

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

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 11<sup>th</sup>, 2007

©2005-2007 Carlos Guestrin

# The general learning problem with missing data

- Marginal likelihood –  $\mathbf{x}$  is observed,  $\mathbf{z}$  is missing:

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log \prod_{j=1}^m P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} | \theta)\end{aligned}$$

*observed parts*

*sum over (marginalize out) hidden vars*

# EM is coordinate ascent

marginal likelihood

$$\ell(\theta : \mathcal{D}) \geq \underset{\uparrow \text{max.}}{F(\theta, Q)} = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}$$

- **M-step:** Fix  $Q$ , maximize  $F$  over  $\theta$  (a lower bound on  $\ell(\theta : \mathcal{D})$ ):

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q^{(t)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m.H(Q^{(t)})$$

Expected counts.

- **E-step:** Fix  $\theta$ , maximize  $F$  over  $Q$ :

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \ell(\theta^{(t)} : \mathcal{D}) - \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))$$

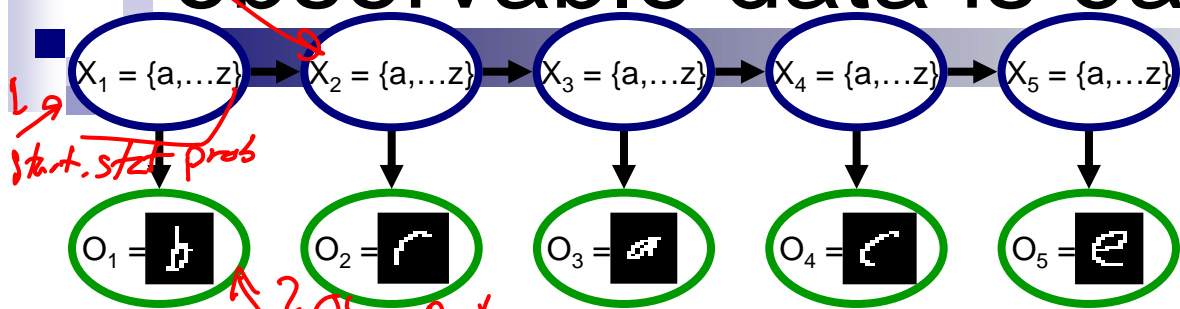
- “Realigns”  $F$  with likelihood:

$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

# What you should know about EM

- K-means for clustering:
  - algorithm
  - converges because it's coordinate ascent
- EM for mixture of Gaussians:
  - How to “learn” maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Be happy with this kind of probabilistic analysis
- Remember, E.M. can get stuck in local minima, and empirically it DOES
- EM is coordinate ascent
- General case for EM

# Learning HMMs from fully observable data is easy



Learn 3 distributions:

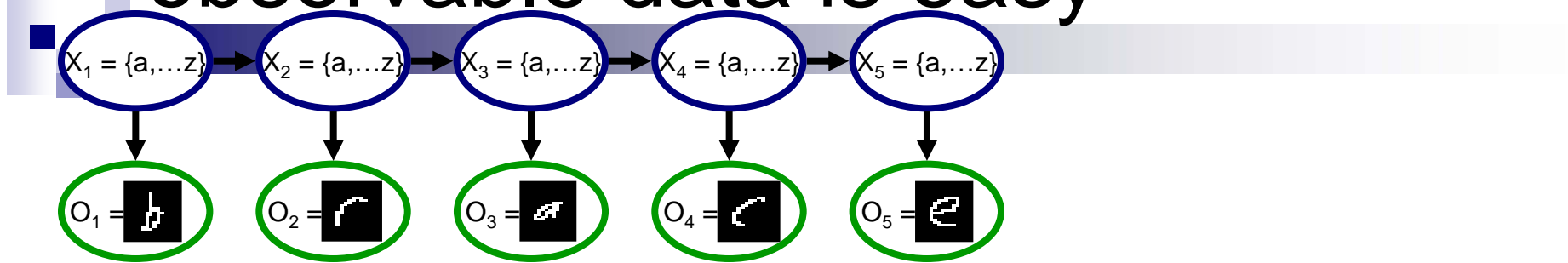
$$1. P(X_1) = \frac{\text{Count}(X_1 = a)}{m = \text{Count}(X_1 = ?)}$$

$$2. P(O_i | X_i) = \frac{\text{Count}(X_i = b, O_i = 'r')}{\text{Count}(X_i = b)}$$

$$3. P(X_i | X_{i-1}) = \frac{\text{Count}(X_{i-1} = g, X_i = c)}{\text{Count}(X_{i-1} = g, X_i = ?)}$$

$X_1 =$   
anything

# Learning HMMs from fully observable data is easy



Learn 3 distributions:

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

$$P(O_i^{\text{pixel 17 is white}} | X_i^a) = \frac{\text{count}(\text{pixel 17 was white, } X_i = a)}{N_i}$$

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

select training data where letter was a

What if O is observed,  
but X is hidden

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

$$\mathbf{O} = (o_1, \dots, o_n)$$

$$\mathbf{X} = (x_1, \dots, x_n)$$

for  $m$  sequences

$$\sum_{j=1}^m \log P(\mathbf{o}^{(j)} | \theta)$$

- Marginal likelihood –  $\mathbf{O}$  is observed,  $\mathbf{X}$  is missing

- For simplicity of notation, training data consists of only one sequence:

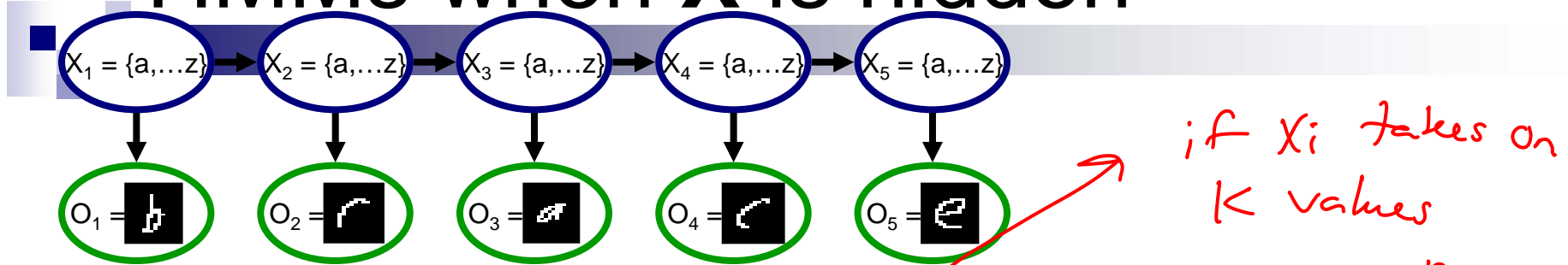
✓ *Observed*

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

- 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



if  $X_i$  takes on  $K$  values

Sum over  $K^n$  assignments

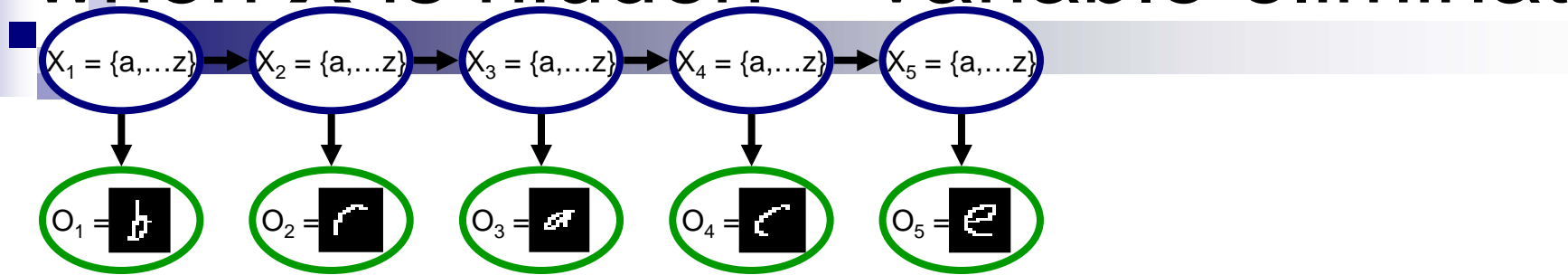
$$\begin{aligned} \ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} | \theta) \\ &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} | \theta) \end{aligned}$$

$$= \log \sum_{x_1} \sum_{x_2} \dots \sum_{x_n} P(x_1) \cdot P(o_1 | x_1) \cdot \prod_{i=2}^n P(x_i | x_{i-1}) \cdot P(o_i | x_i)$$

$$= \log \sum_{x_1} \sum_{x_{n-1}} P(x_1) P(o_1 | x_1) \prod_{i=2}^{n-1} P(x_i | x_{i-1}) P(o_i | x_i) \underbrace{\sum_{x_n} P(x_n | x_{n-1}) \cdot P(o_n | x_n)}_{B_{n-1}(x_{n-1})}$$

use V.E. to  
compute  $\ell(\theta : \mathcal{D})$  in  $O(n)$  time

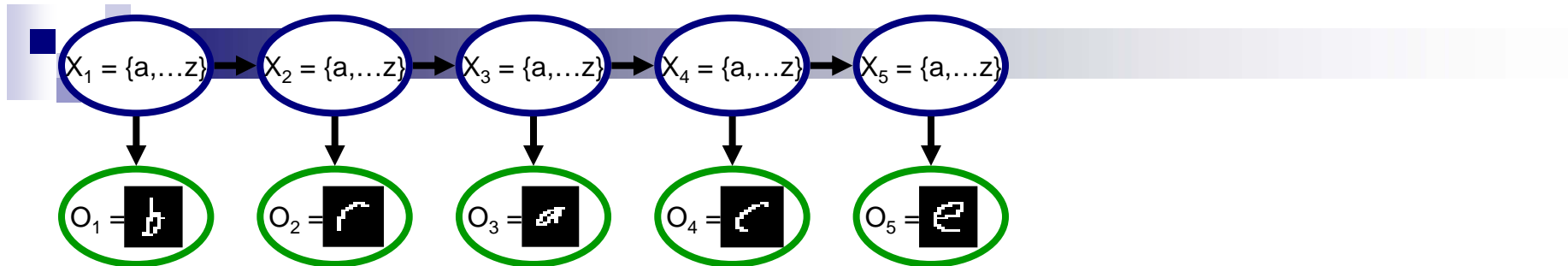
# Computing Log likelihood for HMMs when $\mathbf{X}$ is hidden – variable elimination



- Can compute efficiently with variable elimination:

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

# EM for HMMs when $X$ is hidden



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

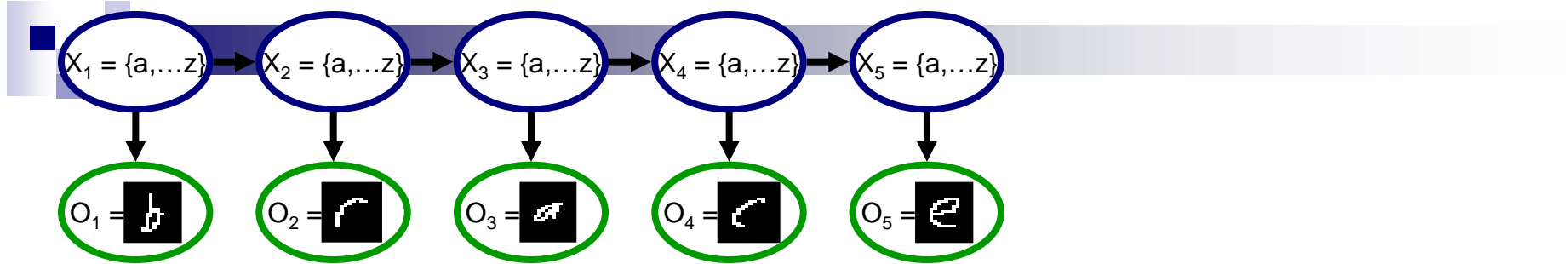
need to compute:  $P(X_i, X_{i+1} | O_1 \dots O_n)$ , use  $\downarrow$  forwards-backwards  
 $= Q(X_i, X_{i+1} | O_1 \dots O_n)$  (faster version of Var. elim.)

- M-step: Recompute parameters with weighted data

if fully observable:  $\hat{P}(X_i = a) = \frac{\text{Count}(X_i = a)}{n}$

if hidden vars:  $\hat{P}(X_i = a) = \frac{\sum_{j=1}^m Q(X_i = a | O^{(j)})}{m}$

# E-step



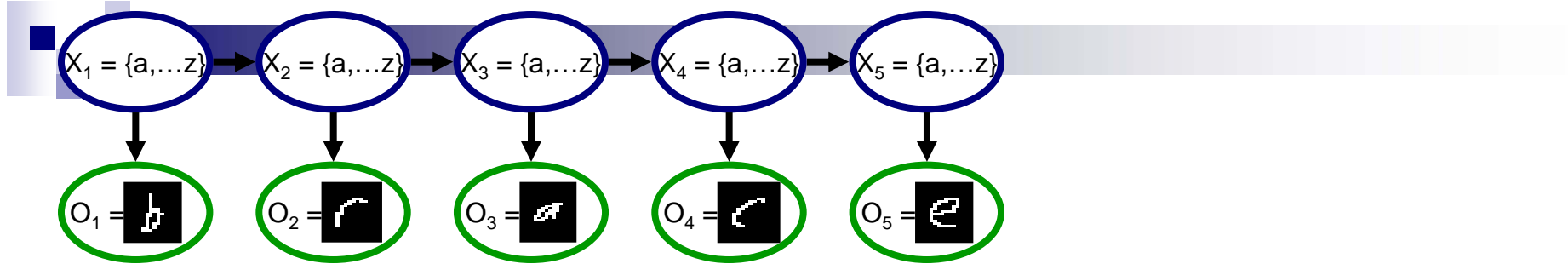
- E-step computes probability of hidden vars  $\mathbf{x}$  given  $\mathbf{o}^{(j)}$

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

*j<sup>th</sup> training example*

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

# The M-step



## ■ Maximization step:

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

*weighted log likelihood*

## ■ Use expected counts instead of counts:

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

$$E_{Q^{(t+1)}}[\text{Count}(\mathbf{x} = \{a, b, c\}, \mathbf{o} = [\boxed{a} \boxed{b}, \boxed{a}])] = \frac{\sum_{j=1}^m Q^{(t+1)}(\mathbf{x} = \{a, b, c\} | \mathbf{o} = [\boxed{a} \boxed{b}, \boxed{a}])}{m}$$

# Decomposition of likelihood $P(X_1) \in \Theta_{X_1}$

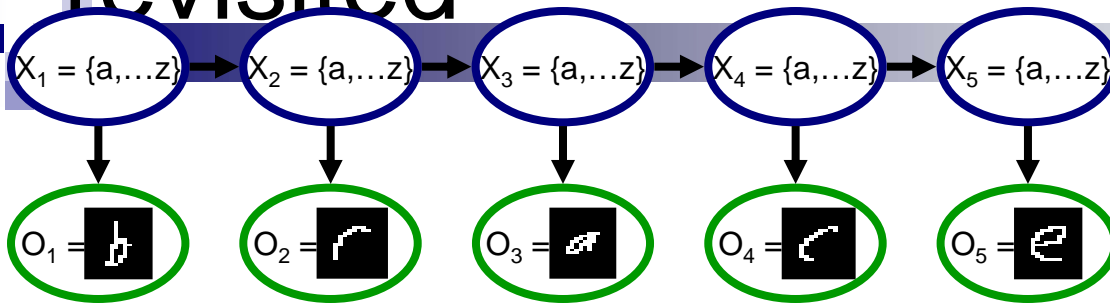
revisited

$$\log a.b = \log a + \log b$$

$$P(O_i | X_i) \in \Theta_{O|X}$$

$$P(X_i | X_{i-1}) \in \Theta_{X_i|X_{i-1}}$$

$$\Theta = (\Theta_{X_1}, \Theta_{X_i|X_{i-1}}, \Theta_{O|X})$$



■ Likelihood optimization decomposes:

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

HMM structure

$$\begin{aligned} & \max_{\theta} \sum_x Q(x | o) \log P(x_1 | \theta_{X_1}) P(o_1 | x_1, \theta_{O|X}) \prod_{t=2}^n P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) P(o_t | x_t, \theta_{O|X}) \\ & \stackrel{\substack{\max_{\theta_{X_1}} \max_{\theta_{O|X}} \max_{\theta_{X_t|X_{t-1}}}}}{=} \max_{\theta} \sum_x Q(x | o) \left[ \log P(x_1 | \theta_{X_1}) + \sum_{t=1}^n \log P(o_t | x_t, \theta_{O|X}) + \sum_{t=2}^n \log P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) \right] \\ & = \underbrace{\left[ \max_{\theta_{X_1}} \sum_x Q(x | o) \log P(x_1 | \theta_{X_1}) \right]}_{\theta_{X_1} \text{ stat. state prior.}} + \underbrace{\left[ \max_{\theta_{O|X}} \sum_x Q(x | o) \sum_{t=1}^n \log P(o_t | x_t, \theta_{O|X}) \right]}_{\theta_{O|X}} + \\ & \quad + \underbrace{\left[ \max_{\theta_{X_t|X_{t-1}}} \sum_x Q(x | o) \sum_{t=2}^n \log P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) \right]}_{\theta_{X_t|X_{t-1}}} \end{aligned}$$

sequences:

# Starting state probability $P(X_1)$

Using expected counts

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

derivation for one sequence

$$\max_{\theta_{X_1}} \sum_{j=1}^m Q(x | o^{(j)}) \log P(x_1 | \theta_{X_1})$$

$Q(X|o)$  explicitly, need  $K^n - 1$  parameters chain rule

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

$$= \max_{\theta_{X_1}} \sum_{j=1}^m \sum_{x_1} Q(x_1 | o^{(j)}) \log P(x_1 | \theta_{X_1}) \cdot \sum_{x_2} \sum_{x_n} Q(x_2 \dots x_n | o^{(j)}, x_1)$$

$$= \max_{\theta_{X_1}} \sum_{j=1}^m \sum_{x_1} Q(x_1 | o^{(j)}) \log P(x_1 | \theta_{X_1})$$

1 (it's a prob.)

only need  $Q(x_1 | o)$

$\rightarrow K-1$  parameters

$\rightarrow \text{argmax}$

$o \in \text{incomes in dataset}$

learning from weighted data

argmax over all sequences:

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

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

## ■ Using expected counts

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

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

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

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

only need  $\forall t \in 2:n, Q(x_{t-1}, x_t | \mathbf{o})$   
 $= (n-1)(K^2-1)$  parameters

$\sum_{x_{t-1}, x_t} \sum_{x_{t-2}, x_{t-1}, x_{t+1}, x_{t+2}} Q(x_{t-2}, x_{t-1}, x_t, x_{t+1}, x_{t+2} | \mathbf{o}^{(j)})$

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

# Observation probability $P(O_t|X_t)$

## ■ Using expected counts

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

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

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

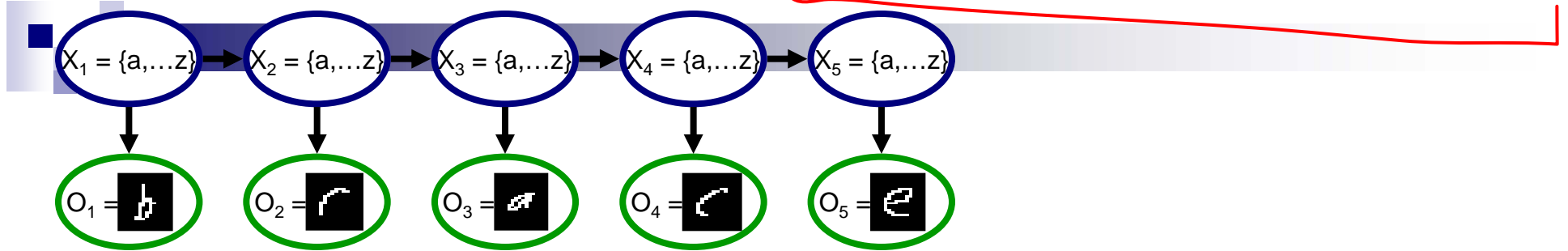
argmax  $\rightarrow$  only need marginals  $Q(x_t | o)$  for a total of  $n \cdot (K-1)$  terms

only count positions where  $O_t=a$

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

# 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  $\mathbf{o}$

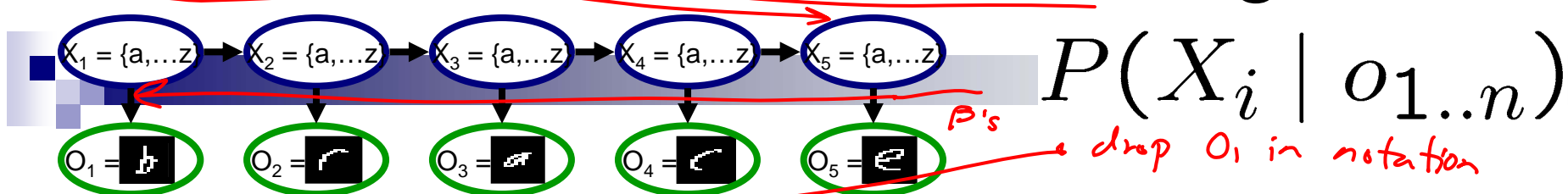
■ Must compute:

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

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

*use forwards-backwards*

# 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* *transition prob*

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

■ Initialization:  $\beta_n(X_n) = 1$

■ For  $i = n-1$  to  $1$

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

*forall xi*

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

*xi*

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

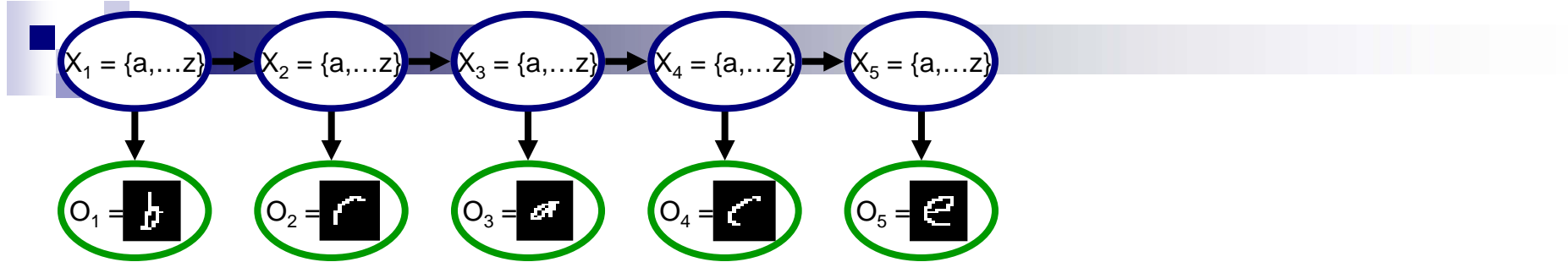
*α<sub>n</sub>(X<sub>n</sub>)*  
*normalized*  
 $= P(X_n | o_{1:n})$

*β<sub>1</sub>(X<sub>1</sub>) α<sub>1</sub>(X<sub>1</sub>)*  
*normalized*  
 $= P(X_1 | o_{1:n})$

*α<sub>5</sub>(a)*  
*α<sub>5</sub>(b)*  
*⋮*  
*α<sub>5</sub>(z)*

# E-step revisited

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

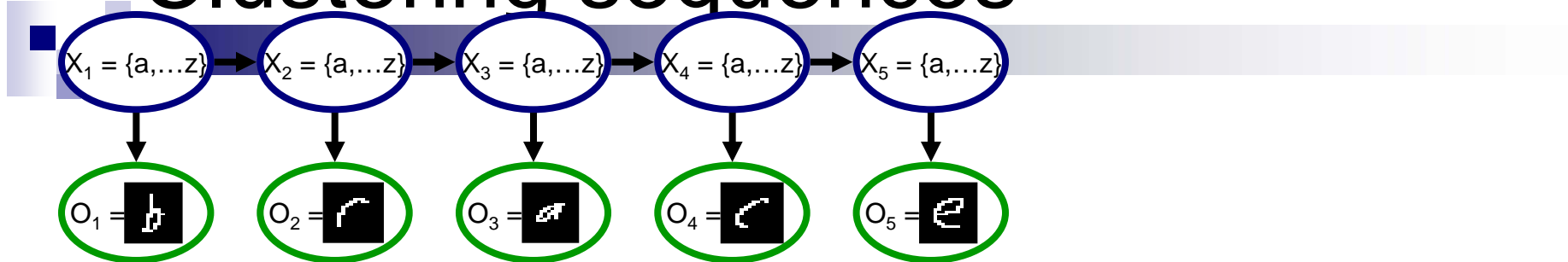


- E-step computes probability of hidden vars  $\mathbf{x}$  given  $\mathbf{o}$
- Must compute:
  - $Q(x_t=a|\mathbf{o})$  – marginal probability of each position
    - Just forwards-backwards!
  - $Q(x_{t+1}=a, x_t=b|\mathbf{o})$  – joint distribution between pairs of positions
 

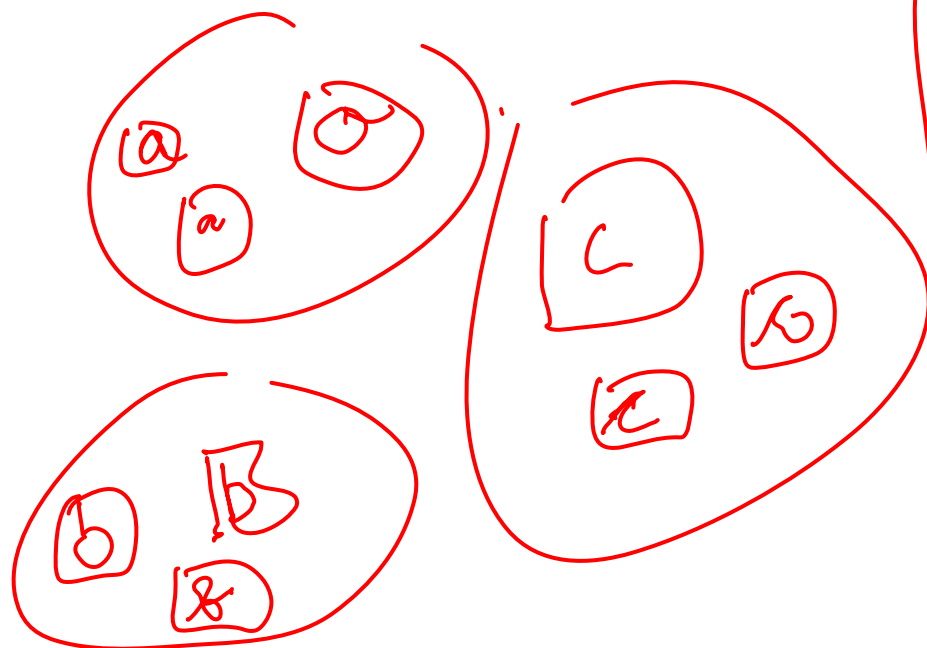
*see reading*      *[simple eqn.]*      *[maybe homework]*

# What can you do with EM for HMMs? 1

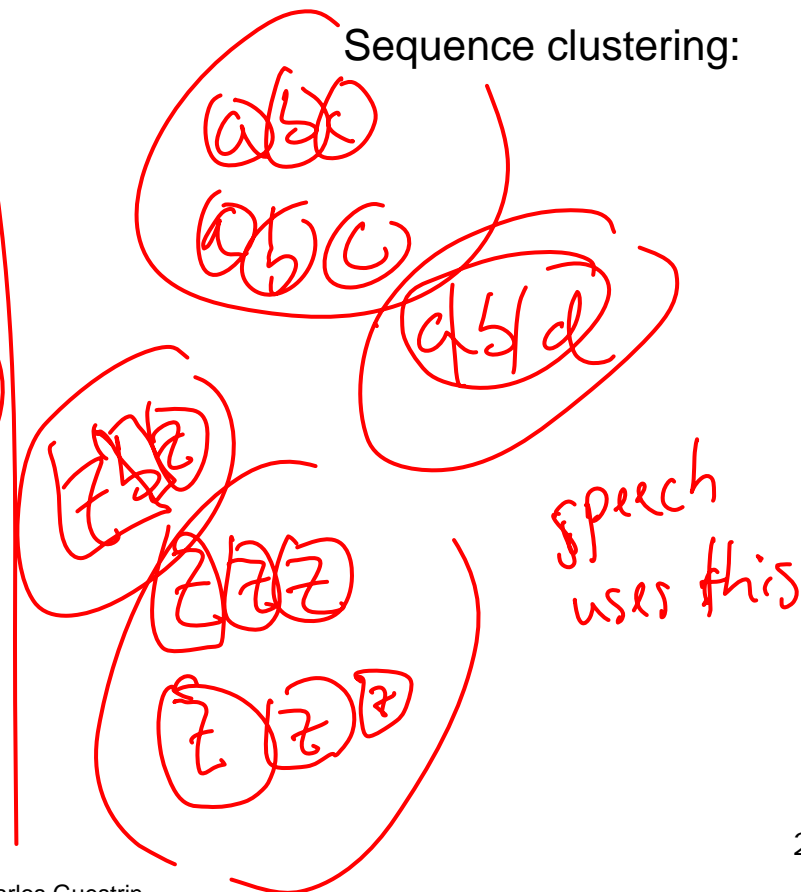
## – Clustering sequences



Independent clustering:

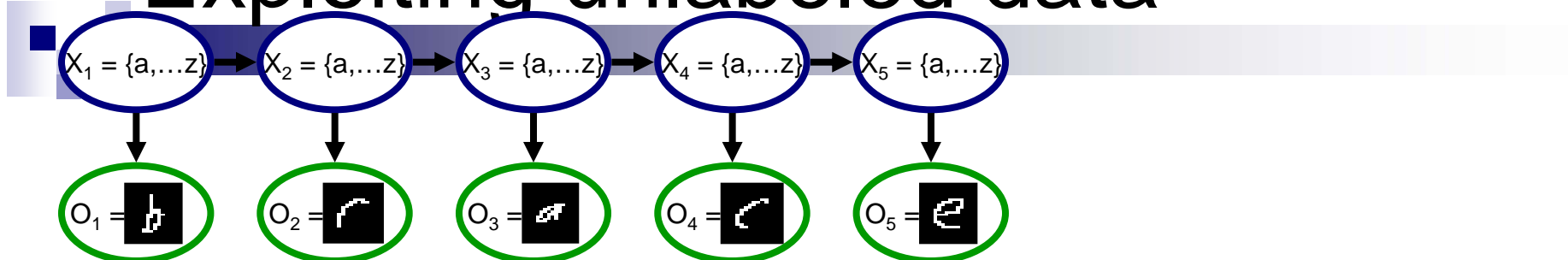


Sequence clustering:



# 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:

- $\langle X = \text{"brace"}, O = \text{B A C D E} \rangle$

- Unlabeled data:

- $\langle X = \text{?????}, O = \text{B A C D E} \rangle$

# Exploiting unlabeled data in clustering

- A few data points are labeled

- $\langle x, o \rangle$

- Most points are unlabeled

- $\langle ?, o \rangle$

- 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~~ <sup>correct label</sup>

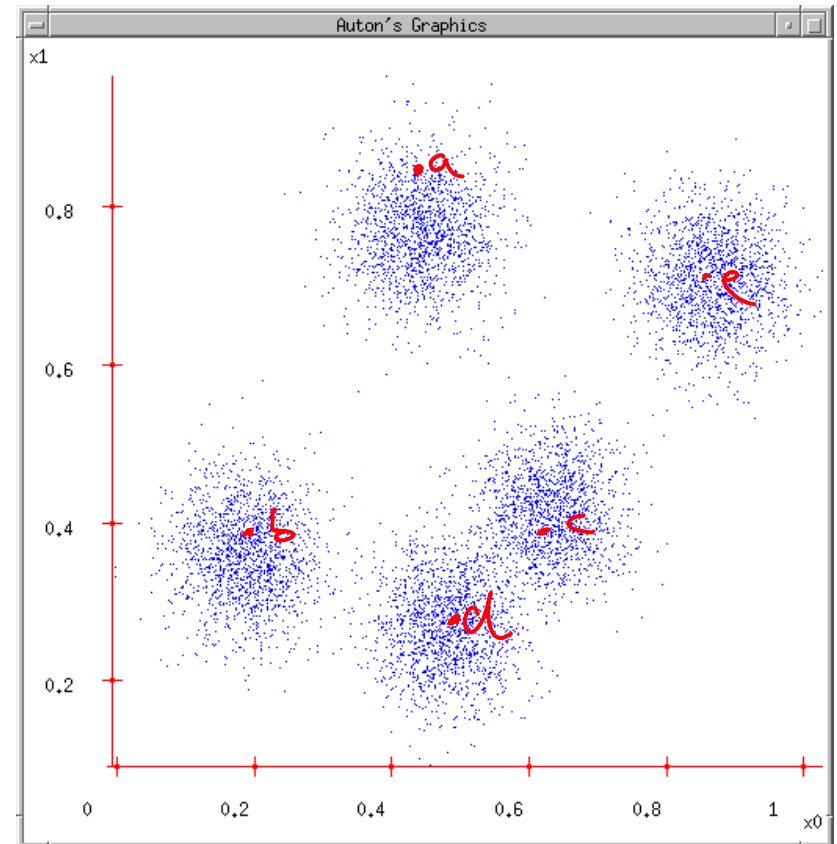
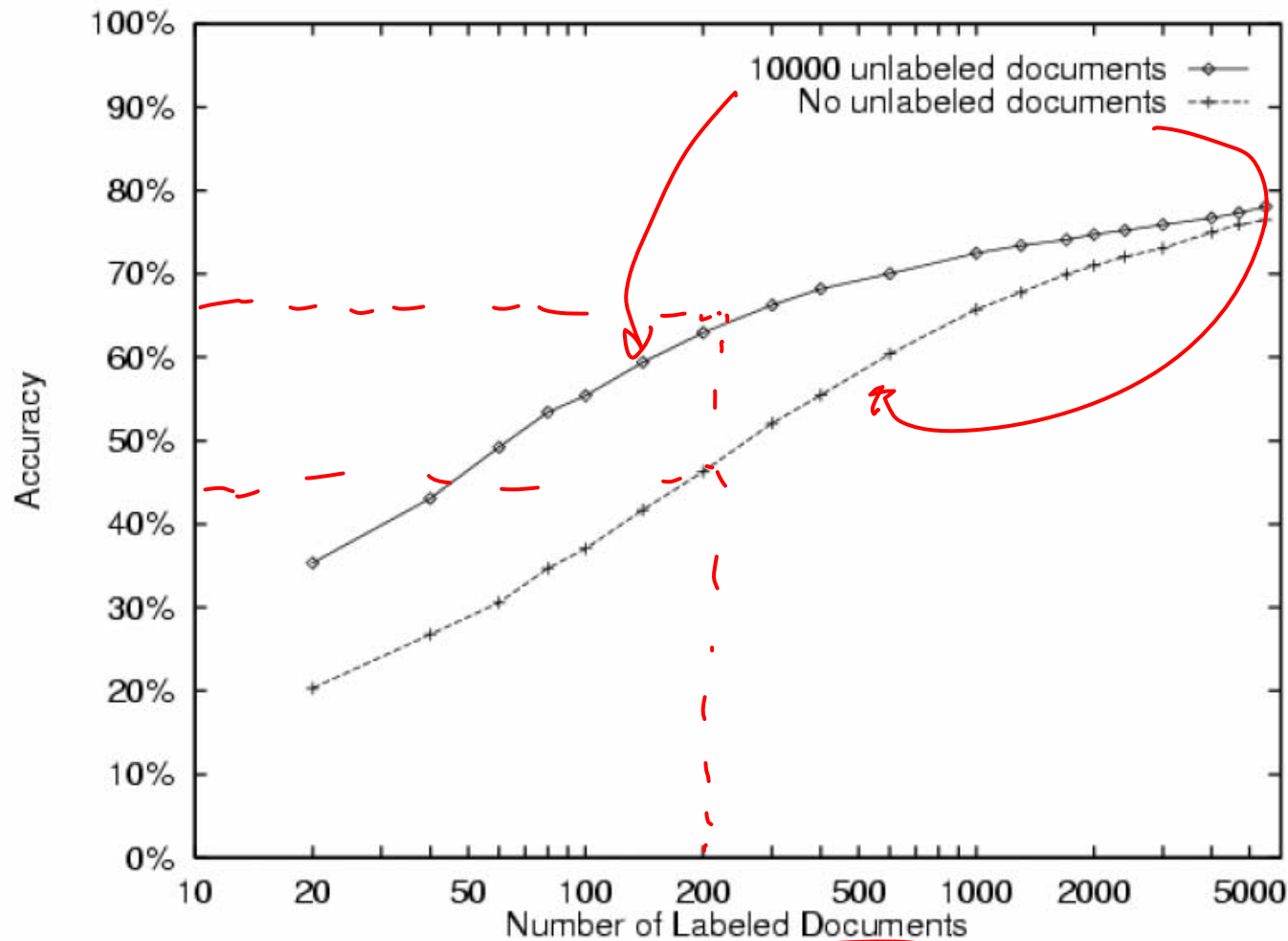


Table 3. Lists of the words most predictive of the course class in the ~~WebKB~~ data set, as they change over iterations of EM for a specific trial. By the second iteration of EM, many common course-related words appear. The symbol *D* indicates an arbitrary digit.

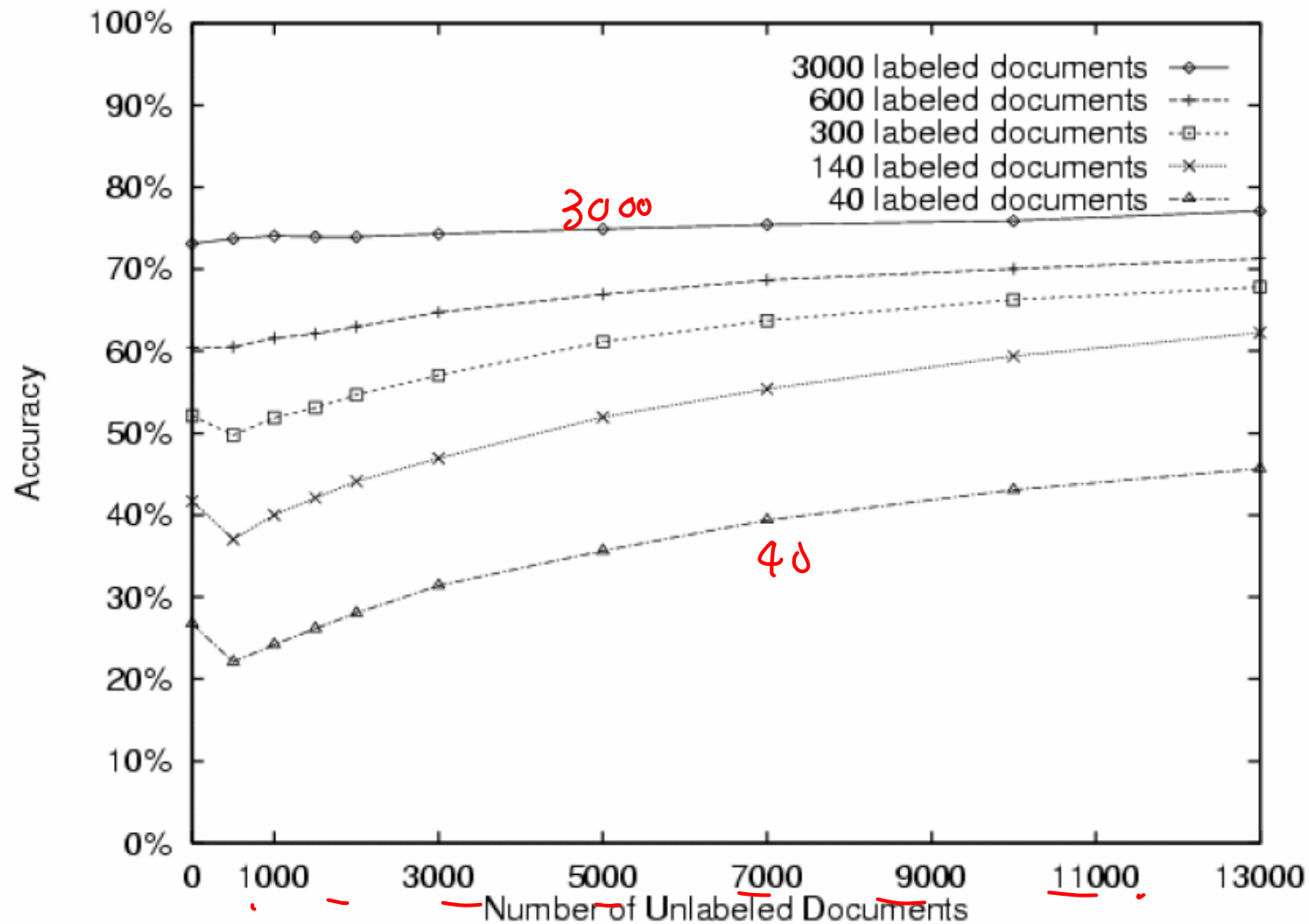
Iteration 0	Iteration 1	Iteration 2
intelligence	<i>DD</i>	<i>D</i>
<i>DD</i>	<i>D</i>	<i>DD</i>
artificial	lecture	<u>lecture</u>
understanding	cc	cc
<i>DDw</i>	<i>D*</i>	<i>DD:DD</i>
dist	<i>DD:DD</i>	<u>due</u>
identical	handout	<i>D*</i>
rus	due	homework
arrange	problem	<u>assignment</u>
games	set	<u>handout</u>
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	<u>quiz</u>
representation	solution	chapter
field	assaf	ascii

Using one  
labeled  
example per  
class

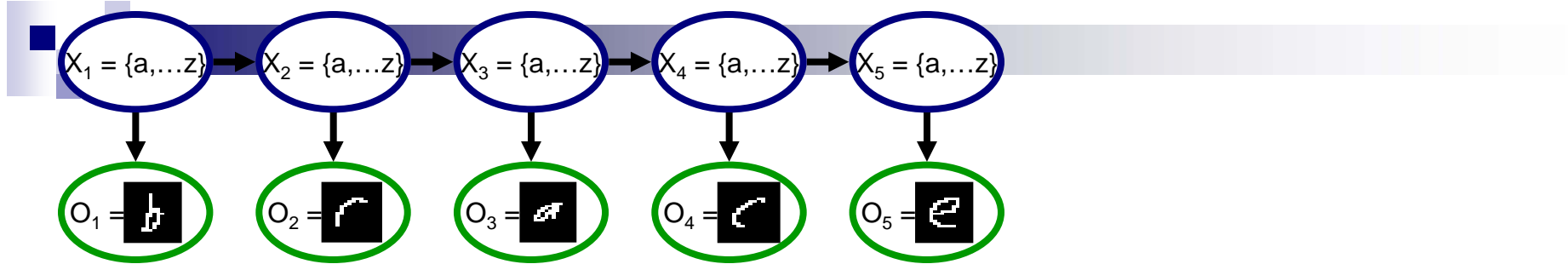
# 20 Newsgroups data – advantage of adding unlabeled data



# 20 Newsgroups data – Effect of additional unlabeled data



# Exploiting unlabeled data in HMMs



- A few data points are labeled
  - $\langle x, o \rangle$
- Most points are unlabeled
  - $\langle ?, o \rangle$
- 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

# 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