11 \times 3 \times 2) =$ X4 indep of X1 given {}?

X and Y are <u>**D-separated**</u> by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is **blocked**

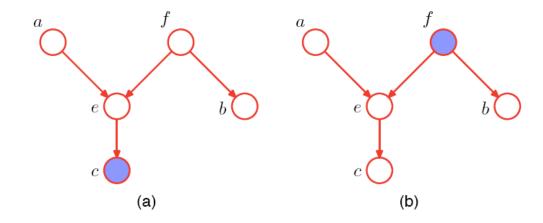
A path from variable A to variable B is **blocked** if it includes a node such that either

- 1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z
- 2. the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

a indep of b given c?

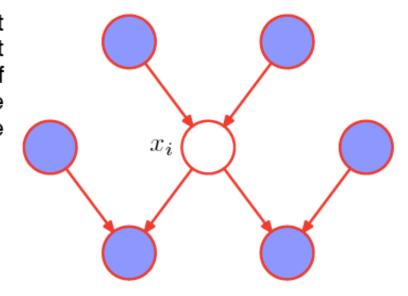
a indep of b given {}?

a indep of b given f?



Markov Blanket

The Markov blanket of a node x_i comprises the set of parents, children and co-parents of the node. It has the property that the conditional distribution of x_i , conditioned on all the remaining variables in the graph, is dependent only on the variables in the Markov blanket.



Java Bayes Net Applet

http://www.pmr.poli.usp.br/ltd/Software/javabayes/Home/applet.html

by Fabio Gagliardi Cozman

What You Should Know

- Bayes nets are convenient representation for encoding dependencies / conditional independence
- BN = Graph plus parameters of CPD's
 - Defines joint distribution over variables
 - Can calculate everything else from that
 - Though inference may be intractable
- Reading conditional independence relations from the graph
 - Each node is cond indep of non-descendents, given only its parents
 - D-separation
 - 'Explaining away'

Learning in Bayes Nets

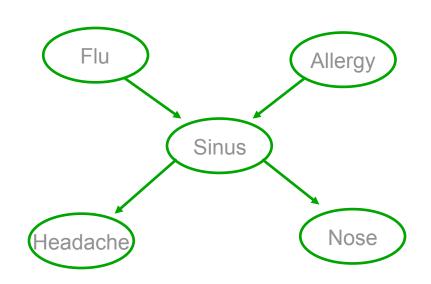
- Four categories of learning problems
 - Graph structure may be known/unknown
 - Variable values may be observed/unobserved
- Easy case: known graph structure, training data is fully observed, learn CPD parameters
- Interesting case: known graph structure, training data is only partly observed, learn CPD parameters

Learning CPTs from Fully Observed Data

 Example: Consider learning the parameter

$$\theta_{s|ij} \equiv P(S=1|F=i,A=j)$$

 MLE (Max Likelihood Estimate) is



$$\theta_{s|ij} = \frac{\sum_{k=1}^K \delta(f_k=i, a_k=j, s_k=1)}{\sum_{k=1}^K \delta(f_k=i, a_k=j)}$$
 kth training

Remember why?

example

MLE estimate of $\theta_{s|ij}$ from fully observed data

Maximum likelihood estimate

$$\theta \leftarrow \arg\max_{\theta} \log P(data|\theta)$$

Our case:

Allerav

Nose

$$P(data|\theta) = \prod_{k=1}^{K} P(f_k, a_k, s_k, h_k, n_k)$$

$$P(data|\theta) = \prod_{k=1}^{K} P(f_k)P(a_k)P(s_k|f_ka_k)P(h_k|s_k)P(n_k|s_k)$$

$$\log P(data|\theta) = \sum_{k=1}^{K} \log P(f_k) + \log P(a_k) + \log P(s_k|f_ka_k) + \log P(h_k|s_k) + \log P(n_k|s_k)$$

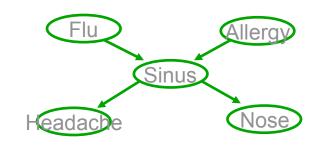
$$\frac{\partial \log P(data|\theta)}{\partial \theta_{s|ij}} = \sum_{k=1}^{K} \frac{\partial \log P(s_k|f_k a_k)}{\partial \theta_{s|ij}}$$

$$\theta_{s|ij} = \frac{\sum_{k=1}^{K} \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^{K} \delta(f_k = i, a_k = j)}$$

Estimate θ from partly observed data

- What if FAHN observed, but not S?
- Can't calculate MLE

$$\theta \leftarrow \arg\max_{\theta} \log \prod_{k} P(f_k, a_k, s_k, h_k, n_k | \theta)$$



- Let X be all observed variable values (over all examples)
- Let Z be all unobserved variable values
- Can't calculate MLE:

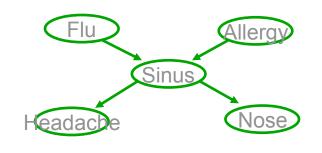
$$\theta \leftarrow \arg\max_{\theta} \log P(X, Z|\theta)$$

• WHAT TO DO? $\theta \leftarrow arsmax los P(x/\theta)$ EM $\Rightarrow arsmax E [los P(x,z|\theta)]$

Estimate θ from partly observed data

- What if FAHN observed, but not S?
- Can't calculate MLE

$$\theta \leftarrow \arg\max_{\theta} \log \prod_{k} P(f_k, a_k, s_k, h_k, n_k | \theta)$$



- Let X be all observed variable values (over all examples)
- Let Z be all unobserved variable values
- Can't calculate MLE:

$$\theta \leftarrow \arg\max_{\theta} \log P(X, Z|\theta)$$

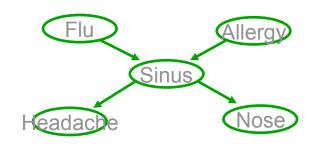
EM seeks* to estimate:

$$\theta \leftarrow \arg\max_{\theta} E_{Z|X,\theta}[\log P(X,Z|\theta)]$$

* EM guaranteed to find local maximum

EM seeks estimate:

$$\theta \leftarrow \arg\max_{\theta} E_{Z|X,\theta}[\log P(X,Z|\theta)]$$



here, observed X={F,A,H,N}, unobserved Z={S}

$$\log P(X, Z | \theta) = \sum_{k=1}^{K} \log P(f_k) + \log P(a_k) + \log P(s_k | f_k a_k) + \log P(h_k | s_k) + \log P(n_k | s_k)$$

$$E_{X|Z,\theta}[logP(X,Z|\theta)] = \sum_{k=1}^{K} \sum_{i=0}^{1} P(s_k = i|f_k, a_k, h_k, n_k)$$
$$[logP(f_k) + log P(a_k) + log P(s_k|f_k a_k) + log P(h_k|s_k) + log P(n_k|s_k)]$$

EM Algorithm

EM is a general procedure for learning from partly observed data

Given observed variables X, unobserved Z (X={F,A,H,N}, Z={S})

Define
$$Q(\theta'|\theta) = E_{P(Z|X,\theta)}[\log P(X,Z|\theta')]$$

Iterate until convergence:

- E Step: Use X and current θ to calculate $P(Z|X,\theta)$
- M Step: Replace current θ by

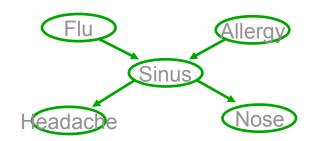
$$\theta \leftarrow \arg \max_{\theta'} Q(\theta'|\theta)$$

Guaranteed to find local maximum.

Each iteration increases $E_{P(Z|X,\theta)}[\log P(X,Z|\theta')]$

E Step: Use X, θ , to Calculate P(Z|X, θ)

observed X={F,A,H,N}, unobserved Z={S}



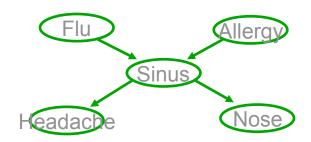
How? Bayes net inference problem.

$$P(S_k = 1 | f_k a_k h_k n_k, \theta) = \frac{P(S_k = 1 | f_k | q_k h_k n_k)}{P(f_k | q_k h_k n_k)}$$

$$P(S_k = 1, f_k | q_k h_k n_k) + P(S_k = 0, f_k | q_k h_k n_k)$$

E Step: Use X, θ , to Calculate P(Z|X, θ)

observed X={F,A,H,N}, unobserved Z={S}



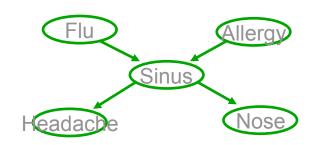
How? Bayes net inference problem.

$$P(S_k = 1 | f_k a_k h_k n_k, \theta) =$$

$$P(S_k = 1 | f_k a_k h_k n_k, \theta) = \frac{P(S_k = 1, f_k a_k h_k n_k | \theta)}{P(S_k = 1, f_k a_k h_k n_k | \theta) + P(S_k = 0, f_k a_k h_k n_k | \theta)}$$

EM and estimating $\theta_{s|ij}$

observed $X = \{F,A,H,N\}$, unobserved $Z=\{S\}$



E step: Calculate for each training example, k

$$P(S_k = 1 | f_k a_k h_k n_k, \theta) = E[s_k] = \frac{P(S_k = 1, f_k a_k h_k n_k | \theta)}{P(S_k = 1, f_k a_k h_k n_k | \theta) + P(S_k = 0, f_k a_k h_k n_k | \theta)}$$

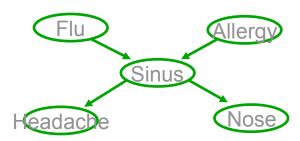
$$P(S_k = 1 | f_k a_k h_k n_k, \theta) = E[s_k] = \frac{P(S_k = 1, f_k a_k h_k n_k | \theta)}{P(S_k = 1, f_k a_k h_k n_k | \theta) + P(S_k = 0, f_k a_k h_k n_k | \theta)}$$

M step: update all relevant parameters. For example:

$$\theta_{s|ij} \leftarrow \frac{\sum_{k=1}^{K} \delta(f_k = i, a_k = j) \ E[s_k]}{\sum_{k=1}^{K} \delta(f_k = i, a_k = j)}$$

Recall MLE was:
$$\theta_{s|ij} = \frac{\sum_{k=1}^{K} \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^{K} \delta(f_k = i, a_k = j)}$$

EM and estimating θ



More generally,

Given observed set X, unobserved set Z of boolean values

E step: Calculate for each training example, k
the expected value of each unobserved variable

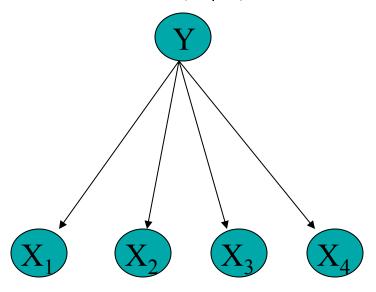
M step:

Calculate estimates similar to MLE, but replacing each count by its expected count

$$\delta(Y=1) \to E_{Z|X,\theta}[Y]$$
 $\delta(Y=0) \to (1 - E_{Z|X,\theta}[Y])$

Using Unlabeled Data to Help Train Naïve Bayes Classifier

Learn P(Y|X)



Υ	X1	X2	Х3	X4
1	0	0	1	1
0	0	1	0	0
0	0	0	1	0
?	0	1	1	0
?	0	1	0	1

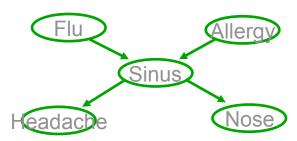
E step: Calculate for each training example, k
the expected value of each unobserved variable

Exp val for Yk Siven Xik Xax ... Xmk

P(Y|X, Xa... Xn) = P(Y) TI P(X; 17)

P(X)

EM and estimating θ



More generally,

Given observed set X, unobserved set Z of boolean values

E step: Calculate for each training example, k
the expected value of each unobserved variable

M step:

Calculate estimates similar to MLE, but replacing each count by its expected count

$$\delta(Y=1) \to E_{Z|X,\theta}[Y]$$
 $\delta(Y=0) \to (1 - E_{Z|X,\theta}[Y])$ $P(\times; | Y=0)$ $P(\times; | Y=0)$