



Learning from Labeled and Unlabeled Data

Machine Learning 10-701

November 7, 2006

Tom M. Mitchell

Machine Learning Department

Carnegie Mellon University



When can Unlabeled Data improve supervised learning?

Important question! In many cases, unlabeled data is plentiful, labeled data expensive

- Medical outcomes ($x = \langle \text{symptoms}, \text{treatment} \rangle$, $y = \text{outcome}$)
- Text classification ($x = \text{document}$, $y = \text{relevance}$)
- Customer modeling ($x = \text{user actions}$, $y = \text{user intent}$)
- Sensor interpretation ($x = \langle \text{video}, \text{audio} \rangle$, $y = \text{who's there}$)

When can Unlabeled Data help supervised learning?

Problem setting:

- Set X of instances drawn from unknown distribution $P(X)$
- Wish to learn target function $f: X \rightarrow Y$ (or, $P(Y|X)$)
- Given a set H of possible hypotheses for f

Given:

- iid labeled examples $L = \{\langle x_1, y_1 \rangle \dots \langle x_m, y_m \rangle\}$
- iid unlabeled examples $U = \{x_{m+1}, \dots, x_{m+n}\}$

Wish to determine:

$$\hat{f} \leftarrow \arg \min_{h \in H} \Pr_{x \in P(X)} [h(x) \neq f(x)]$$



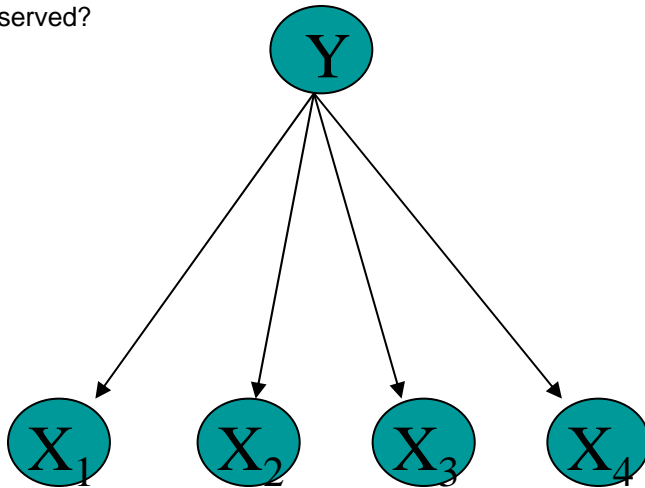
Idea 1: Use Labeled and Unlabeled Data to
Train Bayes Net for $P(X,Y)$

Idea 1: Use Labeled and Unlabeled Data to Train Bayes Net for $P(X,Y)$, then infer $P(Y|X)$

What CPDs are needed?

How do we estimate them from fully observed data?

How do we estimate them from partly observed?



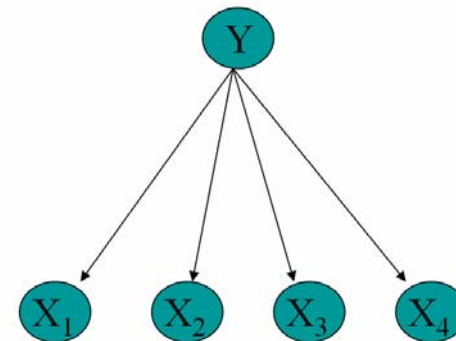
Y	X1	X2	X3	X4
1	0	0	1	1
0	0	1	0	0
0	0	0	1	0
?	0	1	1	0
?	0	1	0	1

Supervised: Naïve Bayes Learner

Train:

For each class y_j of documents

1. Estimate $P(Y=y_j)$
2. For each word w_i estimate $P(X=w_i / Y=y_j)$



Classify (doc):

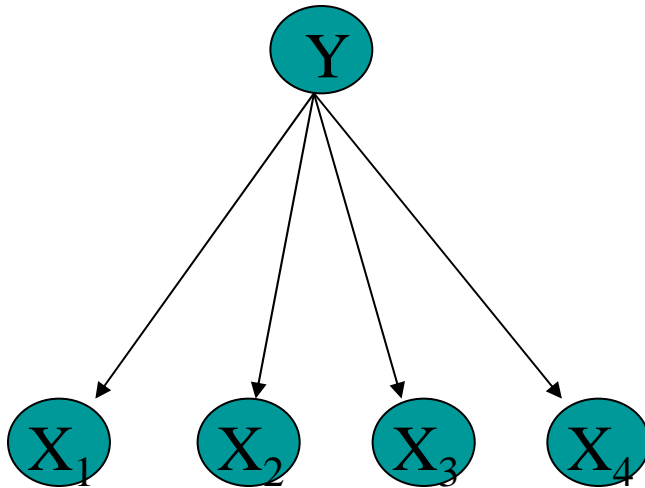
Assign *doc* to most probable* class

$$\hat{P}(y_j | doc) \leftarrow \frac{\hat{P}(y_j) \prod_i \hat{P}(w_i | y_j)}{\sum_k \hat{P}(y_k) \prod_i \hat{P}(w_i | y_k)}$$

* assuming words w_i are conditionally independent, given class

What if we have labels for only *some* documents?

Learn $P(Y|X)$

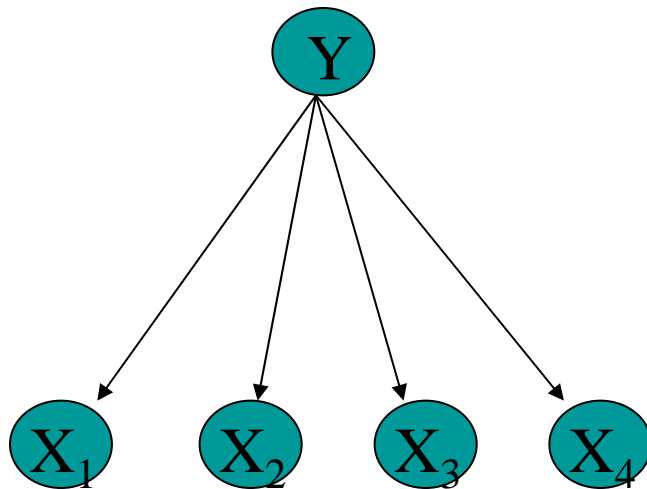


Y	X1	X2	X3	X4
1	0	0	1	1
0	0	1	0	0
0	0	0	1	0
?	0	1	1	0
?	0	1	0	1

What if we have labels for only *some* documents?

[Nigam et al., 2000]

Learn $P(Y|X)$



Y	X1	X2	X3	X4
1	0	0	1	1
0	0	1	0	0
0	0	0	1	0
?	0	1	1	0
?	0	1	0	1

EM: Repeat until convergence

1. Use probabilistic labels to train classifier h
2. Apply h to assign probabilistic labels to unlabeled data

E Step:

$$\begin{aligned} P(y_i = c_j | d_i; \hat{\theta}) &= \frac{P(c_j | \hat{\theta}) P(d_i | c_j; \hat{\theta})}{P(d_i | \hat{\theta})} \\ &= \frac{P(c_j | \hat{\theta}) \prod_{k=1}^{|d_i|} P(w_{d_{i,k}} | c_j; \hat{\theta})}{\sum_{r=1}^{|\mathcal{C}|} P(c_r | \hat{\theta}) \prod_{k=1}^{|d_i|} P(w_{d_{i,k}} | c_r; \hat{\theta})}. \end{aligned}$$

M Step:

w_t is t-th word in vocabulary

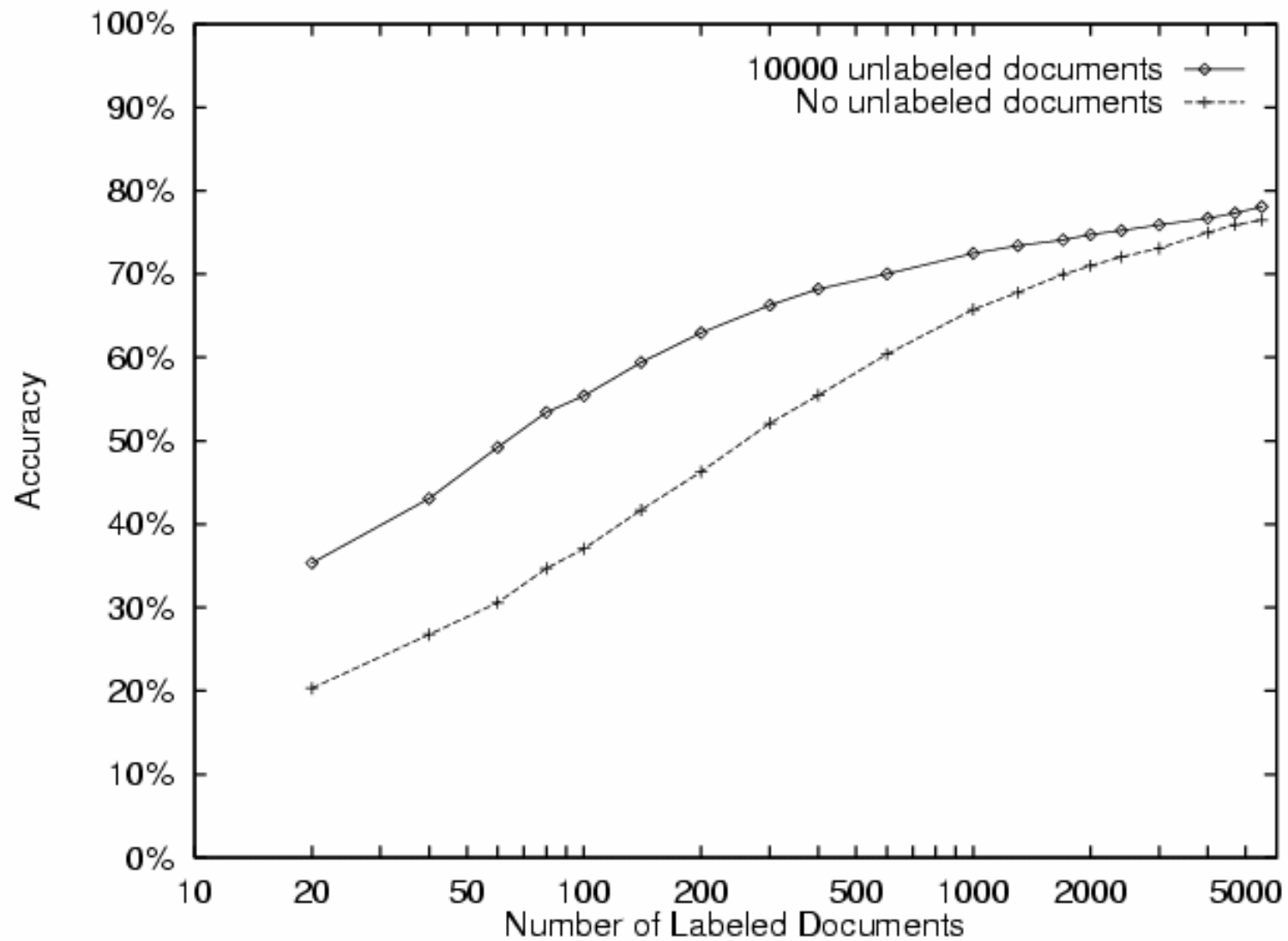
$$\hat{\theta}_{w_t | c_j} \equiv P(w_t | c_j; \hat{\theta}) = \frac{1 + \sum_{i=1}^{|\mathcal{D}|} N(w_t, d_i) P(y_i = c_j | d_i)}{|V| + \sum_{s=1}^{|\mathcal{V}|} \sum_{i=1}^{|\mathcal{D}|} N(w_s, d_i) P(y_i = c_j | d_i)},$$

$$\hat{\theta}_{c_j} \equiv P(c_j | \hat{\theta}) = \frac{1 + \sum_{i=1}^{|\mathcal{D}|} P(y_i = c_j | d_i)}{|\mathcal{C}| + |\mathcal{D}|}.$$

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	Using one labeled example per class	<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	Words sorted by $P(w course) / P(w \neg course)$	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





Why/When will this work?

- What's best case? Worst case? How can we test which we have?

Summary : Semisupervised Learning with EM and Naïve Bayes Model

- If all data is labeled, corresponds to supervised training of Naïve Bayes classifier
- If all data unlabeled, corresponds to unsupervised, mixture-of-multinomial clustering
- If both labeled and unlabeled data, then unlabeled data helps if the Bayes net model is correct (e.g., $P(X)$ is a mixture of class-conditional multinomials with conditionally independent X_i)
- Of course we could use Bayes net models other than Naïve Bayes

Idea 2: Use U to reweight labeled examples

- Most learning algorithms *minimize errors over labeled examples*
- But we really want to *minimize error over future examples* drawn from the same underlying distribution (ie., *true error* of hypothesis)
- If we know the underlying distribution $P(X)$, we could weight each labeled training example $\langle x, y \rangle$ by its probability according to $P(X=x)$
- Unlabeled data allows us to estimate $P(X)$

Idea 2: Use U to reweight labeled examples L

Use $U \rightarrow \hat{P}(X)$ to alter the loss function

- Wish to minimize true error:

$$\hat{f} \leftarrow \arg \min_{h \in H} \sum_{x \in X} \delta(h(x) \neq f(x)) P(x)$$

1 if hypothesis h disagrees with true function f , else 0

- Usually approximate this as:

$$\hat{f} \leftarrow \arg \min_{h \in H} \frac{1}{L} \sum_{\langle x, y \rangle \in L} \delta(h(x) \neq y)$$

Which equals:

$$\hat{f} \leftarrow \arg \min_{h \in H} \sum_{x \in X} \delta(h(x) \neq y) \left[\frac{n(x, L)}{|L|} \right]$$

$n(x, L)$ = number of times x occurs in L

- We can produce a better approximation by incorporating U :

$$\hat{f} \leftarrow \arg \min_{h \in H} \sum_{x \in X} \delta(h(x) \neq f(x)) \left[\frac{n(x, L) + n(x, U)}{|L| + |U|} \delta(n(x, L) > 0) \right]$$

Reweighting Labeled Examples

- Wish to find

$$\hat{f} \leftarrow \arg \min_{h \in H} \sum_{x \in X} \delta(h(x) \neq f(x)) \left[\delta(n(x, L) > 0) \frac{n(x, L) + n(x, U)}{|L| + |U|} \right]$$

- Already have algorithm (e.g., decision tree learner) to find

$$\hat{f} \leftarrow \arg \min_{h \in H} \frac{1}{L} \sum_{\langle x, y \rangle \in L} \delta(h(x) \neq y)$$

- Just reweight examples in L, and have algorithm minimize

$$\hat{f} \leftarrow \arg \min_{h \in H} \frac{1}{L} \sum_{\langle x, y \rangle \in L} \delta(h(x) \neq y) \frac{n(x, L) + n(x, U)}{|L| + |U|}$$

- Or if X is continuous, use L+U to estimate p(X), and minimize

$$\hat{f} \leftarrow \arg \min_{h \in H} \frac{1}{L} \sum_{\langle x, y \rangle \in L} \delta(h(x) \neq y) \hat{p}(x)$$



Idea 3: CoTraining

- In some settings, available data features are redundant and we can train two classifiers based on disjoint features
- In this case, the two classifiers should agree on the classification for each unlabeled example
- Therefore, we can use the unlabeled data to constrain joint training of both classifiers

Redundantly Sufficient Features

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](#))

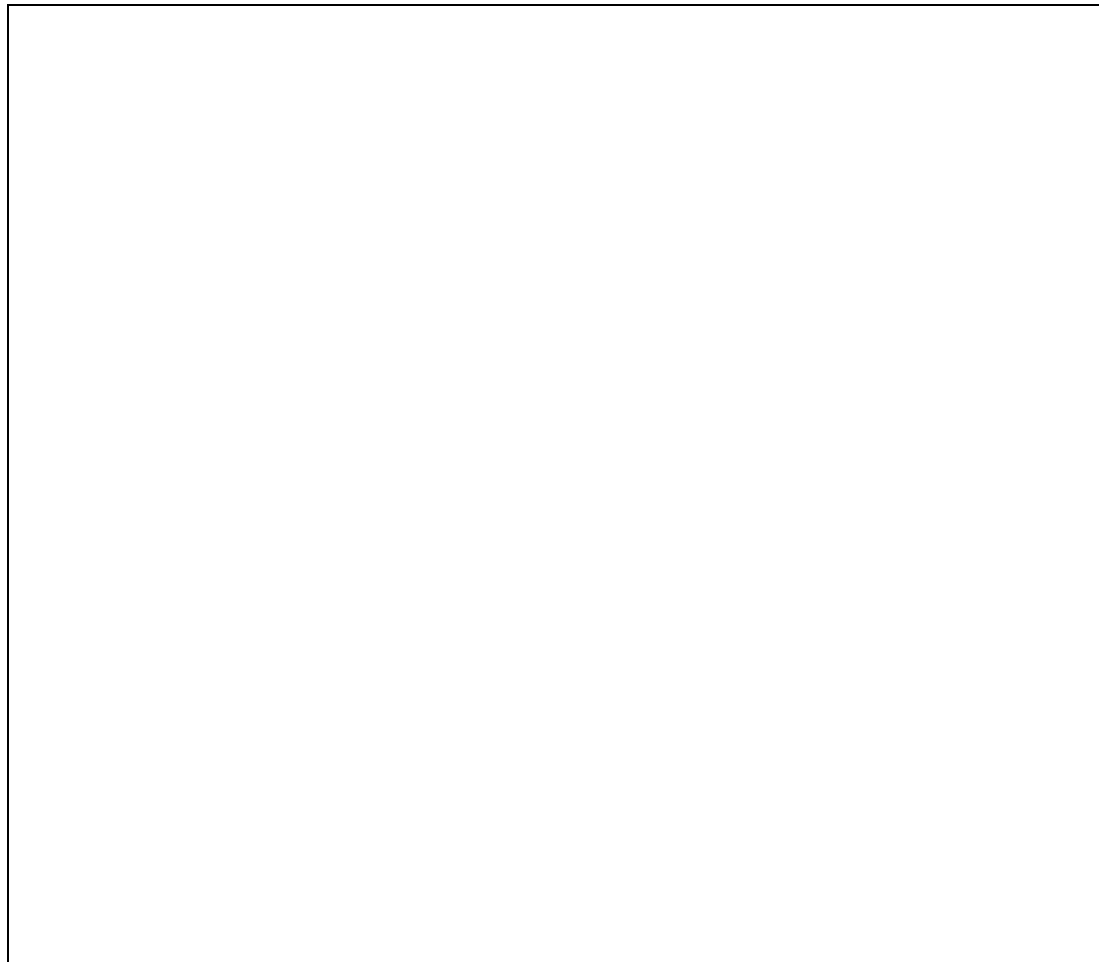
Research Interests:

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

Redundantly Sufficient Features

Professor Faloutsos

my advisor



Redundantly Sufficient Features

**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;

Redundantly Sufficient Features

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;

CoTraining Algorithm #1

[Blum&Mitchell, 1998]

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 exams from U

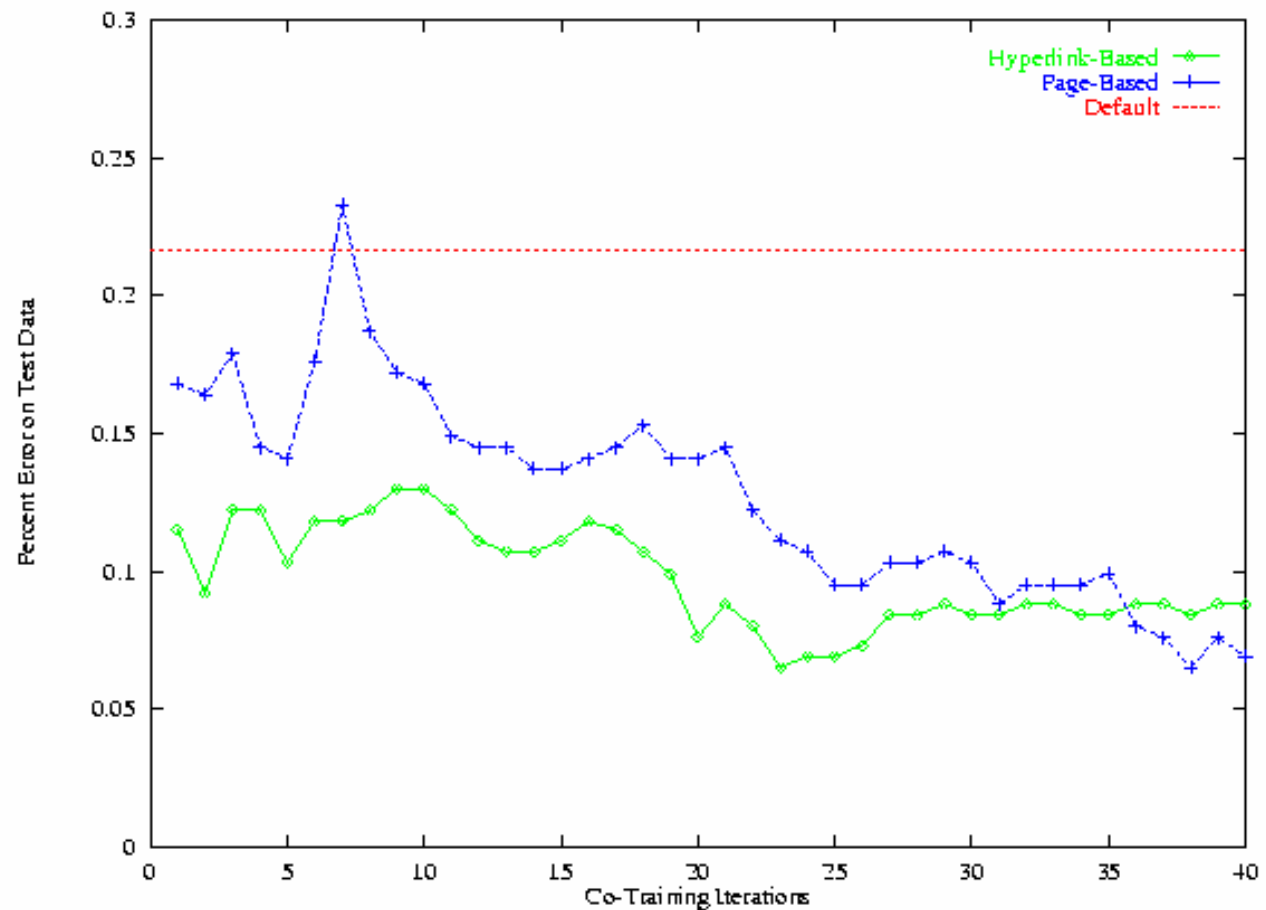
Allow g_2 to label p positive, n negative exams from U

Add these self-labeled examples to L

CoTraining: 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:



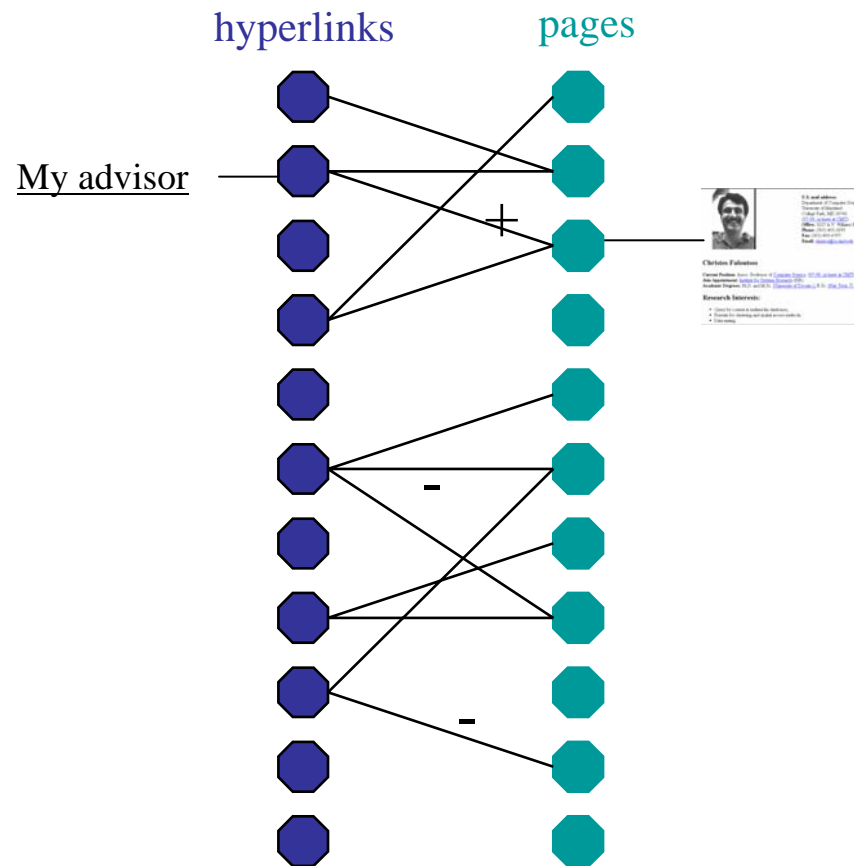
CoTraining setting:

- wish to learn $f: X \rightarrow Y$, given L and U drawn from $P(X)$
- features describing X can be partitioned ($X = X_1 \times X_2$)
such that f can be computed from either X_1 or X_2
$$(\exists g_1, g_2)(\forall x \in X) \quad g_1(x_1) = f(x) = g_2(x_2)$$

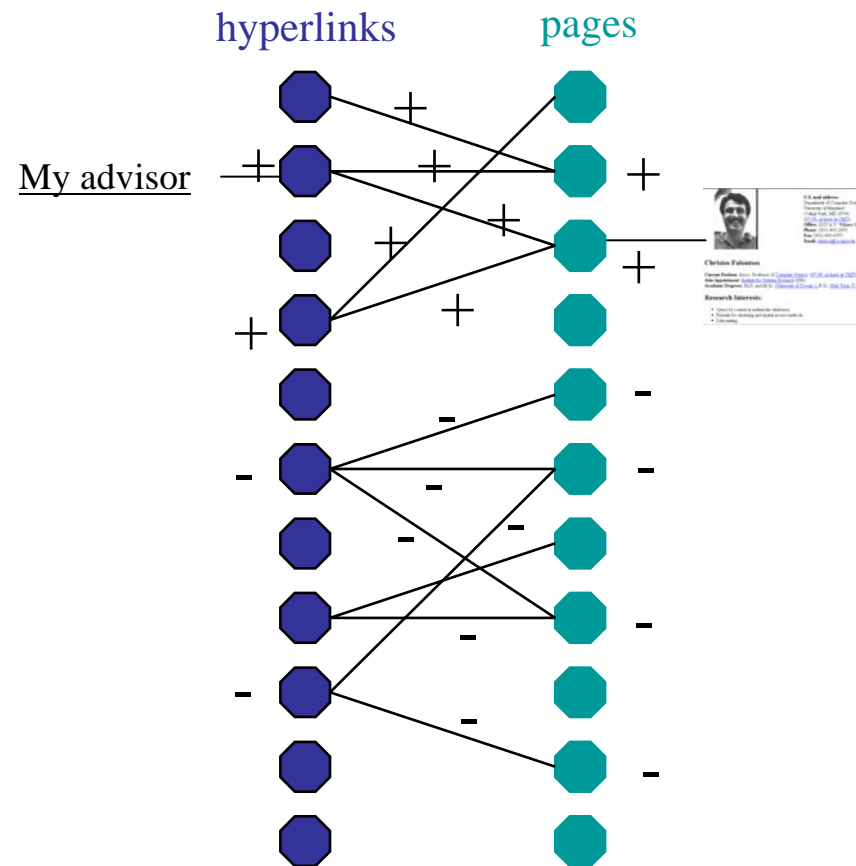
One result [Blum&Mitchell 1998]:

- If
 - X_1 and X_2 are conditionally independent given Y
 - f is PAC learnable from noisy *labeled* data
- Then
 - f is PAC learnable from weak initial classifier plus *unlabeled* data

Co-Training Rote Learner



Co-Training Rote Learner



Expected Rote CoTraining error given m examples

CoTraining setting :

learn $f : X \rightarrow Y$

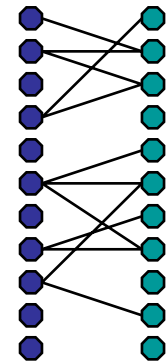
where $X = X_1 \times X_2$

where x drawn from unknown distribution

and $\exists g_1, g_2 \quad (\forall x) g_1(x_1) = g_2(x_2) = f(x)$

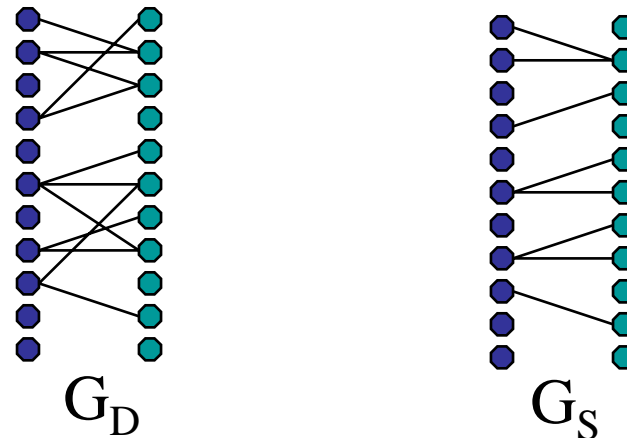
$$E[error] = \sum_j P(x \in g_j)(1 - P(x \in g_j))^m$$

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



How many *unlabeled* examples suffice?

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

PAC Generalization Bounds on CoTraining

[Dasgupta et al., NIPS 2001]

This theorem assumes X_1 and X_2 are conditionally independent given Y

Theorem 1 *With probability at least $1 - \delta$ over the choice of the sample S , we have that for all h_1 and h_2 , if $\gamma_i(h_1, h_2, \delta) > 0$ for $1 \leq i \leq k$ then (a) f is a permutation and (b) for all $1 \leq i \leq k$,*

$$P(h_1 \neq i \mid f(y) = i, h_1 \neq \perp) \leq \frac{\hat{P}(h_1 \neq i \mid h_2 = i, h_1 \neq \perp) + \epsilon_i(h_1, h_2, \delta)}{\gamma_i(h_1, h_2, \delta)}.$$

The theorem states, in essence, that if the sample size is large, and h_1 and h_2 largely agree on the unlabeled data, then $\hat{P}(h_1 \neq i \mid h_2 = i, h_1 \neq \perp)$ is a good estimate of the error rate $P(h_1 \neq i \mid f(y) = i, h_1 \neq \perp)$.

PAC Generalization Bounds on CoTraining

[Dasgupta et al., NIPS 2001]

This theorem assumes X_1 and X_2 are conditionally independent given Y

Theorem 1 *With probability at least $1 - \delta$ over the choice of the sample S , we have that for all h_1 and h_2 , if $\gamma_i(h_1, h_2, \delta) > 0$ for $1 \leq i \leq k$ then (a) f is a permutation and (b) for all $1 \leq i \leq k$,*

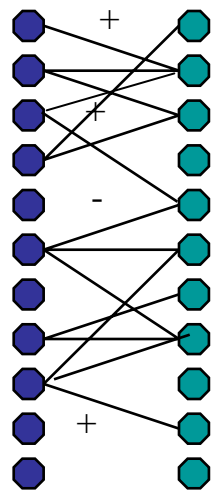
$$P(h_1 \neq i \mid f(y) = i, h_1 \neq \perp) \leq \frac{\hat{P}(h_1 \neq i \mid h_2 = i, h_1 \neq \perp) + \epsilon_i(h_1, h_2, \delta)}{\gamma_i(h_1, h_2, \delta)}.$$

The theorem states, in essence, that if the sample size is large, and h_1 and h_2 largely agree on the unlabeled data, then $\hat{P}(h_1 \neq i \mid h_2 = i, h_1 \neq \perp)$ is a good estimate of the error rate $P(h_1 \neq i \mid f(y) = i, h_1 \neq \perp)$.

$$\gamma_i(h_1, h_2, \delta) = \hat{P}(h_1 = i \mid h_2 = i, h_1 \neq \perp) - \hat{P}(h_1 \neq i \mid h_2 = i, h_1 \neq \perp) - 2\epsilon_i(h_1, h_2, \delta)$$

$$\epsilon_i(h_1, h_2, \delta) = \sqrt{\frac{(\ln 2)(|h_1| + |h_2|) + \ln \frac{2k}{\delta}}{2|S(h_2 = i, h_1 \neq \perp)|}}$$

What if CoTraining Assumption Not Perfectly Satisfied?



- Idea: Want classifiers that produce a *maximally consistent* labeling of the data
- If learning is an optimization problem, what function should we optimize?

Example 2: Learning to extract named entities

location?
↙
I arrived in **Beijing** on Saturday.

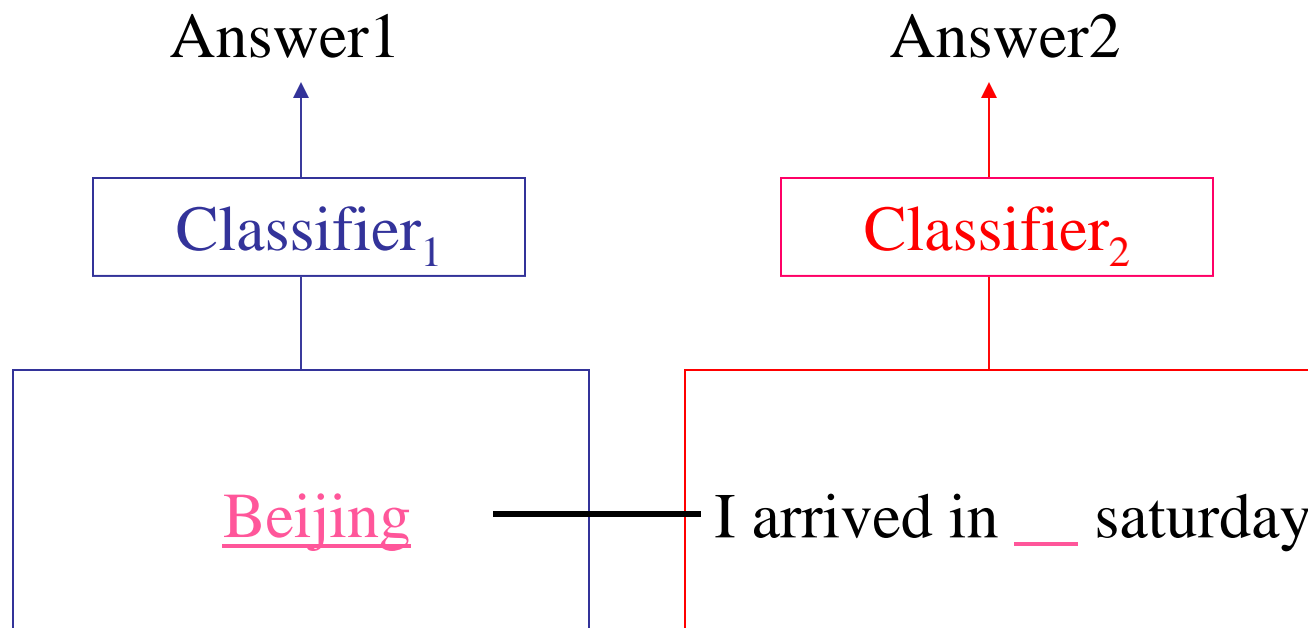
If: “I arrived in <X> on Saturday.”

Then: Location(X)

Co-Training for Named Entity Extraction

(i.e., classifying which strings refer to people, places, dates, etc.)

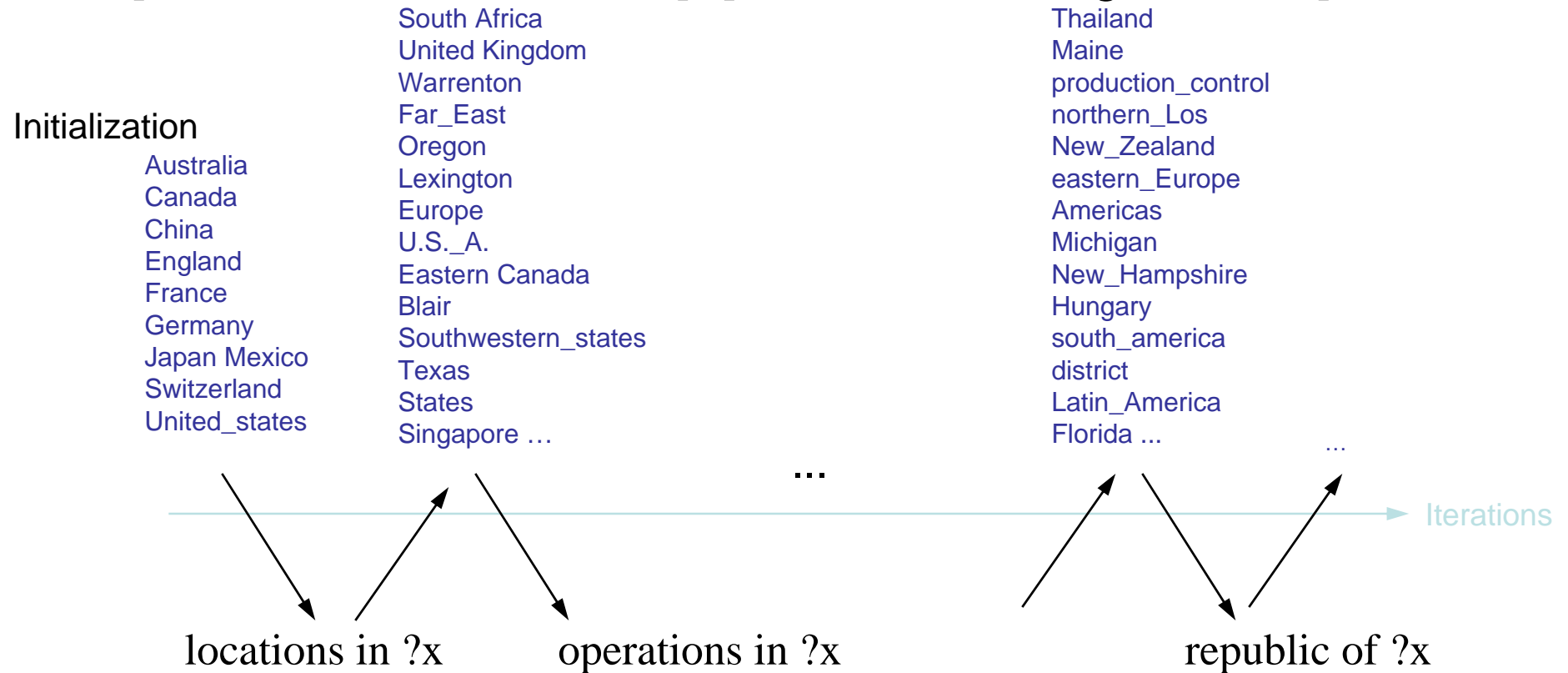
[Riloff&Jones 98; Collins et al., 98; Jones 05]



I arrived in **Beijing** Saturday.

Bootstrap learning to extract named entities

[Riloff and Jones, 1999], [Collins and Singer, 1999], ...



Co-EM [Nigam & Ghani, 2000; Jones 2005]

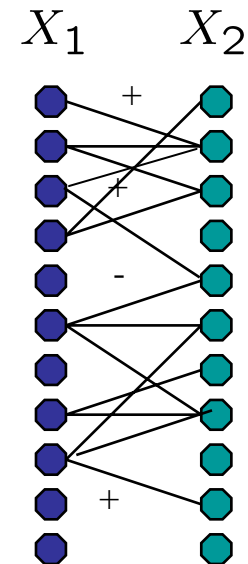
Idea:

- Like co-training, use one set of features to label the other
- Like EM, iterate
 - Assign probabilistic values to unobserved class labels
 - Updating model parameters (= labels of other feature set)

Goal to learn $X_1 \rightarrow Y$, $X_2 \rightarrow Y$, $X_1 \times X_2 \rightarrow Y$

$$P(Y|X_1 = k) = \sum_j P(Y|X_2 = j)P(X_2 = j|X_1 = k)$$

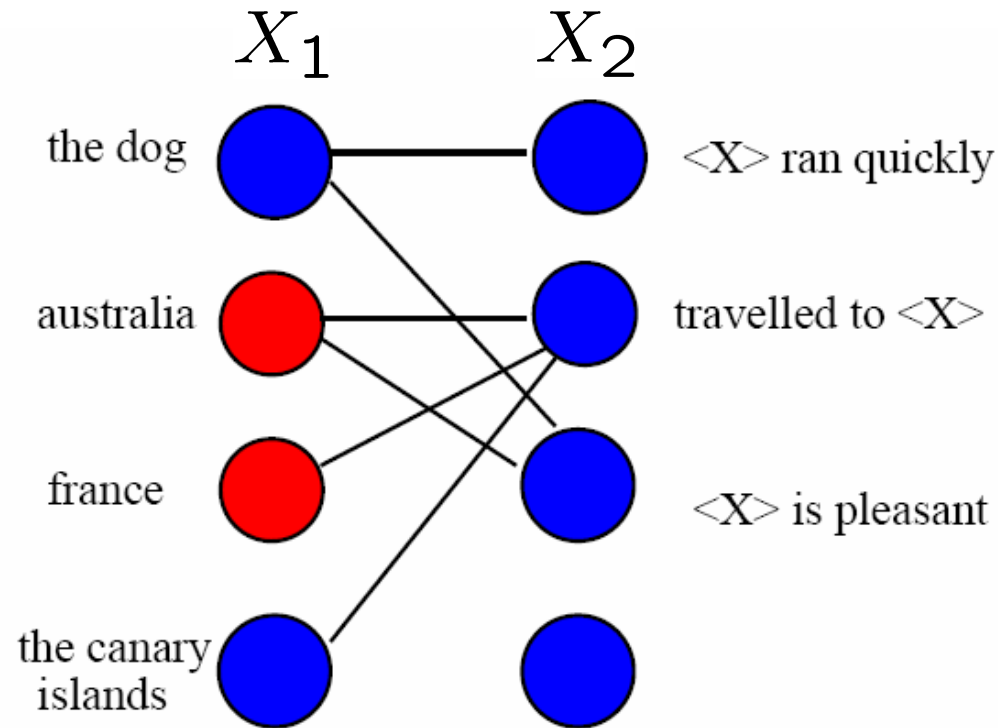
$$P(Y|X_2 = j) = \sum_k P(Y|X_1 = k)P(X_1 = k|X_2 = j)$$



CoEM applied to Named Entity Recognition

[Rosie Jones, 2005], [Ghani & Nigam, 2000]

$$X_1 \times X_2 \rightarrow Y$$



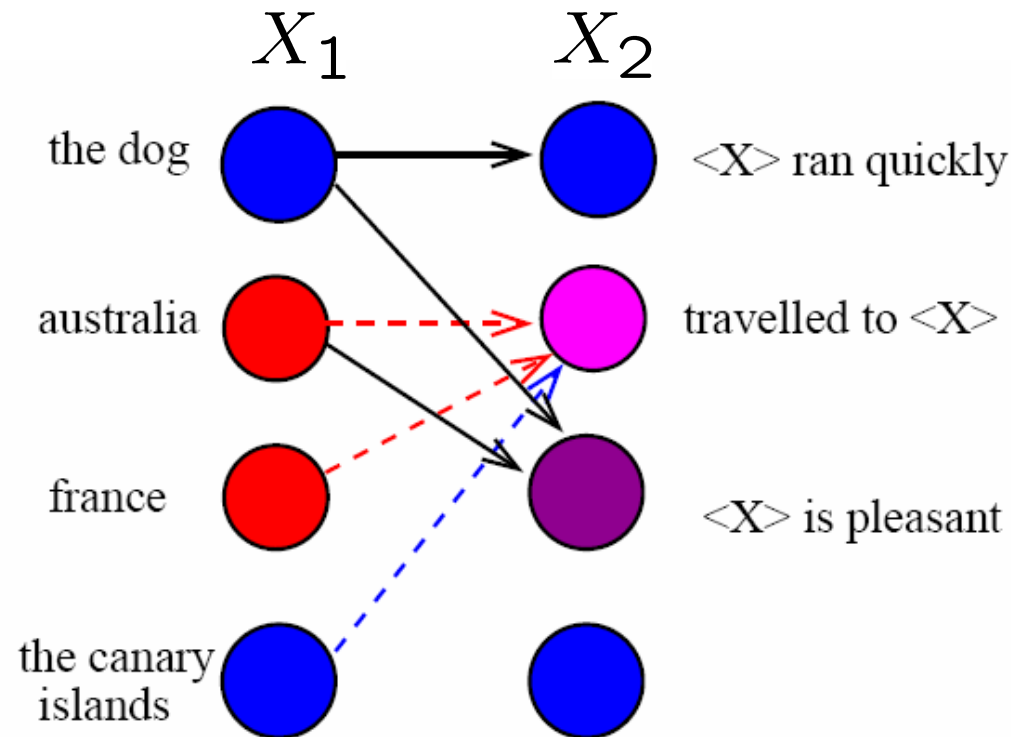
Update
rules:

$$P(Y|X_1 = k) = \sum_j P(Y|X_2 = j)P(X_2 = j|X_1 = k)$$

$$P(Y|X_2 = j) = \sum_k P(Y|X_1 = k)P(X_1 = k|X_2 = j)$$

CoEM applied to Named Entity Recognition

[Rosie Jones, 2005], [Ghani & Nigam, 2000]



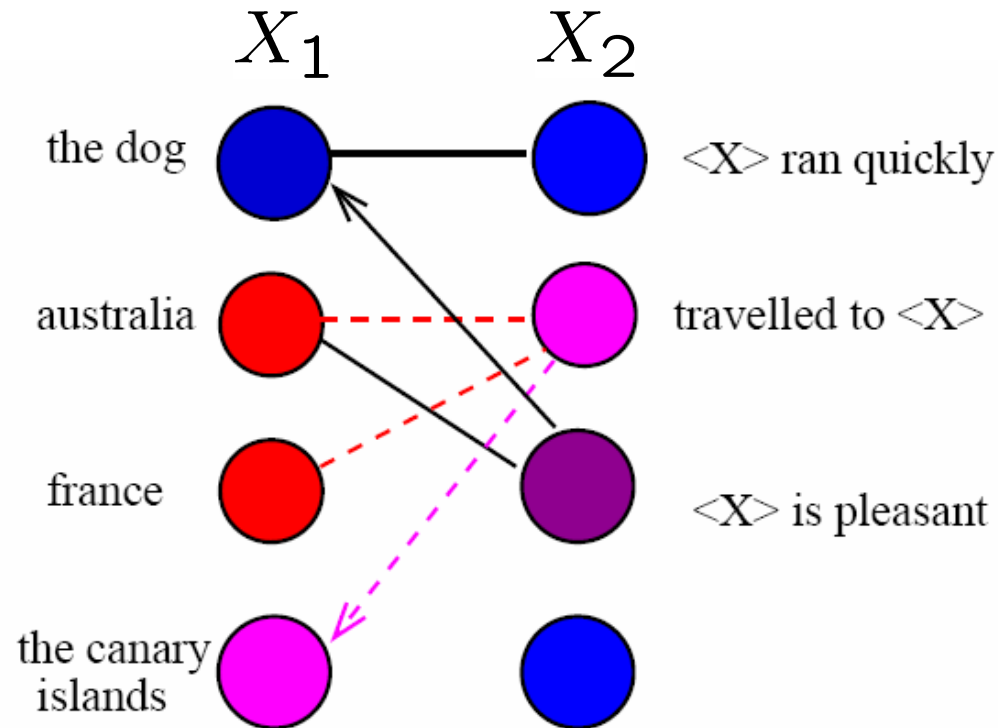
Update
rules:

$$P(Y|X_1 = k) = \sum_j P(Y|X_2 = j)P(X_2 = j|X_1 = k)$$

$$P(Y|X_2 = j) = \sum_k P(Y|X_1 = k)P(X_1 = k|X_2 = j)$$

CoEM applied to Named Entity Recognition

[Rosie Jones, 2005], [Ghani & Nigam, 2000]



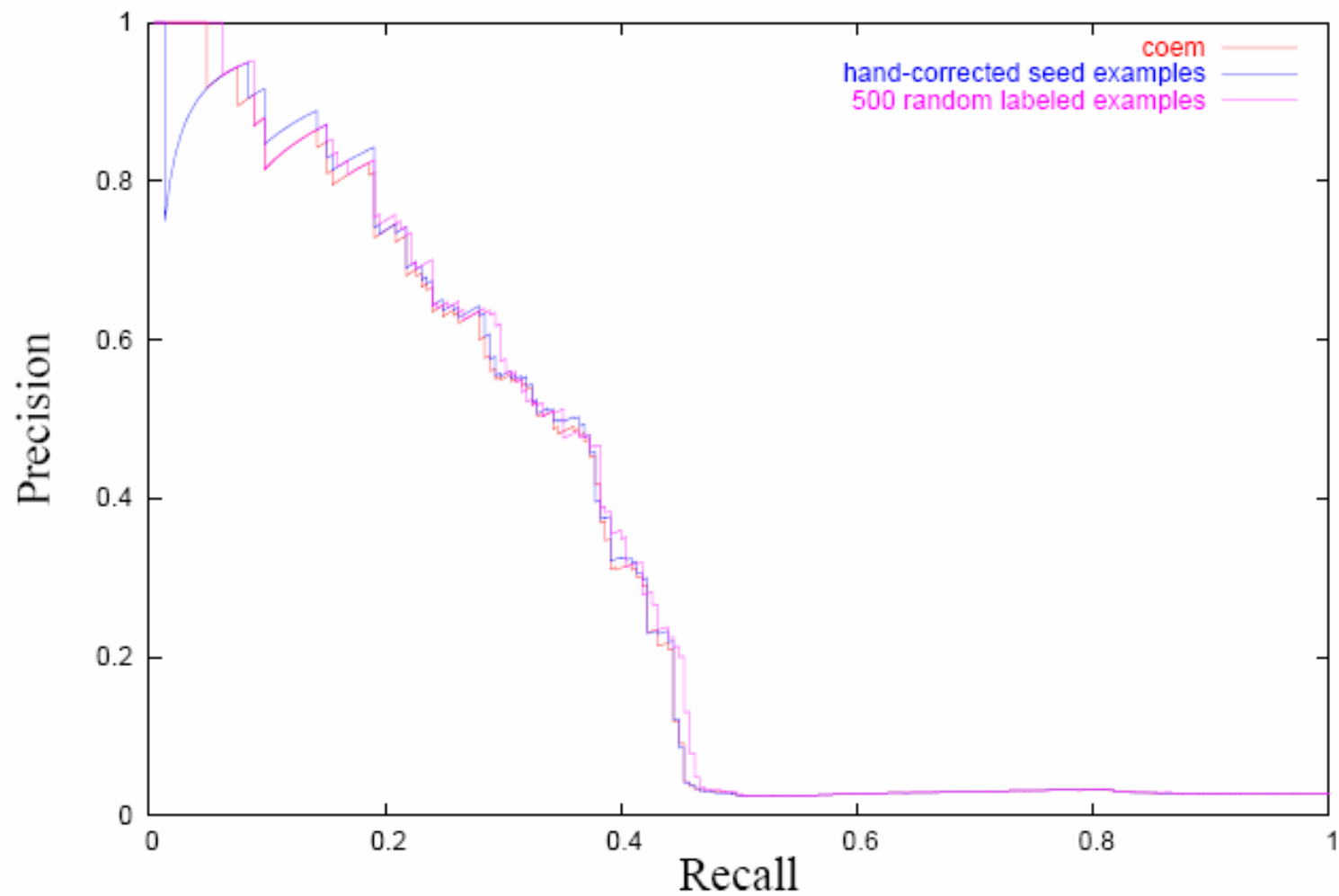
Update
rules:

$$P(Y|X_1 = k) = \sum_j P(Y|X_2 = j)P(X_2 = j|X_1 = k)$$

$$P(Y|X_2 = j) = \sum_k P(Y|X_1 = k)P(X_1 = k|X_2 = j)$$

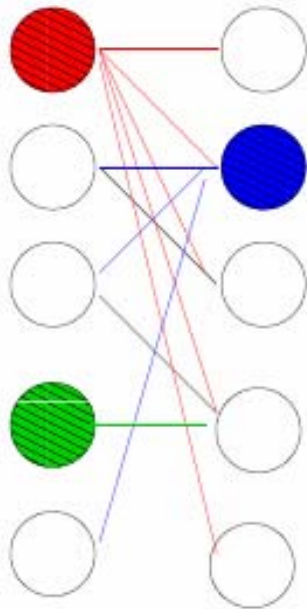
Bootstrapping Results

locations



Some nodes are more important than others [Jones, 2005]

Can use this for active learning...



Noun-phrase	Outdegree
you	1656
we	1479
it	1173
company	1043
this	635
all	520
they	500
information	448
us	367
any	339
products	332
i	319
site	314
one	311
1996	282
he	269
customers	269
these	263
them	263
time	234

Context	Outdegree
<x> including	683
including <x>	612
<x> provides	565
provides <x>	565
provide <x>	390
<x> include	389
include <x>	375
<x> provide	364
one of <x>	354
<x> made	345
<x> offers	338
offers <x>	320
<x> said	287
<x> used	283
includes <x>	279
to provide <x>	266
use <x>	263
like <x>	260
variety of <x>	252
<x> includes	250

CoTraining Summary

- Unlabeled data improves supervised learning when example features are redundantly sufficient
 - Family of algorithms that train multiple classifiers
- Theoretical results
 - Expected error for rote learning
 - If X_1, X_2 conditionally independent given Y , Then
 - PAC learnable from weak initial classifier plus unlabeled data
 - disagreement between $g_1(x_1)$ and $g_2(x_2)$ bounds final classifier error
- Many real-world problems of this type
 - Semantic lexicon generation [Riloff, Jones 99], [Collins, Singer 99]
 - Web page classification [Blum, Mitchell 98]
 - Word sense disambiguation [Yarowsky 95]
 - Speech recognition [de Sa, Ballard 98]
 - Visual classification of cars [Levin, Viola, Freund 03]



4. Use U to Detect/Preempt Overfitting

- Overfitting is a problem for many learning algorithms (e.g., decision trees, neural networks)
- The symptom of overfitting: complex hypothesis h_2 performs better on training data than simpler hypothesis h_1 , but worse on test data
- Unlabeled data can help detect overfitting, by comparing predictions of h_1 and h_2 over the unlabeled examples
 - The rate at which h_1 and h_2 disagree on U should be the same as the rate on L , unless overfitting is occurring

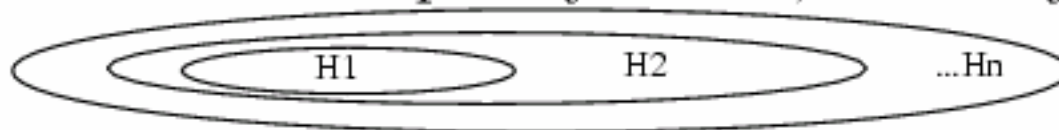
4. Use U to Detect/Preempt Overfitting

Define *metric* over $H \cup \{f\}$

definition $\rightarrow d(h_1, h_2) \equiv \int \delta(h_1(x) \neq h_2(x))p(x)dx$

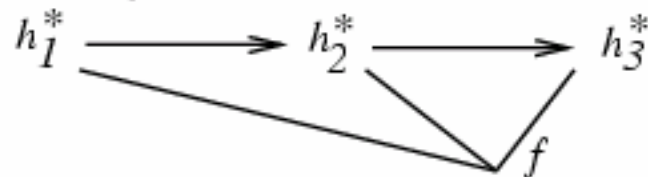
estimates $\rightarrow \hat{d}(h_1, f) = \frac{1}{|L|} \sum_{x_i \in L} \delta(h_1(x_i) \neq y_i)$
 $\rightarrow \hat{d}(h_1, h_2) = \frac{1}{|U|} \sum_{x \in U} \delta(h_1(x) \neq h_2(x))$

Organize H into complexity classes, sorted by $P(h)$



Let h_i^* be hypothesis with lowest $\hat{d}(h, f)$ in H_i

Prefer h_1^* , h_2^* , or h_3^* ?



- Definition of distance metric

- Non-negative $d(f,g) \geq 0$;
- symmetric $d(f,g) = d(g,f)$;
- triangle inequality $d(f,g) \leq d(f,h) + d(h,g)$

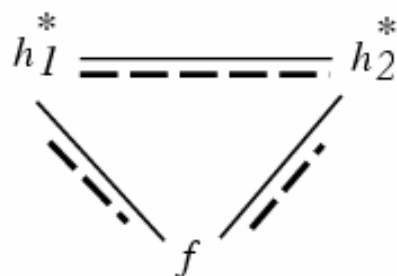
- Classification with zero-one loss:

$$d(h_1, h_2) \equiv \int \delta(h_1(x) \neq h_2(x)) p(x) dx$$

- Regression with squared loss:

$$d(h_1, h_2) \equiv \sqrt{\int (h_1(x) - h_2(x))^2 p(x) dx}$$

Idea: Use U to Avoid Overfitting



Note:

- $\hat{d}(h_i^*, f)$ optimistically biased (too short)
- $\hat{d}(h_i^*, h_j^*)$ unbiased
- Distances must obey triangle inequality!

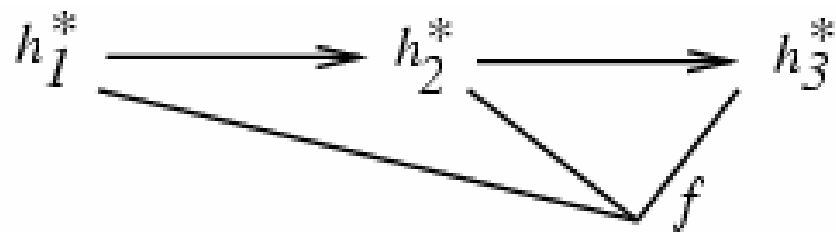
$$d(h_1, h_2) \leq d(h_1, f) + d(f, h_2)$$

→ Heuristic:

- Continue training until $\hat{d}(h_i, h_{i+1})$ fails to satisfy triangle inequality

Procedure TRI

- Given hypothesis sequence h_0, h_1, \dots
- Choose the last hypothesis h_ℓ in the sequence that satisfies the triangle inequality $d(h_k, h_\ell) \leq d(h_k, \widehat{P_{Y|X}}) + d(h_\ell, \widehat{P_{Y|X}})$ with every preceding hypothesis h_k , $0 \leq k < \ell$. (Note that the inter-hypothesis distances $d(h_k, h_\ell)$ are measured on the *unlabeled* training data.)



Experimental Evaluation of TRI

[Schuermans & Southey, MLJ 2002]

- Use it to select degree of polynomial for regression
- Compare to alternatives such as cross validation, structural risk minimization, ...

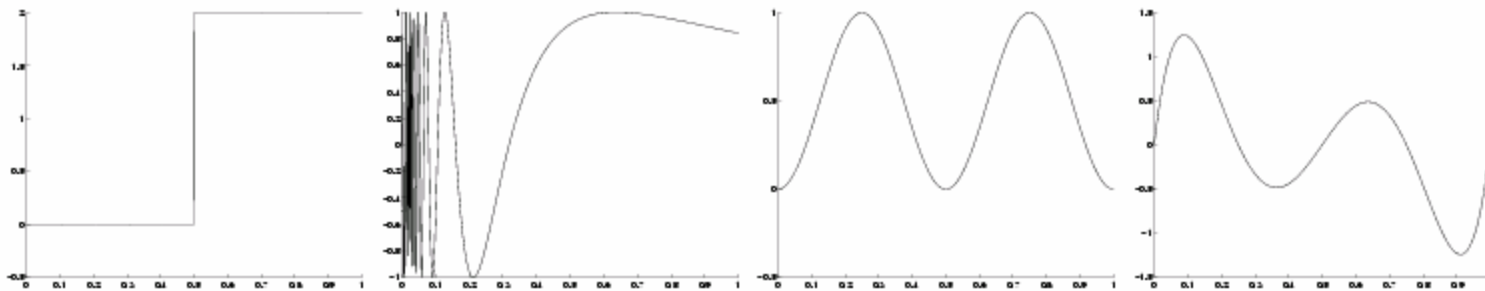


Figure 5: Target functions used in the polynomial curve fitting experiments (in order): $\text{step}(x \geq 0.5)$, $\sin(1/x)$, $\sin^2(2\pi x)$, and a fifth degree polynomial.

Generated y
values contain
zero mean
Gaussian noise ε

$$Y = f(x) + \varepsilon$$

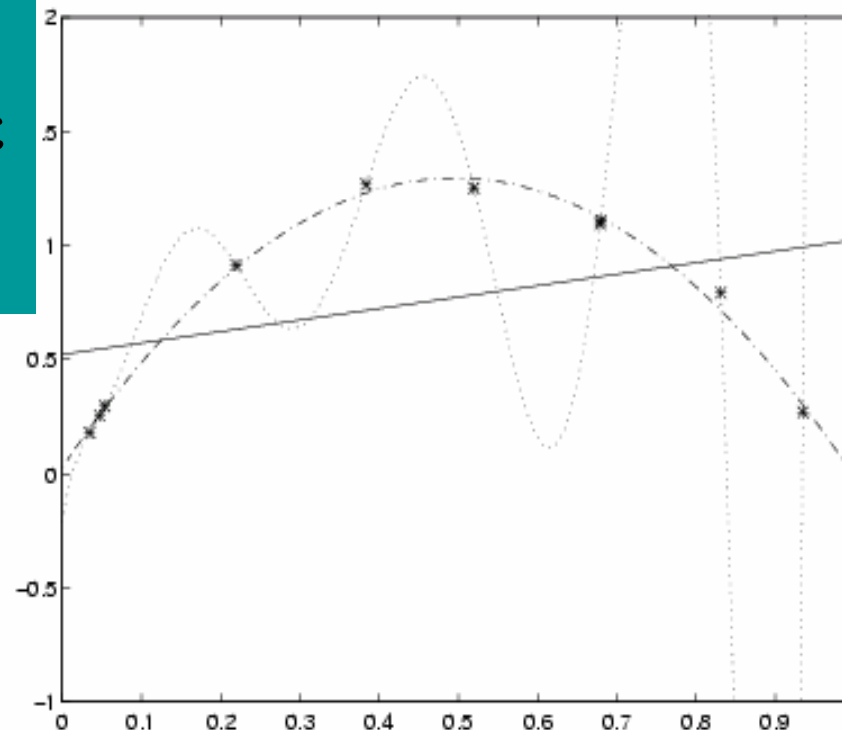


Figure 4: An example of minimum squared error polynomials of degrees 1, 2, and 9 for a set of 10 training points. The large degree polynomial demonstrates erratic behavior off the training set.

Approximation ratio:

true error of selected hypothesis

true error of best hypothesis considered

Results using 200 unlabeled, t labeled

Cross validation (Ten-fold)

Structural risk minimization

Worst
performance
in top .50 of
trials

$t = 20$	TRI	CVT	SRM	RIC	GCV	BIC	AIC	FPE	ADJ
25	1.00	1.06	1.14	7.54	5.47	15.2	22.2	25.8	1.02
50	1.06	1.17	1.39	224	118	394	585	590	1.12
75	1.17	1.42	3.62	5.8e3	3.9e3	9.8e3	1.2e4	1.2e4	1.24
95	1.44	6.75	56.1	6.1e5	3.7e5	7.8e5	9.2e5	8.2e5	1.54
100	2.41	1.1e4	2.2e4	1.5e8	6.5e7	1.5e8	1.5e8	8.2e7	3.02

$t = 30$	TRI	CVT	SRM	RIC	GCV	BIC	AIC	FPE	ADJ
25	1.00	1.08	1.17	4.69	1.51	5.41	5.45	2.72	1.06
50	1.08	1.17	1.54	34.8	9.19	39.6	40.8	19.1	1.14
75	1.19	1.37	9.68	258	91.3	266	266	159	1.25
95	1.45	6.11	419	4.7e3	2.7e3	4.8e3	5.1e3	4.0e3	1.51
100	2.18	643	1.6e7	1.6e7	1.6e7	1.6e7	1.6e7	1.6e7	2.10

Table 1: Fitting $f(x) = \text{step}(x \geq 0.5)$ with $P_x = U(0, 1)$ and $\sigma = 0.05$. Tables give distribution of approximation ratios achieved at training sample size $t = 20$ and $t = 30$, showing percentiles of approximation ratios achieved in 1000 repeated trials.

$t = 20$	TRI	CVT	SRM	RIC	GCV	BIC	AIC	FPE	ADJ
25	2.04	1.03	1.00	1.00	1.06	1.00	1.01	1.58	1.02
50	3.11	1.37	1.33	1.34	1.94	1.35	1.61	18.2	1.32
75	3.87	2.23	2.30	2.13	10.0	2.75	4.14	1.2e3	1.83
95	5.11	9.45	8.84	8.26	5.0e3	11.8	82.9	1.8e5	3.94
100	8.92	105	526	105	2.0e7	2.1e3	2.7e5	2.4e7	6.30

$t = 30$	TRI	CVT	SRM	RIC	GCV	BIC	AIC	FPE	ADJ
25	1.50	1.00	1.00	1.00	1.00	1.00	1.00	1.02	1.01
50	3.51	1.16	1.03	1.05	1.11	1.02	1.08	1.45	1.27
75	4.15	1.64	1.45	1.48	2.02	1.39	1.88	6.44	1.60
95	5.51	5.21	5.06	4.21	26.4	5.01	19.9	295	3.02
100	9.75	124	1.4e3	20.0	9.1e3	28.4	9.4e3	1.0e4	8.35

Table 4: Fitting $f(x) = \sin^2(2\pi x)$ with $P_x = U(0, 1)$ and $\sigma = 0.05$. Tables give distribution of approximation ratios achieved at training sample size $t = 20$ and $t = 30$, showing percentiles of approximation ratios achieved in 1000 repeated trials.

Bound on Error of TRI Relative to Best Hypothesis Considered

Proposition 1 *Let h_m be the optimal hypothesis in the sequence h_0, h_1, \dots (that is, $h_m = \arg \min_{h_k} d(h_k, \widehat{P_{Y|X}})$) and let h_ℓ be the hypothesis selected by TRI. If (i) $m \leq \ell$ and (ii) $d(h_m, \widehat{P_{Y|X}}) \leq d(h_m, P_{Y|X})$ then*

$$d(h_\ell, P_{Y|X}) \leq 3d(h_m, P_{Y|X}) \quad (6)$$

Extension to TRI:

Adjust for expected bias of training data estimates
[Schuermans & Southey, MLJ 2002]

Procedure ADJ

- Given hypothesis sequence h_0, h_1, \dots
- For each hypothesis h_ℓ in the sequence
 - multiply its estimated distance to the target $d(h_\ell, \widehat{P}_{Y|X})$ by the worst ratio of unlabeled and labeled distance to some predecessor h_k to obtain an adjusted distance estimate $d(\widehat{\widehat{h_\ell}}, \widehat{\widehat{P_{Y|X}}}) = d(h_\ell, \widehat{P_{Y|X}}) \frac{d(h_k, h_\ell)}{d(\widehat{\widehat{h_k}}, \widehat{\widehat{P_{Y|X}}})}$.
- Choose the hypothesis h_n with the smallest adjusted distance $d(\widehat{\widehat{h_n}}, \widehat{\widehat{P_{Y|X}}})$.

Experimental results: averaged over multiple target functions,
outperforms TRI

What you should know

1. Unlabeled can help EM learn Bayes nets for $P(X,Y)$
 - If we assume the Bayes net structure is correct
2. Using unlabeled data to reweight labeled examples gives better approximation to true error
 - If we assume examples drawn from stationary $P(X)$
3. CoTraining multiple classifiers, using unlabeled data as constraints
 - If we assume redundantly sufficient features, with different conditional distributions given the class
4. Use unlabeled data to detect/preempt overfitting
 - If we assume priors over H that correctly order hypotheses

Further Reading

- Semi-Supervised Learning, O. Chapelle, B. Sholkopf, and A. Zien (eds.), MIT Press, 2006. (excellent new book)
- EM for Naïve Bayes classifiers: K.Nigam, et al., 2000. "Text Classification from Labeled and Unlabeled Documents using EM", *Machine Learning*, 39, pp.103—134.
- CoTraining: A. Blum and T. Mitchell, 1998. "Combining Labeled and Unlabeled Data with Co-Training," *Proceedings of the 11th Annual Conference on Computational Learning Theory (COLT-98)*.
- S. Dasgupta, et al., "PAC Generalization Bounds for Co-training", *NIPS 2001*
- Model selection: D. Schuurmans and F. Southey, 2002. "Metric-Based methods for Adaptive Model Selection and Regularization," *Machine Learning*, 48, 51—84.