

Machine Learning 10-701

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

February 17, 2011

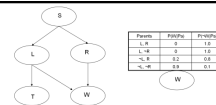
Today:

- Graphical models
- Learning from fully labeled data
- Learning from partly observed data
- EM

Readings:

- Required:
- Bishop chapter 8, through 8.2

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

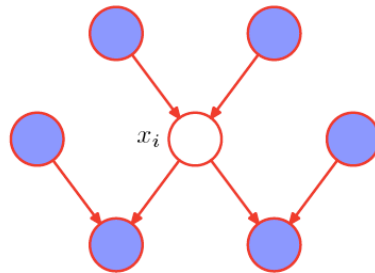
- Each node denotes a random variable
- Edges denote dependencies
- CPD for each node X_i defines $P(X_i | Pa(X_i))$
- The joint distribution over all variables is defined as

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

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

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.



from [Bishop, 8.2]

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-descendants, given its immediate parents
 - D-separation
 - 'Explaining away'

Java Bayes Net Applet

<http://www.pmr.poli.usp.br/ItD/Software/javabayes/Home/applet.html>

by **Fabio Gagliardi Cozman**

Learning of Bayes Nets

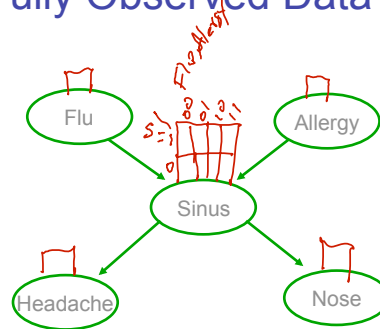
- Four categories of learning problems
 - Graph structure may be known/unknown
 - Variable values may be fully observed / partly unobserved
- Easy case: learn parameters for graph structure is *known*, and data is *fully observed*
- Interesting case: graph *known*, data *partly known*
- Gruesome case: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

δ if argument is true, else 0

- Remember why?

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:

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

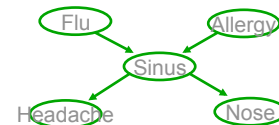
k^{th} train. example

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

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

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

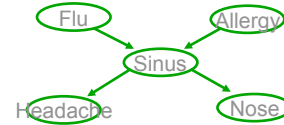
$$\theta_{s|ij} = \frac{\sum_{k=1}^K \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^K \delta(f_k = i, a_k = j)}$$



Estimate 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:

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

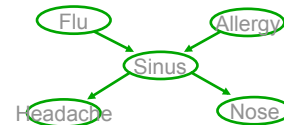
- WHAT TO DO?

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

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:

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

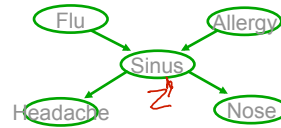
- EM seeks* to estimate:

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

* 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 $Z=\{S\}$

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

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

E step provides this!

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

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

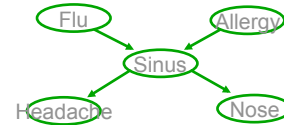
value of each var $\in Z$ for each train example

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 $P(Z|X,\theta)$

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

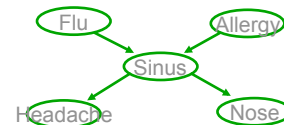


- How? Bayes net inference problem.

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

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

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



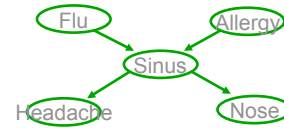
- How? Bayes net inference problem.

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

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

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:

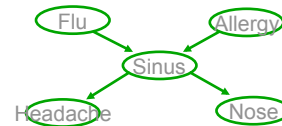
$$\theta_{s|ij} \leftarrow \frac{\sum_{k=1}^K \delta(f_k = i, a_k = j) E[s_k]}{\sum_{k=1}^K \delta(f_k = i, a_k = j)}$$

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 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

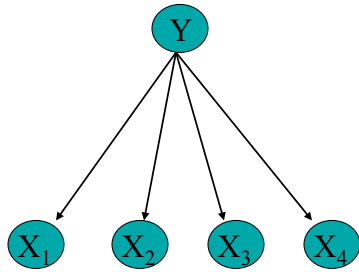
M step:

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

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

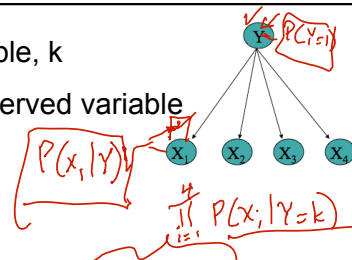
Using Unlabeled Data to Help Train Naïve Bayes Classifier

Learn $P(Y|X)$



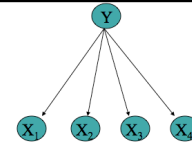
Y	X1	X2	X3	X4
✓ 1	0	0	1	1
✓ 0	0	1	0	0
✓ 0	0	0	1	0
E[Y] ?	0	1	1	0
E[Y] ?	0	1	0	1

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



$$E[Y] = \frac{P(Y=1|x_1, \dots, x_4, \theta)}{P(Y|x_1, \dots, x_4, \theta)} = \frac{P(Y=1|x_1, \dots, x_4, \theta)}{\sum_j P(X_1, \dots, X_4|Y=j, \theta) P(Y=j, \theta)}$$

EM and estimating θ



Given observed set X, unobserved set Y of boolean values

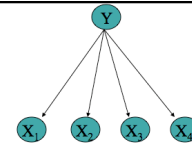
E step: Calculate for each training example, k
the expected value of each unobserved variable Y

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

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 θ



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

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

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

- **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]

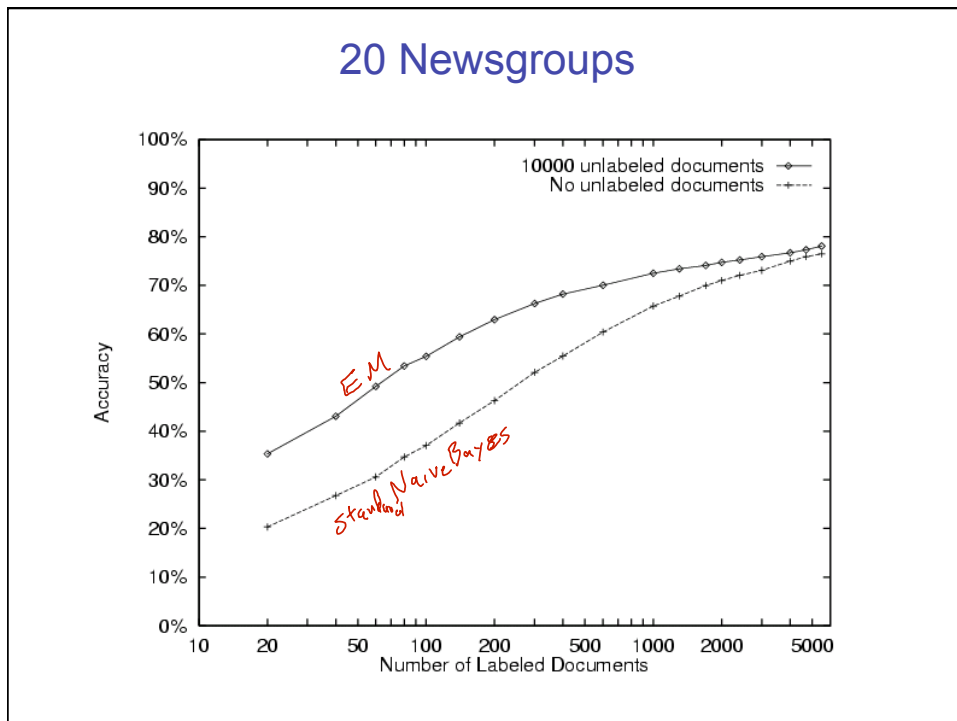


Experimental Evaluation

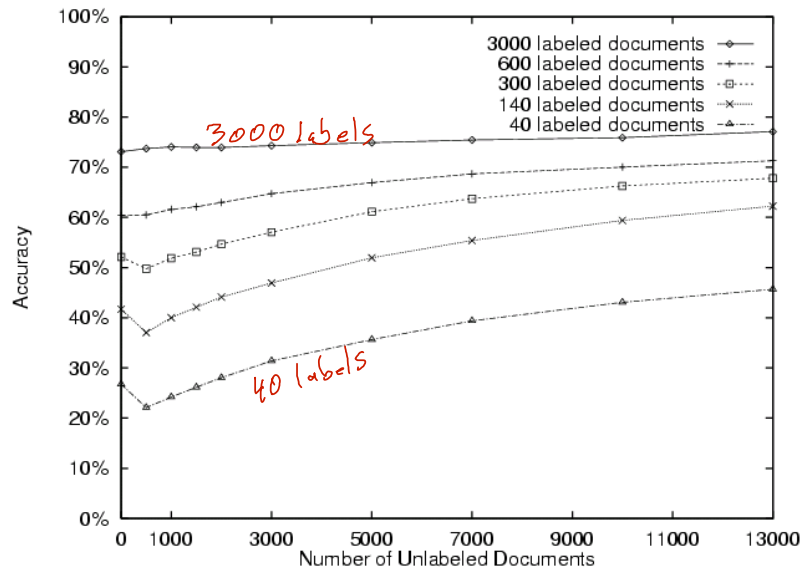
- Newsgroup postings
 - 20 newsgroups, 1000/group
- Web page classification
 - student, faculty, course, project
 - 4199 web pages
- Reuters newswire articles
 - 12,902 articles
 - 90 topics categories

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



20 Newsgroups



Unsupervised clustering

Just extreme case for EM with
zero labeled examples...