# 10601
# Machine Learning

Semi supervised learning

# Can Unlabeled Data improve supervised learning?

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

- Medical outcomes (x=<patient,treatment>, y=outcome)

- Text classification (x=document, y=relevance)

- Customer modeling (x=user actions, y=user intent)

- …

# When can Unlabeled Data help supervised learning?

Consider setting:

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

Given:

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

Determine:

$$\hat{f} \leftarrow \arg\min_{h \in H} \Pr_{x \in P(X)} [h(x) \neq f(x)]$$
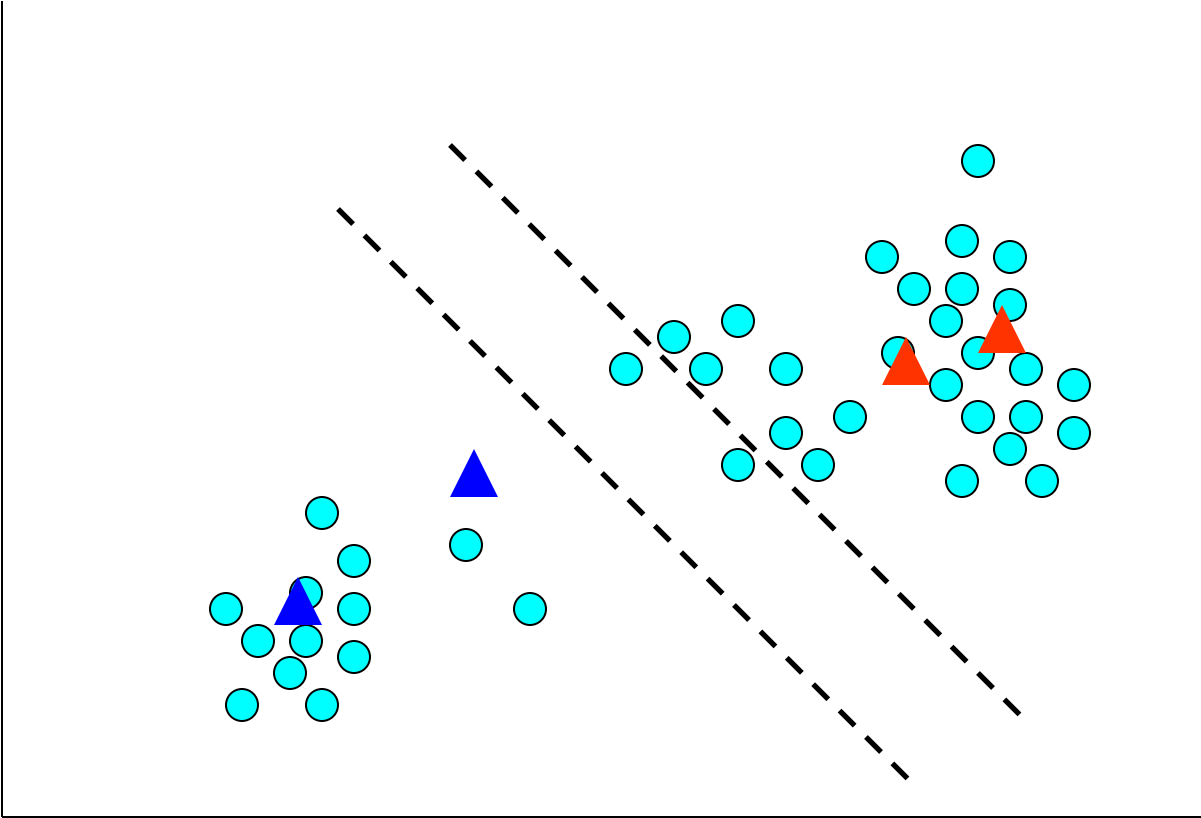
# Four Ways to Use Unlabeled Data for Supervised Learning

1. Use to re-weight labeled examples

2. Use to help EM learn class-specific generative models

3. If problem has redundantly sufficient features, use CoTraining.

4. Use to detect/preempt overfitting

# 1. Use unlabeled data to reweight labeled examples

- Most machine learning algorithms (neural nets, decision trees, SVMs) attempt to *minimize errors over labeled examples*

- But our ultimate goal is to *minimize error over future examples* drawn from the same underlying distribution

- If we know the underlying distribution, we should weight each training example by its probability according to this distribution

- Unlabeled data allows us to estimate this distribution more accurately, and to reweight our labeled examples accordingly

# Example

# 1. reweight labeled examples

Can use $U \to \hat{P}(X)$ to alter optimization problem

- Wish to find

$$\hat{f} \leftarrow \underset{h \in H}{\operatorname{argmin}} \sum_{x \in X} \delta(h(x) \neq f(x)) P(x)$$

- Often approximate as

$$\hat{f} \leftarrow \underset{h \in H}{\operatorname{argmin}} \frac{1}{|L|} \sum_{\langle x, y \rangle \in L} \delta(h(x) \neq y)$$

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

# 1. reweight labeled examples

Can use $U \to \hat{P}(X)$ to alter optimization problem

- Wish to find

$$\hat{f} \leftarrow \operatorname*{argmin}_{h \in H} \sum_{x \in X} \delta(h(x) \neq f(x)) P(x)$$

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

- Often approximate as

$$\hat{f} \leftarrow \operatorname*{argmin}_{h \in H} \frac{1}{|L|} \sum_{\langle x,y \rangle \in L} \delta(h(x) \neq y)$$
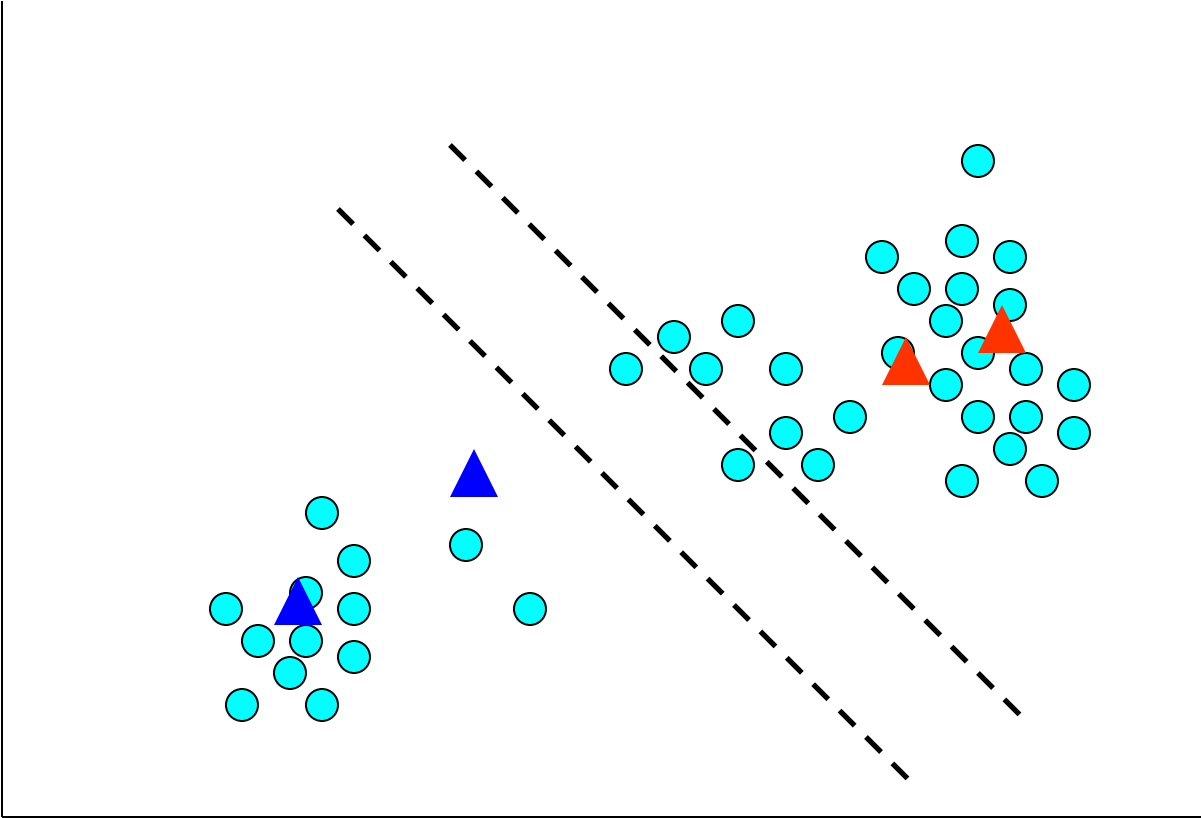
# of times we have x in the labeled set

$$\hat{f} \leftarrow \operatorname*{argmin}_{h \in H} \sum_{x \in X} \delta(h(x) \neq f(x)) \frac{n(x, L)}{|L|}$$

# 1. reweight labeled examples

Can use $U \rightarrow \hat{P}(X)$ to alter optimization problem

- Wish to find

$$\hat{f} \leftarrow \underset{h \in H}{\operatorname{argmin}} \sum_{x \in X} \delta(h(x) \neq f(x)) P(x)$$

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

- Often approximate as

$$\hat{f} \leftarrow \underset{h \in H}{\operatorname{argmin}} \frac{1}{|L|} \sum_{\langle x, y \rangle \in L} \delta(h(x) \neq y)$$

$$\hat{f} \leftarrow \underset{h \in H}{\operatorname{argmin}} \sum_{x \in X} \delta(h(x) \neq f(x)) \frac{n(x, L)}{|L|}$$

\# of times we have x in the labeled set

- Can use $U$ for improved approximation:

$$\hat{f} \leftarrow \underset{h \in H}{\operatorname{argmin}} \sum_{x \in X} \delta(h(x) \neq f(x)) \frac{n(x, L) + n(x, U)}{|L| + |U|}$$

\# of times we have x in the unlabeled set

# Example

# 2. Improve EM clustering algorithms

- Consider completely unsupervised clustering, where we assume data X is generated by a mixture of probability distributions, one for each cluster
  - For example, Gaussian mixtures
- Some classifier learning algorithms such as Gaussian Bayes classifiers also assumes the data X is generated by a mixture of distributions, one for each class Y
- Supervised learning: estimate $P(X|Y)$ from labeled data
- Opportunity: estimate $P(X|Y)$ from labeled and unlabeled data, using EM as in clustering

# Bag of Words Text Classification



| | |
|---|---|
| aardvark | 0 |
| about | 2 |
| all | 2 |
| Africa | 1 |
| apple | 0 |
| anxious | 0 |
| ... | |
| gas | 1 |
| ... | |
| oil | 1 |
| … | |
| Zaire | 0 |

# Baseline: Naïve Bayes Learner

***Train***:

For each class $c_j$ of documents

1. Estimate $P(c_j)$

2. For each word $w_i$ estimate $P(w_i / c_j)$

***Classify (doc):***

Assign *doc* to most probable class

$$\arg \max_j P(c_j) \prod_{w_i \in doc} P(w_i \mid c_j)$$

Naïve Bayes assumption: words are conditionally independent, given class

## Faculty

| | |
|---|---|
| associate | 0.00417 |
| chair | 0.00303 |
| member | 0.00288 |
| ph | 0.00287 |
| director | 0.00282 |
| fax | 0.00279 |
| journal | 0.00271 |
| recent | 0.00260 |
| received | 0.00258 |
| award | 0.00250 |

## Students

| | |
|---|---|
| resume | 0.00516 |
| advisor | 0.00456 |
| student | 0.00387 |
| working | 0.00361 |
| stuff | 0.00359 |
| links | 0.00355 |
| homepage | 0.00345 |
| interests | 0.00332 |
| personal | 0.00332 |
| favorite | 0.00310 |

## Courses

| | |
|---|---|
| homework | 0.00413 |
| syllabus | 0.00399 |
| assignments | 0.00388 |
| exam | 0.00385 |
| grading | 0.00381 |
| midterm | 0.00374 |
| pm | 0.00371 |
| instructor | 0.00370 |
| due | 0.00364 |
| final | 0.00355 |

## Departments

| | |
|---|---|
| departmental | 0.01246 |
| colloquia | 0.01076 |
| epartment | 0.01045 |
| seminars | 0.00997 |
| schedules | 0.00879 |
| webmaster | 0.00879 |
| events | 0.00826 |
| facilities | 0.00807 |
| eople | 0.00772 |
| postgraduate | 0.00764 |

## Research Projects

| | |
|---|---|
| investigators | 0.00256 |
| group | 0.00250 |
| members | 0.00242 |
| researchers | 0.00241 |
| laboratory | 0.00238 |
| develop | 0.00201 |
| related | 0.00200 |
| arpa | 0.00187 |
| affiliated | 0.00184 |
| project | 0.00183 |

## Others

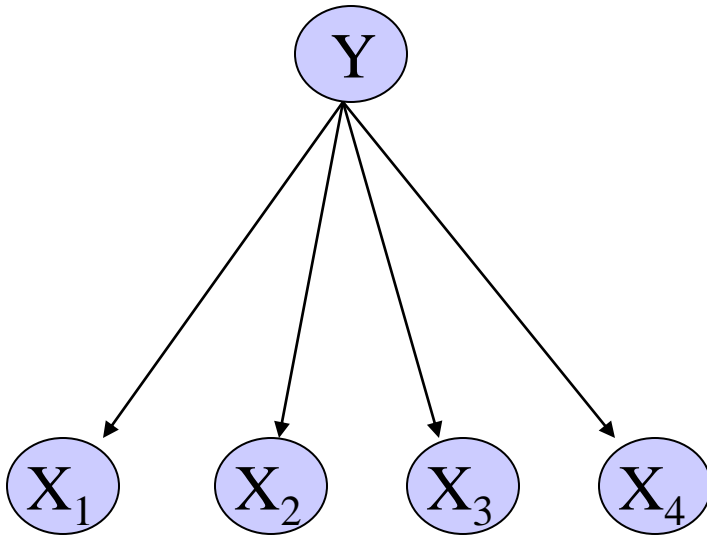| | |
|---|---|
| type | 0.00164 |
| jan | 0.00148 |
| enter | 0.00145 |
| random | 0.00142 |
| program | 0.00136 |
| net | 0.00128 |
| time | 0.00128 |
| format | 0.00124 |
| access | 0.00117 |
| begin | 0.00116 |

# Expectation Maximization (EM) Algorithm

- Use labeled data $L$ to learn initial classifier $h$

Loop:

- E Step:
  - Assign probabilistic labels to $U$, based on $h$

- M Step:
  - Retrain classifier $h$ using both $L$ (with fixed membership) and assigned labels to $U$ (soft membership)

- Under certain conditions, guaranteed to converge to locally maximum likelihood $h$

# 2. Generative Bayes model

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  |

**E Step:**

$$\mathrm{P}(y_i = c_j | d_i; \hat{\theta}) = \frac{\mathrm{P}(c_j | \hat{\theta}) \mathrm{P}(d_i | c_j; \hat{\theta})}{\mathrm{P}(d_i | \hat{\theta})}$$

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

$w_t$ is t-th word in vocabulary

**M Step:**

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

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

**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.
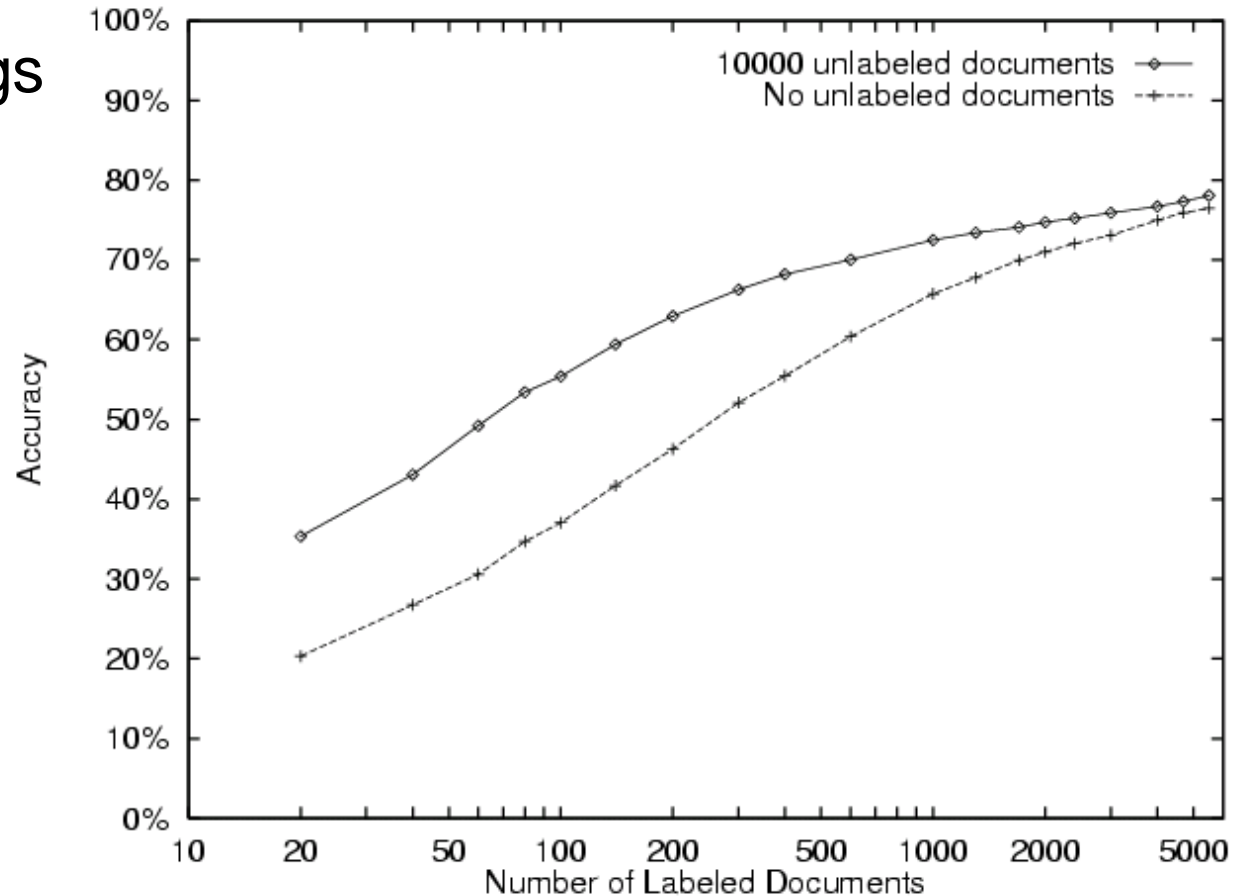
| Iteration 0 | Iteration 1 | Iteration 2 |
|:---:|:---:|:---:|
| intelligence | $DD$ | $D$ |
| $DD$ | $D$ | $DD$ |
| artificial | lecture | lecture |
| understanding | cc | cc |
| $DD$w | $D^\star$ | $DD{:}DD$ |
| dist | $DD{:}DD$ | due |
| identical | handout | $D^\star$ |
| rus | due | homework |
| arrange | problem | assignment |
| games | set | handout |
| dartmouth | tay | set |
| natural | $DD$am | hw |
| cognitive | yurttas | exam |
| logic | homework | problem |
| proving | kfoury | $DD$am |
| prolog | sec | postscript |
| knowledge | postscript | solution |
| human | exam | quiz |
| representation | solution | chapter |
| field | assaf | ascii |

Using one labeled example per class

# Experimental Evaluation

Newsgrop postings
– 20 newsgroups,
1000/group

# 3. If Problem Setting Provides Redundantly Sufficient Features, use CoTraining

- In some settings, available data features are so redundant that we can train two classifiers using different 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 training of both classifiers, forcing them to agree

# 3. CoTraining

$$learn \quad f : X \to Y$$

$$where \quad X = X_1 \times X_2$$

$$where \quad x \quad drawn \quad from \quad unknown \quad distribution$$

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

# 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
### [Blum&Mitchell, 1998]

Given: labeled data L,

unlabeled data U

Loop:

Train g1 (hyperlink classifier) using L

Train g2 (page classifier) using L

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

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

Add the intersection of the self-labeled examples to L
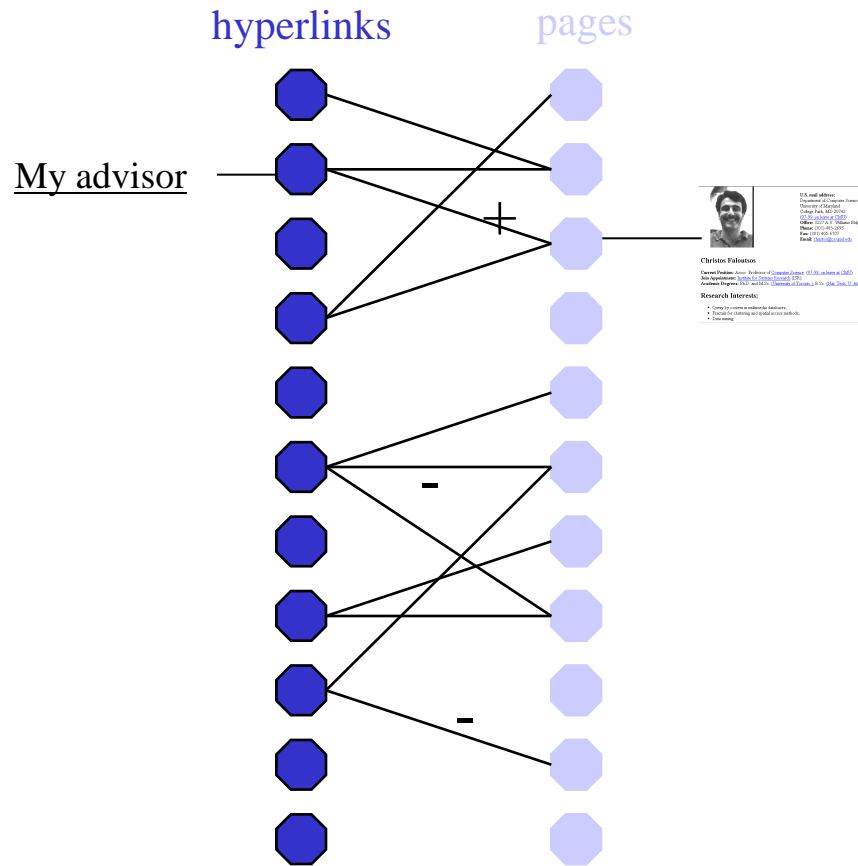
# CoTraining: Experimental Results

- 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% (when both agree)

Typical run:

# Co-Training Rote Learner

hyperlinks          pages

My advisor

# Classifying Jobs for FlipDog

# 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 h2 performs better on training data than simpler hypothesis h1, but worse on test data

- Unlabeled data can help detect overfitting, by comparing predictions of h1 and h2 over the unlabeled examples
  - The rate at which h1 and h2 disagree on U should be the same as the rate on L, unless overfitting is occuring

# Defining a distance metric

- Definition of distance metric
  - Non-negative *d(f,g)≥0;*
  - symmetric *d(f,g)=d(g,f);*
  - triangle inequality *d(f,g) · 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:

$$d(h_1, h_2) \equiv \sqrt{\int (h_1(x) - h_2(x))^2 p(x)dx}$$
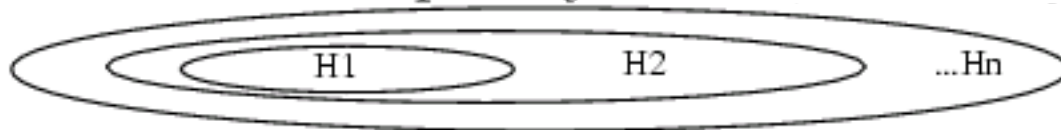
# Using the distance metric

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

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

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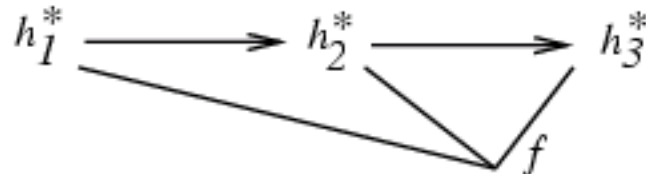
$$\hat{d}(h_1, h_2) = \frac{1}{|U|} \sum_{x \in U} \delta(h_1(x) \neq h_2(x))$$

Organize $H$ into complexity classes



Let $h_i^*$ be hypothesis with lowest $\hat{d}(h, f)$ in $H_i$
Prefer $h_1^*$, $h_2^*$, or $h_3^*$?

# 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!

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

$\rightarrow$ Heuristic:

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

Generated y values contain zero mean Gaussian noise $\varepsilon$

$Y=f(x)+\varepsilon$



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.

# Experimental Evaluation of TRI
## [Schuurmans & Southey, MLJ 2002]

• Use it to select degree of polynomial for regression

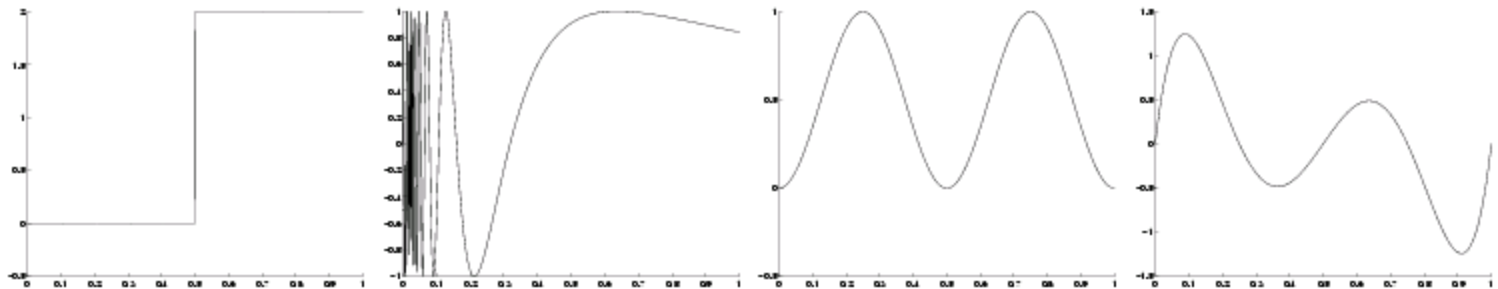• Compare to alternatives such as cross validation, structural risk minimization, …



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.

**Approximation ratio:**

$$\frac{\text{true error of selected hypothesis}}{\text{true error of best hypothesis considered}}$$

**Results using 200 unlabeled, t labeled**

Cross validation (Ten-fold)

Structural risk minimization

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.

# Summary

Several ways to use unlabeled data in supervised learning

1. Use to reweight labeled examples

2. Use to help EM learn class-specific generative models

3. If problem has redundantly sufficient features, use CoTraining

4. Use to detect/preempt overfitting

Ongoing research area

# Further Reading

- <u>EM approach</u>: K.Nigam, et al., 2000. "Text Classification from Labeled and Unlabeled Documents using EM", *Machine Learning,* 39, pp.103—134.

- <u>CoTraining</u>: 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*

- <u>Model selection</u>: D. Schuurmans and F. Southey, 2002. "Metric-Based methods for Adaptive Model Selection and Regularizaiton," Machine Learning, 48, 51—84.

# Acknowledgment

Some of these slides are based in on slides from Tom Mitchell.