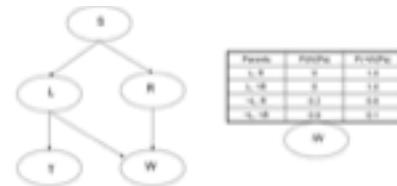


# Bayes Net – Learning

## Bayes Network Definition



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

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

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

# Learning of Bayes Nets

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

# Learning CPTs from Fully Observed Data

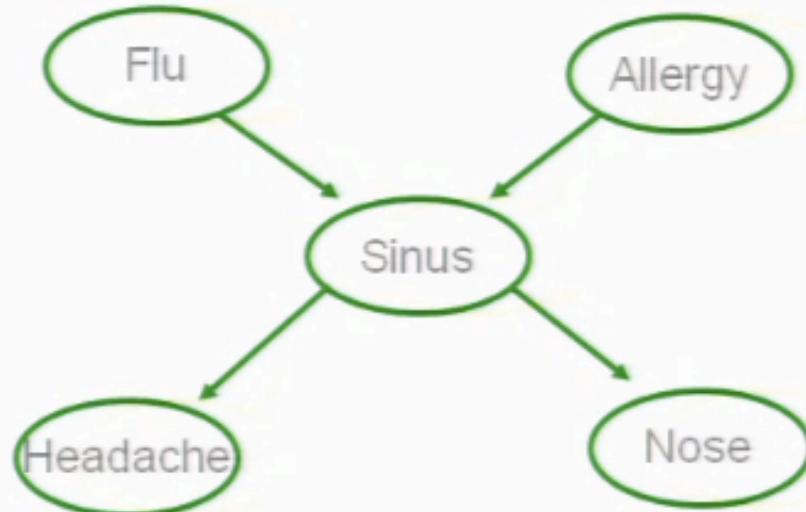
- Example: Consider learning the parameter

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

- MLE (Max Likelihood Estimate) is

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

k<sup>th</sup> training example



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

- Maximum likelihood estimate

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

- Our case:

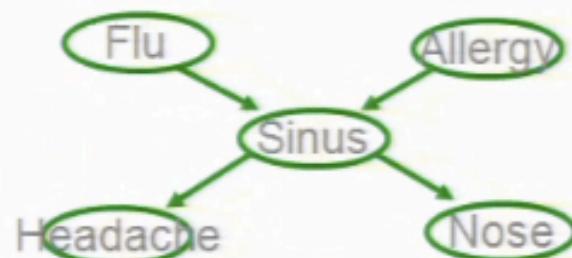
$$P(\text{data}|\theta) = \prod_{k=1}^K P(f_k, a_k, s_k, h_k, n_k)$$

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

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

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

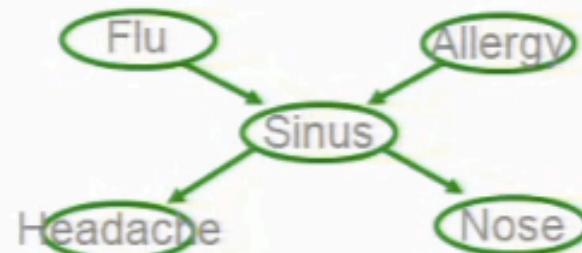
$$\theta_{s|ij} = \frac{\sum_{k=1}^K \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^K \delta(f_k = i, a_k = j)}$$



## Estimate $\theta$ 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:

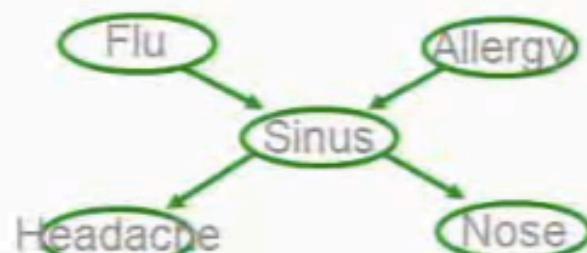
$$\theta \leftarrow \arg \max_{\theta} \log P(X, Z | \theta)$$

- WHAT TO DO?

## Estimate $\theta$ 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:

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

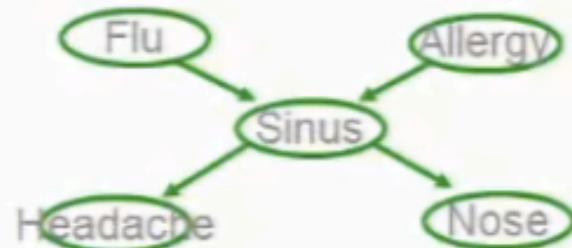
- EM seeks\* to estimate:

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

\* 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  $Z = \{S\}$

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

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

## 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  $\theta$  to calculate  $P(Z|X,\theta)$
- M Step: Replace current  $\theta$  by

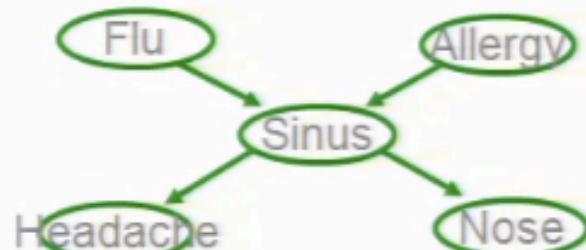
$$\theta \leftarrow \arg \max_{\theta'} Q(\theta'|\theta)$$

Guaranteed to find local maximum.

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

## E Step: Use $X, \theta$ , to Calculate $P(Z|X, \theta)$

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

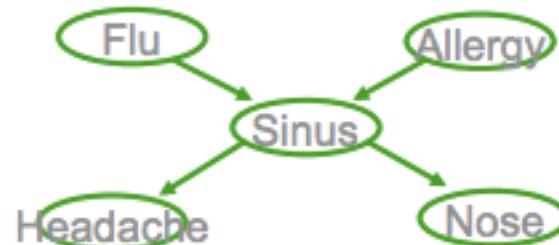


- How? Bayes net inference problem.

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

## E Step: Use $X$ , $\theta$ , to Calculate $P(Z|X, \theta)$

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



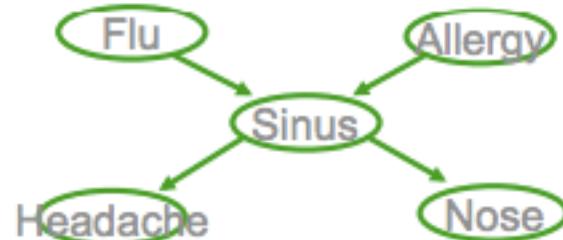
- How? Bayes net inference problem.

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

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

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

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



E step: Calculate  $P(Z_k|X_k; \theta)$  for each training example, k

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

M step: update all relevant parameters. For example:

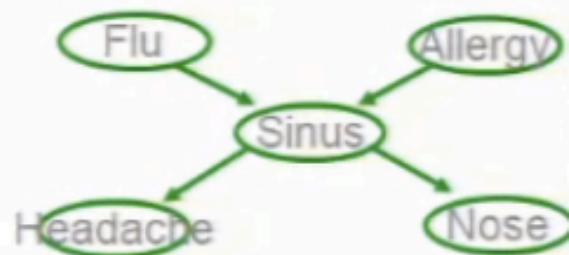
$$\theta_{s|ij} \leftarrow \frac{\sum_{k=1}^K \delta(f_k = i, a_k = j) E[s_k]}{\sum_{k=1}^K \delta(f_k = i, a_k = j)}$$

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 $\theta$

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:

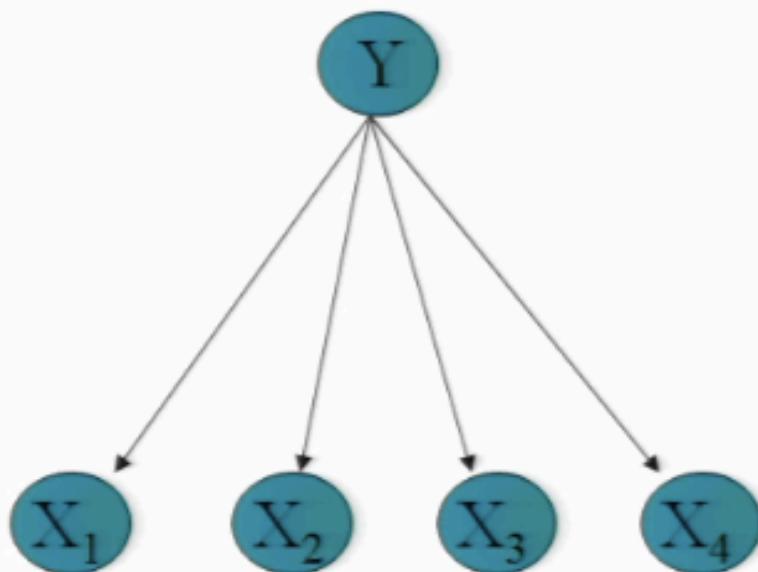
Calculate estimates similar to MLE, but  
replacing each count by its expected count

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

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

# Using Unlabeled Data to Help Train Naïve Bayes Classifier

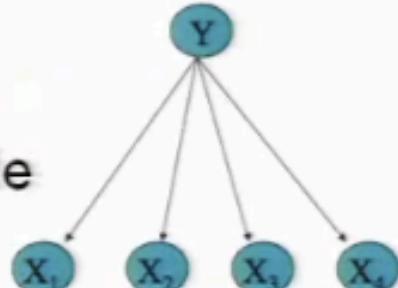
Learn  $P(Y|X)$



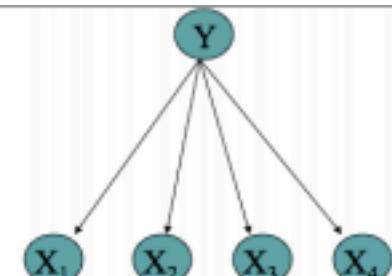
Y	X1	X2	X3	X4
1	0	0	1	1
0	0	1	0	0
0	0	0	1	0
?	0	1	1	0
?	0	1	0	1

E step: Calculate for each training example, k

the expected value of each unobserved variable



## EM and estimating $\theta$



Given observed set X, unobserved set Y of boolean values

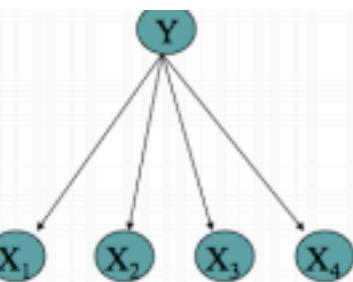
E step: Calculate for each training example, k

the expected value of each unobserved variable Y

$$E_{P(Y|X_1 \dots X_N)}[y(k)] = P(y(k) = 1 | x_1(k), \dots, x_N(k); \theta) = \frac{P(y(k) = 1) \prod_i P(x_i(k) | y(k) = 1)}{\sum_{j=0}^1 P(y(k) = j) \prod_i P(x_i(k) | y(k) = j)}$$

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

let's use  $y(k)$  to indicate value of Y on kth example



## EM and estimating $\theta$

Given observed set  $X$ , unobserved set  $Y$  of boolean values

E step: Calculate for each training example,  $k$

the expected value of each unobserved variable  $Y$

$$E_{P(Y|X_1 \dots X_N)}[y(k)] = P(y(k) = 1|x_1(k), \dots, x_N(k); \theta) = \frac{P(y(k) = 1) \prod_i P(x_i(k)|y(k) = 1)}{\sum_{j=0}^1 P(y(k) = j) \prod_i P(x_i(k)|y(k) = j)}$$

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

$$\theta_{ij|m} = \hat{P}(X_i = j|Y = m) = \frac{\sum_k P(y(k) = m|x_1(k) \dots x_N(k)) \delta(x_i(k) = j)}{\sum_k P(y(k) = m|x_1(k) \dots x_N(k))}$$

MLE would be:  $\hat{P}(X_i = j|Y = m) = \frac{\sum_k \delta((y(k) = m) \wedge (x_i(k) = j))}{\sum_k \delta(y(k) = m)}$

- **Inputs:** Collections  $\mathcal{D}^l$  of labeled documents and  $\mathcal{D}^u$  of unlabeled documents.
- Build an initial naive Bayes classifier,  $\hat{\theta}$ , from the labeled documents,  $\mathcal{D}^l$ , only. Use maximum a posteriori parameter estimation to find  $\hat{\theta} = \arg \max_{\theta} P(\mathcal{D}|\theta)P(\theta)$  (see Equations 5 and 6).
- Loop while classifier parameters improve, as measured by the change in  $l_c(\theta|\mathcal{D}; \mathbf{z})$  (the complete log probability of the labeled and unlabeled data
  - **(E-step)** Use the current classifier,  $\hat{\theta}$ , to estimate component membership of each unlabeled document, i.e., the probability that each mixture component (and class) generated each document,  $P(c_j|d_i; \hat{\theta})$  (see Equation 7).
  - **(M-step)** Re-estimate the classifier,  $\hat{\theta}$ , given the estimated component membership of each document. Use maximum a posteriori parameter estimation to find  $\hat{\theta} = \arg \max_{\theta} P(\mathcal{D}|\theta)P(\theta)$  (see Equations 5 and 6).
- **Output:** A classifier,  $\hat{\theta}$ , that takes an unlabeled document and predicts a class label.

From [Nigam et al., 2000]



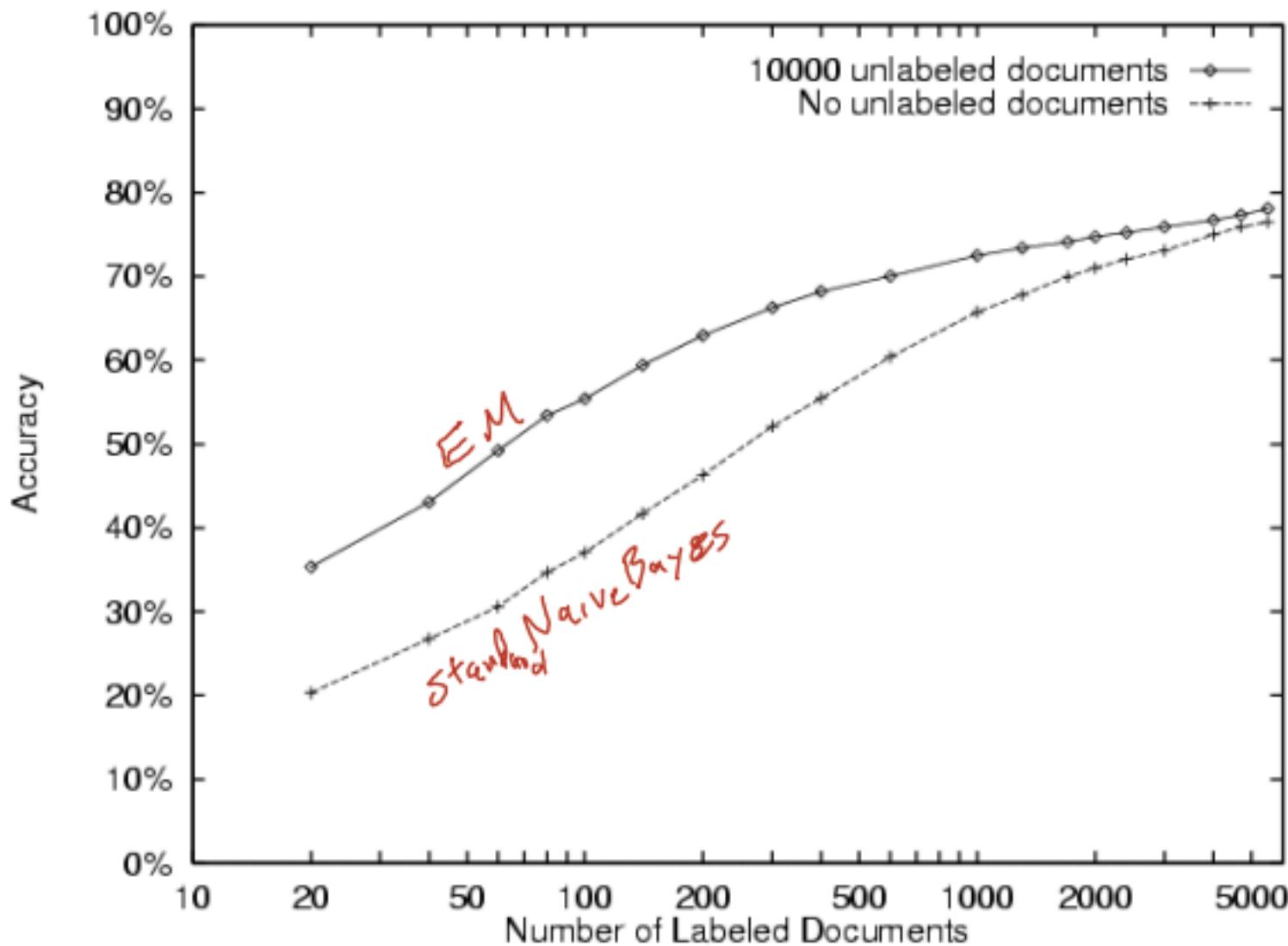
## Experimental Evaluation

- Newsgroup postings
  - 20 newsgroups, 1000/group
- Web page classification
  - student, faculty, course, project
  - 4199 web pages
- Reuters newswire articles
  - 12,902 articles
  - 90 topics categories

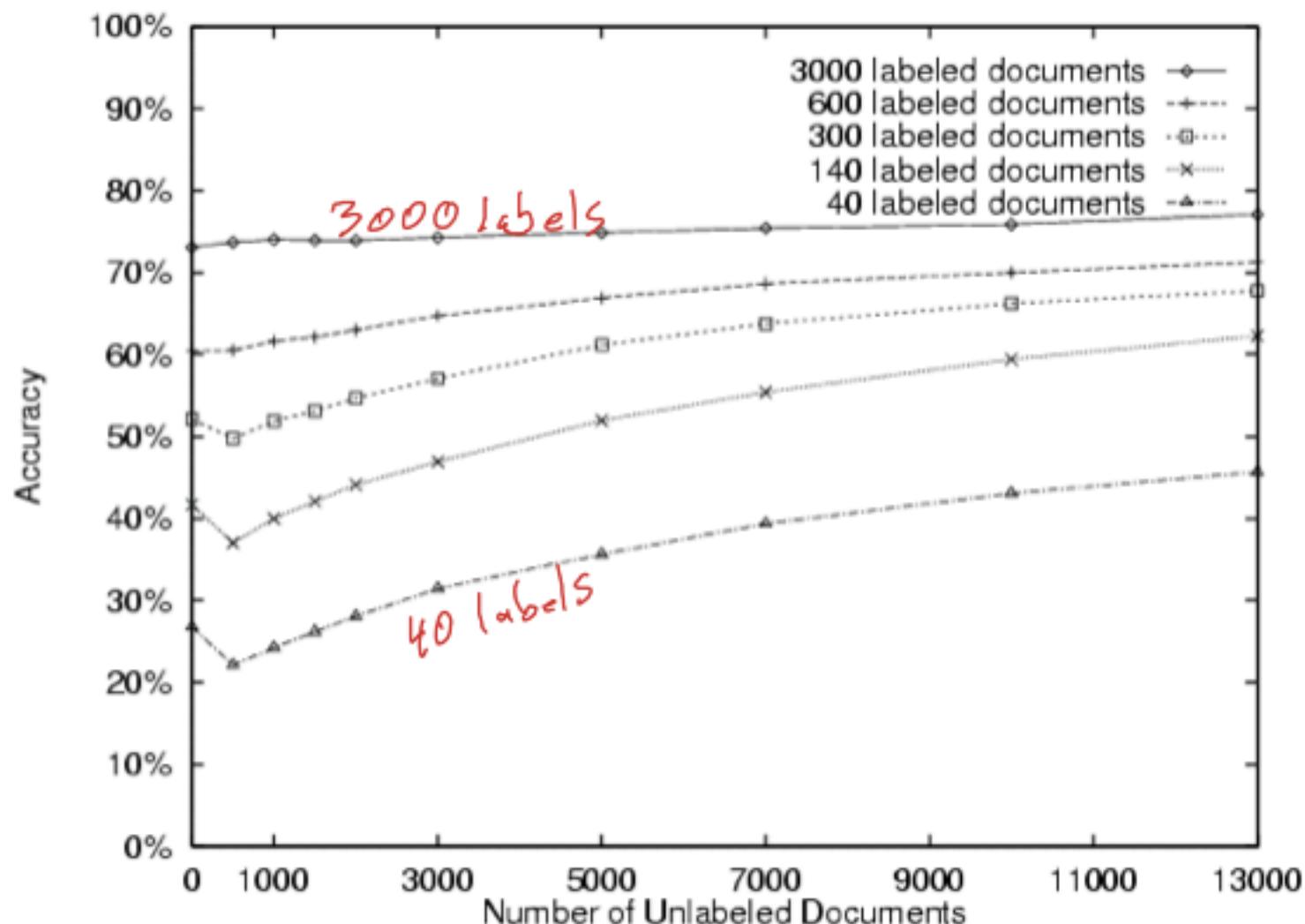
Table 3. Lists of the words most predictive of the course class in the WebKB data set, as they change over iterations of EM for a specific trial. By the second iteration of EM, many common course-related words appear. The symbol  $D$  indicates an arbitrary digit.

Iteration 0		Iteration 1	Iteration 2
intelligence	word w ranked by $P(w Y=\text{course}) / P(w Y \neq \text{course})$	$DD$	$D$
$DD$		$D$	$DD$
artificial		lecture	lecture
understanding		cc	cc
$DDw$		$D^*$	$DD:DD$
dist		$DD:DD$	due
identical		handout	$D^*$
rus		due	homework
arrange		problem	assignment
games		set	handout
dartmouth		tay	set
natural		$DDam$	hw
cognitive		yurttas	exam
logic		homework	problem
proving		kfoury	$DDam$
prolog		sec	postscript
knowledge		postscript	solution
human		exam	quiz
representation		solution	chapter
field		assaf	ascii
	Using one labeled example per class		

# 20 Newsgroups



# 20 Newsgroups



# Learning Bayes Net Structure

# How can we learn Bayes Net graph structure?

In general case, open problem

- can require lots of data (else high risk of overfitting)
- can use Bayesian methods to constrain search

One key result:

- Chow-Liu algorithm: finds “best” tree-structured network
- What’s best?
  - suppose  $P(\mathbf{X})$  is true distribution,  $T(\mathbf{X})$  is our tree-structured network, where  $\mathbf{X} = \langle X_1, \dots, X_n \rangle$
  - Chow-Liu minimizes Kullback-Leibler divergence:

$$KL(P(\mathbf{X}) \parallel T(\mathbf{X})) \equiv \sum_k P(\mathbf{X} = k) \log \frac{P(\mathbf{X} = k)}{T(\mathbf{X} = k)}$$

## Chow-Liu Algorithm

Key result: To minimize  $KL(P \parallel T)$ , it suffices to find the tree network  $T$  that maximizes the sum of mutual informations over its edges

Mutual information for an edge between variable A and B:

$$I(A, B) = \sum_a \sum_b P(a, b) \log \frac{P(a, b)}{P(a)P(b)}$$

This works because for tree networks with nodes  $\mathbf{X} \equiv \langle X_1 \dots X_n \rangle$

$$\begin{aligned} KL(P(\mathbf{X}) \parallel T(\mathbf{X})) &\equiv \sum_k P(\mathbf{X} = k) \log \frac{P(\mathbf{X} = k)}{T(\mathbf{X} = k)} \\ &= - \sum_i I(X_i, Pa(X_i)) + \sum_i H(X_i) - H(X_1 \dots X_n) \end{aligned}$$

## Chow-Liu Algorithm

1. for each pair of vars A,B, use data to estimate  $P(A,B)$ ,  $P(A)$ ,  $P(B)$

2. for each pair of vars A,B calculate mutual information

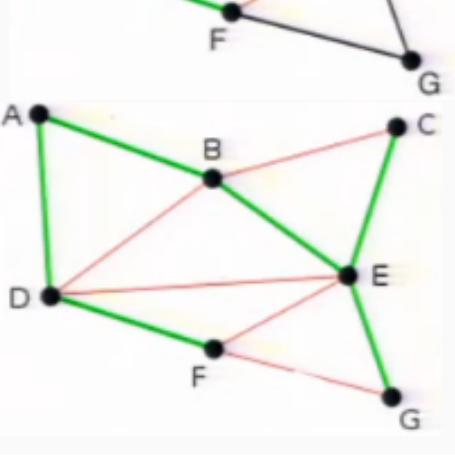
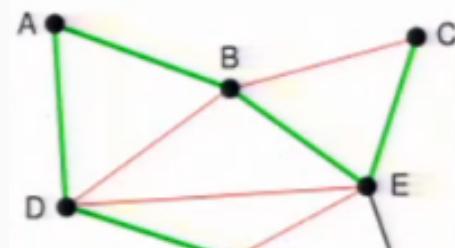
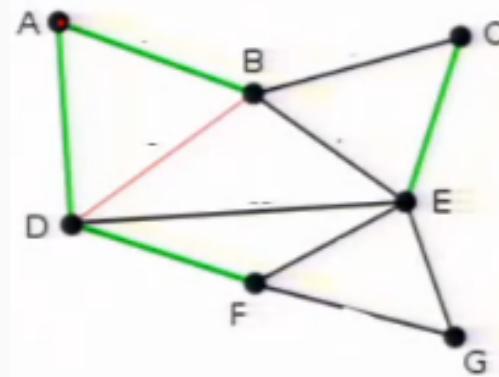
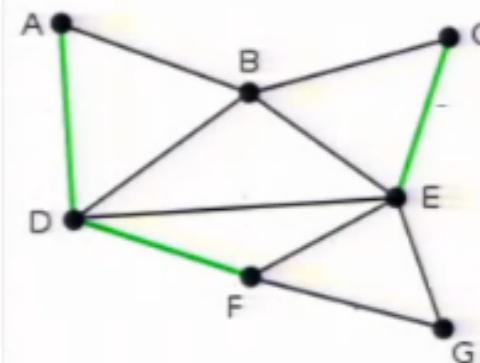
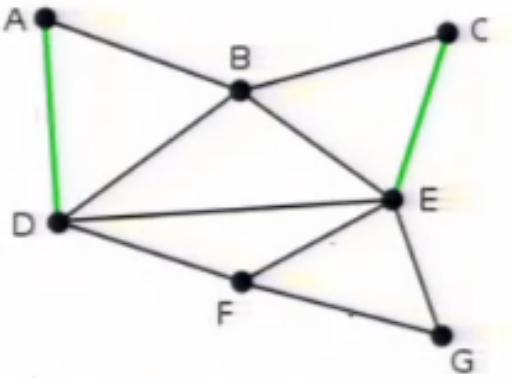
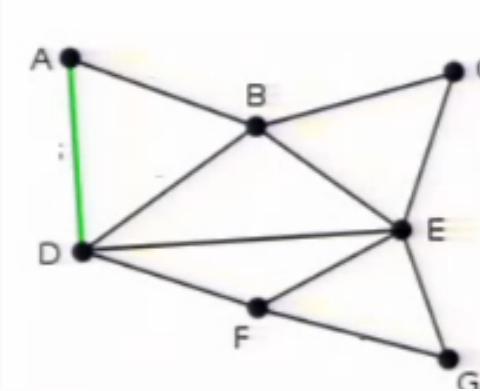
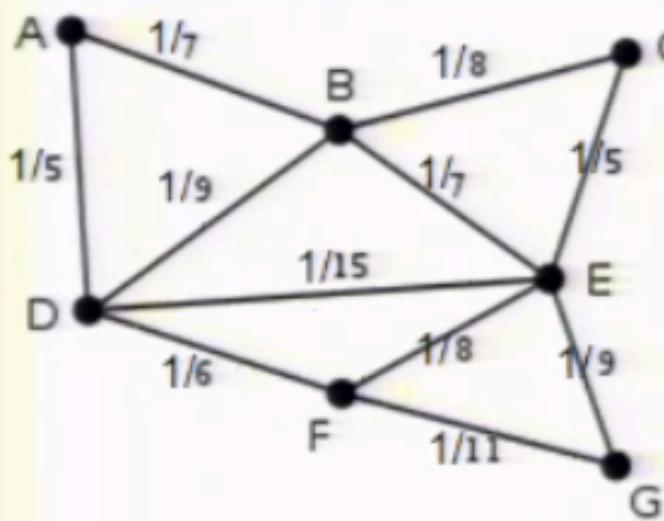
$$I(A, B) = \sum_a \sum_b P(a, b) \log \frac{P(a, b)}{P(a)P(b)}$$

3. calculate the maximum spanning tree over the set of variables, using edge weights  $I(A,B)$   
(given N vars, this costs only  $O(N^2)$  time)

4. add arrows to edges to form a directed-acyclic graph
5. learn the CPD's for this graph

# Chow-Liu algorithm example

## Greedy Algorithm to find Max-Spanning Tree



# Bayes Nets – What You Should Know

- **Representation**
  - Bayes nets represent joint distribution as a DAG + Conditional Distributions
  - D-separation lets us decode conditional independence assumptions
- **Inference**
  - NP-hard in general
  - For some graphs, closed form inference is feasible
  - Approximate methods too, e.g., Monte Carlo methods, ...
- **Learning**
  - Easy for known graph, fully observed data (MLE's, MAP est.)
  - EM for partly observed data
  - Learning graph structure: Chow-Liu for tree-structured networks