Learning from Labeled and Unlabeled Data part 2: coupled training

Machine Learning 10-601

April 1, 2009

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

Last Time...

- 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. Use unlabeled data to detect/preempt overfitting
 - Based on distance metric, triangle inequality
 - If we assume priors over H that correctly order hypotheses

When can Unlabeled Data help supervised learning?

Problem setting (the PAC learning 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:

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

Wish to find hypothesis with lowest true error:

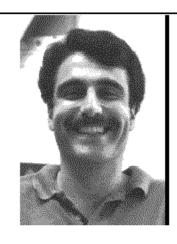
$$\widehat{f} \leftarrow \arg\min_{h \in H} \Pr_{x \in P(X)} [h(x) \neq f(x)]$$

Idea 4: CoTraining, Coupled Training

- 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

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

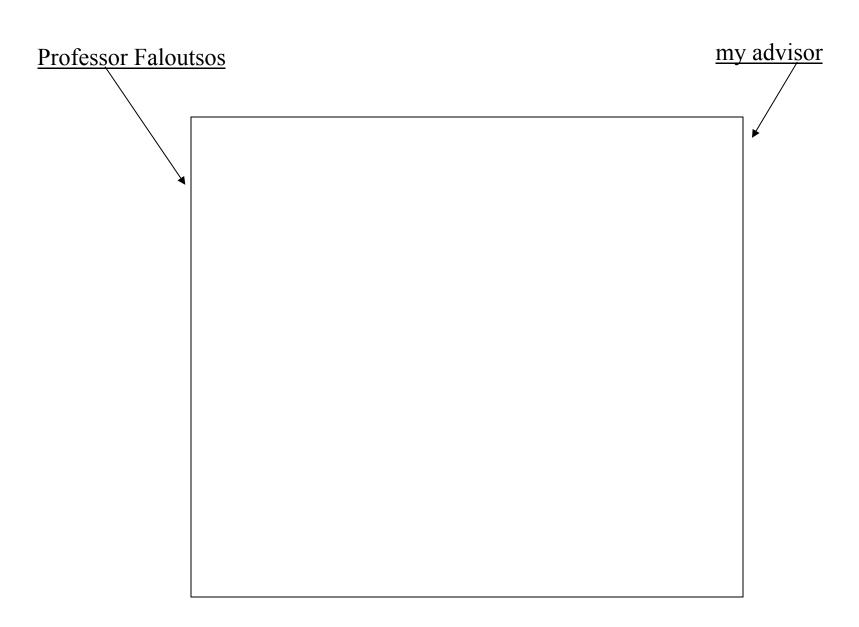
Current Position: Assoc. Professor of Computer Science. (97-98: on leave at CMU)

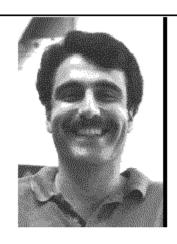
Join Appointment: Institute for Systems Research (ISR).

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;





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).

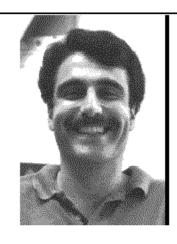
Academic Degrees: Ph.D. and M.Sc. (University of Toronto.); B.Sc. (Nat. Tech. U. Atherical Control of the Contr

Research Interests:

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

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 g1 (hyperlink classifier) using L

Train g2 (page classifier) using L

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

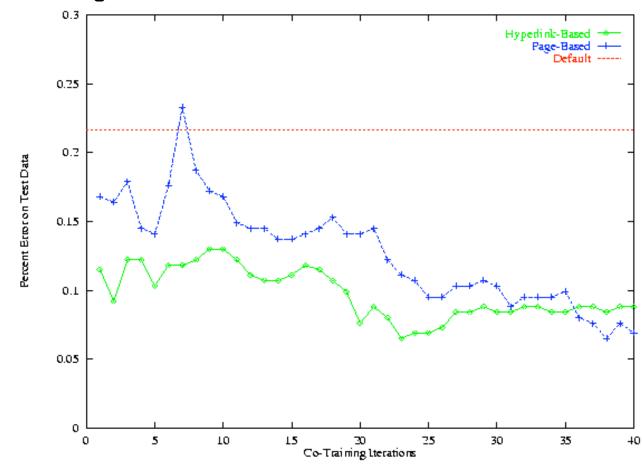
Allow g2 to label p positive, n negative examps from U

Add 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 → Y, given L and U drawn from P(X)
- features describing X can be partitioned (X = X1 x X2) such that f can be computed from either X1 or X2 $(\exists g_1, g_2)(\forall x \in X)$ $g_1(x_1) = f(x) = g_2(x_2)$

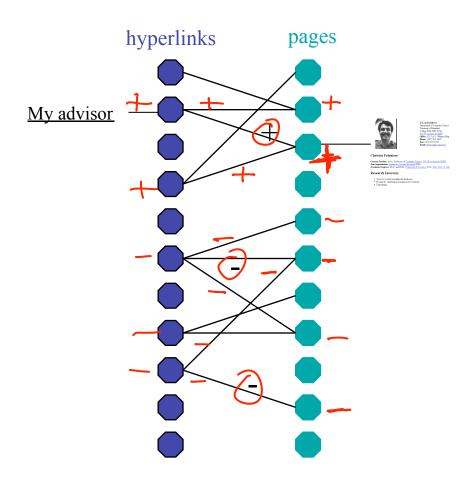
One result [Blum&Mitchell 1998]:

- If
 - X1 and X2 are conditionally independent given Y
- accuracy > 0.5

Classifier with

- f is PAC learnable from noisy labeled data
- Then
 - f is PAC learnable from weak initial classifier plus polynomial number of *unlabeled* examples

Example: Co-Training Rote Learners f1:hyperlink → Y, f2: page → Y



Example: Co-Training Rote Learner

M labeled examps. E[err] hyperlinks My advisor -P(X)

(X)

(X)

Expected Rote CoTraining error given *m* examples

CoTraining setting:

learn $f: X \to Y$

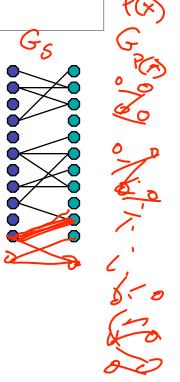
where $X = X_1 \times X_2$

where x drawn from unknown distribution

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

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

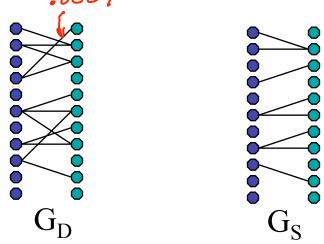
Where g_j is the *j*th connected component of graph of L+U, m is number of labeled examples



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 X1 and X2 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 $1 \le i \le k$ then (a) f is a permutation and (b) for all $1 \le i \le k$,

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

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 $\widehat{P}(h_1 \neq i \mid h_2 = i, h_1 \neq \bot)$ is a good estimate of the error rate $P(h_1 \neq i \mid f(y) = i, h_1 \neq \bot)$.

$$\gamma_{i}(h_{1}, h_{2}, \delta) = \widehat{P}(h_{1} = i \mid h_{2} = i, h_{1} \neq \bot) - \widehat{P}(h_{1} \neq i \mid h_{2} = i, h_{1} \neq \bot) - 2\epsilon_{i}(h_{1}, h_{2}, \delta)$$

$$\epsilon_{i}(h_{1}, h_{2}, \delta) = \sqrt{\frac{(\ln 2)(|h_{1}| + |h_{2}|) + \ln \frac{2k}{\delta}}{2|S(h_{2} = i, h_{1} \neq \bot)|}}$$

PAC Generalization Bounds on CoTraining

[Dasgupta et al., NIPS 2001]

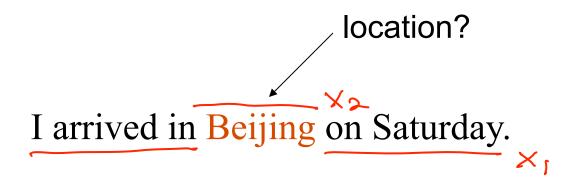
This theorem assumes X1 and X2 are conditionally independent given Y

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 $1 \le i \le k$ then (a) f is a permutation and (b) for all $1 \le i \le k$,

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

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 $\widehat{P}(h_1 \neq i \mid h_2 = i, h_1 \neq \bot)$ is a good estimate of the error rate $P(h_1 \neq i \mid f(y) = i, h_1 \neq \bot)$.

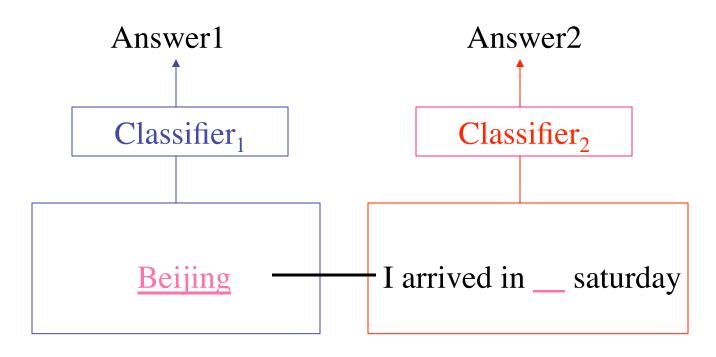
Example 2: Learning to extract named entities



If: "I arrived in <X> on Saturday."

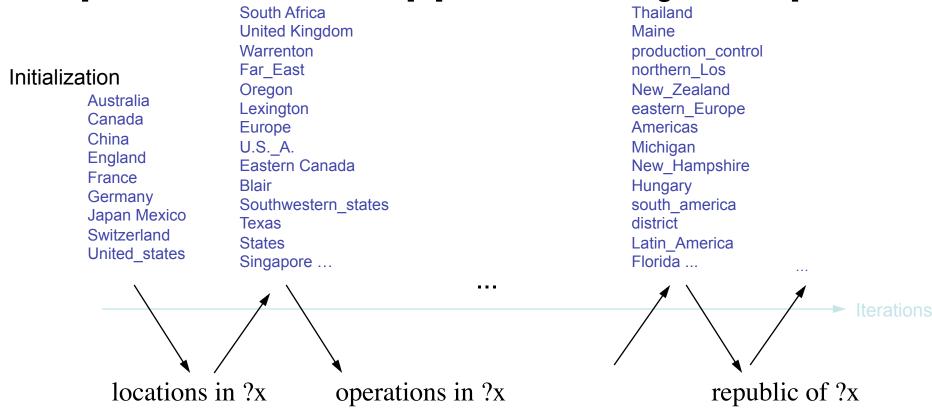
Then: Location(X)

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], ...



The Problem with Semi-Supervised Bootstrap Learning

it's underconstrained!!

Paris
Pittsburgh
Seattle
Cupertino

San Francisco Austin denial

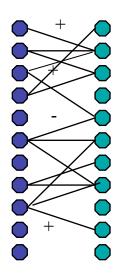


mayor of arg1 live in arg1



arg1 is home of traits such as arg1

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?

What Objective Function?

$$E = E1 + E2$$

$$E1 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_1(x_1))^2$$

$$E2 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_2(x_2))^2$$

$$E2 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_2(x_2))^2$$

Error on labeled examples

What Objective Function?

$$E = E1 + E2 + c_3 E3$$

$$E1 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_1(x_1))^2$$

$$E2 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_2(x_2))^2$$

$$E3 = \sum_{x \in U} (\hat{g}_1(x_1) - \hat{g}_2(x_2))^2$$

$$E2 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_2(x_2))^2$$

$$E3 = \sum_{\mathbf{x} \in U} (\hat{g}_1(x_1) - \hat{g}_2(x_2))^2$$

Error on labeled examples

Disagreement over unlabeled

What Objective Function?

$$E = E1 + E2 + c_3 E3 + c_4 E4$$

$$E1 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_1(x_1))^2$$

$$E2 = \sum_{\langle x, y \rangle \in L} (y - \hat{g}_2(x_2))^2$$

$$E3 = \sum_{x \in U} (\hat{g}_1(x_1) - \hat{g}_2(x_2))^2$$
 Misfit to estimated class priors

Error on labeled examples

Disagreement over unlabeled

$$E4 = \left(\left(\frac{1}{|L|} \sum_{\langle x, y \rangle \in L} y \right) - \left(\frac{1}{|L| + |U|} \sum_{x \in L \cup U} \frac{\hat{g}_1(x_1) + \hat{g}_2(x_2)}{2} \right) \right)^2$$

What Function Approximators?

What Function Approximators?

$$\hat{g}_1(x) = \frac{1}{1 + e^{\sum_{j=1}^{N} w_{j,1} x_j}}$$

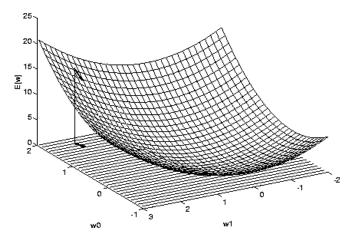
$$\hat{g}_1(x) = \frac{1}{1 + e^{\sum_{j=1}^{w_{j,1}x_j}}} \qquad \hat{g}_2(x) = \frac{1}{1 + e^{\sum_{j=1}^{w_{j,2}x_j}}}$$

- Same fn form as Naïve Bayes, Max Entropy
- Use gradient descent to simultaneously learn g1 and g2, directly minimizing E = E1 + E2 + E3 + E4
- No word independence assumption, use both labeled and unlabeled data

Gradient CoTraining

$$\hat{g}_1(x) = \frac{1}{1 + e^{\sum_{j=1}^{w_{j,1}x_j}}}$$

$$\hat{g}_1(x) = \frac{1}{1 + e^{\sum_{j=1}^{w_{j,1}x_j}}} \qquad \hat{g}_2(x) = \frac{1}{1 + e^{\sum_{j=1}^{w_{j,2}x_j}}}$$



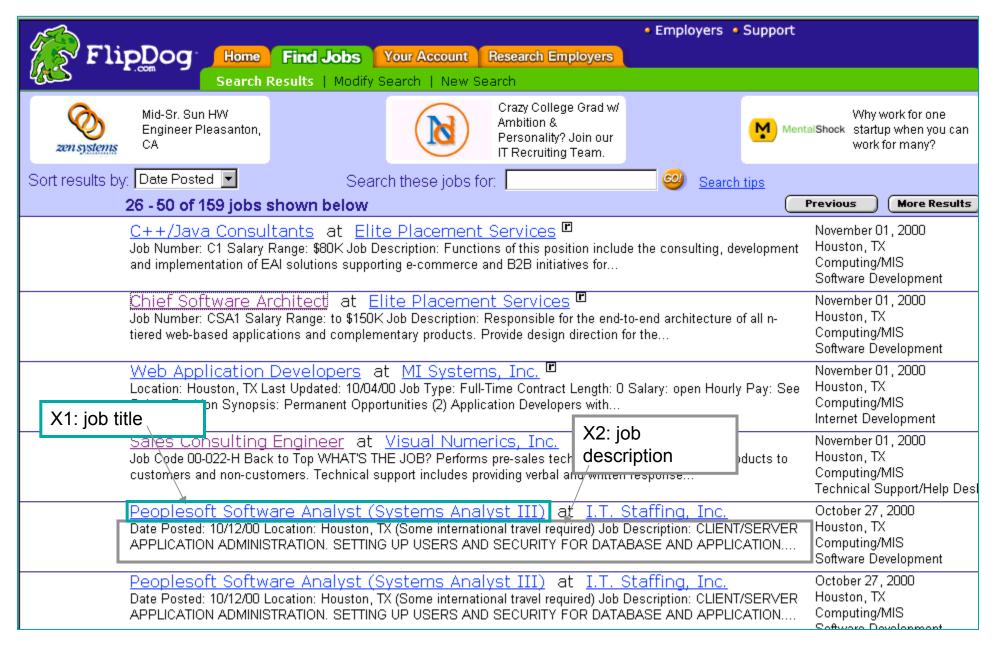
Gradient

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \cdots \frac{\partial E}{\partial w_n} \right]$$

Training rule:

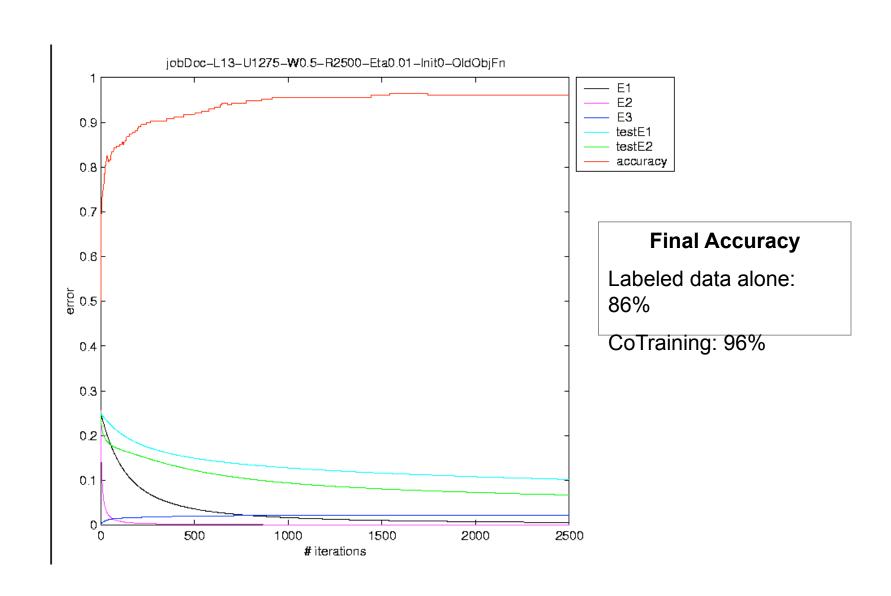
$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

Classifying Jobs for FlipDog



Gradient CoTraining

Classifying FlipDog job descriptions: SysAdmin vs. WebProgrammer



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 X1,X2 conditionally independent given Y, Then
 - PAC learnable from weak initial classifier plus unlabeled data
 - disagreement between g1(x1) and g2(x2) 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]

What you should know

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

Further Reading

- <u>Semi-Supervised Learning</u>, O. Chapelle, B. Sholkopf, and A. Zien (eds.), MIT Press, 2006. (excellent 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.
- <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
- Model selection: D. Schuurmans and F. Southey, 2002. "Metric-Based methods for Adaptive Model Selection and Regularization," Machine Learning, 48, 51—84.

Toward Never-Ending Learning of Semantic Knowledge

Justin Betteridge, Andrew Carlson, Estevam R. Hruschka Jr., Tom M. Mitchell

(with help from Sue Ann Hong, Sophie Wang, Richard Wang)

Carnegie Mellon University

March 2009

Our Goal: Never-Ending Language Learning

Goal:

- run 24x7, forever
- each day:
 - 1. extract more facts from the web to populate and extend initial ontology
 - 2. learn to read better than yesterday

Our <u>Goal</u>: Never-Ending Language Learning

Goal:

- run 24x7, forever
- each day:
 - 1. extract more facts from the web to populate initial ontology
 - 2. learn to read better than yesterday

Today...

Given:

- initial ontology defining dozens of classes and relations
- 10-20 seed examples of each

Task:

- learn to extract / extract to learn
- running over 200M web pages, for a few days

Browse the KB

- ~ 18,000+ entities, ~ 30,000 extracted beliefs
- learned from 10-20 seed examples, 200M unlabeled web pages
- ~ 2 days computation on M45 cluster (thanks Yahoo!)

on the web:

<u>initial</u> populated

The Problem with Semi-Supervised Bootstrap Learning

it's underconstrained!!

Paris
Pittsburgh
Seattle
Cupertino

San Francisco Austin denial



mayor of arg1 live in arg1



arg1 is home of traits such as arg1

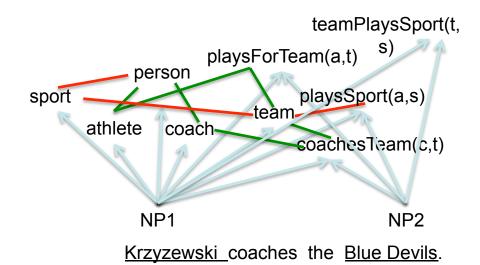
The Key to Accurate Semi-Supervised Learning

coach(NP)

NP context

Krzyzewski coaches the Blue Devils.

hard (underconstrained) semi-supervised learning problem



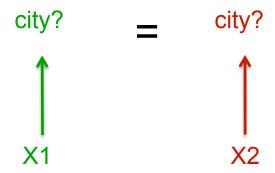
much easier (more constrained) semi-supervised learning problem

Wish to learn $f: X \rightarrow Y$

e.g., city: NounPhraseInSentence \rightarrow {0,1}

Constraint type 1 (co-training):

if X can be split into redundantly sufficient X1, X2 then learn both f1: X1 \rightarrow Y, and f2: X2 \rightarrow Y

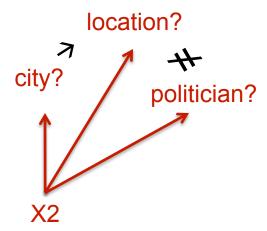


X: Luke is mayor of Pittsburgh.

Wish to learn f: $X \rightarrow Y$

e.g., city: NounPhraseInSentence \rightarrow {0,1}

Constraint type 2: couple training of multiple classes
Ontology provides coupling constraints

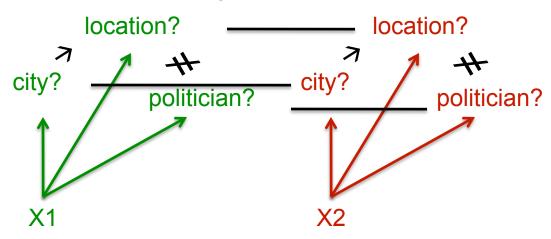


Luke is mayor of Pittsburgh.

Wish to learn f: $X \rightarrow Y$

e.g., city: NounPhraseInSentence \rightarrow {0,1}

Constraint type 2: couple training of multiple classes Ontology provides coupling constraints



Luke is mayor of Pittsburgh.

Coupling unsupervised training of multiple functions

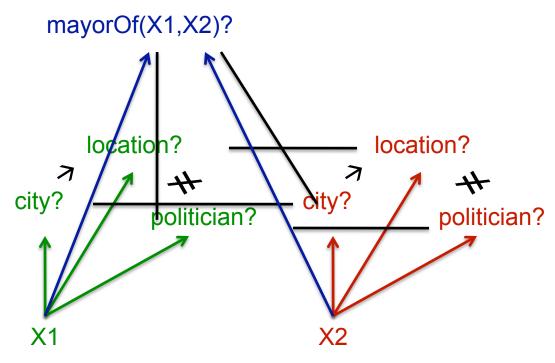
Cotraining: learn f: X > Y, by coupling training of

- f1: X1 → Y
- f2: X2 → Y

Coupled functions: learn f1: $X \rightarrow Y1$, f2: $X \rightarrow Y2$ where coupling is based on a constraint C(Y1,Y2)

- e.g., mutallyExclusive(Person, City)
- e.g., subset(Athlete, Person)

Constraint type 3 (couple training of <u>classes</u> and <u>relations</u>)



Luke is mayor of Pittsburgh.

Coupled Bootstrap Learner algorithm

Algorithm 1: CBL Algorithm

Input: An ontology \mathcal{O} , and text corpus C

Output: Trusted instances/patterns for each

predicate

SHARE initial instances/patterns among predicates;

for $i = 1, 2, \ldots, \infty$ do

foreach $predicate \ p \in \mathcal{O}$ **do**

EXTRACT new candidate

instances/patterns;

FILTER candidates;

TRAIN instance/pattern classifiers;

Assess candidates using trained

classifiers;

PROMOTE highest-confidence candidates;

end

SHARE promoted items among predicates;

end

In the **ontology**: categories, relations, seed instances and patterns, type information, mutual **SXELVATION** Extraction (M45) relations, and type checking Arg1 HQ in Arg2 → (CBC || Filtering (M45) e || San Jose), ... Ase simento Not enough विभिन्न कि strength of and patterns. Use type-checking. Score patterns with estimate of

precision

learned extraction patterns: Company

```
retailers like such clients as an operating business_of__ being_acquired_by__
   firms_such_as__ a_flight_attendant_for__ chains_such_as__ industry_leaders_such_as__
   advertisers like social networking sites such as a senior manager at
   competitors like stores like is an ebay company discounters like
   a_distribution_deal_with__ popular_sites_like__ a_company_such_as__ vendors_such_as__
   rivals_such_as__ competitors_such_as__ has_been_quoted_in_the__ providers_such_as__
   company_research_for__ providers_like__ giants_such_as__ a_social_network_like__
   popular_websites_like__ multinationals_like__ social_networks_such_as__
   the_former_ceo_of__ a_software_engineer_at__ a_store_like__ video_sites_like__
   a_social_networking_site_like__ giants_like__ a_company_like__ premieres_on__
   corporations such as corporations like professional profile on outlets like
   the executives at stores such as is the only carrier a big company like
   social media sites such as has an article today manufacturers such as
   companies like social media sites like companies including firms like
   networking_websites_such_as__ networks like carriers like
   social networking websites like an executive at insured via
   provides dialup access a patent infringement lawsuit against
   social networking sites like social network sites like carriers such as
   are_shipped_via__ social_sites_like__ a_licensing_deal_with__ portals_like__
   vendors like the accounting firm of industry leaders like retailers such as
   chains_like__ prior_fiscal_years_for__ such_firms_as__ provided_free_by__
   manufacturers like airlines_like__ airlines_such_as__
```

learned extraction patterns: playsSport(arg1,arg2)

```
arg1 was playing arg2 arg2 megastar arg1 arg2 icons arg1 arg2 player named arg1
   arg2 prodigy arg1 arg1 is the tiger woods of arg2 arg2 career of arg1
   arg2 greats as arg1 arg1 plays arg2 arg2 player is arg1 arg2 legends arg1
   arg1 announced his retirement from arg2 arg2 operations chief arg1
   arg2 player like arg1 arg2 and golfing personalities including arg1 arg2 players like arg1
   arg2 greats like arg1 arg2 players are steffi graf and arg1 arg2 great arg1
   arg2 champ arg1 arg2 greats such as arg1 arg2 professionals such as arg1
   arg2 course designed by arg1 arg2 hit by arg1 arg2 course architects including arg1
   arg2 greats arg1 arg2 icon arg1 arg2 stars like arg1 arg2 pros like arg1
   arg1 retires from arg2 arg2 phenom arg1 arg2 lesson from arg1
   arg2 architects robert trent jones and arg1 arg2 sensation arg1 arg2 architects like arg1
   arg2 pros arg1 arg2 stars venus and arg1 arg2 legends arnold palmer and arg1
   arg2 hall of famer arg1 arg2 racket in arg1 arg2 superstar arg1 arg2 legend arg1
   arg2_legends_such_as_arg1 arg2_players_is_arg1 arg2_pro_arg1 arg2_player_was_arg1
   arg2 god arg1 arg2 idol arg1 arg1 was born to play arg2 arg2 star arg1
   arg2_hero_arg1_arg2_course_architect_arg1_arg2_players_are_arg1
   arg1 retired from professional arg2 arg2 legends as arg1 arg2 autographed by arg1
   arg2 related quotations spoken by arg1 arg2 courses were designed by arg1
   arg2 player since arg1 arg2 match between arg1 arg2 course was designed by arg1
   arg1 has retired from arg2 arg2 player arg1 arg1 can hit a arg2
   arg2_legends_including_arg1_arg2_player_than_arg1_arg2_legends_like_arg1
   arg2 courses designed by legends arg1 arg2 player of all time is arg1
   arg2 fan knows arg1 arg1 learned to play arg2 arg1 is the best player in arg2
   arg2 signed by arg1 arg2 champion arg1
```

Experimental Evaluation

- 31 predicates
 - 15 relations, 16 categories
- Domains:
 - Companies
 - Sports
- Run for 15 iterations:
 - Full system
 - No Sharing of promoted items
 - No Relation/Category coupling
- Evaluated a sample of promoted items

	5 iterations			10 iterations			15 iterations		
Predicate	Full	NS	NCR	Full	NS	NCR	Full	NS	NCR
Actor	93	100	100	93	97	100	100	97	100
Athlete	100	100	100	100	93	100	100	73	100
Board Game	93	76	93	89	27	93	89	30	93
City	100	100	100	100	97	100	100	100	100
Coach	100	63	73	97	53	43	97	47	47
Company	100	100	100	97	90	97	100	90	100
Country	60	40	60	30	43	27	40	23	40
Economic Sector	77	63	73	57	67	67	50	63	40
Hobby	67	63	67	40	40	57	20	23	30
Person	97	97	90	97	93	97	93	97	93
Politician	93	93	97	73	53	90	90	53	87
Product	97	87	90	90	87	100	97	90	77
Product Type	93	93	90	70	73	97	77	80	67
Scientist	100	90	97	97	63	97	93	60	100
Sport	100	90	100	93	67	83	97	27	90
Sports Team	100	97	100	97	70	100	90	50	100
Category Average	92	84	89	82	70	84	83	63	79
Acquired(Company, Company)	77	77	80	67	80	47	70	63	47
CeoOf(Person, Company)	97	87	100	90	87	97	90	80	83
CoachesTeam(Coach, Sports Team)	100	100	100	100	100	97	100	100	90
CompetesIn(Company, Econ. Sector)	97	97	80	100	93	67	97	63	60
CompetesWith(Company, Company)	93	80	60	77	70	37	70	60	43
HasOfficesIn(Company, City)	97	93	40	93	90	27	93	57	30
HasOperationsIn(Company, Country)	100	95	50	100	97	40	90	83	13
HeadquarteredIn(Company, City)	77	90	20	70	77	27	70	60	7
LocatedIn(City, Country)	90	67	57	63	50	43	73	50	30
PlaysFor(Athlete, Sports Team)	100	100	0	100	97	7	100	43	0
PlaysSport(Athlete, Sport)	100	100	27	93	80	10	100	40	30
TeamPlaysSport(Sports Team, Sport)	100	100	77	100	97	80	93	83	67
Produces(Company, Product)	91	83	90	83	93	67	93	80	57
HasType(Product, Product Type)	73	63	17	33	67	33	40	57	27
Relation Average	92	88	57	84	84	48	84	66	42
All	92	86	74	83	76	68	84	64	62

Table 1: Precision (%) for each predicate. Results are presented after 5, 10, and 15 iterations, for the Full, No Sharing (NS), and No Category/Relation Coupling (NCR) configurations of CBL.

Extending Freebase

	Freebase	CBL	Est.	Est. New
Category	Matches	Instances	Prec.	Instances
Actor	465	522	100	57
Athlete	54	117	100	63
Board Game	6	18	89	10
City	1665	1799	100	134
Company	995	1937	100	942
Econ. Sector	137	1541	50	634
Politician	74	962	90	792
Product	0	1259	97	1221
Sports Team	139	414	90	234
Sport	134	613	97	461

Table 3: Estimated numbers of 'New Instances', which are correct instances promoted by CBL in the Full 15 iteration run which do not have a match in Freebase, and the values used in calculating them (number of Freebase/CBL matches, number of CBL instances, and the estimated precision of CBL for the predicate).

Summary

For never-ending language learning, the key is achieving accurate semi-supervised training

- → Constrain learning by coupling the training of many types of knowledge (functions)
 - sample complexity <u>decreases</u> as ontology size increases
- → Want an architecture in which current learning makes future learning even more accurate
 - -- learn symbolic rules which become new probabilistic constraints
- → Want architecture where self-consistency ≈ correctness