# Reading the Web: Advanced Statistical Language Processing

www.cs.cmu.edu/~tom/rtw09/

Machine Learning 10-709

September 10, 2009

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

### The Plan

- What will you get out of this class?
  - knowledge of state of the art in semi-supervised learning, statistical NLP, never-ending-learning
  - an opportunity to advance it
  - a research infrastructure you might use in the future
- What will you be expected to do?
  - read, discuss, critique research papers
  - design and perform a research project in this area
- What will we build on?
  - the RTW project data and knowledge base

### The Goals\*

- Build the first cumulative <u>never-ending learner</u>
- Advance state of Natural Language Understanding
- Build and publish the world's largest structured knowledge base

<sup>\*</sup> choose research problems wisely -2/3 of success in research is (re)choosing the problem

# The Problem Specification

#### Inputs:

- initial ontology
- handful of examples of each predicate in ontology
- the web
- occasional access to human trainer

#### The task:

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

# What We Have Today

#### 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 per predicate, 200M unlabeled web pages
- ~ 2 days computation on M45 cluster

Initial ontology: <u>Initial ontology</u>

After a few days of self-supervised learning:

http://rtw.ml.cmu.edu/sslnlp09/index.html http://rtw.ml.cmu.edu/wsdm10 online/

# Semi-Supervised Bootstrap Learning

it's underconstrained!!

Paris
Pittsburgh
Seattle
Cupertino

San Francisco Austin denial

anxiety selfishness Berlin





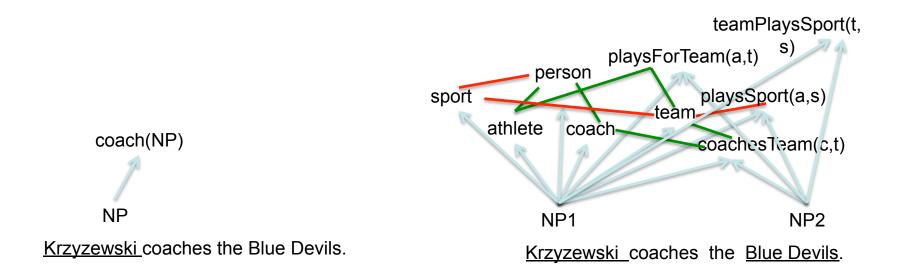




mayor of arg1 live in arg1

arg1 is home of traits such as arg1

### The Key to Accurate Semi-Supervised Learning



hard (underconstrained) semi-supervised learning problem **much easier** (more constrained) semi-supervised learning problem

The Key: Couple the training of many functions to make unlabeled data more informative

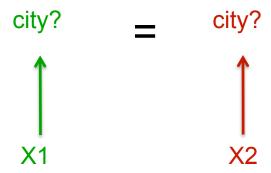
# Coupled training type 1

Wish to learn  $f: X \rightarrow Y$ 

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

Coupling type 1 (co-training): Learn 2 functions with different input features  $f1: X1 \rightarrow Y$ , and  $f2: X2 \rightarrow Y$ 

Coupling: force their outputs to agree over unlabeled examples



X: Luke is mayor of Pittsburgh.

# Coupled training type 2

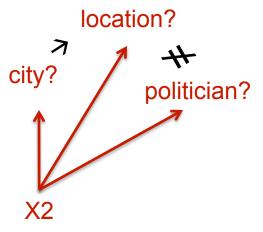
Wish to learn f:  $X \rightarrow Y1$ , f:  $X \rightarrow Y2$ , where g(y1,y2)

city: NounPhraseInSentence  $\rightarrow$  {0,1}

politician: NounPhraseInSentence → {0,1}

Constraint type 2: force outputs to satisfy g(y1,y2)

Ontology provides coupling constraints



Luke is mayor of Pittsburgh.

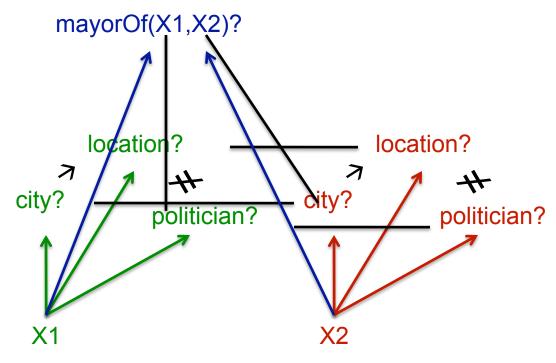
# Coupled training type 3

### Constraint type 3 (argument type consistency)

mayorOf: NP1 x NP2  $\rightarrow$  {0,1}

city: NP1  $\rightarrow$  {0,1}

politician: NP2  $\rightarrow$  {0,1}



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

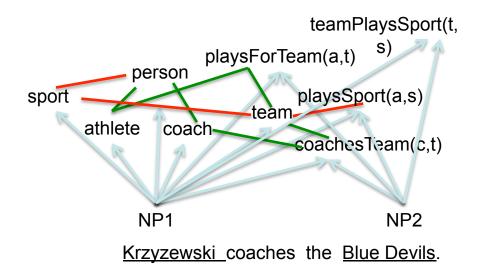
# learned extraction patterns: Company

```
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
```

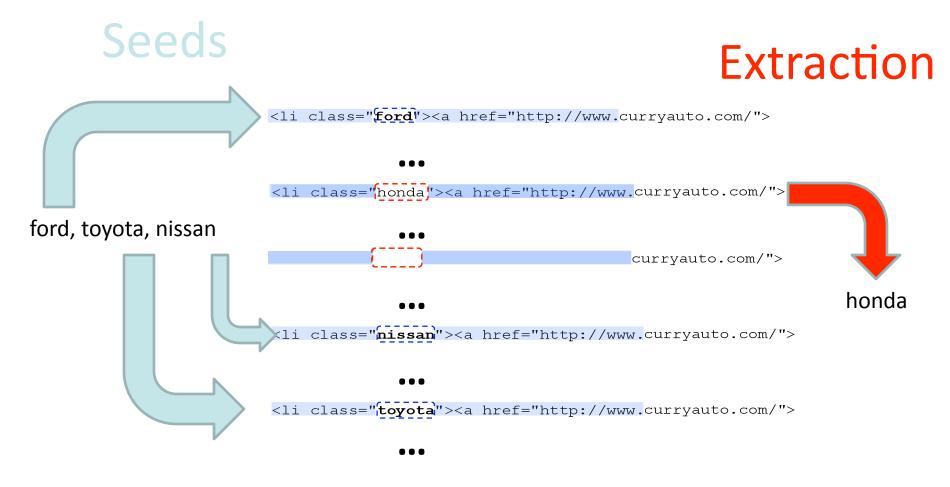
If the key to accurate self-supervised learning is coupling the training of many functions, then how can we create even more coupling?



1. introduce additional coupling by adding a learner that uses HTML features instead of free text features

# SEAL Set Expander for Any Language

\*Richard C. Wang and William W. Cohen: Language-Independent Set Expansion of Named Entities using the Web. In *Proceedings of IEEE International Conference on Data Mining* (ICDM 2007), Omaha, NE, USA. 2007.



### SEAL

For each class being learned,

On each iteration

Retrain CBL from current KB, allow it to add to KB

Retrain SEAL from current KB, allow it to add to KB

#### Typical learned SEAL extractors:

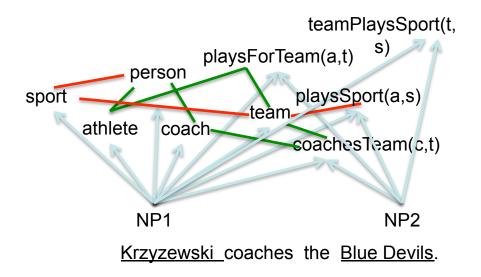
```
URL:
          http://www.shopcarparts.com/
Wrapper:
          .html" CLASS="shopcp">[...] Parts</A> <br>
          acura, audi, bmw, buick, cadillac, chevrolet, chevy, chrysler, daewoo, daihatsu, dodge, ea
Content:
          http://www.allautoreviews.com/
   URL:
Wrapper:
          </a><br> <a href="auto_reviews/[...]/
Content:
          acura, audi, bmw, buick, cadillac, chevrolet, chrysler, dodge, ford, gmc, honda, hyundai, i
          http://www.hertrichs.com/
   URL:
          class="franchise [...]"> <h4><a href="#">
Wrapper:
          buick, chevrolet, chrysler, dodge, ford, gmc, isuzu, jeep, lincoln, mazda, mercury, nissan,
Content:
```

	Precision (%)				1	Promo				
Predicate	$\operatorname{CPL}$	$\operatorname{UPL}$	CSEAL	SEAL	MBL	CPL	$\operatorname{UPL}$	CSEAL	SEAL	MBL
AcademicField	70	83	90	97	100	46	903	203	1000	181
Actor	100	33	100	97	100	199	1000	1000	1000	380
Animal	80	50	90	70	97	741	1000	144	974	307
Athlete	87	17	100	87	100	132	930	276	1000	555
AwardTrophyTournament	57	7	53	7	77	86	902	146	1000	79
BoardGame	80	13	70	77	90	10	907	126	1000	31
BodyPart	77	17	97	63	93	176	922	80	1000	61
Building	33	50	30	0	93	597	1000	57	1000	14
Celebrity	100	90	100	100	97	347	1000	72	747	514
CEO	33	30	100	77	100	3	902	322	1000	30
City	97	100	97	87	97	1000	1000	368	1000	603
Clothing	97	20	43	27	97	83	973	167	1000	102
Coach	93	63	100	83	100	188	838	619	1000	242
Company	97	83	100	100	97	1000	1000	245	1000	784
Conference	93	53	97	90	100	95	990	437	928	92
Country	57	33	97	37	93	1000	1000	130	1000	207
EconomicSector	60	23	100	10	77	1000	1000	34	1000	138
Emotion	77	53	87	60	83	483	992	183	1000	211
Food	90	70	97	80	100	811	1000	89	1000	272
Furniture	100	0	57	57	90	55	963	215	1000	95
Hobby	77	33	77	50	90	357	936	77	1000	127
KitchenItem	73	3	88	13	100	11	900	8	960	2
Mammal	83	50	93	50	90	224	1000	154	1000	169
Movie	97	57	97	100	100	718	1000	566	1000	183
NewspaperCompany	90	60	60	97	100	179	1000	1000	1000	241
Politician	80	60	97	37	100	178	990	30	1000	101
Product	90	83	-	77	70	1000	1000	0	999	127
	73	63	27	63	50	712	1000	31	1000	159
ProductType Profession	73	53		57	93	916	973	0	1000	171
Profession ProfessionalOrganization	93	63	100	77	93 87	104	943	58	1000	163
	95 95	3	90	27	100	104	912	149	1000	54
Reptile Room	64	0	33	7	100	25	913	12		3
Scientist	97	30	100	17	100	83	913	928	643 1000	130
	77	7	7	7	85	43	985	28	733	26
Shape	77	13	63	83	73	283	1000	225	1000	284
Sport SportsEquipment	20	10	57	23	23	203 58	902	52	1000	174
	100	7	80	23 27	23 86		902	10	1000	174
SportsLeague	90	30	80 87		87	11				
SportsTeam		57		87	90	301	903	864	944	506
StateOrProvince	93		53 83	63	77	102	767	944	1000	343
StateOrProvince	77	63		93 90		202	1000	114	1000	161 59
Tool	40	13	93 52		97 97	561	1000	713	1000	
Trait	53 93	40 97	100	47		234	1000	21	1000	44
University				90	93	1000	1000	961	1000	516
Vehicle	67	30	50	13	77	460	1000	50	1000	98
Average	78	41	78	59	90	360	960	271	976	199
Weighted average	79	42	86	59	91					

Table 2: Precision (%) and counts of promoted instances for each category using CPL, UPL, CSEAL, SEAL MBL.

# If the key to accurate self-supervised learning is <u>coupling</u> the training of many functions,

then how can we create even more coupling?



2. allow learner to discover new coupling constraints (by mining its extracted beliefs)

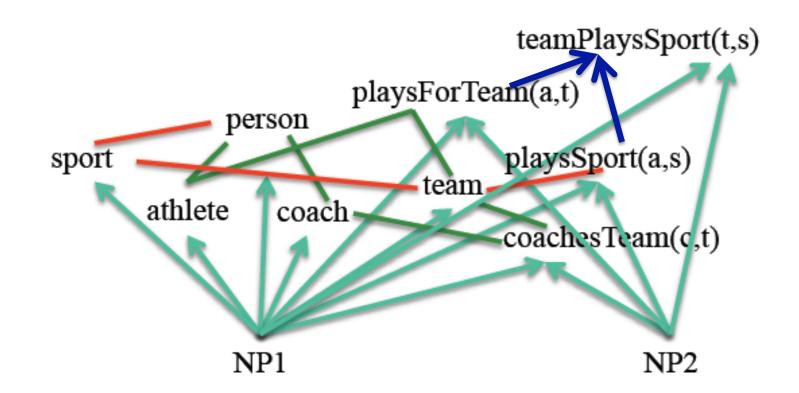
#### Learned Probabilistic Horn Clause Rules

- 40 learned rules for teamPlaySport, playSport,
- when applied, inferred over 800 new beliefs
  - e.g., teamPlaysSport(Caps,hockey),
  - playSport(JasonGiambi,baseball)

```
0.84 playsSport(?x,?y) ← playsFor(?x,?z), teamPlaysSport(?z,?y)
0.70 playsSport(?x,baseball) ← playsFor(?x,cubs)
...
0.81 teamPlaysSport(?x,?y) ← playsForTeam(?x,?z), playsSport(?z,?y)
0.70 teamPlaysSport(?x,basketball) ← playsAgainst(?x,pistons)
0.64 teamPlaysSport(?x,?y) ← playsAgainst(?x ?z), teamPlaysSport(?z,?y)
```

#### Learned Probabilistic Horn Clause Rules

0.81 teamPlaysSport(?x,?y) ← playsForTeam(?x,?z), playSport(?z,?y)



# Summary: what we have to work from

### Data/Knowledge:

- KB with 10<sup>4</sup> extracted beliefs, .85-.90 accurate
- statistics on co-occurrence frequency of
  - 10<sup>5</sup> NPs x 10<sup>5</sup> contexts

### Learning algorithms

- coupled semi-supervised learning of name entity extractors, relation extractors
- learning probabilistic first-order horn clauses

### Homework 1: due next thursday

- get co-occurrence data: 10<sup>5</sup> NP's x 10<sup>5</sup> Contexts
- look at <a href="http://rtw.ml.cmu.edu/wsdm10">http://rtw.ml.cmu.edu/wsdm10</a> online/
  - labeled data: "Instances Promoted by Meta-Bootstrap Learner"
- do something interesting, prepare a 2 slide, 3 min presentation

#### Example:

learn to classify <NP,Context>, or just NP's as city, person, emotion, ...

- supervised, semi-supervised, unsupervised, ...

# Projects Ideas 1

- Add a morphology-based entity extractor
  - Omalinski is probably a person
  - yet another redundant information source
- Ultra High-dimensional training
  - each noun phrase as bag of 10<sup>6</sup> contexts
  - each context as bag of 10<sup>5</sup> noun phrases
- Add first self-reflection capability
  - where is my performance weakest?
  - what should I do next?

# Projects Ideas 2

- Better rule learning algorithm
  - learning from positive data only?
  - accuracy estimates based on resampling?
- Prep phrase attachment
  - How can KB and background statistics be used?
- Co-reference resolution
  - IBM versus Int.Bus.Mach versus it

# Projects Ideas 3

### Active learning

- what questions should I ask in today's 5 min session with a human?
- what new data should I download?

### Temporal scoping

- How can we determine when in time a fact holds?
- Consider earliest web page containing the fact?
- Date web pages?
- Read the temporal scope?