

Recommended reading:

“An Introduction to HMMs and Bayesian Networks,”
Z. Ghahramani, *Int. Journal of Pattern Recognition and AI*,
15(1):9-42, (2001)

Especially Section 4

EM for HMMs a.k.a. The Baum-Welch Algorithm

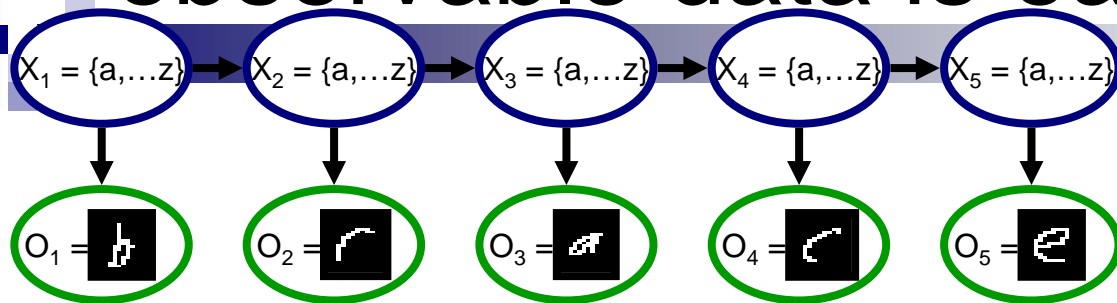
Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 12th, 2006

Learning HMMs from fully observable data is easy



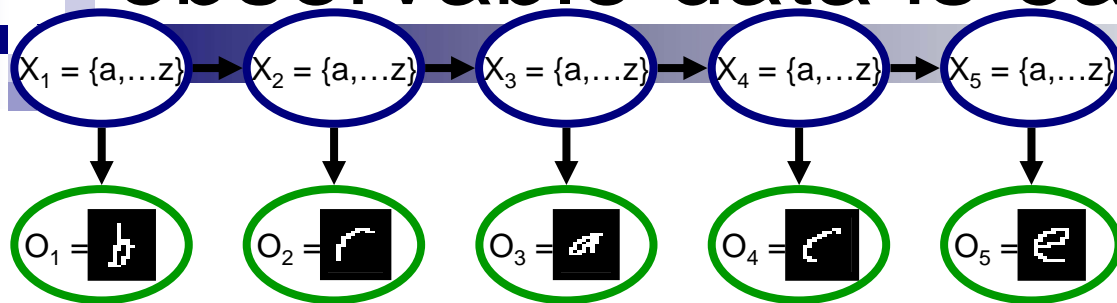
Learn 3 distributions:

$$P(X_i^a) = \frac{\text{Count}(X_i = a)}{m}$$

$$P(O_i^{\text{Pixel 17=on}} | X_i^a) = \frac{\text{Count}(A=a, \text{Pixel 17=on})}{\text{Count}(A=a)}$$

$$P(X_i^a | X_{i-1}^b) = \frac{\text{Count}(X_{i-1}=b, X_i=a)}{\text{Count}(X_{i-1}=b, X_i=?)}$$

Learning HMMs from fully observable data is easy



Learn 3 distributions:

$$P(X_1^a) = \frac{\text{count}(\# \text{ first letter was } a)}{N = \text{dataset size}}$$

$$P(O_i^{\text{pixel } l \text{ is white}} | X_i^a) = \frac{\text{count}(\text{pixel } l \text{ was white, } X_i = a)}{n_i}$$

$$P(X_i^a | X_{i-1}^b)$$

What if O is observed,
but X is hidden

select training data where letter was a

Log likelihood for HMMs when \mathbf{X} is hidden

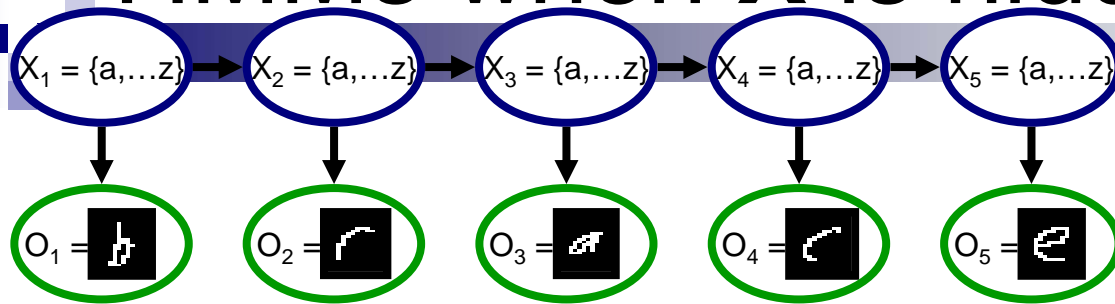
- Marginal likelihood – \mathbf{O} is observed, \mathbf{X} is missing
 - For simplicity of notation, training data consists of only one sequence:

$$\begin{aligned} \ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} | \theta) \quad \leftarrow \text{marginal likelihood} \\ &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} | \theta) \end{aligned}$$

- If there were m sequences:

$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o}^{(j)} | \theta)$$

Computing Log likelihood for HMMs when \mathbf{X} is hidden



$$\begin{aligned} \ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} \mid \theta) \\ &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} \mid \theta) \end{aligned}$$

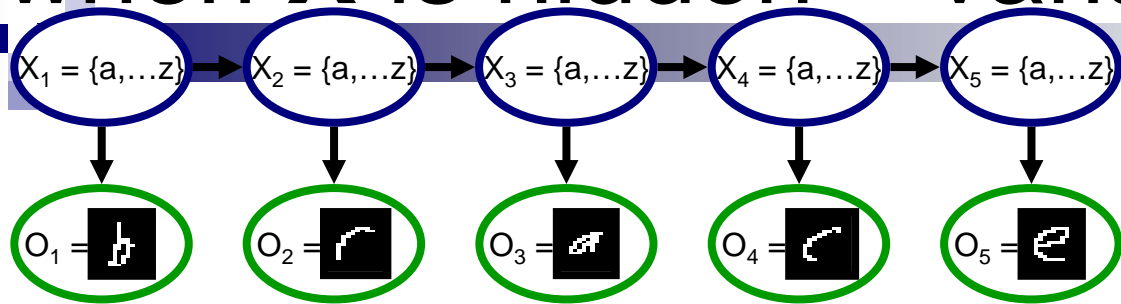
sum naively

K - letters

n - positions (word length)

K^n - terms

Computing Log likelihood for HMMs when \mathbf{X} is hidden – variable elimination

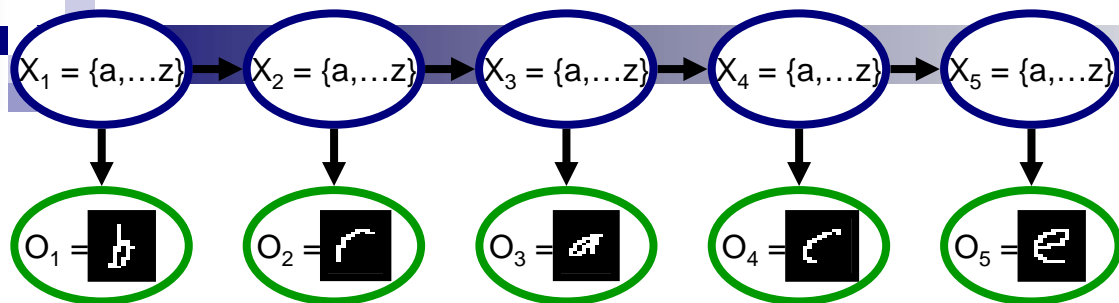


- Can compute efficiently with variable elimination:

$$\begin{aligned}
 \ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} | \theta) \\
 &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} | \theta) = \log \sum_{x_1, \dots, x_n} P(x_1) \cdot P(o_1 | x_1) \prod_{t=2}^n P(x_t | x_{t-1}) P(o_t | x_t) \\
 &= \log \sum_{x_1, \dots, x_{n-1}} P(x_1) P(o_1 | x_1) \prod_{t=2}^{n-1} P(x_t | x_{t-1}) P(o_t | x_t) \underbrace{\sum_{x_n} P(x_n | x_{n-1}) P(o_n | x_n)}_{g(x_{n-1})}
 \end{aligned}$$

eliminate x_{n-1}
 x_{n-2}
 \vdots

EM for HMMs when X is hidden



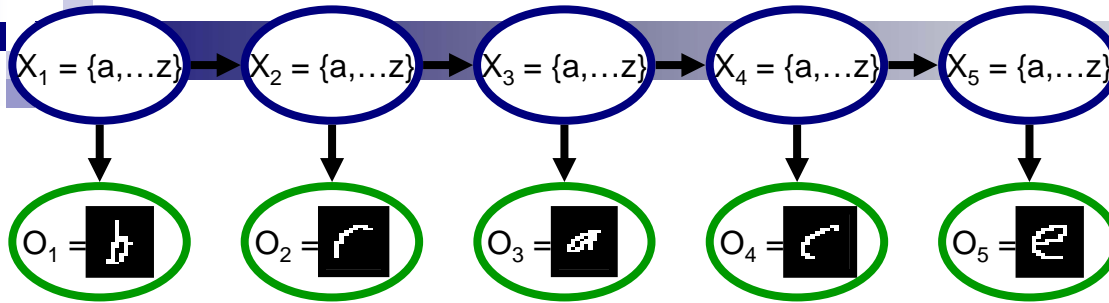
- E-step: Use inference (forwards-backwards algorithm)

$$P(X_3 = a \mid O = \boxed{\text{black}})$$

- M-step: Recompute parameters with weighted data

learn weighted data!!
😊

E-step



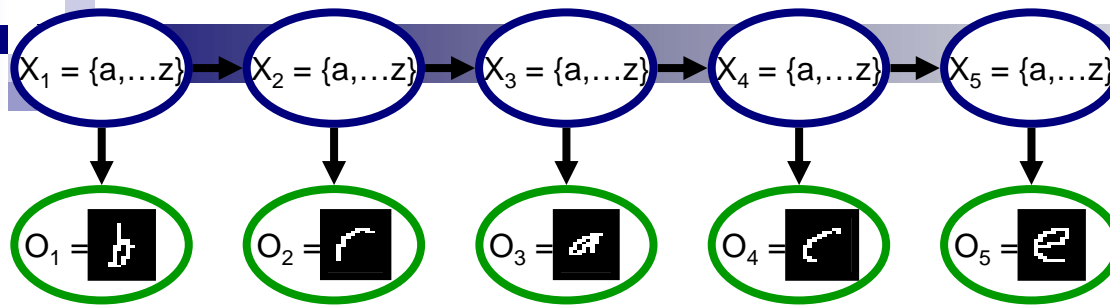
- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}

$$Q^{(t+1)}(\underline{\mathbf{x}} \mid \underline{\mathbf{o}}) = P(\mathbf{x} \mid \mathbf{o}, \theta^{(t)})$$

iteration t+1

- Will correspond to inference
 - use forward-backward algorithm!

The M-step



Starting: $P(x_1) \leftarrow \theta_{x_1=a}, \theta_{x_1=b} \dots$
 transition: $P(x_i | x_{i-1}) \leftarrow \theta_{x_i=a | x_{i-1}=b}$
 model:

obs. : $P(o_t | x_t) \leftarrow \theta_{\text{pixel 7=0} | x_t=0}$
 model:

Maximization step:

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{\mathbf{x}} Q^{(t+1)}(\mathbf{x} | \mathbf{o}) \log P(\mathbf{x}, \mathbf{o} | \theta)$$

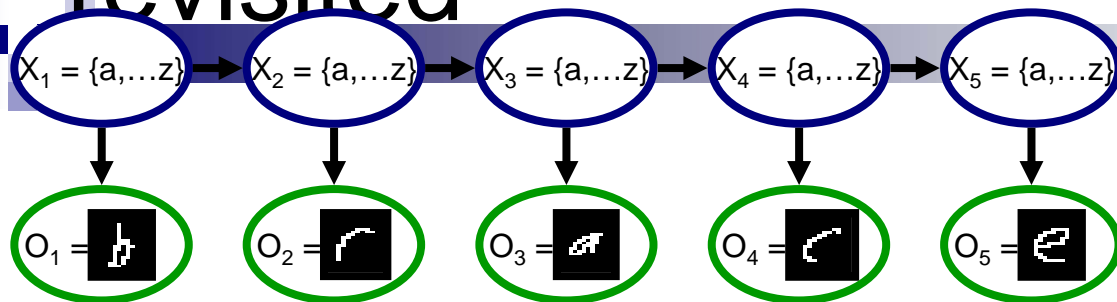
likelihood completed data
 weighted

Use expected counts instead of counts:

- If learning requires $\text{Count}(\mathbf{x}, \mathbf{o})$
- Use $E_{Q^{(t+1)}}[\text{Count}(\mathbf{x}, \mathbf{o})]$

Decomposition of likelihood $P(X_1) \leftarrow \theta_{X_1}$

revisited



$$P(O_i | X_i) \leftarrow \theta_{O|X}$$

$$P(X_i | X_{i-1}) \leftarrow \theta_{X_i | X_{i-1}}$$

$\log a \cdot b = \log a + \log b$

■ Likelihood optimization decomposes:

$$\max_{\theta} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log P(\mathbf{x}, \mathbf{o} | \theta) =$$

$$\max_{\theta} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log P(x_1 | \theta_{X_1}) P(o_1 | x_1, \theta_{O|X}) \prod_{t=2}^n P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) P(o_t | x_t, \theta_{O|X})$$

$$= \max_{\theta} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \left[\log P(x_1 | \theta_{X_1}) + \sum_{t=1}^n \log P(o_t | x_t, \theta_{O|X}) + \sum_{t=2}^n \log P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) \right]$$

$\theta_{X_1}, \theta_{O|X}, \theta_{X_t|X_{t-1}}$

learn start state dist.

learn obs. model

$$= \left[\max_{\theta_{X_1}} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log P(x_1 | \theta_{X_1}) \right] + \left[\max_{\theta_{O|X}} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \sum_{t=1}^n \log P(o_t | x_t, \theta_{O|X}) \right] +$$

learn transition model

$$+ \left[\max_{\theta_{X_t|X_{t-1}}} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \sum_{t=2}^n \log P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) \right]$$

Starting state probability $P(X_1)$

$P(a,b) = P(a) \cdot P(b|a)$

(chain rule)

Using expected counts

$Q \in$ prob. dist. ✓

$$\square P(X_1=a) = \theta_{X_1=a}$$

$$Q(x_1 \dots x_n | o) = Q(x_1 | o) \cdot Q(x_2 \dots x_n | x_1, o)$$

$$\max_{\theta_{X_1}} \sum_{\mathbf{x}} Q(\mathbf{x} | o) \log P(x_1 | \theta_{X_1}) = \max_{\theta_{X_1}} \sum_{x_1, \dots, x_n} Q(x_1 \dots x_n | o) \log P(x_1 | \theta_{X_1})$$

$$= \max_{\theta_{X_1}} \sum_{x_1, \dots, x_n} Q(x_1 | o) \cdot Q(x_2 \dots x_n | x_1, o) \log P(x_1 | \theta_{X_1}) =$$

$$= \max_{\theta_{X_1}} \sum_{x_1} Q(x_1 | o) \log P(x_1 | \theta_{X_1}) \cdot \sum_{x_2 \dots x_n} Q(x_2 \dots x_n | o)$$

$$= \max_{\theta_{X_1}} \sum_{x_1} Q(x_1 | o) \log P(x_1 | \theta_{X_1})$$



$$\theta_{X_1=a} = \frac{\sum_{j=1}^m Q(X_1 = a | o^{(j)})}{m}$$

$$P(x_1 = x_i | \theta_{x_1}) = \theta_{x_1 = x_i}$$

$$\frac{\partial}{\partial \theta_{x_1}} \sum_{j=1}^m \sum_{x_1} Q(x_1 | o^{(j)}) \log P(x_1 | \theta_{x_1}) = 0$$

$$= \sum_{j=1}^m \sum_{x_1} Q(x_1 | o^{(j)}) \frac{\partial}{\partial \theta_{x_1}} \log \theta_{x_1}$$

$$= \sum_{j=1}^m \left[Q(x_1 = t | o^{(j)}) \frac{1}{\theta_{x_1 = t}} + Q(x_1 = f | o^{(j)}) \frac{1}{1 - \theta_{x_1 = t}} \right] = 0$$

re arrange

$$\theta_{x_1 = t} = \frac{\sum_{j=1}^m Q(x_1 = t | o^{(j)})}{\sum_{j=1}^m \left[Q(x_1 = t | o^{(j)}) + Q(x_1 = f | o^{(j)}) \right]}$$

$$m = \sum_{j=1}^m \left[Q(x_1 = t | o^{(j)}) + Q(x_1 = f | o^{(j)}) \right]$$

$$\begin{aligned} P(x_1 = t | \theta_{x_1}) &= \theta_{x_1 = t} \\ P(x_1 = f | \theta_{x_1}) &= 1 - \theta_{x_1 = t} \\ \frac{\partial}{\partial x} \log x &= \frac{1}{x} \end{aligned}$$

Transition probability $P(X_t|X_{t-1})$

$$\log \Pi = \sum \log$$

- Using expected counts

- $P(X_t=a|X_{t-1}=b) = \theta_{X_t=a|X_{t-1}=b}$

$$\max_{\theta_{X_t|X_{t-1}}} \sum_{j=1}^m Q(\mathbf{x} | \mathbf{o}) \log \prod_{t=2}^n P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) = \max_{\theta_{X_t|X_{t-1}}} \sum_{j=1}^m \sum_{t=2}^n \sum_x Q(x | \mathbf{o}) \log P(x_t | x_{t-1}, \theta)$$

$$= \max_{\theta_{X_t|X_{t-1}}} \sum_{j=1}^m \sum_{t=2}^n \sum_{x_1 \dots x_n} Q(x_t, x_{t-1} | \mathbf{o}) \cdot Q(x_1 \dots x_{t-2}, x_{t+1} \dots x_n | x_t, x_{t-1}, \theta) \log P(x_t | x_{t-1}, \theta)$$

$$= \max_{\theta_{X_t|X_{t-1}}} \sum_{j=1}^m \sum_{t=2}^n \sum_{x_t, x_{t-1}} Q(x_t, x_{t-1} | \mathbf{o}) \log P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}})$$

$$\theta_{X_t=a|X_{t-1}=b} = \frac{\sum_{j=1}^m \sum_{t=2}^n Q(X_t = a, X_{t-1} = b | \mathbf{o}^{(j)})}{\sum_{j=1}^m \sum_{t=2}^n \sum_{i=1}^k Q(X_t = i, X_{t-1} = b | \mathbf{o}^{(j)})}$$

Observation probability $P(O_t|X_t)$

- Using expected counts

- $P(O_t=a|X_t=b) = \theta_{O_t=a|X_t=b}$

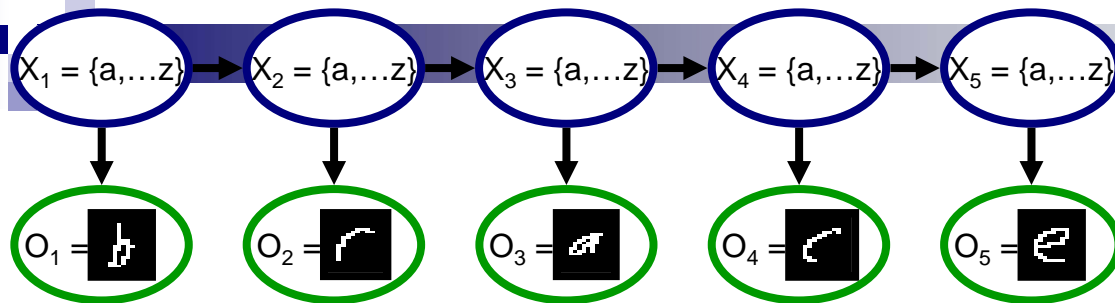
$$\max_{\theta_{O|X}} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log \prod_{t=1}^n P(o_t | x_t, \theta_{O|X}) = \max_{\theta_{O|X}} \sum_{t=1}^n \sum_{x_t} Q(x_t | \mathbf{o}) \log P(o_t | x_t, \theta_{O|X})$$

$\partial(O_t^{(j)} = a)$ ← $\begin{cases} 1: & \text{if } t\text{'th obs. of} \\ & j\text{'th training} \\ & \text{example} = a \\ 0; & \text{otherwise} \end{cases}$

$$\theta_{O_t=a|X_t=b} = \frac{\sum_{j=1}^m \sum_{t=1}^n \delta(o_t^{(j)} = a) Q(X_t = b | \mathbf{o}^{(j)})}{\sum_{j=1}^m \sum_{t=1}^n Q(X_t = b | \mathbf{o}^{(j)})}$$

E-step revisited

$$Q^{(t+1)}(\mathbf{x} | \mathbf{o}) = P(\mathbf{x} | \mathbf{o}, \theta^{(t)})$$



- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}

- Must compute:

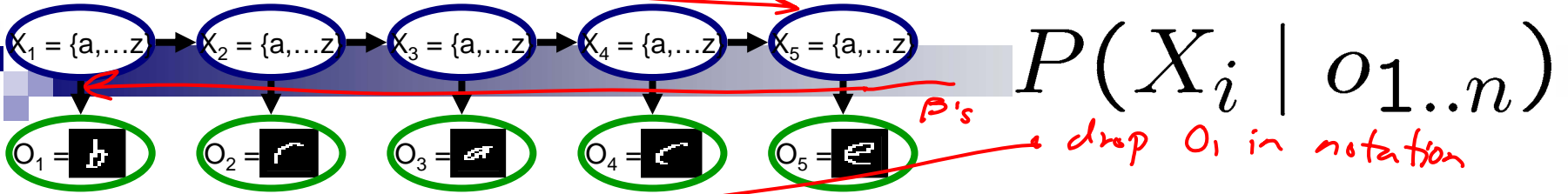
$$P(X_t = a | \mathbf{o})$$

□ $Q(x_t = a | \mathbf{o})$ – marginal probability of each position

$$P(X_{t+1} = a, X_t = b | \mathbf{o})$$

□ $Q(x_{t+1} = a, x_t = b | \mathbf{o})$ – joint distribution between pairs of positions

α's The forwards-backwards algorithm



■ Initialization: $\alpha_1(X_1) = P(X_1)P(o_1 | X_1)$

■ For $i = 2$ to n

□ Generate a forwards factor by eliminating X_{i-1}

sum out previous var prob obs

$$\alpha_i(X_i) = \sum_{x_{i-1}} P(o_i | X_i) P(X_i | X_{i-1} = x_{i-1}) \alpha_{i-1}(x_{i-1})$$

transition prob

■ Initialization: $\beta_n(X_n) = 1$

■ For $i = n-1$ to 1

□ Generate a backwards factor by eliminating X_{i+1}

$\alpha_5(a)$
 $\alpha_5(b)$
 \vdots
 $\alpha_5(z)$

$\alpha_n(x_n)$
 normalized
 $= P(X_n | o_{1:n})$

$\beta_i(x_i) \alpha_i(x_i)$
 normalized
 $= P(X_i | o_{1:n})$

$\forall x_i$

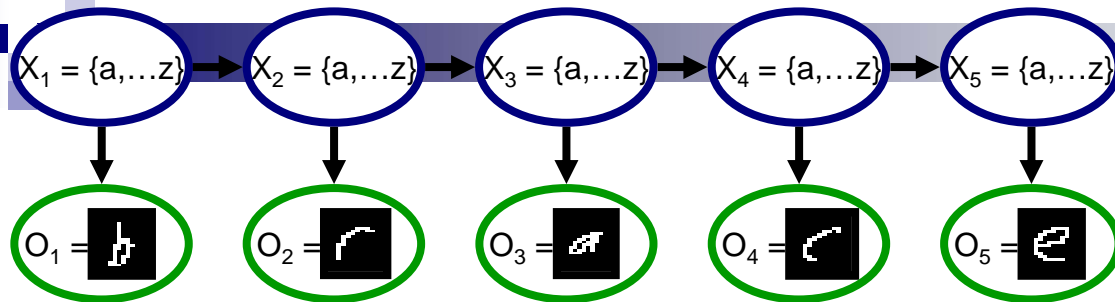
$$\beta_i(X_i) = \sum_{x_{i+1}} P(o_{i+1} | x_{i+1}) P(x_{i+1} | X_i) \beta_{i+1}(x_{i+1})$$

x_i

■ $\forall i$, probability is: $P(X_i | o_{1..n}) = \alpha_i(X_i) \beta_i(X_i)$

E-step revisited

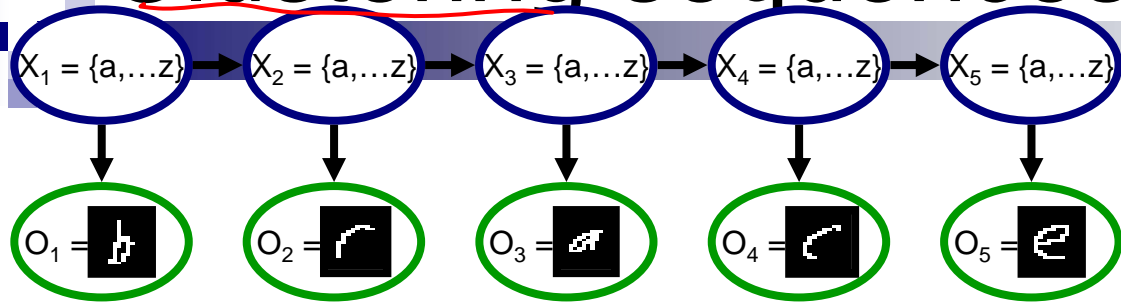
$$Q^{(t+1)}(\mathbf{x} | \mathbf{o}) = P(\mathbf{x} | \mathbf{o}, \theta^{(t)})$$



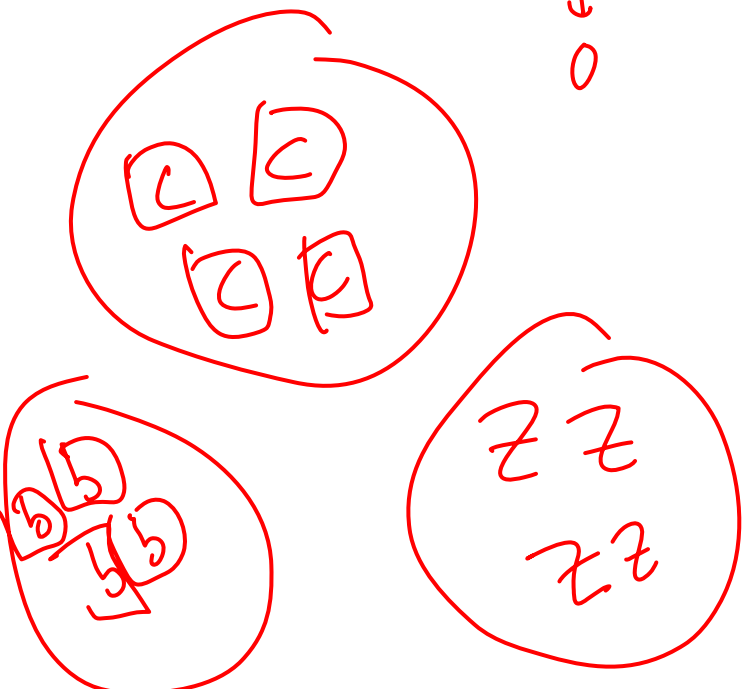
- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}
- Must compute:
 - $Q(x_t = a | \mathbf{o})$ – marginal probability of each position
 - Just forwards-backwards! $P(x_t = a | o_1 \dots o_n)$
 - $Q(x_{t+1} = a, x_t = b | \mathbf{o})$ – joint distribution between pairs of positions
 - Homework! 😊

What can you do with EM for HMMs? 1

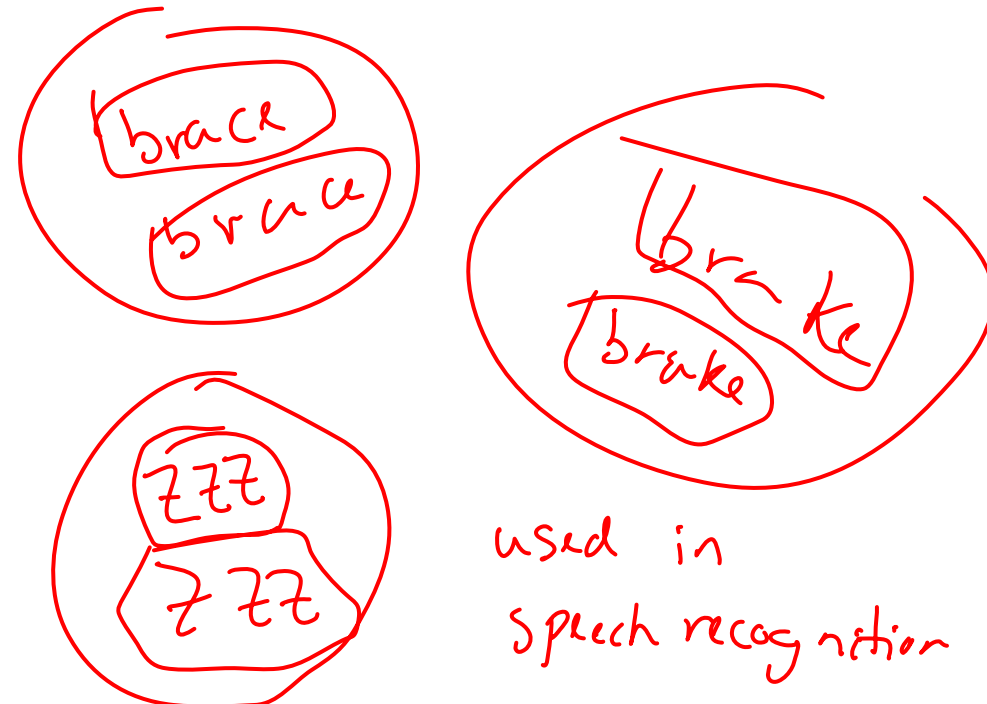
- Clustering sequences



Independent clustering:

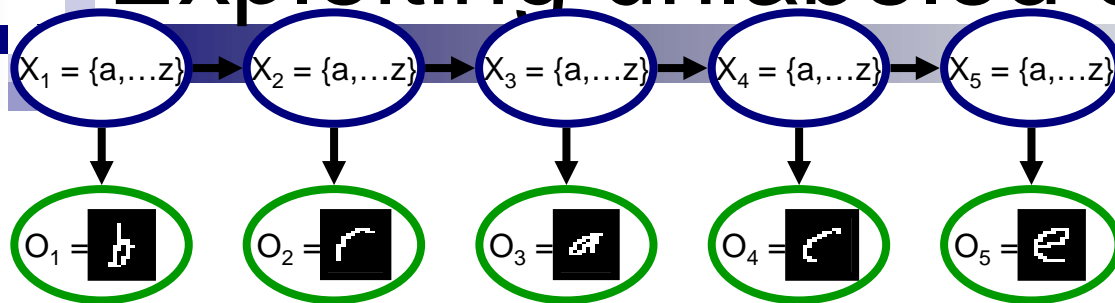


Sequence clustering:



What can you do with EM for HMMs? 2

– Exploiting unlabeled data



- Labeling data is hard work → save (graduate student) time by using both labeled and unlabeled data

- Labeled data:

- $\langle X = \text{"brace"}, O = \text{b r a d z} \rangle$
⋮

- Unlabeled data:

- $\langle X = \text{?????}, O = \text{b r a d z} \rangle$

Exploiting unlabeled data in clustering

- A few data points are labeled

- $\langle x, o \rangle$ $p_1 = \langle x=1, o = \{0.5, 0.8\} \rangle$
 $o_1 \quad o_2$

- Most points are unlabeled

- $\langle ?, o \rangle$ $p_2 = \langle x=?, o = \{0.4, 0.9\} \rangle$
 $o_1 \quad o_2$

- In the E-step of EM:

- If i'th point is unlabeled: $p(x=j|o_i)$

- compute $Q(X|o_i)$ as usual
- point generate by cluster j*

- If i'th point is labeled:

- set $Q(X=x|o_i)=1$ and $Q(X \neq x|o_i)=0$

- M-step as usual

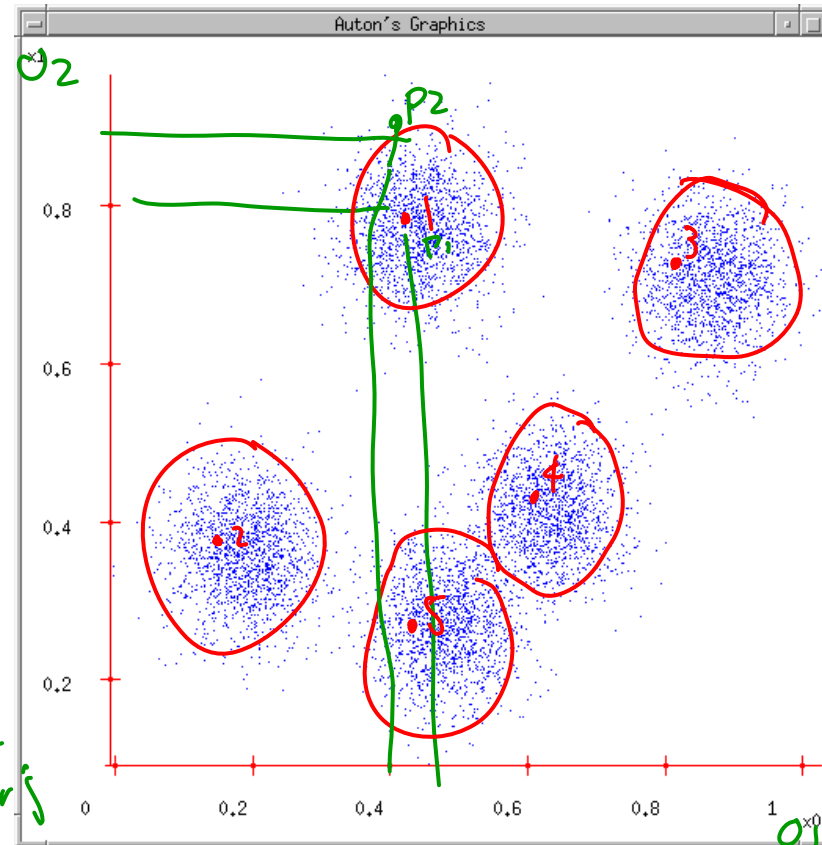


Table 3. Lists of the words most predictive of the course class in the WebKB data set, as they change over iterations of EM for a specific trial. By the second iteration of EM, many common course-related words appear. The symbol *D* indicates an arbitrary digit.

Iteration 0	Iteration 1	Iteration 2
intelligence	<i>DD</i>	<i>D</i>
<i>DD</i>	<i>D</i>	<i>DD</i>
artificial	lecture	lecture
understanding	cc	cc
<i>DDw</i>	<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

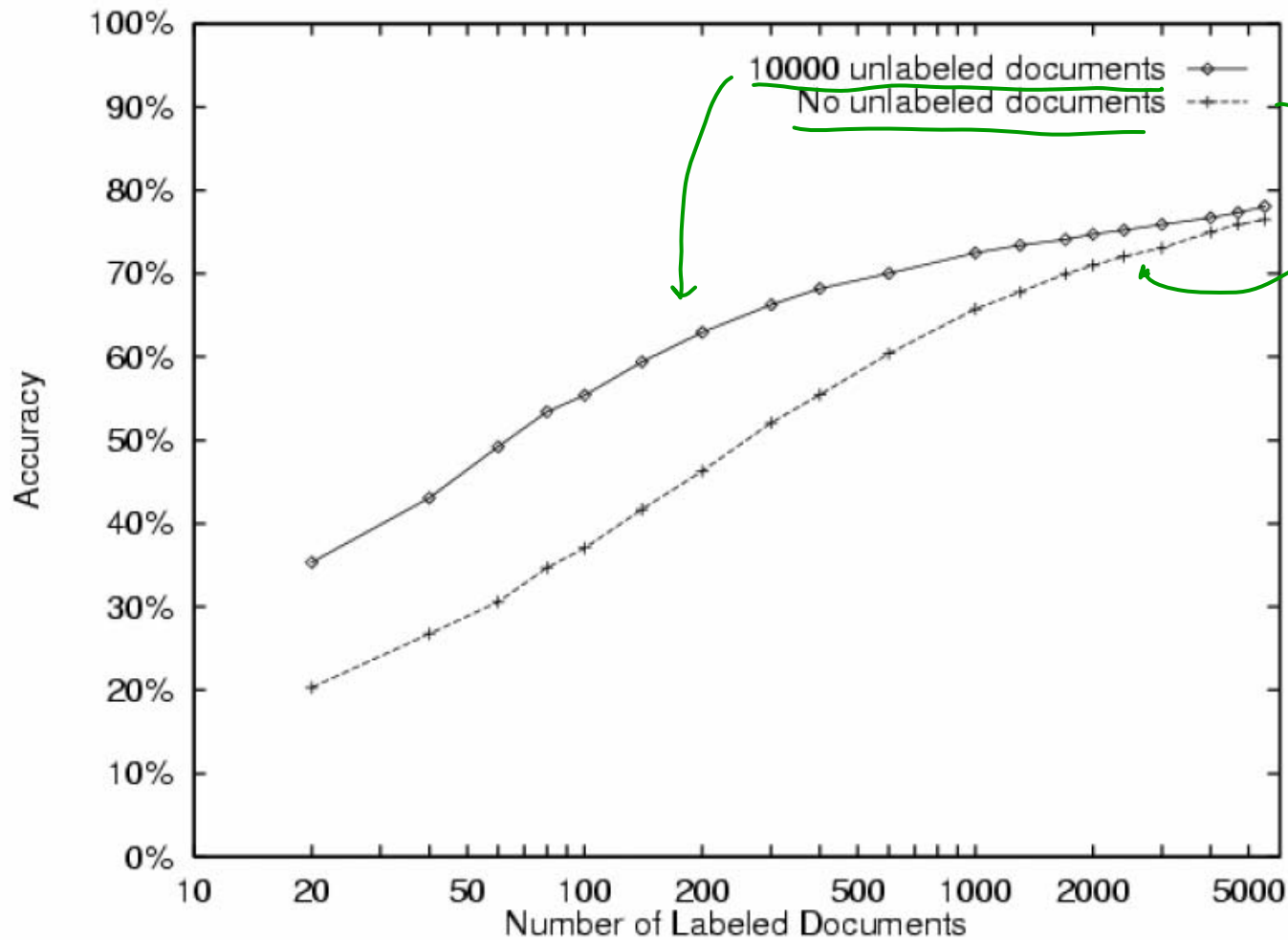
Using one labeled example per class

words likely associated with class = course

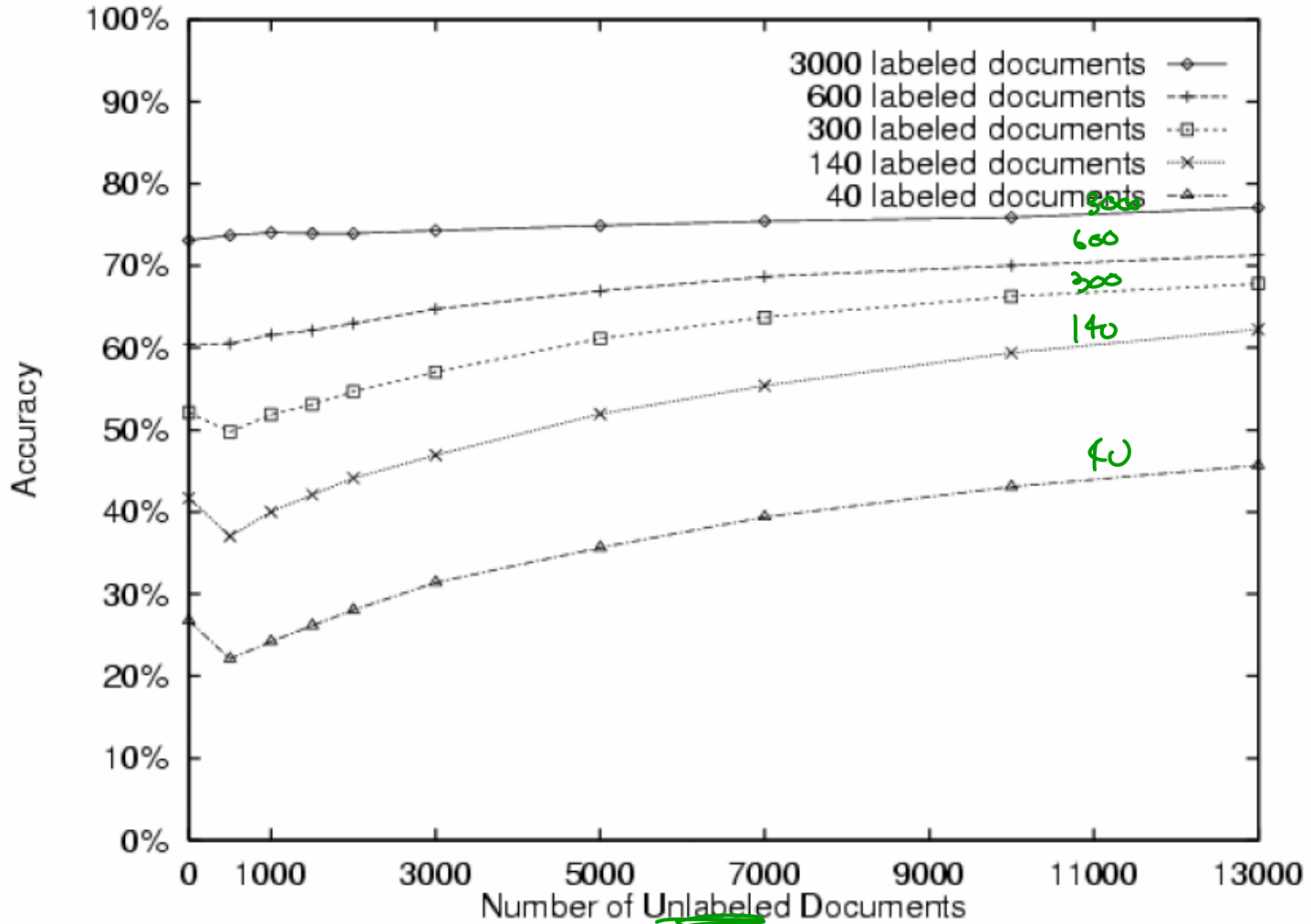
$P(\text{word} | X = \text{course})$ is high

quiz

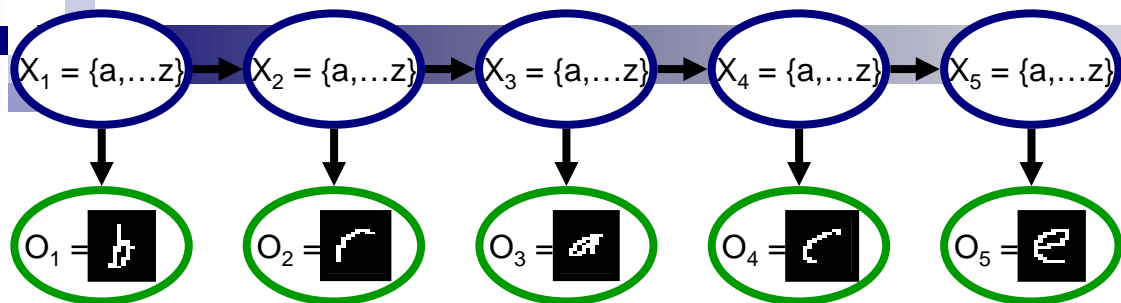
20 Newsgroups data – advantage of adding unlabeled data



20 Newsgroups data – Effect of additional unlabeled data



Exploiting unlabeled data in HMMs



- A few data points are labeled

□ $\langle x, o \rangle$ $\langle x = \text{"brace"}, o = \{b\} \{r\} \{a\} \{c\} \{e\} \rangle$

- Most points are unlabeled

□ $\langle ?, o \rangle$ $\langle x = \text{"?????"}, o = \{b\} \{r\} \{a\} \{c\} \{e\} \rangle$

- In the E-step of EM:

- If i'th point is unlabeled:
 - compute $Q(X|o_i)$ as usual
- If i'th point is labeled:
 - set $Q(X=x|o_i)=1$ and $Q(X \neq x|o_i)=0$

- M-step as usual

- Speed up by remembering counts for labeled data

$Q(X_3 = a | \{b\} \{r\} \{a\} \{c\} \{e\}) = 1$
 $Q(X_3 \neq a | \{b\} \{r\} \{a\} \{c\} \{e\}) = 0$

(pseudo-counts)₂₄

What you need to know

- Baum-Welch = EM for HMMs
- E-step:
 - Inference using forwards-backwards
- M-step:
 - Use weighted counts
- Exploiting unlabeled data:
 - Some unlabeled data can help classification
 - Small change to EM algorithm
 - In E-step, only use inference for unlabeled data

Acknowledgements



- Experiments combining labeled and unlabeled data provided by Tom Mitchell



EM for Bayes Nets

Machine Learning – 10701/15781

Carlos Guestrin

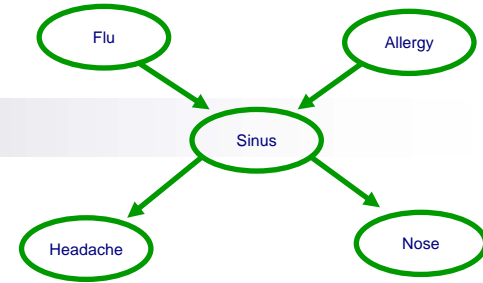
Carnegie Mellon University

April 12th, 2006

Data likelihood for BNs

- Given structure, log likelihood of fully observed data:

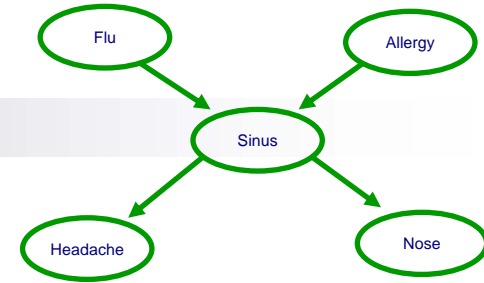
$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$



Marginal likelihood

- What if S is hidden?

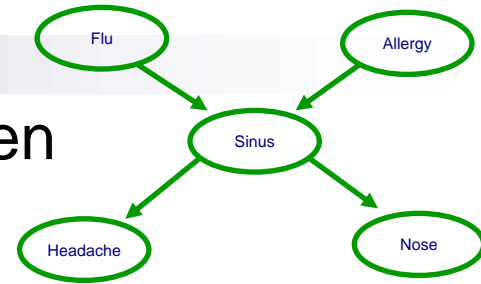
$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$



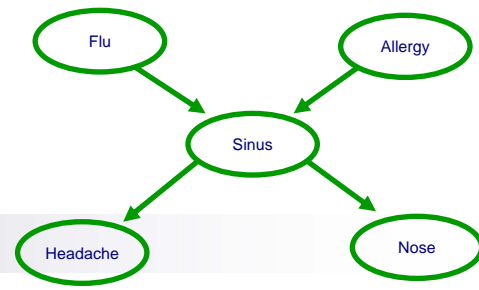
Log likelihood for BNs with hidden data

- Marginal likelihood – \mathbf{O} is observed, \mathbf{H} is hidden

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \sum_{j=1}^m \log P(\mathbf{o}^{(j)} \mid \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{h}} P(\mathbf{h}, \mathbf{o}^{(j)} \mid \theta)\end{aligned}$$



E-step for BNs

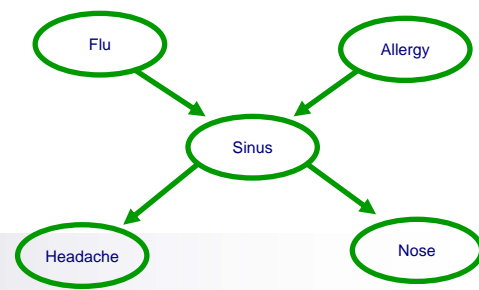


- E-step computes probability of hidden vars \mathbf{h} given \mathbf{o}

$$Q^{(t+1)}(\mathbf{x} \mid \mathbf{o}) = P(\mathbf{x} \mid \mathbf{o}, \theta^{(t)})$$

- Corresponds to inference in BN

The M-step for BNs



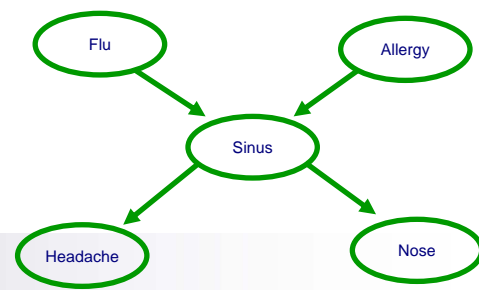
- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{\mathbf{x}} Q^{(t+1)}(\mathbf{h} | \mathbf{o}) \log P(\mathbf{h}, \mathbf{o} | \theta)$$

- Use expected counts instead of counts:

- If learning requires $\text{Count}(\mathbf{h}, \mathbf{o})$
- Use $E_{Q^{(t+1)}}[\text{Count}(\mathbf{h}, \mathbf{o})]$

M-step for each CPT



- M-step decomposes per CPT

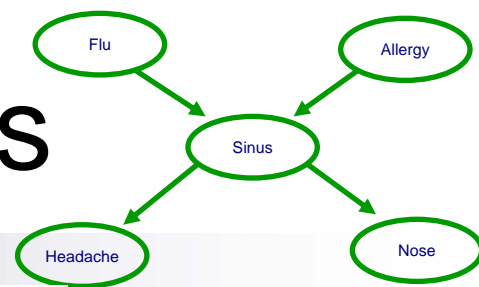
- Standard MLE:

$$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\text{Count}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\text{Count}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

- M-step uses expected counts:

$$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\text{ExCount}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\text{ExCount}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

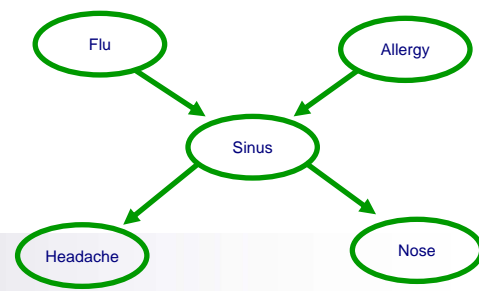
Computing expected counts



$$P(X_i = x_i | \text{Pa}_{X_i} = \mathbf{z}) = \frac{\text{ExCount}(X_i = x_i, \text{Pa}_{X_i} = \mathbf{z})}{\text{ExCount}(\text{Pa}_{X_i} = \mathbf{z})}$$

- M-step requires expected counts:
 - For a set of vars \mathbf{A} , must compute $\text{ExCount}(\mathbf{A}=\mathbf{a})$
 - Some of \mathbf{A} in example j will be observed
 - denote by $\mathbf{A}_O = \mathbf{a}_O^{(j)}$
 - Some of \mathbf{A} will be hidden
 - denote by \mathbf{A}_H
- Use inference (E-step computes expected counts):
 - $\text{ExCount}^{(t+1)}(\mathbf{A}_O = \mathbf{a}_O^{(j)}, \mathbf{A}_H = \mathbf{a}_H) \leftarrow P(\mathbf{A}_H = \mathbf{a}_H | \mathbf{A}_O = \mathbf{a}_O^{(j)}, \theta^{(t)})$

Data need not be hidden in the same way



- When data is fully observed
 - A data point is
- When data is partially observed
 - A data point is
- But unobserved variables can be different for different data points
 - e.g.,
- Same framework, just change definition of expected counts
 - $\text{ExCount}^{(t+1)}(\mathbf{A}_O = \mathbf{a}_O^{(i)}, \mathbf{A}_H = \mathbf{a}_H) \leftarrow P(\mathbf{A}_H = \mathbf{a}_H \mid \mathbf{A}_O = \mathbf{a}_O^{(i)}, \theta^{(t)})$

What you need to know

- EM for Bayes Nets
- E-step: inference computes expected counts
 - Only need expected counts over X_i and \mathbf{Pa}_{X_i}
- M-step: expected counts used to estimate parameters
- Hidden variables can change per datapoint

- Use labeled and unlabeled data → some data points are complete, some include hidden variables