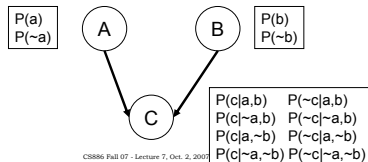


Bayesian Networks

aka belief networks, probabilistic networks

- A BN over variables $\{X_1, X_2, \dots, X_n\}$ consists of:
 - a DAG whose nodes are the variables
 - a set of CPTs $(Pr(X_i | \text{Parents}(X_i)))$ for each X_i

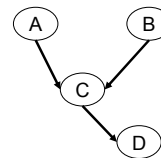


1

Bayesian Networks

aka belief networks, probabilistic networks

- Key notions
 - parents of a node: $\text{Par}(X_i)$
 - children of node
 - descendants of a node
 - ancestors of a node
 - family: set of nodes consisting of X_i and its parents
 - CPTs are defined over families in the BN

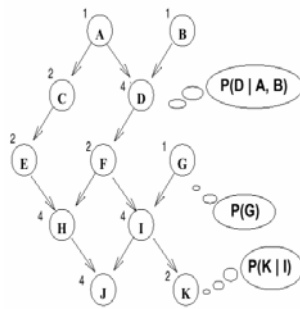


$\text{Parents}(C) = \{A, B\}$
 $\text{Children}(A) = \{C\}$
 $\text{Descendants}(B) = \{C, D\}$
 $\text{Ancestors}(D) = \{A, B, C\}$
 $\text{Family}(C) = \{C, A, B\}$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

2

An Example Bayes Net



- A few CPTs are "shown"
- Explicit joint requires $2^{11} - 1 = 2047$ params
- BN requires only 27 params (the number of entries for each CPT is listed)

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

3

Semantics of a Bayes Net

- The structure of the BN means: every X_i is *conditionally independent of all of its nondescendants given its parents*.

$$Pr(X_i | S \cup \text{Par}(X_i)) = Pr(X_i | \text{Par}(X_i))$$

for any subset $S \subseteq \text{NonDescendants}(X_i)$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

4

Semantics of Bayes Nets

- If we ask for $P(X_1, X_2, \dots, X_n)$ we obtain
 - assuming an ordering consistent with network
- By the chain rule, we have:

$$P(X_1, X_2, \dots, X_n)$$

$$= P(X_n | X_{n-1}, \dots, X_1) P(X_{n-1} | X_{n-2}, \dots, X_1) \dots P(X_1)$$

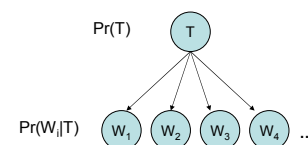
$$= P(X_n | \text{Par}(X_n)) P(X_{n-1} | \text{Par}(X_{n-1})) \dots P(X_1)$$
- Thus, the joint is recoverable using the parameters (CPTs) specified in an arbitrary BN

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

5

Bayes net example I

- Naïve Bayes model
 - Naïve: because words do not depend on each other.
 - Joint = $Pr(T) \prod_i Pr(W_i | T)$

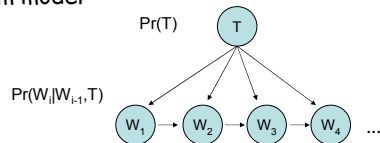


CS886 Fall 07 - Lecture 7, Oct. 2, 2007

6

Bayes net example II

- Tree augmented naïve Bayes classifier
 - Allow arcs between the leaves as long as they form a tree
 - E.g., augment naïve Bayes model with bigram model

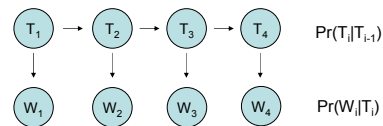


CS886 Fall 07 - Lecture 7, Oct. 2, 2007

7

Bayes net example III

- Hidden Markov model
 - T_i : hidden variable (e.g., tag)
 - W_i : observable variable (e.g., word)



CS886 Fall 07 - Lecture 7, Oct. 2, 2007

8

Maximum Likelihood Learning

- ML learning of Bayes net parameters:
 - For $\theta_{V=\text{true}, \text{pa}(V)=v} = \Pr(V=\text{true} | \text{pa}(V)=v)$
 - $\theta_{V=\text{true}, \text{pa}(V)=v} = \frac{\#[V=\text{true}, \text{pa}(V)=v]}{\#[V=\text{true}, \text{pa}(V)=v] + \#[V=\text{false}, \text{pa}(V)=v]}$
 - Assumes all attributes have values...
- What if values of some attributes are missing?

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

9

Incomplete data

- But many real-world problems have **hidden variables** (a.k.a **latent variables**)
 - Values of some attributes missing
 - Incomplete data \rightarrow unsupervised learning
- Examples:
 - Part of speech tagging
 - Topic modeling
 - Market segmentation for marketing
 - Medical diagnosis

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

10

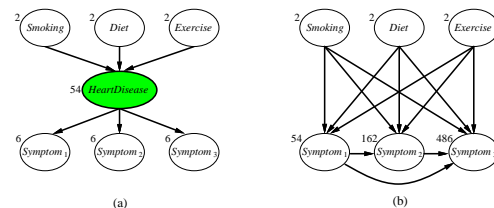
"Naive" solutions for incomplete data

- Solution #1: **Ignore records with missing values**
 - But what if all records are missing values (i.e., when a variable is hidden, none of the records have any value for that variable)
- Solution #2: **Ignore hidden variables**
 - Model may become significantly more complex!

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

11

Heart disease example



(a)

(b)

- a) simpler (i.e., fewer CPT parameters)
- b) complex (i.e., lots of CPT parameters)

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

12

"Direct" maximum likelihood

- Solution 3: **maximize likelihood directly**
 - Let \mathbf{Z} be hidden and \mathbf{E} observable
 - $h_{ML} = \operatorname{argmax}_h P(\mathbf{e}|\mathbf{h})$

$$= \operatorname{argmax}_h \sum_{\mathbf{Z}} P(\mathbf{e}, \mathbf{Z}|\mathbf{h})$$

$$= \operatorname{argmax}_h \sum_{\mathbf{Z}} \prod_i CPT(V_i)$$

$$= \operatorname{argmax}_h \log \sum_{\mathbf{Z}} \prod_i CPT(V_i)$$
 - Problem: can't push log past sum to linearize product

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

13

Expectation-Maximization (EM)

- Solution #4: EM algorithm
 - Intuition: if we knew the missing values, computing h_{ML} would be trivial
- Guess h_{ML}
- Iterate
 - **Expectation**: based on h_{ML} , compute expectation of the missing values
 - **Maximization**: based on expected missing values, compute new estimate of h_{ML}

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

14

Expectation-Maximization (EM)

- More formally:
 - Approximate maximum likelihood
 - Iteratively compute:

$$h_{i+1} = \operatorname{argmax}_h \underbrace{\sum_{\mathbf{Z}} P(\mathbf{Z}|\mathbf{h}_i, \mathbf{e}) \log P(\mathbf{e}, \mathbf{Z}|\mathbf{h})}_{\text{Expectation}}$$

$$\underbrace{\hspace{10em}}_{\text{Maximization}}$$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

15

Expectation-Maximization (EM)

- Derivation
 - $\log P(\mathbf{e}|\mathbf{h}) = \log [P(\mathbf{e}, \mathbf{Z}|\mathbf{h}) / P(\mathbf{Z}|\mathbf{e}, \mathbf{h})]$

$$= \log P(\mathbf{e}, \mathbf{Z}|\mathbf{h}) - \log P(\mathbf{Z}|\mathbf{e}, \mathbf{h})$$

$$= \sum_{\mathbf{Z}} P(\mathbf{Z}|\mathbf{e}, \mathbf{h}) \log P(\mathbf{e}, \mathbf{Z}|\mathbf{h})$$

$$= \sum_{\mathbf{Z}} P(\mathbf{Z}|\mathbf{e}, \mathbf{h}) \log P(\mathbf{Z}|\mathbf{e}, \mathbf{h})$$

$$\geq \sum_{\mathbf{Z}} \tilde{P}(\mathbf{Z}|\mathbf{e}, \mathbf{h}) \log P(\mathbf{e}, \mathbf{Z}|\mathbf{h})$$
- EM finds a **local maximum** of $\sum_{\mathbf{Z}} P(\mathbf{Z}|\mathbf{e}, \mathbf{h}) \log P(\mathbf{e}, \mathbf{Z}|\mathbf{h})$ which is a **lower bound** of $\log P(\mathbf{e}|\mathbf{h})$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

16

Expectation-Maximization (EM)

- **Log inside sum can linearize product**
 - $h_{i+1} = \operatorname{argmax}_h \sum_{\mathbf{Z}} P(\mathbf{Z}|\mathbf{h}_i, \mathbf{e}) \log P(\mathbf{e}, \mathbf{Z}|\mathbf{h})$

$$= \operatorname{argmax}_h \sum_{\mathbf{Z}} P(\mathbf{Z}|\mathbf{h}_i, \mathbf{e}) \log \prod_j CPT_j$$

$$= \operatorname{argmax}_h \sum_{\mathbf{Z}} P(\mathbf{Z}|\mathbf{h}_i, \mathbf{e}) \sum_j \log CPT_j$$
- **Monotonic improvement of likelihood**
 - $P(\mathbf{e}|\mathbf{h}_{i+1}) \geq P(\mathbf{e}|\mathbf{h}_i)$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

17

Candy Example

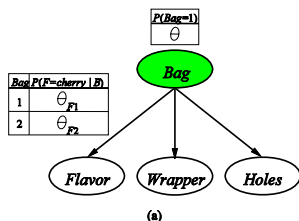
- Suppose you buy two bags of candies of unknown type (e.g. flavour ratios)
- You plan to eat sufficiently many candies of each bag to learn their type
- Ignoring your plan, your roommate mixes both bags...
- How can you learn the type of each bag despite being mixed?

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

18

Candy Example

- "Bag" variable is hidden

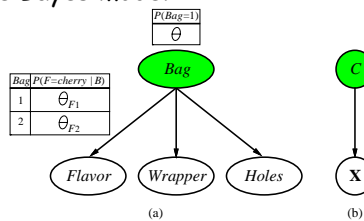


CS886 Fall 07 - Lecture 7, Oct. 2, 2007

19

Unsupervised Clustering

- "Class" variable is hidden
- Naïve Bayes model



CS886 Fall 07 - Lecture 7, Oct. 2, 2007

20

Candy Example

- Unknown Parameters:
 - $\theta_i = P(\text{Bag}=i)$
 - $\theta_{Fi} = P(\text{Flavour}=\text{cherry}|\text{Bag}=i)$
 - $\theta_{Wi} = P(\text{Wrapper}=\text{red}|\text{Bag}=i)$
 - $\theta_{Hi} = P(\text{Hole}=\text{yes}|\text{Bag}=i)$
- When eating a candy:
 - F, W and H are observable
 - B is hidden

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

21

Candy Example

- Let true parameters be:
 - $\theta=0.5, \theta_{F1}=\theta_{W1}=\theta_{H1}=0.8, \theta_{F2}=\theta_{W2}=\theta_{H2}=0.3$
- After eating 1000 candies:

	W=red		W=green	
	H=1	H=0	H=1	H=0
F=cherry	273	93	104	90
F=lime	79	100	94	167

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

22

Candy Example

- EM algorithm
- Guess h_0 :
 - $\theta=0.6, \theta_{F1}=\theta_{W1}=\theta_{H1}=0.6, \theta_{F2}=\theta_{W2}=\theta_{H2}=0.4$
- Alternate:
 - Expectation: expected # of candies in each bag
 - Maximization: new parameter estimates

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

23

Candy Example

- Expectation: expected # of candies in each bag
 - $\#[\text{Bag}=i] = \sum_j P(B=i|f_j, w_j, h_j)$
 - Compute $P(B=i|f_j, w_j, h_j)$ by variable elimination (or any other inference alg.)
- Example:
 - $\#[\text{Bag}=1] = 612$
 - $\#[\text{Bag}=2] = 388$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

24

Candy Example

- Maximization: relative frequency of each bag
 - $\theta_1 = 612/1000 = 0.612$
 - $\theta_2 = 388/1000 = 0.388$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

25

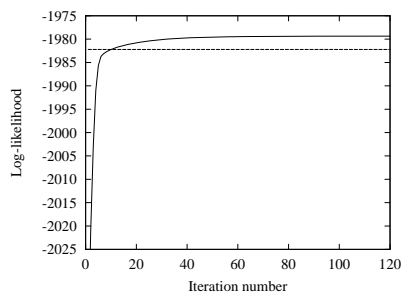
Candy Example

- Expectation: expected # of cherry candies in each bag
 - $\#[B=i, F=\text{cherry}] = \sum_j P(B=i | f_j=\text{cherry}, w_j, h_j)$
 - Compute $P(B=i | f_j=\text{cherry}, w_j, h_j)$ by variable elimination (or any other inference alg.)
- Maximization:
 - $\theta_{F1} = \#[B=1, F=\text{cherry}] / \#[B=1] = 0.668$
 - $\theta_{F2} = \#[B=2, F=\text{cherry}] / \#[B=2] = 0.389$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

26

Candy Example



CS886 Fall 07 - Lecture 7, Oct. 2, 2007

27

Bayesian networks

- EM algorithm for general Bayes nets
- Expectation:
 - $\#[V_i=v_{ij}, Pa(V_i)=pa_{ik}] = \text{expected frequency}$
- Maximization:
 - $\theta_{v_{ij}, pa_{ik}} = \#[V_i=v_{ij}, Pa(V_i)=pa_{ik}] / \#[Pa(V_i)=pa_{ik}]$

CS886 Fall 07 - Lecture 7, Oct. 2, 2007

28