



EM for Bayes Nets

Machine Learning – 10701/15781

Carlos Guestrin

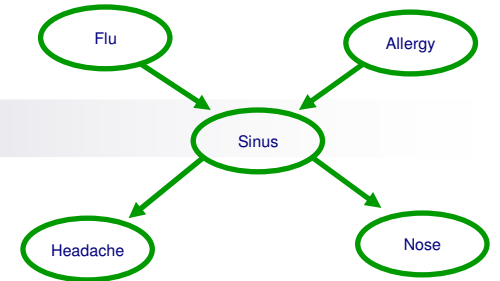
Carnegie Mellon University

April 17th, 2006

Data likelihood for BNs

- Given structure, log likelihood of fully observed data:

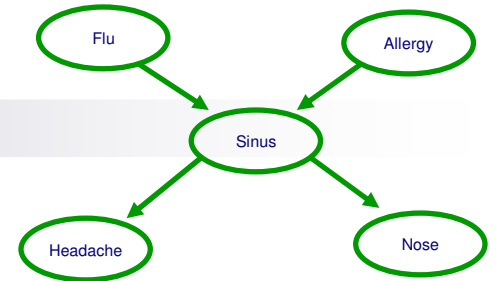
$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$



Marginal likelihood

- What if S is hidden?

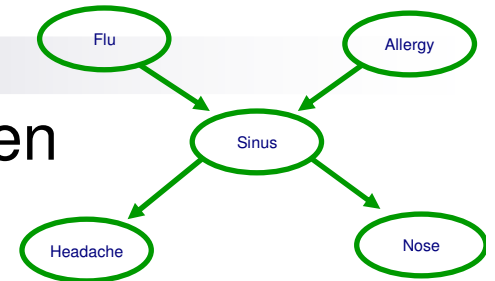
$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$



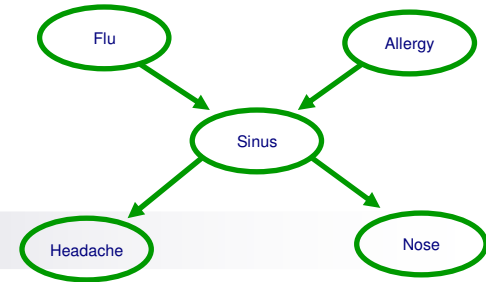
Log likelihood for BNs with hidden data

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

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



E-step for BNs

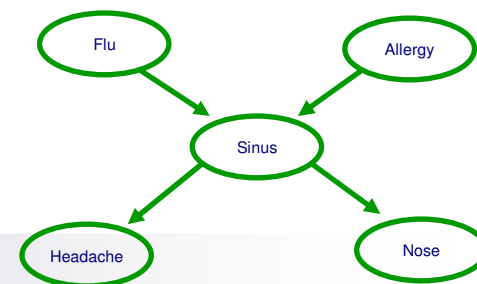


- E-step computes probability of hidden vars **h** given **o**

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

- Corresponds to inference in BN

The M-step for BNs



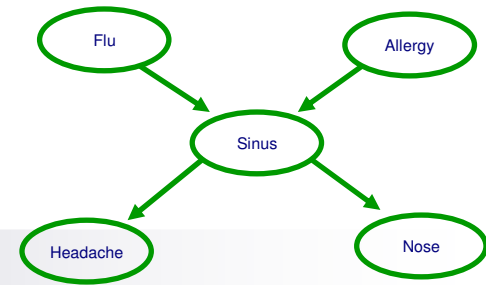
- Maximization step:

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

- 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

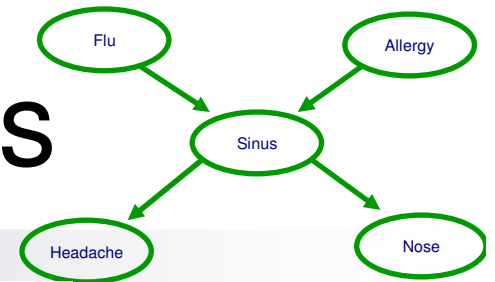
- Standard MLE:

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

- M-step uses expected counts:

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

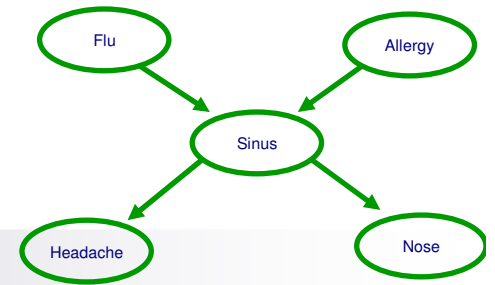
Computing expected counts



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

- 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):
 - $\text{ExCount}^{(t+1)}(\mathbf{A}_O = \mathbf{a}_O^{(j)}, \mathbf{A}_H = \mathbf{a}_H) \leftarrow P(\mathbf{A}_H = \mathbf{a}_H \mid \mathbf{A}_O = \mathbf{a}_O^{(j)}, \theta^{(t)})$

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
 - $\text{ExCount}^{(t+1)}(\mathbf{A}_O = \mathbf{a}_O^{(j)}, \mathbf{A}_H = \mathbf{a}_H) \leftarrow P(\mathbf{A}_H = \mathbf{a}_H \mid \mathbf{A}_O = \mathbf{a}_O^{(j)}, \theta^{(t)})$

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



Reading:
Blum & Mitchell 1999
(see class website)

Co-Training for Semi-supervised learning

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 17th, 2006

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;

Redundant information – webpage text

**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. Athens](#))

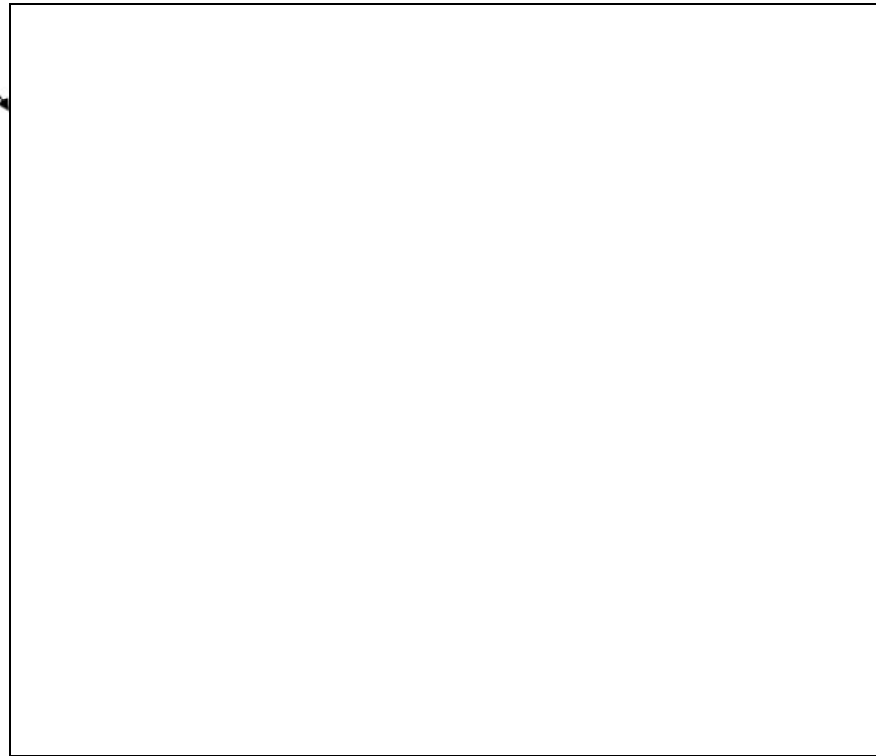
Research Interests:

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

Redundant information – anchor text for hyperlinks

Professor Faloutsos

my advisor




Exploiting redundant information in semi-supervised learning

- Want to predict Y from features \mathbf{X}
 - $f(\mathbf{X}) \mapsto Y$
 - have some labeled data \mathbf{L}
 - lots of unlabeled data \mathbf{U}
- Co-training assumption: \mathbf{X} is very expressive
 - $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$
 - can learn
 - $g_1(\mathbf{X}_1) \mapsto Y$
 - $g_2(\mathbf{X}_2) \mapsto Y$

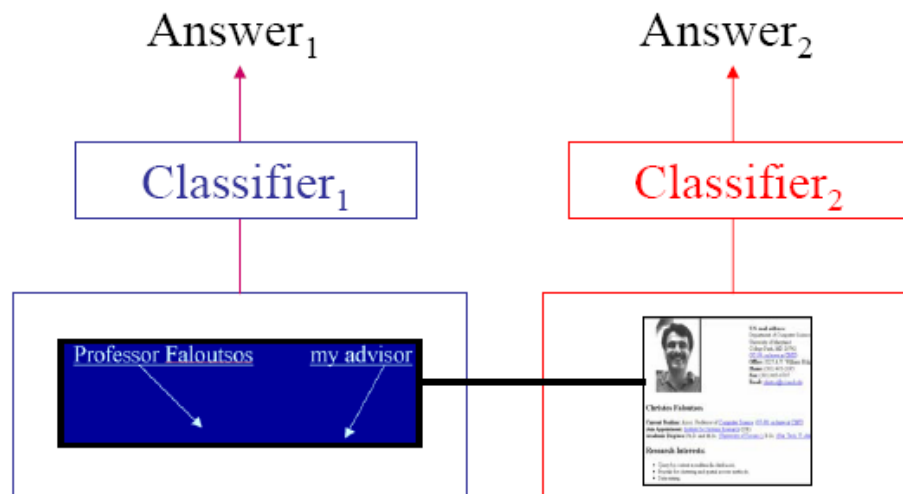
Professor Faloutsos

my advisor

	<p>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</p>
<p>Christos Faloutsos</p> <p>Current Position: Assoc. Professor of Computer Science, (97-99, 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. Athens)</p> <p>Research Interests:</p> <ul style="list-style-type: none">• Query by content in multimedia databases;• Fractals for clustering and spatial access methods;• Data mining.	

Co-Training

- Key idea: Classifier₁ and Classifier₂ must:
 - Correctly classify labeled data
 - **Agree** on unlabeled data



Co-Training Algorithm

[Blum & Mitchell '99]

Given: labeled data L ,

unlabeled data U

Loop:

Train g_1 (hyperlink classifier) using L

Train g_2 (page classifier) using L

Allow g_1 to label p positive, n negative examps from U

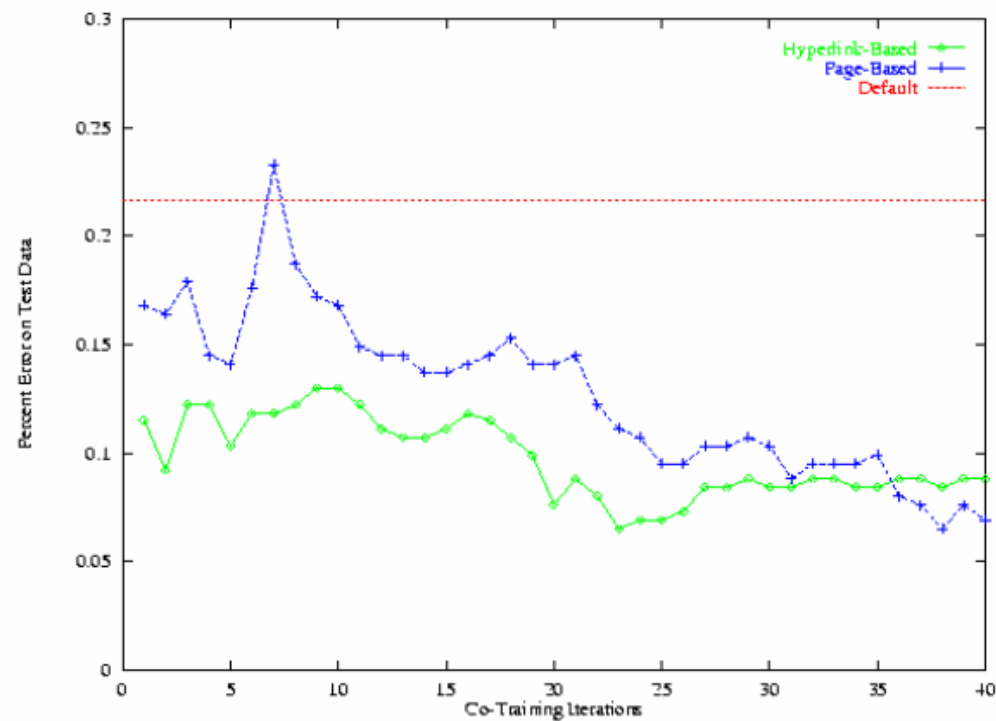
Allow g_2 to label p positive, n negative examps from U

Add these self-labeled examples to L

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:



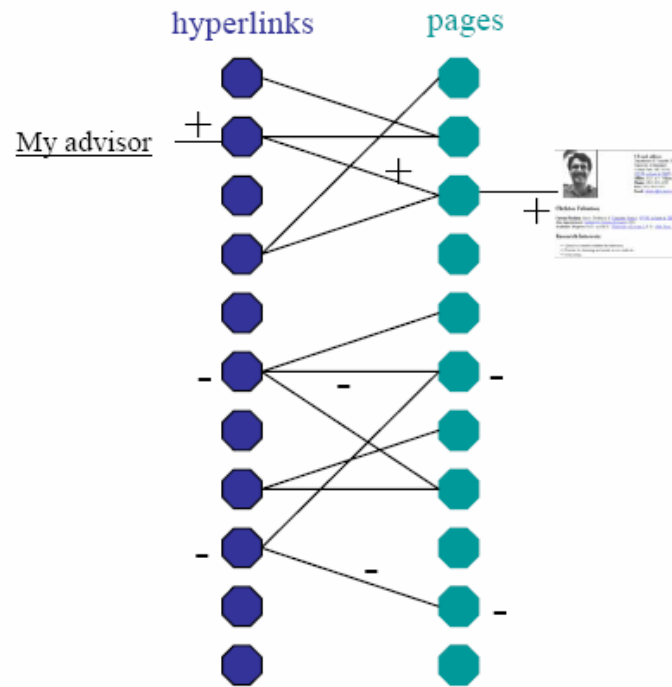
Co-Training theory

- Want to predict Y from features \mathbf{X}
 - $f(\mathbf{X}) \mapsto Y$
- Co-training assumption: \mathbf{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?

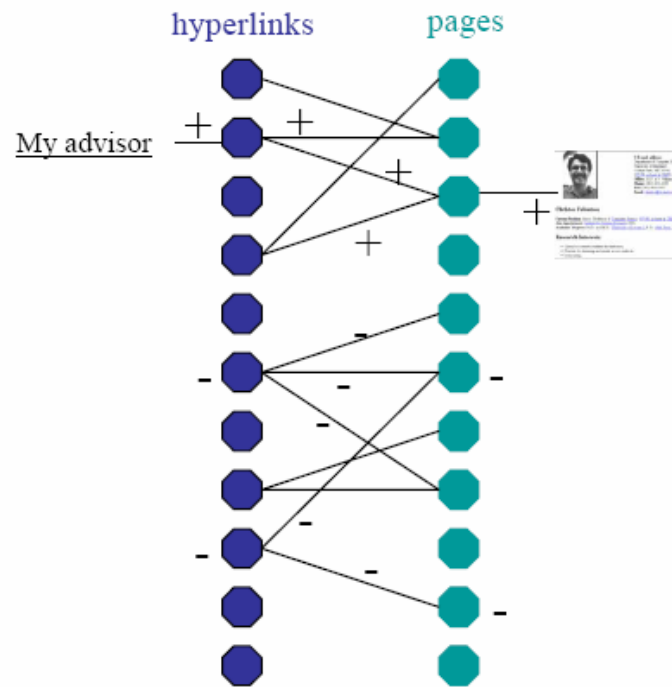
Understanding Co-Training: A simple setting

- Suppose \mathbf{X}_1 and \mathbf{X}_2 are discrete
 - $|\mathbf{X}_1| = |\mathbf{X}_2| = N$
- No label noise
- Without unlabeled data, how hard is it to learn g_1 (or g_2)?

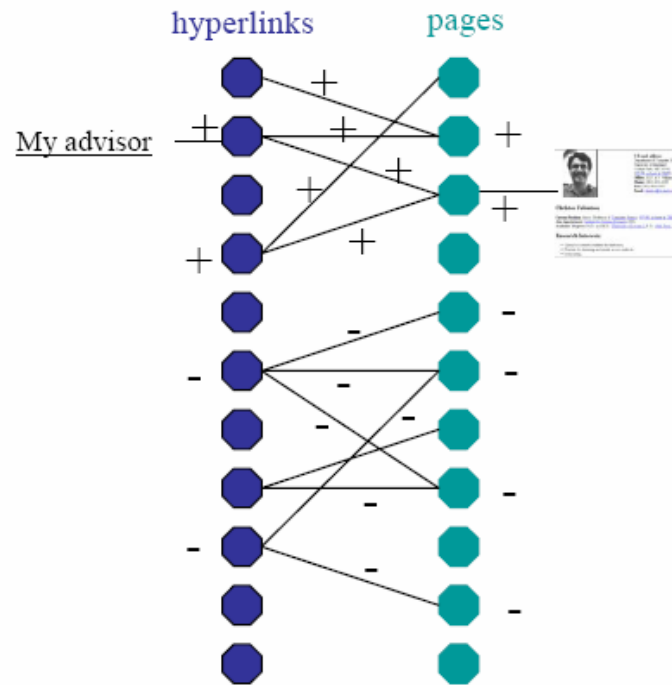
Co-Training in simple setting – Iteration 0



Co-Training in simple setting – Iteration 1

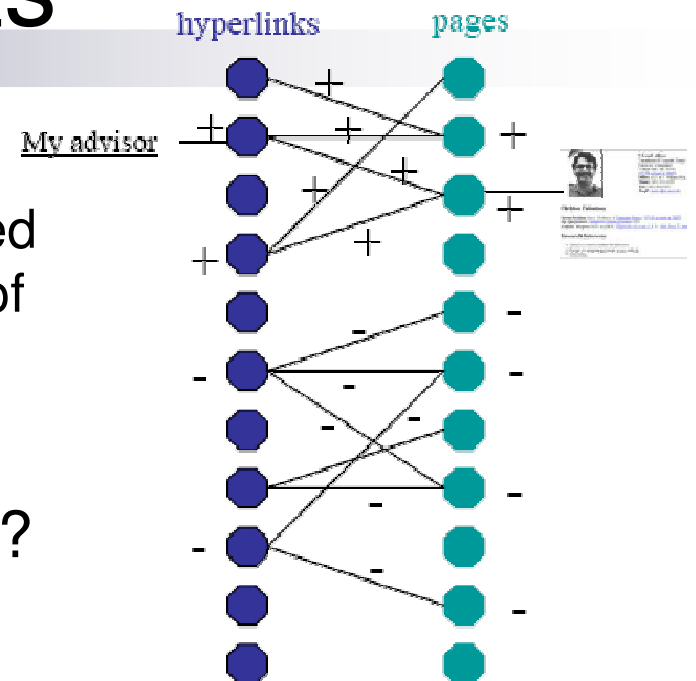


Co-Training in simple setting – after convergence



Co-Training in simple setting – Connected components

- Suppose infinite **unlabeled** data
 - Co-training must have at least one labeled example in each connected component of L+U graph
- What's probability of making an error?



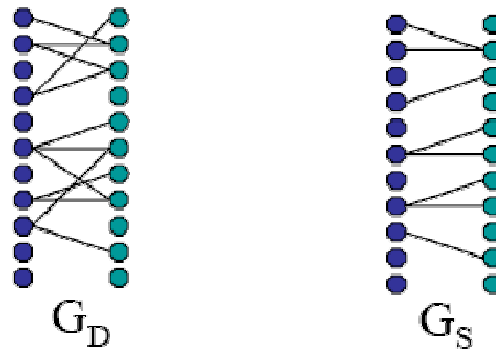
$$E[error] = \sum_j P(x \in g_j)(1 - P(x \in g_j))^m$$

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

- For k Connected components, how much labeled data?

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 , α is min cut over all connected components of G_D

Co-Training theory

- Want to predict Y from features \mathbf{X}
 - $f(\mathbf{X}) \mapsto Y$
- Co-training assumption: \mathbf{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*: $\exists g_1, g_2, \forall \mathbf{x} \ g_1(\mathbf{x}_1) = f(\mathbf{x}), g_2(\mathbf{x}_2) = f(\mathbf{x})$
- One co-training result [Blum & Mitchell '99]
 - If
 - $(\mathbf{X}_1 \perp \mathbf{X}_2 \mid Y)$
 - 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

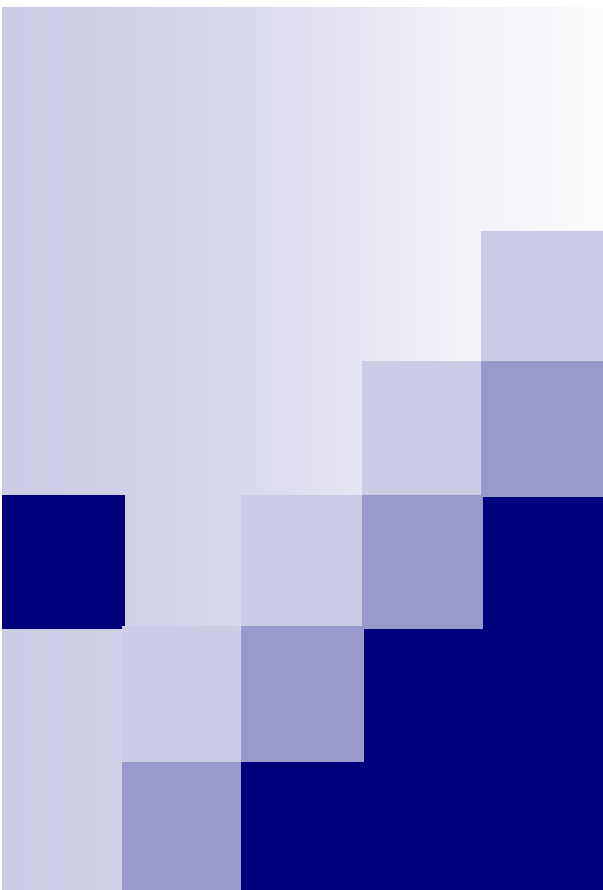
What you need to know about co-training

- Unlabeled data can help supervised learning (a lot) when there are (mostly) independent redundant features
- One theoretical result:
 - If $(\mathbf{X}_1 \perp \mathbf{X}_2 \mid Y)$ 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



Reading:
Vapnik 1998
Joachims 1999 (see class website)

Transductive SVMs

Machine Learning – 10701/15781

Carlos Guestrin

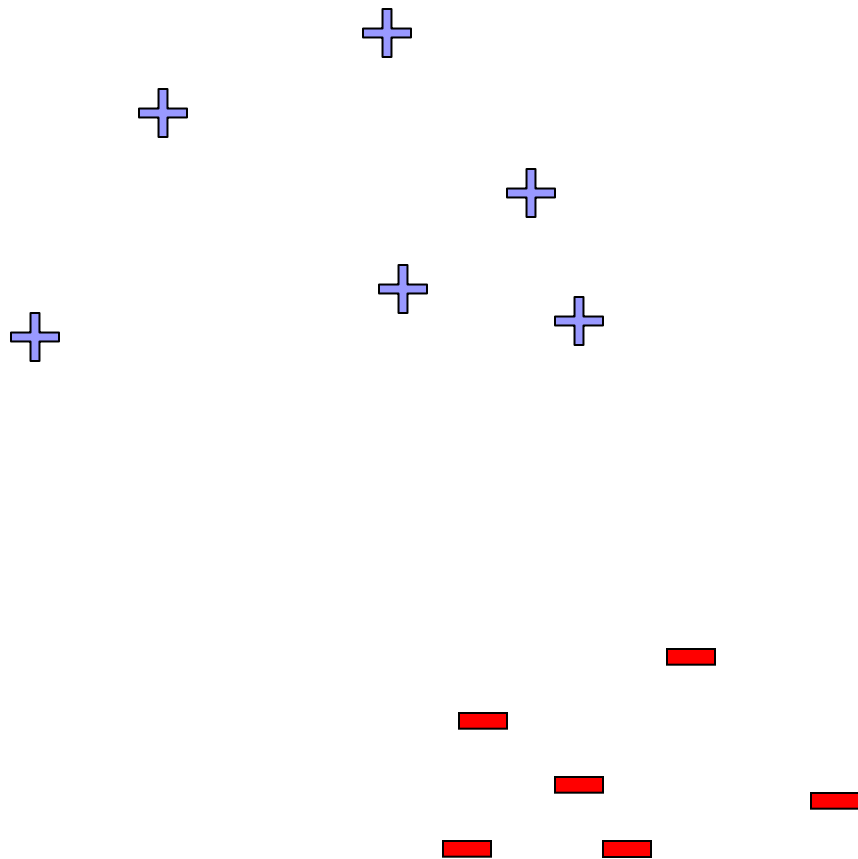
Carnegie Mellon University

April 17th, 2006

Semi-supervised learning and discriminative models

- We have seen semi-supervised learning for generative models
 - EM
- What can we do for discriminative models
 - Not regular EM
 - we can't compute $P(x)$
 - But there are discriminative versions of EM
 - Co-Training!
 - Many other tricks... let's see an example

Linear classifiers – Which line is better?



Data:

$$\begin{aligned} &\langle x_1^{(1)}, \dots, x_1^{(m)}, y_1 \rangle \\ &\quad \vdots \\ &\langle x_n^{(1)}, \dots, x_n^{(m)}, y_n \rangle \end{aligned}$$

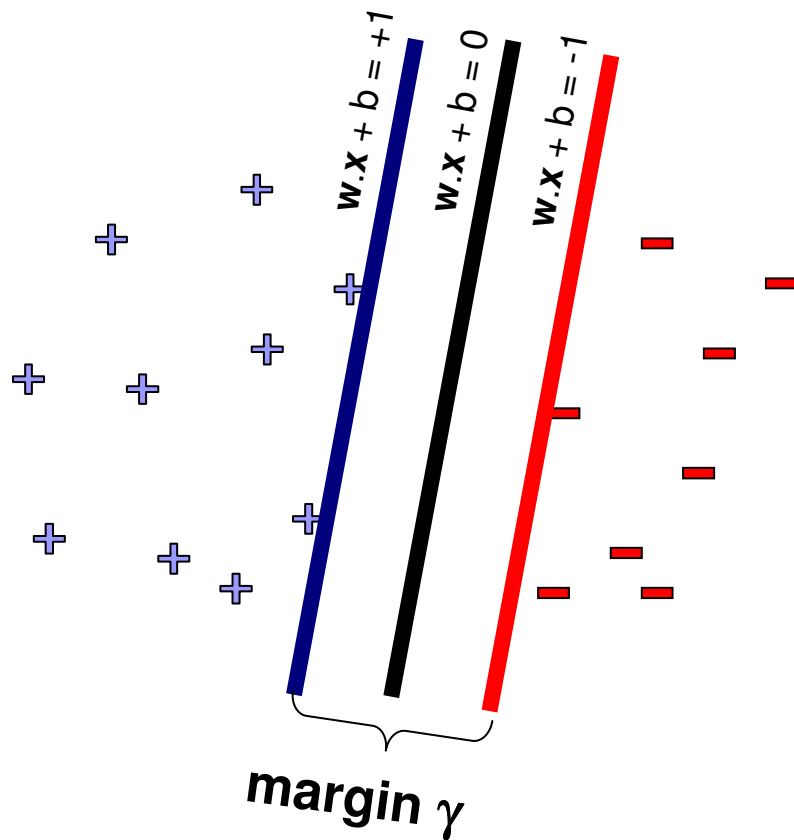
Example i:

$$\langle x_i^{(1)}, \dots, x_i^{(m)} \rangle \text{ — } m \text{ features}$$

$$y_i \in \{-1, +1\} \text{ — class}$$

$$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$$

Support vector machines (SVMs)



$$\text{minimize}_w \quad w \cdot w$$
$$(w \cdot x_j + b) y_j \geq 1, \quad \forall j$$

- Solve efficiently by quadratic programming (QP)
 - Well-studied solution algorithms
- Hyperplane defined by support vectors

What if we have unlabeled data?



n_L Labeled Data:

$$\langle x_1^{(1)}, \dots, x_1^{(m)}, y_1 \rangle$$

\vdots

$$\langle x_n^{(1)}, \dots, x_n^{(m)}, y_{n_L} \rangle$$

Example i:

$$\langle x_i^{(1)}, \dots, x_i^{(m)} \rangle \text{ — } m \text{ features}$$

$$y_i \in \{-1, +1\} \text{ — class}$$

n_U Unlabeled Data:

$$\langle x_1^{(1)}, \dots, x_1^{(m)}, ? \rangle$$

\vdots

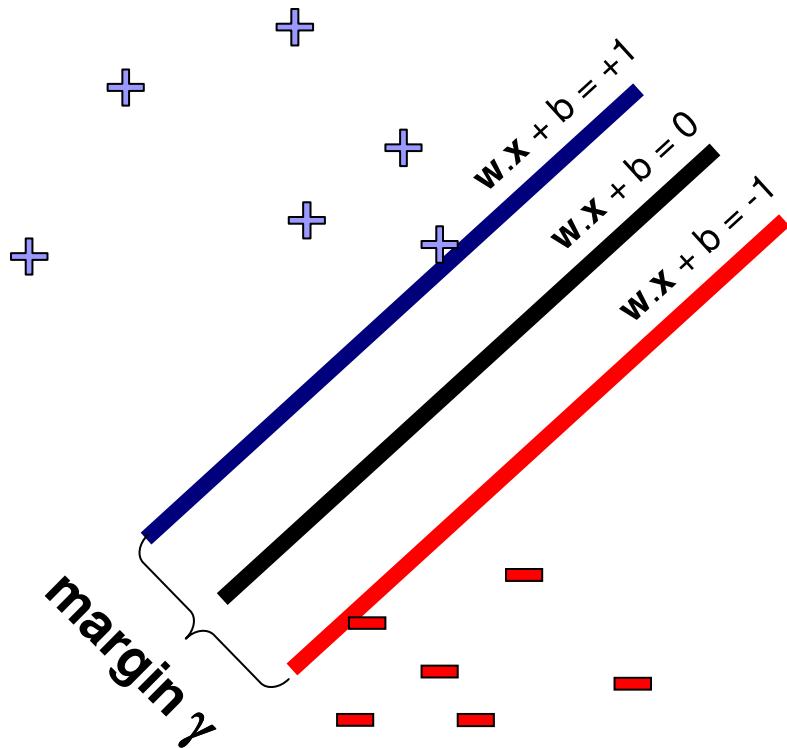
$$\langle x_n^{(1)}, \dots, x_n^{(m)}, ? \rangle$$

33

$$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$$

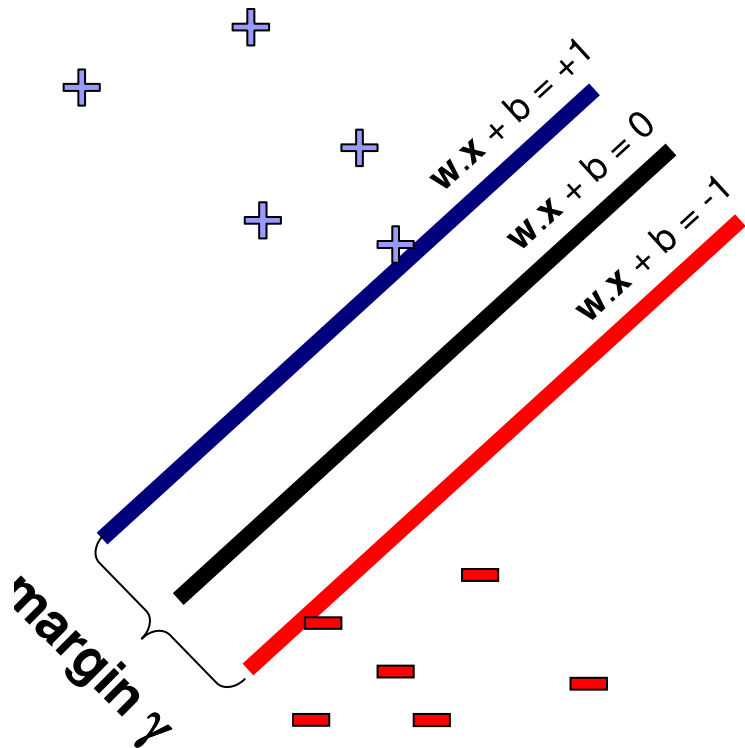
Transductive support vector machines (TSVMs)

minimize_w $w \cdot w$



$$(w \cdot x_j + b) y_j \geq 1, \quad \forall j$$

Transductive support vector machines (TSVMs)



$$\text{minimize}_{\mathbf{w}, \{\hat{y}_1, \dots, \hat{y}_{n_U}\}} \quad \mathbf{w} \cdot \mathbf{w}$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \quad \forall j = 1, \dots, n_L$$

$$(\mathbf{w} \cdot \mathbf{x}_u + b) \hat{y}_u \geq 1, \quad \forall u = 1, \dots, n_U$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_U$$

What's the difference between transductive learning and semi-supervised learning?

- Not much, and
- A lot!!!
- Semi-supervised learning:
 - labeled and unlabeled data → learn \mathbf{w}
 - use \mathbf{w} on test data
- Transductive learning
 - same algorithms for labeled and unlabeled data, but...
 - unlabeled data is test data!!!
- You are learning on the test data!!!
 - OK, because you never look at the labels of the test data
 - can get better classification
 - but be very very very very very very very very careful!!!
 - never use test data prediction accuracy to tune parameters, select kernels, etc.

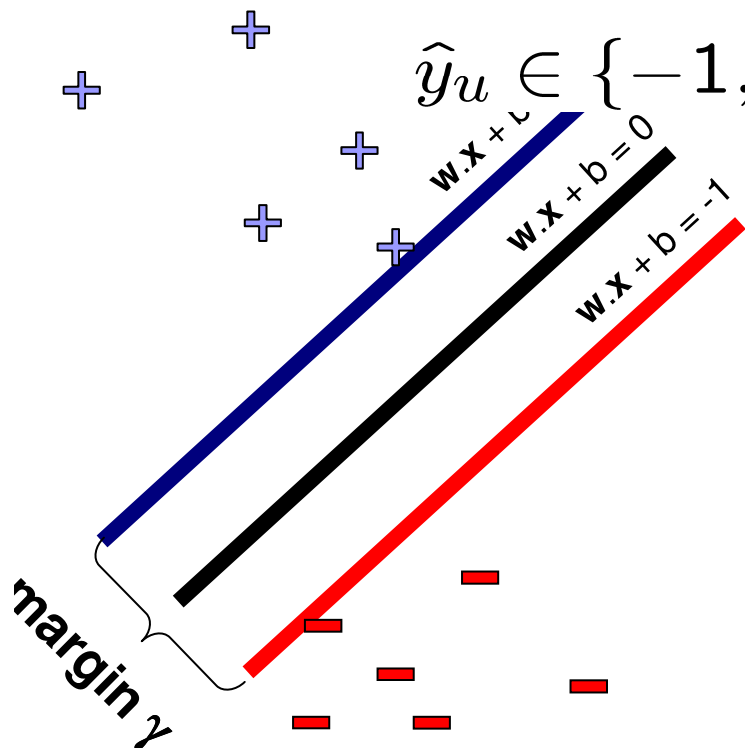
Adding slack variables

$$\text{minimize}_{\mathbf{w}, \{\hat{y}_1, \dots, \hat{y}_{n_U}\}} \quad \mathbf{w} \cdot \mathbf{w}$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 \quad \forall j = 1, \dots, n_L$$

$$(\mathbf{w} \cdot \mathbf{x}_u + b) \hat{y}_u \geq 1 \quad \forall u = 1, \dots, n_U$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_U$$



Transductive SVMs – now with slack variables! [Vapnik 98]

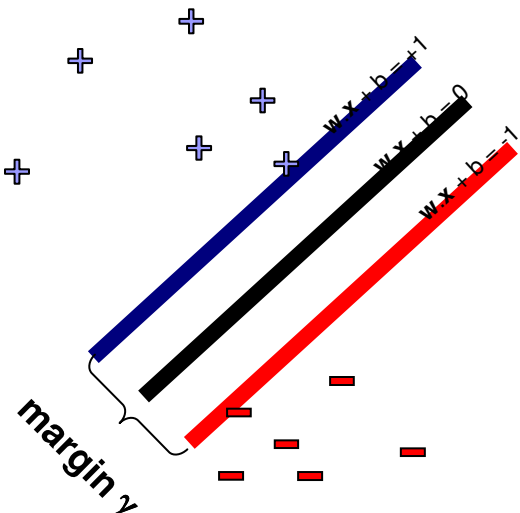
■ Optimize $w, \{\xi_1, \dots, \xi_{n_L}\}, \{\hat{y}_1, \dots, \hat{y}_{n_U}\}, \{\hat{\xi}_1, \dots, \hat{\xi}_{n_U}\}$

$$\text{minimize } w \cdot w + C \sum_j \xi_j + \hat{C} \sum_u \hat{\xi}_u$$

$$(w \cdot x_j + b) y_j \geq 1 - \xi_j, \quad \forall j = 1, \dots, n_L$$

$$(w \cdot x_u + b) \hat{y}_u \geq 1 - \hat{\xi}_u, \quad \forall u = 1, \dots, n_u$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_u$$



Learning Transductive SVMs is hard!

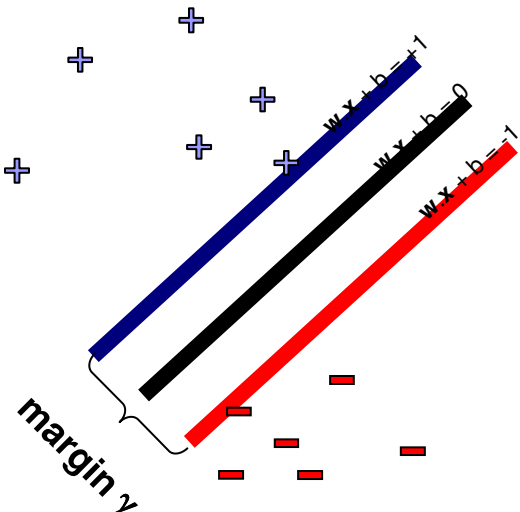
■ Optimize $w, \{\xi_1, \dots, \xi_{n_L}\}, \{\hat{y}_1, \dots, \hat{y}_{n_U}\}, \{\hat{\xi}_1, \dots, \hat{\xi}_{n_U}\}$

$$\text{minimize} \quad w \cdot w + C \sum_j \xi_j + \hat{C} \sum_u \hat{\xi}_u$$

$$(w \cdot x_j + b) y_j \geq 1 - \xi_j, \quad \forall j = 1, \dots, n_L$$

$$(w \cdot x_u + b) \hat{y}_u \geq 1 - \hat{\xi}_u, \quad \forall u = 1, \dots, n_u$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_u$$



■ Integer Program

□ NP-hard!!!

□ Well-studied solution algorithms, but will not scale up to very large problems

A (heuristic) learning algorithm for Transductive SVMs [Joachims 99]

$$\text{minimize} \quad \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j + \hat{C} \sum_u \hat{\xi}_u$$

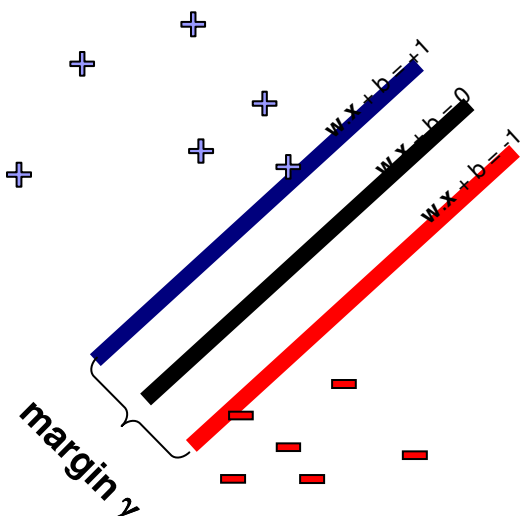
$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j = 1, \dots, n_L$$

$$(\mathbf{w} \cdot \mathbf{x}_u + b) \hat{y}_u \geq 1 - \hat{\xi}_u, \quad \forall u = 1, \dots, n_u$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_u$$

- If you set \hat{C} to zero \rightarrow ignore unlabeled data
- Intuition of algorithm:

- start with small \hat{C}
- add labels to some unlabeled data based on classifier prediction
- slowly increase \hat{C}
- keep on labeling unlabeled data and re-running classifier



Some results classifying news articles – from [Joachims 99]

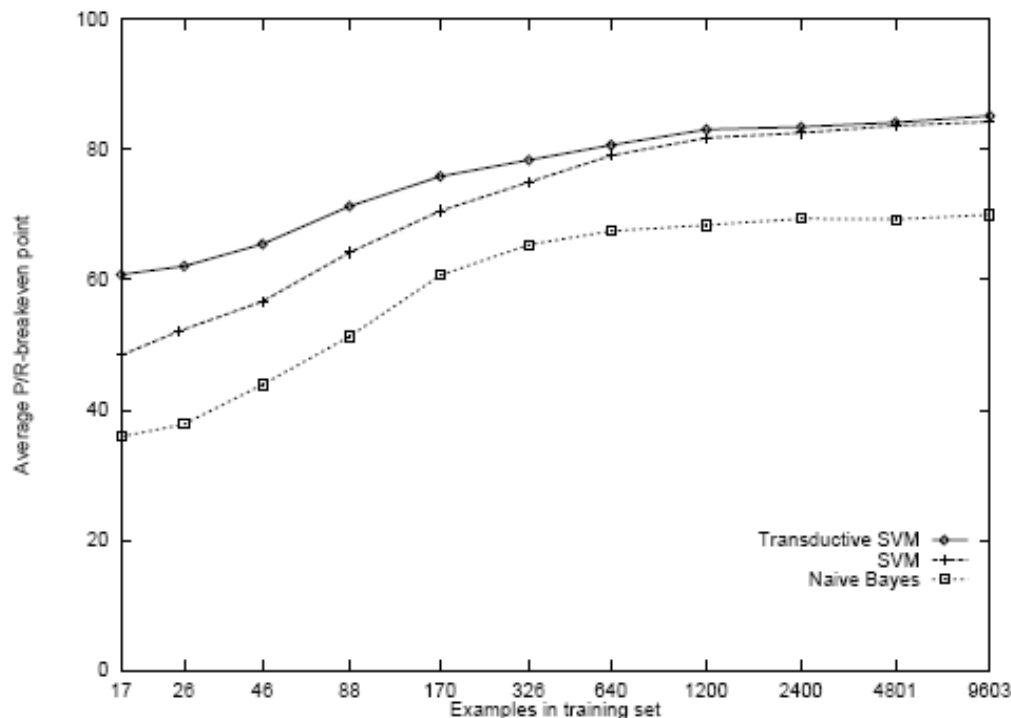


Figure 6: Average P/R-breakeven point on the Reuters dataset for different training set sizes and a test set size of 3,299.

What you need to know about transductive SVMs

- What is transductive v. semi-supervised learning
- Formulation for transductive SVM
 - can also be used for semi-supervised learning
- Optimization is hard!
 - Integer program
- There are simple heuristic solution methods that work well here