# EM for Bayes Nets

Machine Learning – 10701/15781
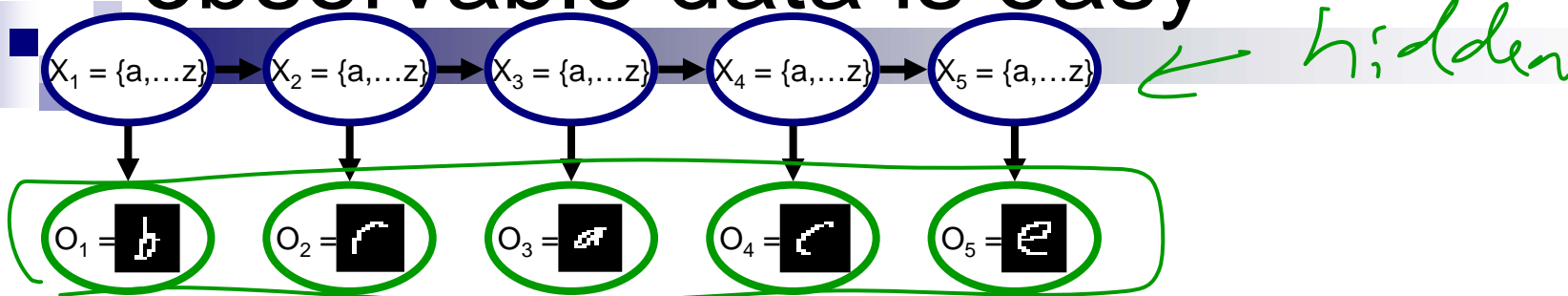
Carlos Guestrin

Carnegie Mellon University

April 16th, 2007

1

# Learning HMMs from fully observable data is easy

$X_1 = \{a,...z\}$ → $X_2 = \{a,...z\}$ → $X_3 = \{a,...z\}$ → $X_4 = \{a,...z\}$ → $X_5 = \{a,...z\}$  ← hidden

$O_1 =$ [img]   $O_2 =$ [img]   $O_3 =$ [img]   $O_4 =$ [img]   $O_5 =$ [img]

**Learn 3 distributions:**

$$P(X_1^{=a}) = \frac{\text{Count (\# first letter was a)}}{N = \text{data set size}}$$

select training data where letter was a

$$P(O_i^{=\text{pixel 17 is white}} \mid X_i^{=a}) = \frac{\text{Count ( pixel 17 was white, } X_i = a)}{\text{any...in}}$$

$$P(X_i^{=a} \mid X_{i-}^{=b}$$

> # What if **O** is observed, but **X** is hidden

# Log likelihood for HMMs when **X** is hidden

$\mathbf{O} = (O_1, \ldots, O_n)$

$X = (x_1, \ldots, x_n)$

for m sequences

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

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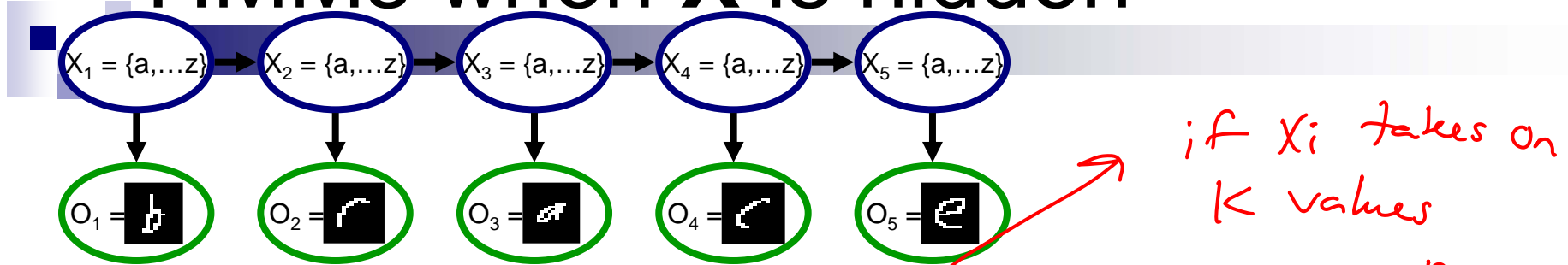  - For simplicity of notation, training data consists of only one sequence:

  ✓ Observed

$$\ell(\theta : \mathcal{D}) = \log P(\mathbf{o} \mid \theta)$$

$$= \log \sum_{\mathbf{X}} P(\mathbf{x}, \mathbf{o} \mid \theta)$$

$$P(X_1 | \theta) \cdot P(O_1 | X_1, \theta) \cdot \prod_{t=2}^{n} P(X_t | X_{t-1}, \theta)$$

$$P(O_t | X_t, \theta)$$

  - If there were m sequences:

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

3

# Computing Log likelihood for HMMs when **X** is hidden



$X_1 = \{a,...z\}$ → $X_2 = \{a,...z\}$ → $X_3 = \{a,...z\}$ → $X_4 = \{a,...z\}$ → $X_5 = \{a,...z\}$

$O_1 = $ ... $O_2 = $ ... $O_3 = $ ... $O_4 = $ ... $O_5 = $ ...

if $X_i$ takes on $K$ values

Sum over $K^n$ assignments

$$\ell(\theta : \mathcal{D}) = \log P(\mathbf{o} \mid \theta)$$

$$= \log \sum_{\mathbf{X}} P(\mathbf{x}, \mathbf{o} \mid \theta)$$

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

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

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

**4**

# The M-step



- **Maximization step:**

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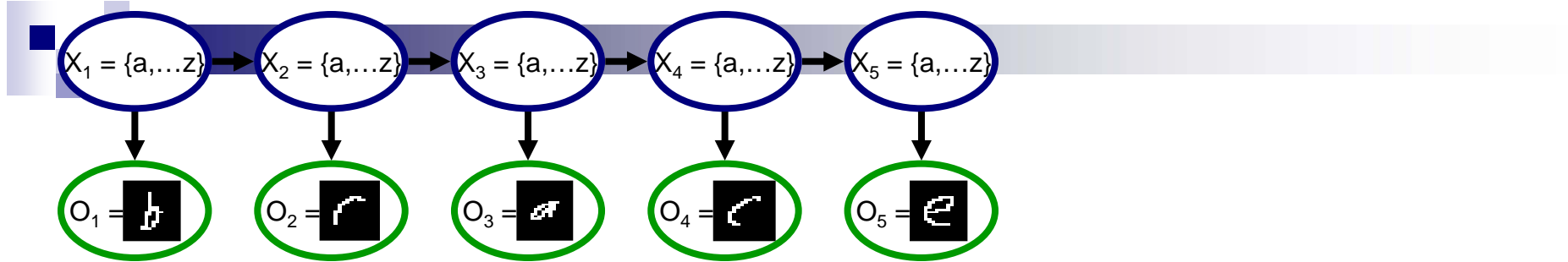
weighted log likelihood

- **Use expected counts instead of counts:**
  - □ If learning requires Count(**x**,**o**)
  - □ Use $E_{Q(t+1)}[\text{Count}(\mathbf{x},\mathbf{o})]$

$$E_{Q^{(t+1)}}\left[\text{Count}\left(X = \{a,b,c\}, O[\boxed{a}][b], [d]\right)\right] = \sum_{j=1}^{m} Q^{(t+1)}(x = \{a,b,c\}[\boxed{a}][b], [d])$$

# E-step revisited

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



- E-step computes probability of hidden vars **x** given **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 ]

# Exploiting unlabeled data in clustering

- A few data points are labeled
  - $<x,o>$

- Most points are unlabeled
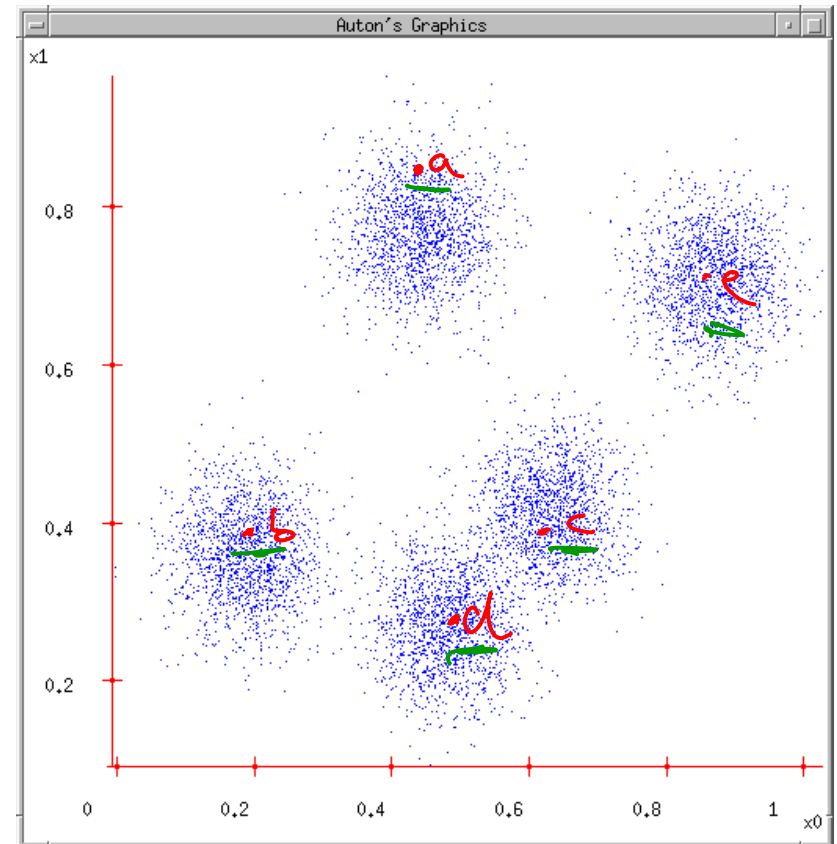  - $<?,o>$

- In the E-step of EM:
  - If i'th point is unlabeled:
    - compute $Q(X|o_i)$ as usual
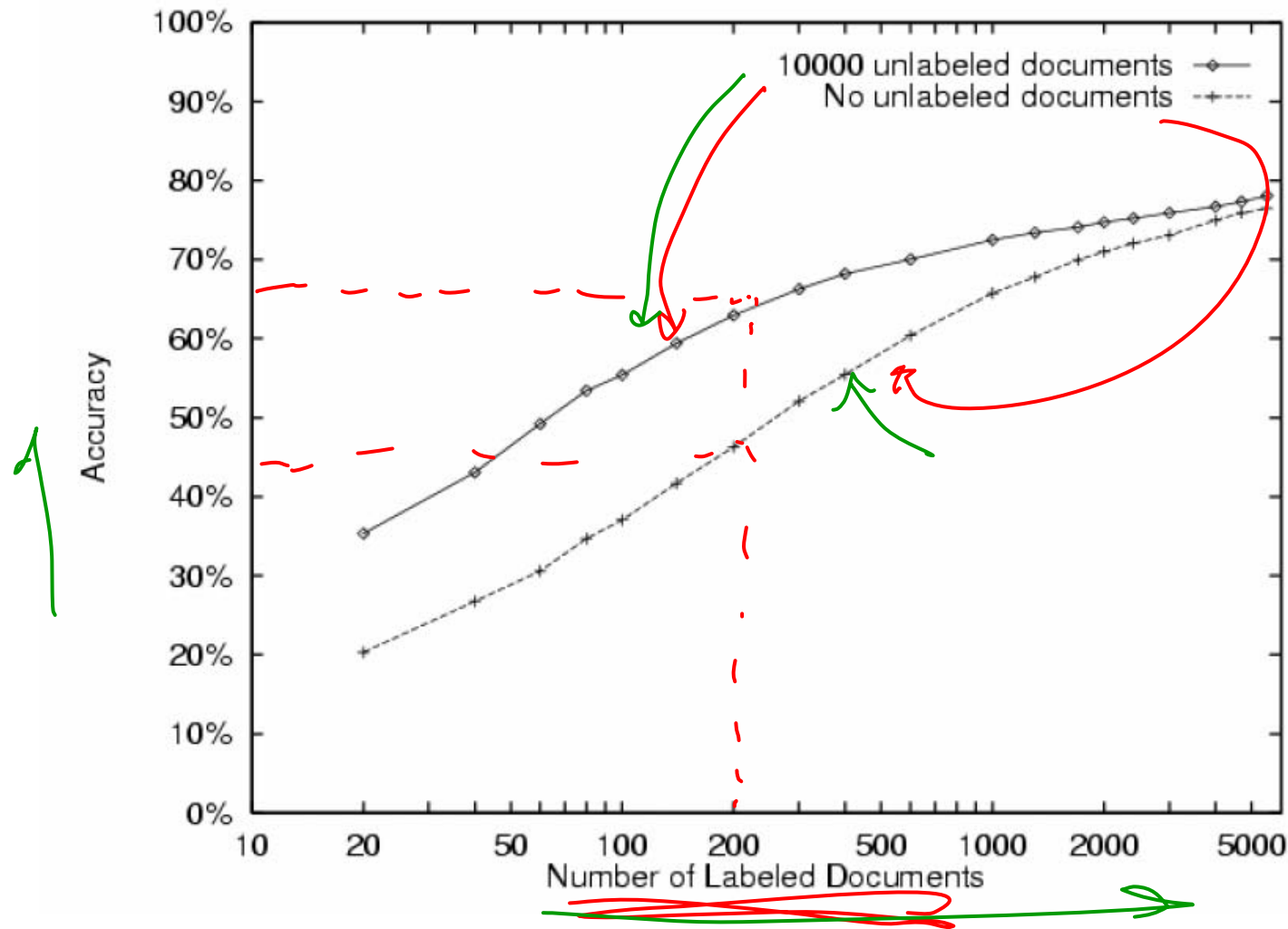  - If i'th point is labeled:
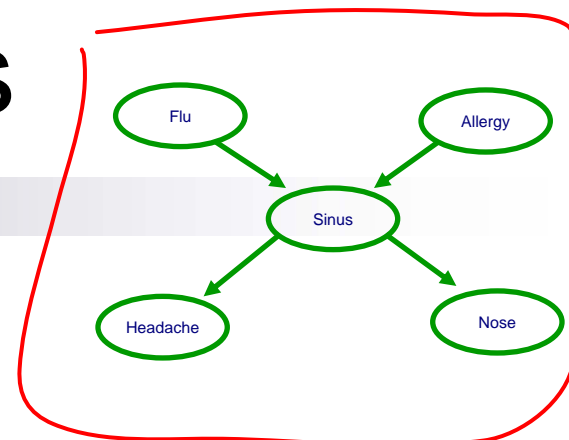    - set $Q(X=x|o_i)=1$ and $Q(X \neq x|o_i)=0$
      
      *correct label*          *not label*

- M-step as usual

# 20 Newsgroups data – advantage of adding unlabeled data

# Data likelihood for BNs

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

- Given structure, log likelihood of fully observed data:

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

$$\log \prod_{j=1}^{m} P(f^{(j)} \mid \theta_F) \cdot P(a^{(j)} \mid \theta_A) \cdot P(s^{(j)} \mid a^{(j)}, f^{(j)}, \theta_{S|FA})$$

$$\cdot P(h^{(j)} \mid s^{(j)}, \theta_{H|S}) \cdot P(n^{(j)} \mid s^{(j)}, \theta_{N|S})$$
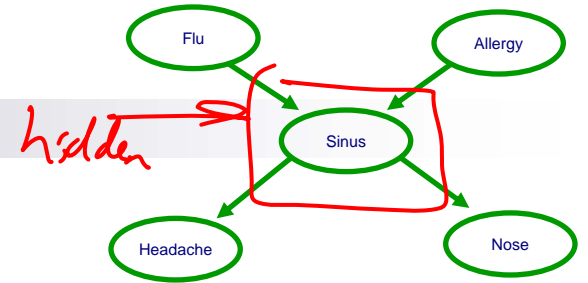
learning Flu CPT      A CPT      S|FA CPT

$$= \left[ \sum_{j=1}^{m} \log P(f^{(j)} \mid \theta_F) \right] + \left[ \sum_{j=1}^{m} \log P(a^{(j)} \mid \theta_A) \right] + \left[ \sum_{j=1}^{m} \log P(s^{(j)} \mid f^{(j)}, a^{(j)}, \theta_{S|FA}) \right]$$

$$+ \ldots$$

become independent learning problems

©2005-2007 Carlos Guestrin

# Marginal likelihood



- What if S is hidden?

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

$$= \sum_{j=1}^{m} \log \sum_{S} P(a^{(j)} \mid \theta_A) \, P(f^{(j)} \mid \theta_F) \cdot P(s \mid f^{(j)}, a^{(j)}, \theta_{S \mid FA})$$
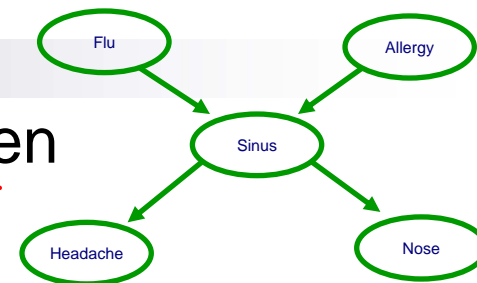
$$\cdot P(h^{(j)} \mid s, \theta_{H \mid S}) \cdot P(n^{(j)} \mid s, \theta_{N \mid S})$$

$$\log \sum$$

doesn't decompose
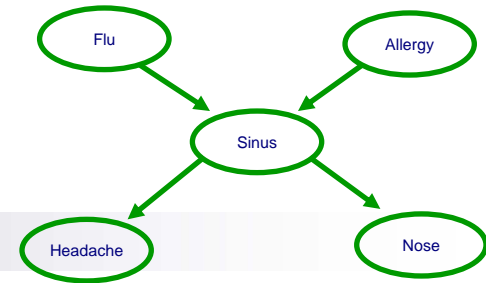
EM for BNS
same derivation (Jensen's, etc)

# Log likelihood for BNs with hidden data

■ Marginal likelihood – **O** is observed, **H** is hidden

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

Flu     Allergy

Sinus

Headache     Nose

# E-step for BNs

Flu → Sinus ← Allergy
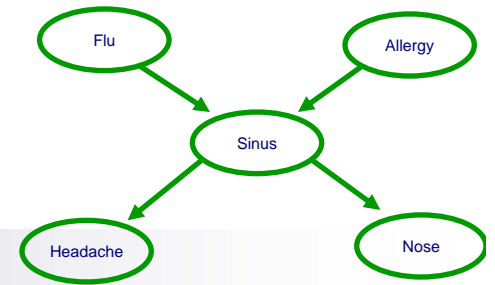Sinus → Headache
Sinus → Nose

- E-step computes probability of hidden vars **h** given **o**

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

*if* $|H| = 100$

$K^{100} - 1$ *params*
*(very large)*

- Corresponds to inference in BN

V.E.

# The M-step for BNs



- **Maximization step:**
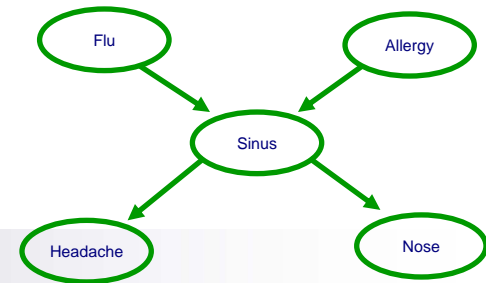
$$\theta^{(t+1)} \leftarrow \arg\max_\theta \sum_h Q^{(t+1)}(\mathbf{h} \mid \mathbf{o}^{(j)}) \log P(\mathbf{h}, \mathbf{o}^{(j)} \mid \theta)$$

- **Use expected counts instead of counts:**
  - ☐ If learning requires Count(**h**,**o**)
  - ☐ Use $E_{Q(t+1)}[\text{Count}(\mathbf{h},\mathbf{o})] = \sum_{j=1}^{m} \delta(O^{(j)} = o) \cdot Q^{(t+1)}(h \mid O^{(j)})$

# M-step for each CPT



- ## M-step decomposes per CPT

  - □ Standard MLE:

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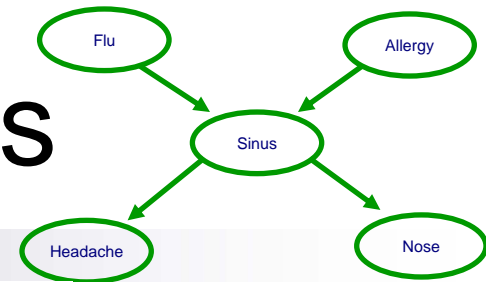
  *MLE CPT*

  $$\hat{P}(S=t \mid A=f, F=t)$$
  $$= \frac{\text{Count}(S=t, A=f, F=t)}{\text{Count}(A=f, F=t)}$$

  - □ M-step uses expected counts:

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

**14**

# Computing expected counts

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

- **M-step requires expected counts:**
  - For a set of vars **A**, must compute ExCount(**A**=**a**)
  - Some of **A** in example *j* will be observed
    - denote by **A**$_O$ = **a**$_O$
  - Some of **A** will be hidden
    - denote by **A**$_H$ = $a_W$

- **Use inference (E-step computes expected counts):**
  - ExCount$^{(t+1)}$(**A**$_O$ = **a**$_O$, **A**$_H$ = **a**$_H$) ← ~~P(A$_H$ = a$_H$, A$_O$ = a$_O$ |~~

$$= \sum_{j=1}^{m} \delta(A_O^{(j)} = a_O) \cdot \underbrace{P(A_H = a_H \mid O^{(j)}, \theta^{(t)})}_{\text{inference (VE)}}$$

15

# Data need not be hidden in the same way



- When data is fully observed
  - A data point is $\langle F=t,\ A=f,\ S=t,\ H=t,\ N=f \rangle$

- When data is partially observed
  - A data point is $\langle F=t, A=?, S=?, H=t, N=f \rangle$

- But unobserved variables can be different for different data points
  - e.g., $\langle F=t, A=t, S=t, H=?, N=? \rangle$,
    $\langle F=?, A=f, S=t, H=?, N=f \rangle$

- Same framework, just change definition of expected counts
  - ExCount$^{(t+1)}(\mathbf{A_O} = \mathbf{a_O}^{(i)}, \mathbf{A_H} = \mathbf{a_H}) \leftarrow$ ~~P(A_H = a_H, A_O = a_O^{(i)}|θ^{(t)})~~

    *set of hidden vars are a function of j..*

# What you need to know

- EM for Bayes Nets

- E-step: inference computes expected counts
  - Only need expected counts over $X_i$ and $\mathbf{Pa}_{xi}$

- M-step: expected counts used to estimate parameters

- Hidden variables can change per datapoint

- Use labeled and unlabeled data $\rightarrow$ some data points are complete, some include hidden variables

# Announcements

- No recitation this week

  *Spring Carnival*

- On Wednesday, Special lecture on learning with text data by Prof. Noah Smith (LTI)

# Co-Training for Semi-supervised learning

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 16th, 2007

**19**

# Redundant information



Professor Faloutsos    my advisor

U.S. mail address:
Department of Computer Science
University of Maryland
College Park, MD 20742
(97-99: on leave at CMU)
Office: 3227 A.V. Williams Bldg.
Phone: (301) 405-2695
Fax: (301) 405-6707
Email: christos@cs.umd.edu

**Christos Faloutsos**

Current Position: Assoc. Professor of Computer Science. (97-98: on leave at CMU)
Join Appointment: Institute for Systems Research (ISR).
Academic Degrees: Ph.D. and M.Sc. (University of Toronto.); B.Sc. (Nat. Tech. U. Ath...

**Research Interests:**

- Query by content in multimedia databases;
- Fractals for clustering and spatial access methods;
- Data mining;

class
$\in \{$ Faculty, Student, project, ... $\}$

# Redundant information – webpage text



Christos Faloutsos

**U.S. mail address:**
Department of Computer Science
University of Maryland
College Park, MD 20742
(97-99: on leave at CMU)
**Office:** 3227 A.V. Williams Bldg.
**Phone:** (301) 405-2695
**Fax:** (301) 405-6707
**Email:** christos@cs.umd.edu

**Current Position:** Assoc. Professor of Computer Science. (97-98: on leave at CMU)
**Join Appointment:** Institute for Systems Research (ISR).
**Academic Degrees:** Ph.D. and M.Sc. (University of Toronto.); B.Sc. (Nat. Tech. U. Athe

## Research Interests:

- Query by content in multimedia databases;
- Fractals for clustering and spatial access methods;
- Data mining;

21

# Redundant information – anchor text for hyperlinks

# Exploiting redundant information in semi-supervised learning

- **Want to predict Y from features X**
  - $f(\mathbf{X}) \mapsto Y$
  - have some labeled data **L**
  - lots of unlabeled data **U**

- **Co-training assumption: X is very expressive**
  - $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$
  - can learn
    - $g_1(\mathbf{X}_1) \mapsto Y$
    - $g_2(\mathbf{X}_2) \mapsto Y$

Can do alot with unlabeled data, especially if $X_1 \perp X_2 \mid Y$

Professor Faloutsos                                    my advisor

**U.S. mail address:**
Department of Computer Science
University of Maryland
College Park, MD 20742
(97-99: on leave at CMU)
Office: 3227 A.V. Williams Bldg.
Phone: (301) 405-2695
Fax: (301) 405-6707
Email: christos@cs.umd.edu

**Christos Faloutsos**

**Current Position:** Assoc. Professor of Computer Science. (97-98: on leave at CMU)
**Join Appointment:** Institute for Systems Research (ISR).
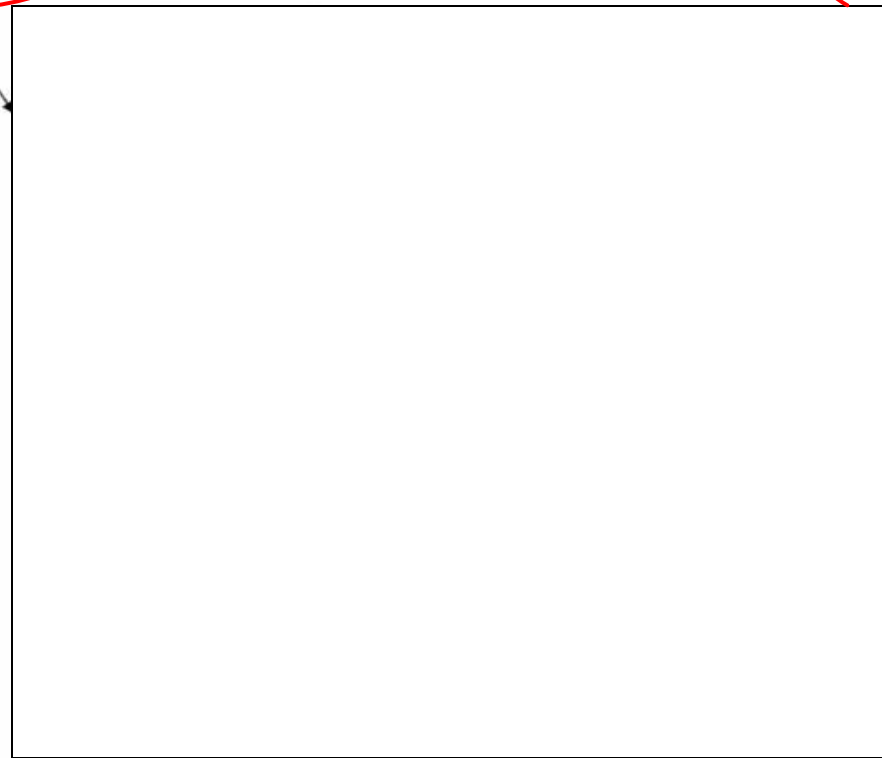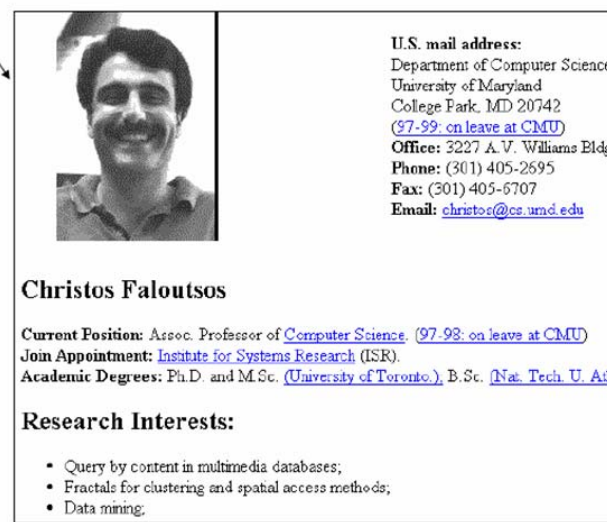**Academic Degrees:** Ph.D. and M.Sc. (University of Toronto), B.Sc. (Nat. Tech. U. Ath

**Research Interests:**

- Query by content in multimedia databases;
- Fractals for clustering and spatial access methods;
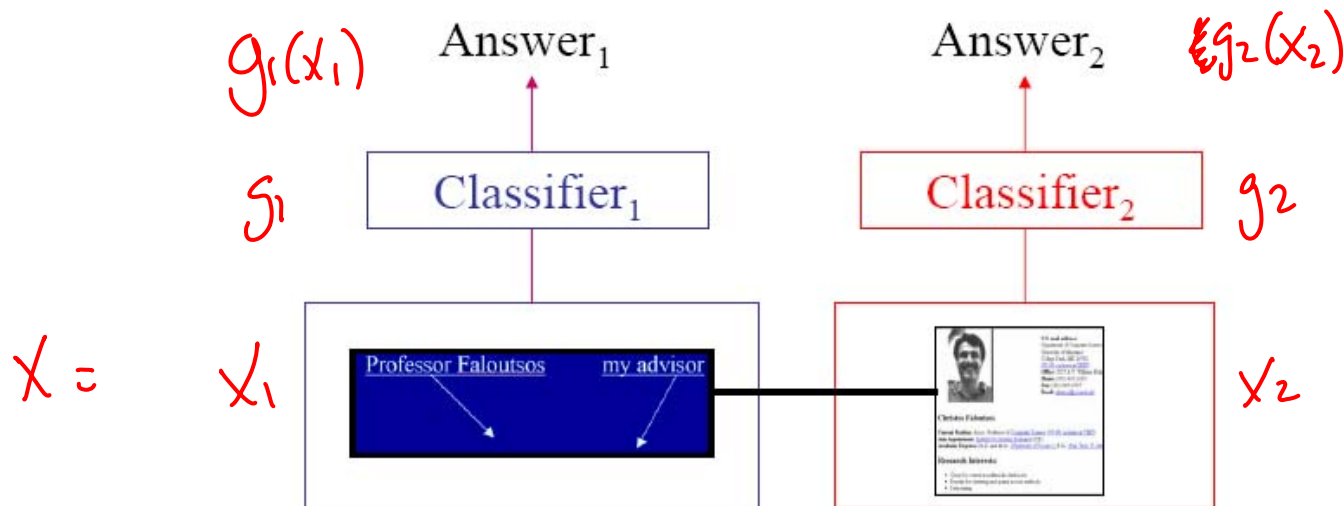- Data mining;

# Co-Training

- Key idea: Classifier$_1$ ($g_1$) and Classifier$_2$ ($g_2$) must:
    - ☐ Correctly classify labeled data
    - ☐ **Agree** on unlabeled data

*if x is labeled as Y, I want*
$$g_1(x_1) = Y$$
$$g_2(x_2) = Y$$

*if x is unlabeled*
*want* $g_1(x_1) = g_2(x_2)$

$g_1(x_1)$    Answer$_1$        Answer$_2$   $g_2(x_2)$

$g_1$    | Classifier$_1$ |        | Classifier$_2$ |   $g_2$

$x = $   $x_1$

Professor Faloutsos    my advisor              $x_2$

24

# Co-Training Algorithm
## [Blum & Mitchell '99]

*(example of the Co-training principle)*

Given: labeled data L,

      unlabeled data U

Loop:

    Train g1 (hyperlink classifier) using L    $X_1$

    Train g2 (page classifier) using L    $X_2$

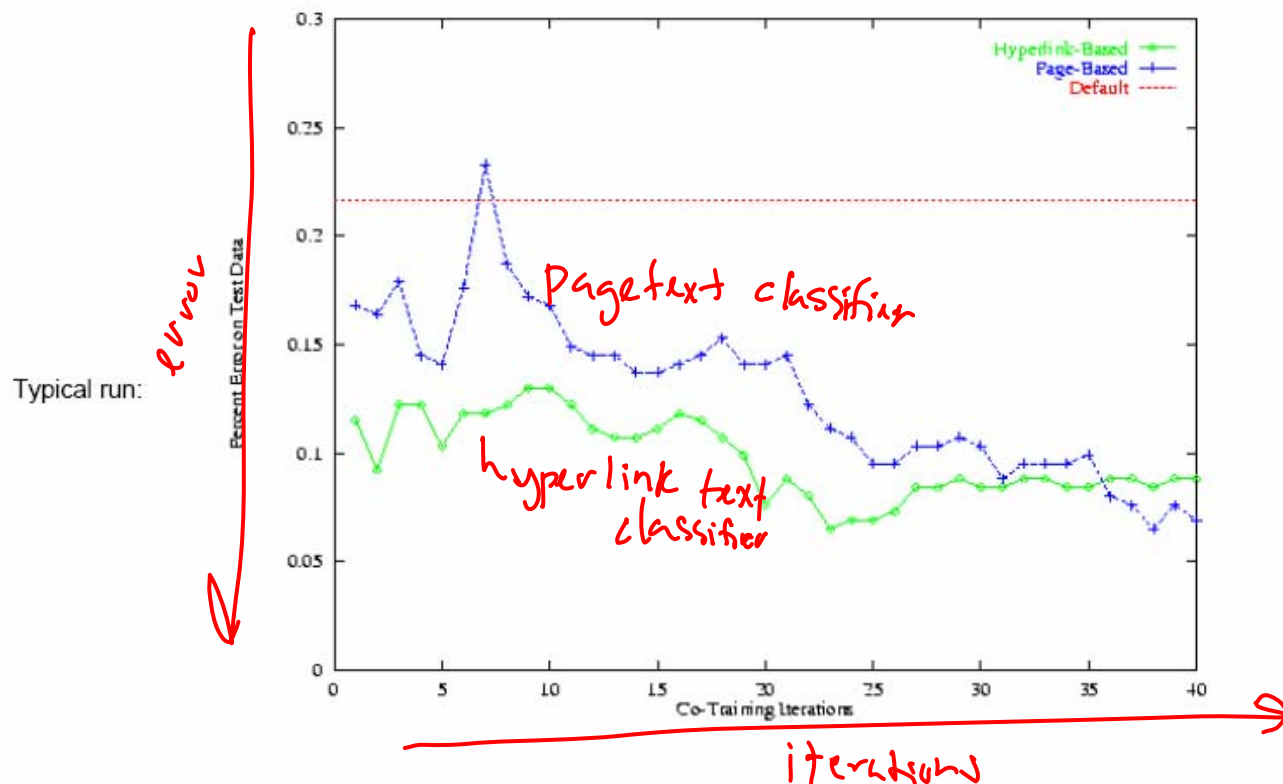    Allow g1 to label $p$ positive, $n$ negative examps from U

    Allow g2 to label $p$ positive, $n$ negative examps from U

    ~~Add~~ *Move* these self-labeled examples to L

**25**

# Co-Training experimental results

- begin with 12 labeled web pages (academic course)
- provide 1,000 additional unlabeled web pages
- average error: learning from labeled data 11.1%;
- average error: cotraining 5.0%

Typical run:

# Co-Training theory

- Want to predict Y from features **X**
  - $f(\mathbf{X}) \mapsto Y$
- Co-training assumption: **X** is very expressive
  - $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$
  - want to learn $g_1(\mathbf{X}_1) \mapsto Y$ and $g_2(\mathbf{X}_2) \mapsto Y$

- *Assumption*: $\exists\, g_1, g_2, \forall\, \mathbf{x}\; g_1(\mathbf{x}_1) = f(\mathbf{x}),\; g_2(\mathbf{x}_2) = f(\mathbf{x})$
- Questions:
  - Does unlabeled data always help?
  - How many labeled examples do I need?
  - How many unlabeled examples do I need?

27

# Understanding Co-Training: A simple setting

- Suppose $X_1$ and $X_2$ are discrete
  - $|X_1| = |X_2| = N$

  possible values

  _if $X_1$ is described by $n$ binary features, $N = 2^n$_

- No label noise

- Without unlabeled data, how hard is it to learn $g_1$ (or $g_2$)?

$|H| = 2^N$

$X_1$    ↶ hypothesis space

\# training examples is dependent on

1   $\{+, -\}$      $g_1 \in H$      $\ln|H| = N \cdot \ln 2$

2   $\{+, -\}$

$\vdots$    $\vdots$

$\vdots$    $\vdots$

$n$   $\{+, -\}$

# Co-Training in simple setting – Iteration 0



you get
if a web page
with $X_1 = 12$...
& $X_2 = 18$

$X_1$

text of hyperlinks

set of webpages form pages registexton

labeled data

$X_2$

edge $X_1 = x_1$
to $X_2 = x_2$
means
$X_1$ & $X_2$
co occurred
on a
webpage

My advisor

unlabeled webpage
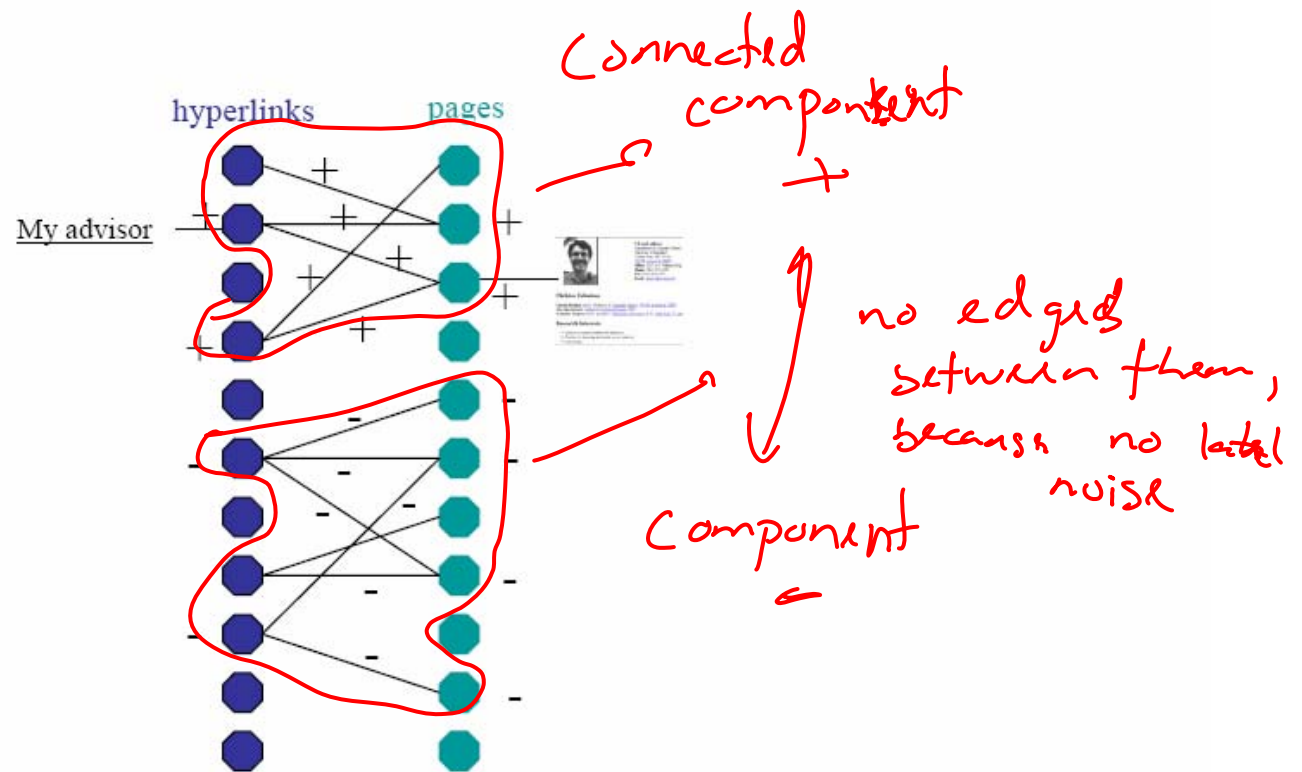
NO LABEL NOISE

1

12
16

N

N

one webpage
$X_1 = 16$ & $X_2 = 17$

-17

**29**

# Co-Training in simple setting – Iteration 1

# Co-Training in simple setting – after convergence

# Co-Training in simple setting – Connected components

**hyperlinks**   **pages**

- Suppose infinite **unlabeled** data
  - Co-training must have at least one labeled example in each connected component of L+U graph

  *component $g_j$*

- What's probability of making an error?

  *with m datapoints*

  *∃ connected component, where no data was labeled*

  *test point $x$*

  $$E[\text{error}] = \sum_{g_j \in \text{components}} P(x \in g_j) \left(1 - P(x \in g_j)\right)^m$$

  *no training data in $g_j$*

  $$E[error] = \sum_{j} P(x \in g_j)(1 - P(x \in g_j))^m$$

  Where $g_j$ is the $j$th connected component of graph of L+U, $m$ is number of labeled examples

- For k Connected components, how much labeled data?

  *about K data points instead of N*

# How much unlabeled data?

Want to assure that connected components in the underlying distribution, $G_D$, are connected components in the observed sample, $G_S$



$O(\log(N)/\alpha)$ examples assure that with high probability, $G_S$ has same connected components as $G_D$ [Karger, 94]

N is size of $G_D$, $\alpha$ is min cut over all connected components of $G_D$

# Co-Training theory

- Want to predict Y from features **X**
  - □ f(**X**)   Y
- Co-training assumption: **X** is very expressive
  - □ **X** = (**X**$_1$,**X**$_2$)
  - □ want to learn g$_1$(**X**$_1$)   Y and g$_2$(**X**$_2$)   Y

- *Assumption*: $\exists$ g$_1$, g$_2$, $\forall$ **x** g$_1$(**x**$_1$) = f(**x**), g$_2$(**x**$_2$) = f(**x**)

- One co-training result [Blum & Mitchell '99]
  - □ If
    - (**X**$_1$ $\perp$ **X**$_2$ | Y)
    - g$_1$ & g$_2$ are PAC learnable from noisy data (and thus f)
  - □ Then
    - f is PAC learnable from weak initial classifier plus unlabeled data

# What you need to know about co-training

- Unlabeled data can help supervised learning (a lot) when there are (mostly) independent redundant features

- One theoretical result:
    - If ($X_1 \perp X_2 \mid Y$) and $g_1$ & $g_2$ are PAC learnable from noisy data (and thus f)
    - Then f is PAC learnable from weak initial classifier plus unlabeled data
    - Disagreement between $g_1$ and $g_2$ provides bound on error of final classifier

- Applied in many real-world settings:
    - Semantic lexicon generation [Riloff, Jones 99] [Collins, Singer 99], [Jones 05]
    - Web page classification [Blum, Mitchell 99]
    - Word sense disambiguation [Yarowsky 95]
    - Speech recognition [de Sa, Ballard 98]
    - Visual classification of cars [Levin, Viola, Freund 03]

# Acknowledgement

- I would like to thank Tom Mitchell for some of the material used in this presentation of co-training