

Relation Extraction II



Bill MacCartney

CS224U

31 January 2013

Many thanks for lots of slides from many people,
including Dan Jurafsky, Jim Martin, Oren Etzioni,
Michele Banko, Rion Snow, Mike Mintz, and Steven Bills



Reminder!

- Lit Review paper due in 2 weeks!
- Start forming your project groups
- We strongly encourage working in groups!
- Bring your project ideas to office hours
- The ACL Anthology Searchbench may help find relevant literature: <http://aclasb.dfgi.de/>



Extracting structured knowledge

Each article can contain hundreds or thousands of items of knowledge

WIKIPEDIA The Free Encyclopedia

article discussion edit this page history Log in / create account

Lawrence Livermore National Laboratory

From Wikipedia, the free encyclopedia

The Lawrence Livermore National Laboratory (LLNL) in Livermore, California is a scientific research laboratory founded by the University of California in 1952. It is funded by the United States Department of Energy (DOE) and managed by Lawrence Livermore National Security, LLC (LNS), a partnership of the University of California, Bechtel Corporation, Babcock and Wilcox, the URS Corporation, and Battelle Memorial Institute. On October 1, 2007 LNS assumed management of LLNL from the University of California, which had exclusively managed and operated the Laboratory since its inception 55 years before.

Contents [hide]

- 1 Background
- 2 Origins
- 3 Weapons projects
- 4 Plutonium research
- 5 National Ignition Facility and photon science
- 6 Global security program
- 7 Other programs
- 8 Key accomplishments
- 9 Unique facilities
- 10 World-class computers
- 11 Sponsors
- 12 Directors
- 13 Organization
- 14 Footnotes
- 15 References
- 16 External links and sources

Background

LLNL is self-described as "a premier research and development institution for science and technology applied to national security".^[1] Its principal responsibility is ensuring the safety, security and reliability of the nation's nuclear weapons through the application of advanced science, engineering and technology. The Laboratory also applies its special expertise and multidisciplinary capabilities to preventing the proliferation and use of weapons of mass destruction, bolstering homeland security and solving other nationally important problems, including energy and environmental security, basic science and economic competitiveness.

LLNL is home to many unique facilities and a number of the most powerful computer systems in the world, according to the TOP500 list, including Blue Gene/L, the world's fastest computer from 2004 until Los Alamos National Laboratory's Roadrunner supercomputer surpassed it in 2008. The Lab is a leader in technical innovation: since 1978, LLNL has received a total of 118 prestigious R&D 100 Awards, including

Coordinates: 37.686024°N 121.709547°W

Lawrence Livermore National Laboratory

University of California Lawrence Livermore National Laboratory

Motto	"Science in the national interest"
Established	1952 by the University of California
Research Type	National security, nuclear science
Budget	US\$1.6 billion
Director	George H. Miller
Staff	6,800
Location	Livermore, California
Campus	3.2 km ² (800 acres)
Operating Agency	Lawrence Livermore National Security, LLC
Website	www.llnl.gov

Aerial view of Lawrence Livermore National Laboratory

"The **Lawrence Livermore National Laboratory (LLNL)** in **Livermore, California** is a **scientific research laboratory** founded by the **University of California** in **1952**."



LLNL EQ Lawrence Livermore National Laboratory
LLNL LOC-IN California
Livermore LOC-IN California
LLNL IS-A scientific research laboratory
LLNL FOUNDED-BY University of California
LLNL FOUNDED-IN 1952



Relation extraction: 5 easy methods

1. Hand-built patterns
2. Bootstrapping methods
3. Supervised methods
4. **Distant supervision**
5. Unsupervised methods

Distant supervision paradigm

Snow, Jurafsky, Ng. 2005. Learning syntactic patterns for automatic hypernym discovery. NIPS 17



Mintz, Bills, Snow, Jurafsky. 2009. Distant supervision for relation extraction without labeled data. ACL-2009.

- Hypothesis: If two entities belong to a certain relation, any sentence containing those two entities is likely to express that relation
- Key idea: use a *database* of relations to get lots of training examples
 - instead of hand-creating a few seed tuples (bootstrapping)
 - instead of using hand-labeled corpus (supervised)



Distant supervision approach

- For each pair of entities in a large database:
 - Grab sentences containing these entities from a corpus
 - Extract lots of noisy features from the sentences
 - Lexical features, syntactic features, named entity tags
 - Combine in a classifier

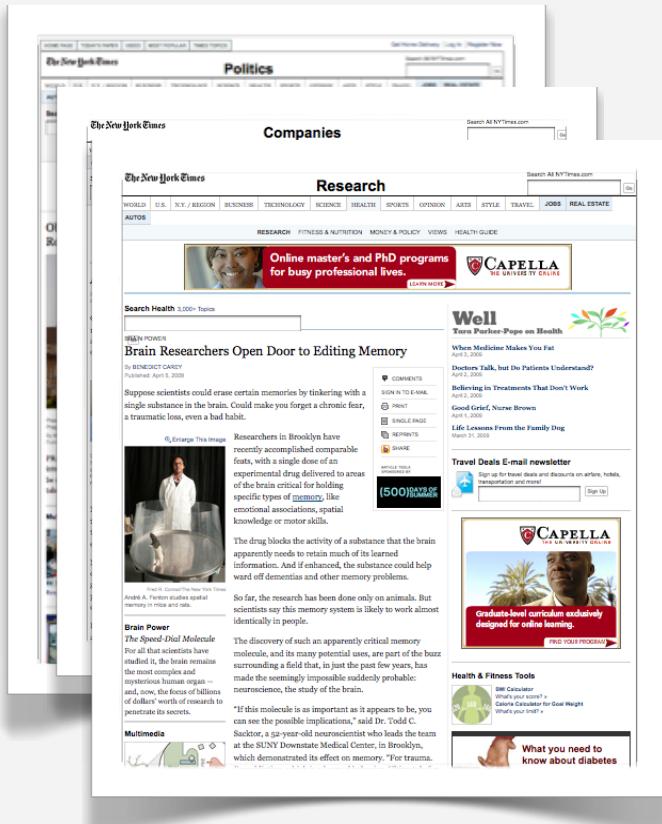


Benefits of distant supervision

- Has advantages of supervised approach
 - leverage rich, reliable hand-created knowledge
 - relations have canonical names
 - can use rich features (e.g. syntactic features)
- Has advantages of unsupervised approach
 - leverage unlimited amounts of text data
 - allows for very large number of weak features
 - not sensitive to training corpus: genre-independent

Hypernyms via distant supervision

We construct a noisy training set consisting of occurrences from our corpus that contain a hyponym-hypernym pair from WordNet.

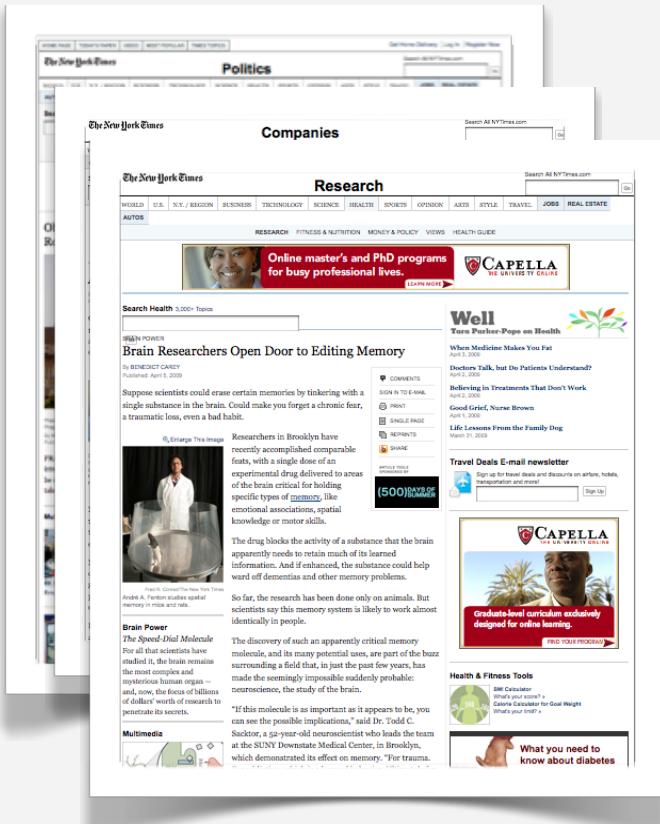


This yields high-signal examples like:

“...consider **authors** like **Shakespeare**...”
 “Some **authors** (including **Shakespeare**)...”
 “**Shakespeare** was the **author** of several...”
 “**Shakespeare**, **author** of *The Tempest*...”

Hypernyms via distant supervision

We construct a noisy training set consisting of occurrences from our corpus that contain a hyponym-hypernym pair from WordNet.



This yields high-signal examples like:

“...consider **authors** like **Shakespeare**...”
 “Some **authors** (including **Shakespeare**)...”
 “**Shakespeare** was the **author** of several...”
 “**Shakespeare**, **author** of *The Tempest*...”

But also noisy examples like:

“The **author** of **Shakespeare in Love**...”
 “...**authors** at the **Shakespeare Festival**...”

Learning hypernym patterns

1. Take corpus sentences

... doubly heavy hydrogen atom called deuterium ...

2. Collect noun pairs

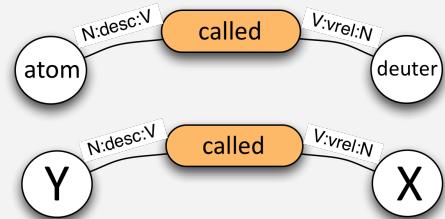
e.g. (atom, deuterium)

752,311 pairs from 6M sentences of newswire

3. Is pair an IS-A in WordNet?

14,387 yes; 737,924 no

4. Parse the sentences



69,592 dependency paths with >5 pairs

5. Extract patterns

logistic regression with 70K features
(converted to 974,288 bucketed binary features)



One of 70,000 patterns

Pattern: <superordinate> called <subordinate>

Learned from cases such as:

- | | |
|-------------------|---|
| (sarcoma, cancer) | ...an uncommon bone cancer called osteogenic sarcoma and to... |
| (deuterium, atom) | ...heavy water rich in the doubly heavy hydrogen atom called deuterium . |

New pairs discovered:

- | | |
|-----------------------------|---|
| (efflorescence, condition) | ...and a condition called efflorescence are other reasons for... |
| (O'Neal_inc, company) | ...The company , now called O'Neal Inc , was sole distributor of... |
| (hat_creek_outfit, ranch) | ...run a small ranch called the Hat Creek Outfit . |
| (hiv-1, aids_virus) | ...infected by the AIDS virus , called HIV-1 . |
| (bateau_mouche, attraction) | ...local sightseeing attraction called the Bateau Mouche ... |

Syntactic dependency paths

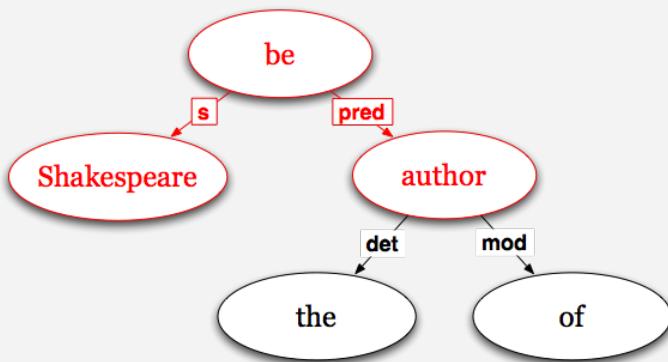
Patterns are based on paths through dependency parses generated by MINIPAR (Lin, 1998)



Example word pair: (Shakespeare, author)

Example sentence: “Shakespeare was the author of several plays...”

Minipar parse:



Extract shortest path:
-N:s:VBE, be, VBE:pred:N

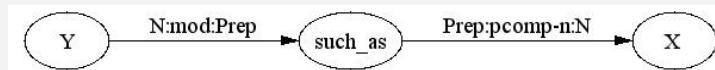
Hearst patterns to dependency paths

Hearst Pattern



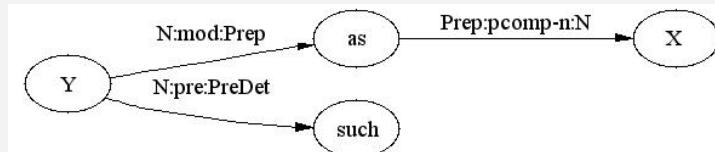
Y such as X ...

MINIPAR Representation



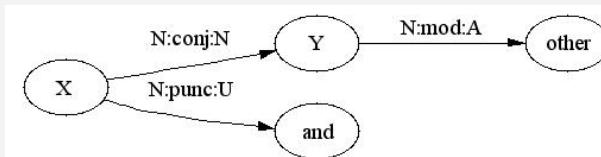
-N:pcomp-n:Prep,such_as,such_as,-Prep:mod:N

Such Y as X ...



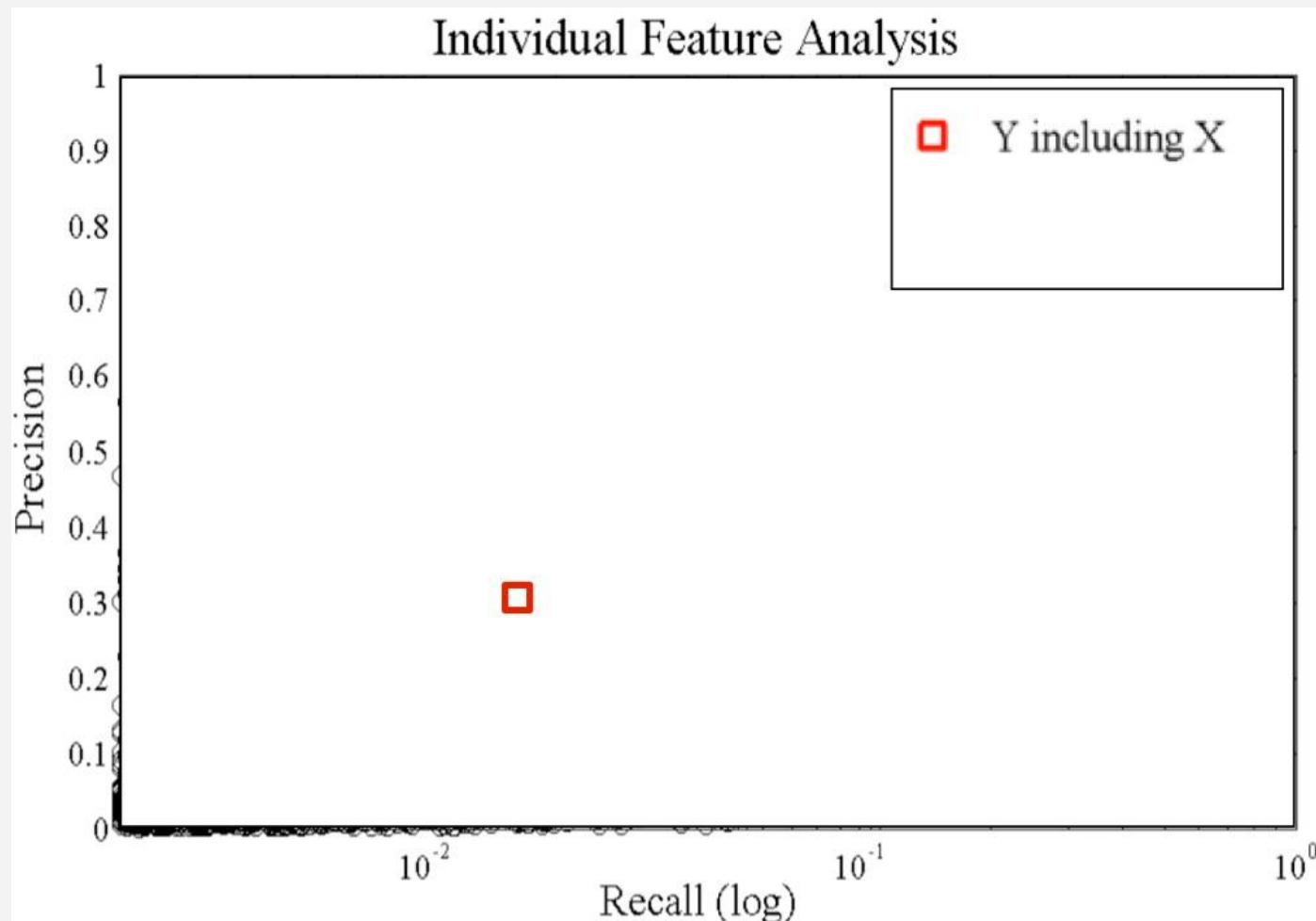
-N:pcomp-n:Prep,as,as,-Prep:mod:N,(such,PreDet:pre:N)}

X ... and other Y

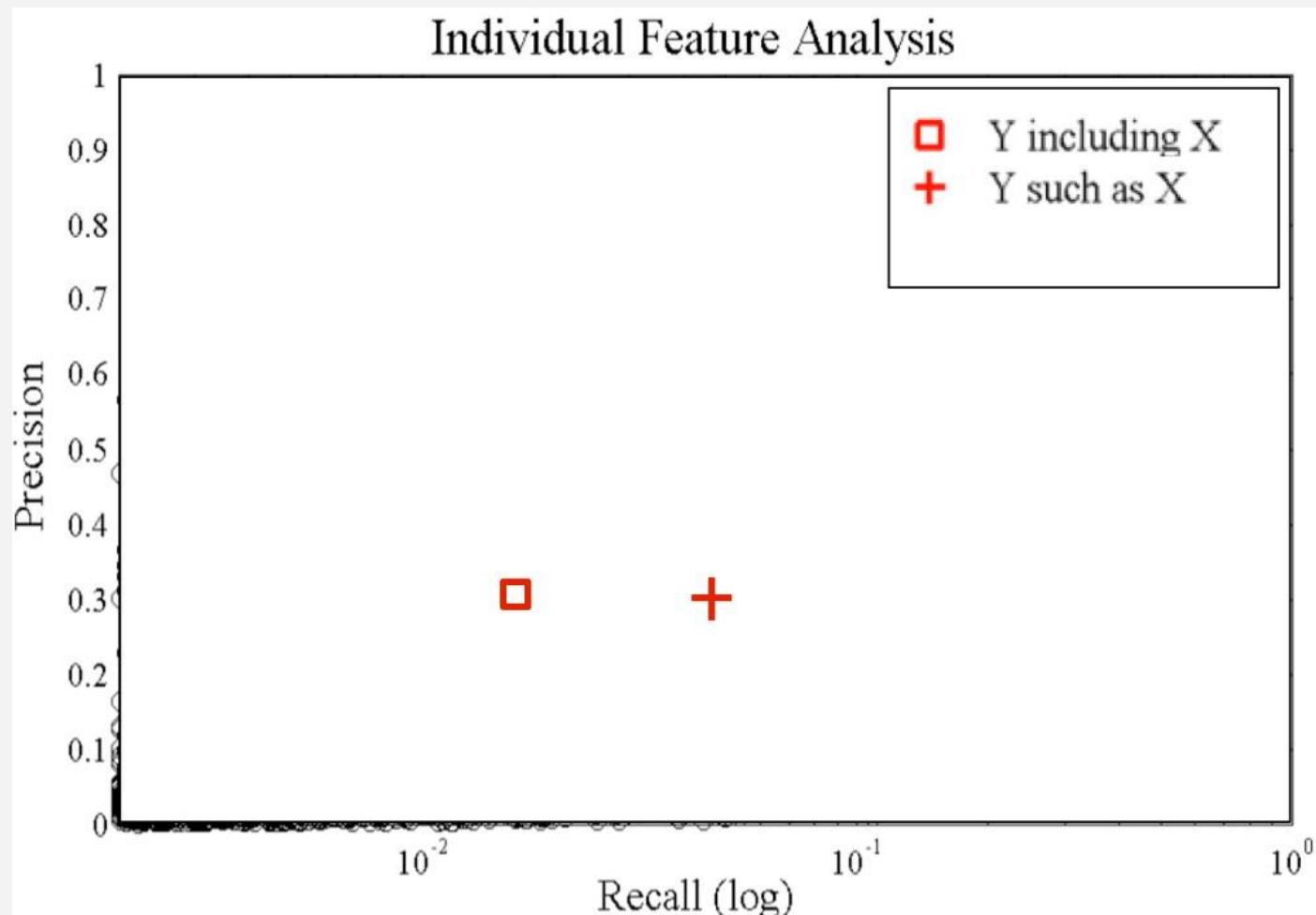


(and,U:punc:N),N:conj:N, (other,A:mod:N)

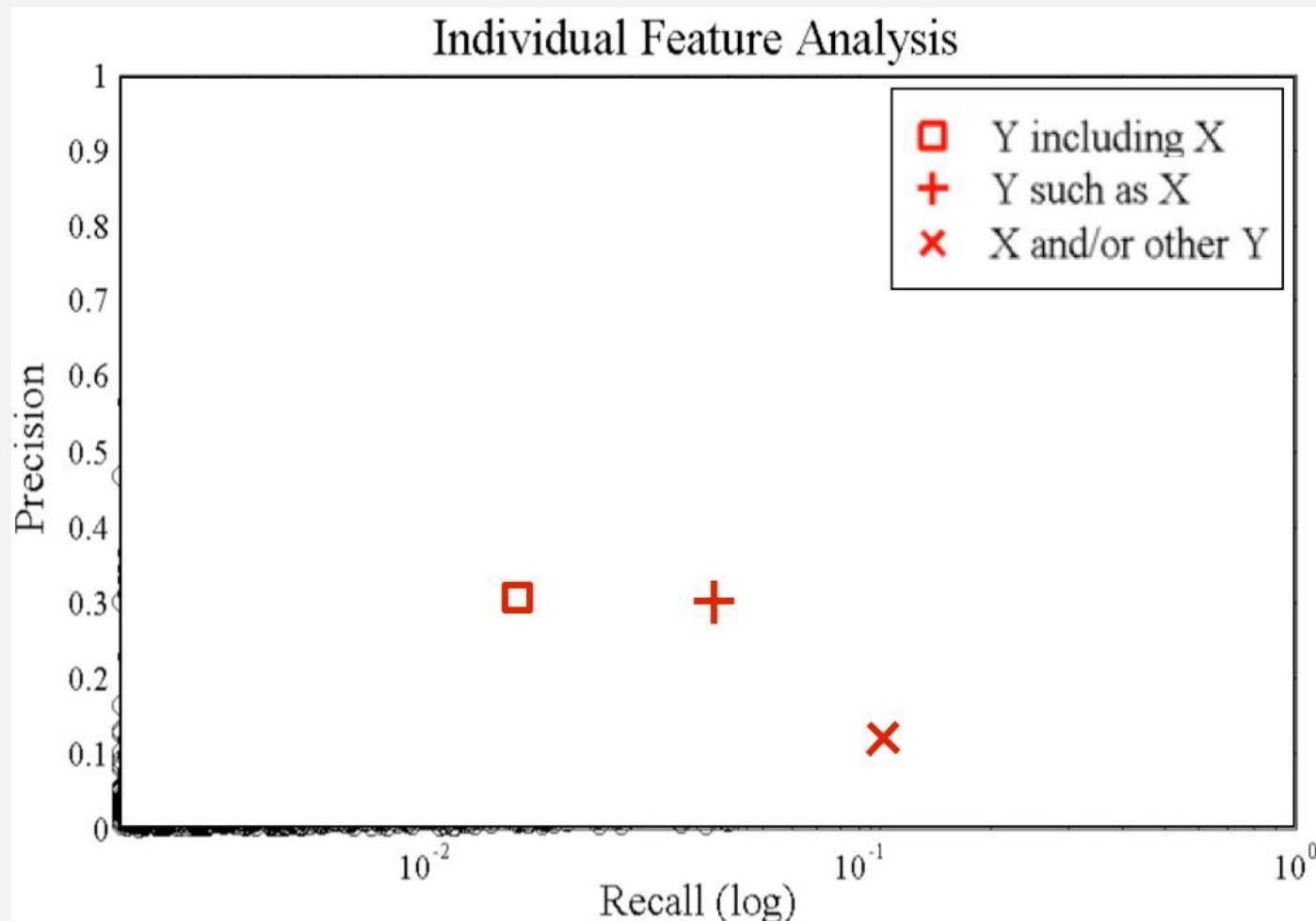
P/R of hypernym extraction patterns



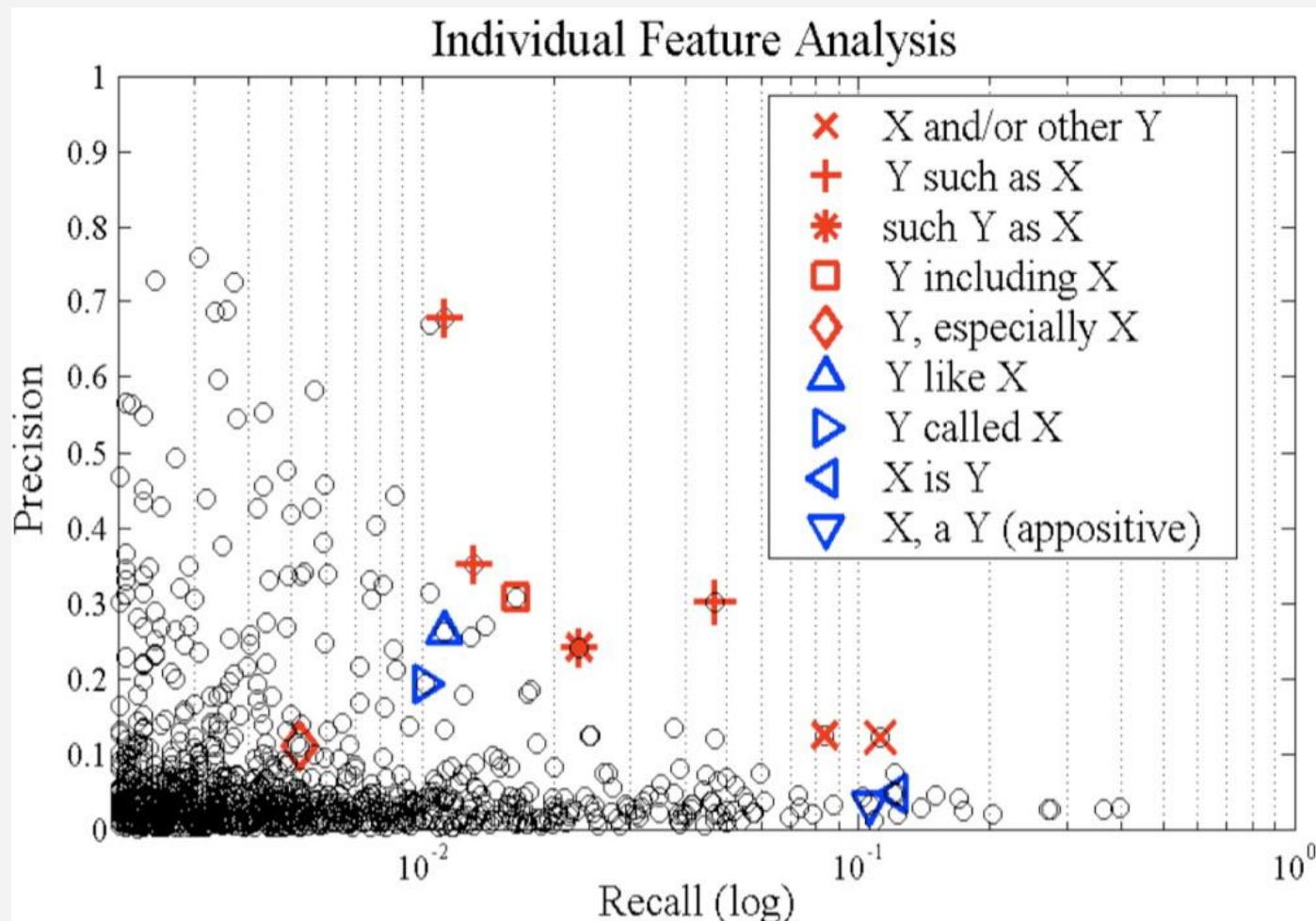
P/R of hypernym extraction patterns



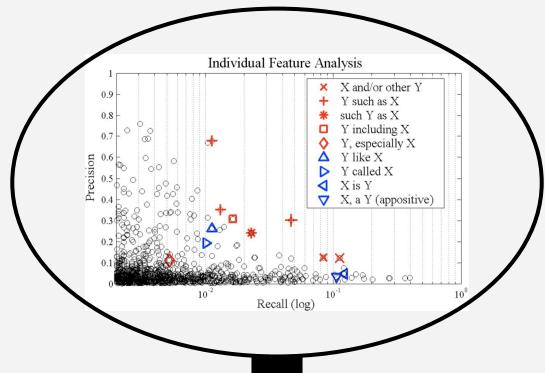
P/R of hypernym extraction patterns



P/R of hypernym extraction patterns

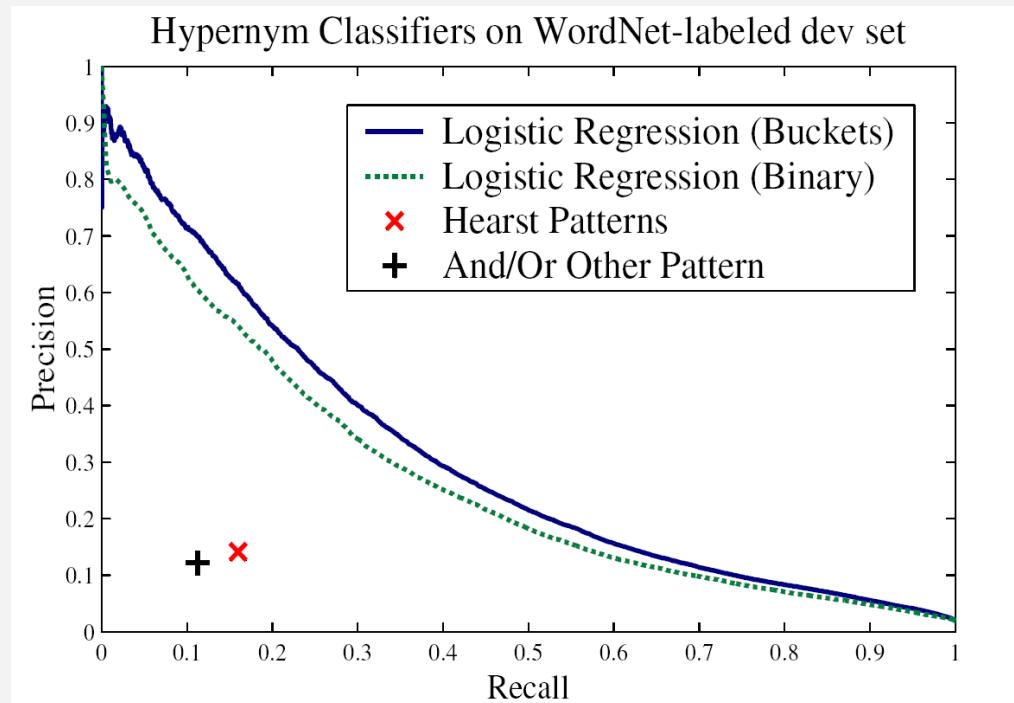


P/R of hypernym classifier



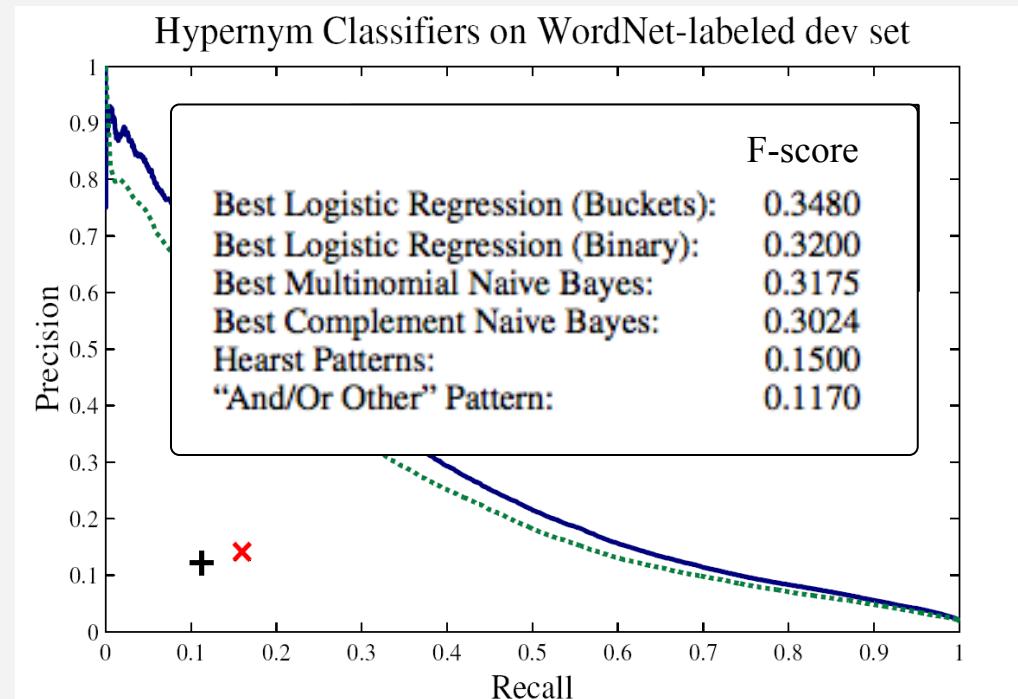
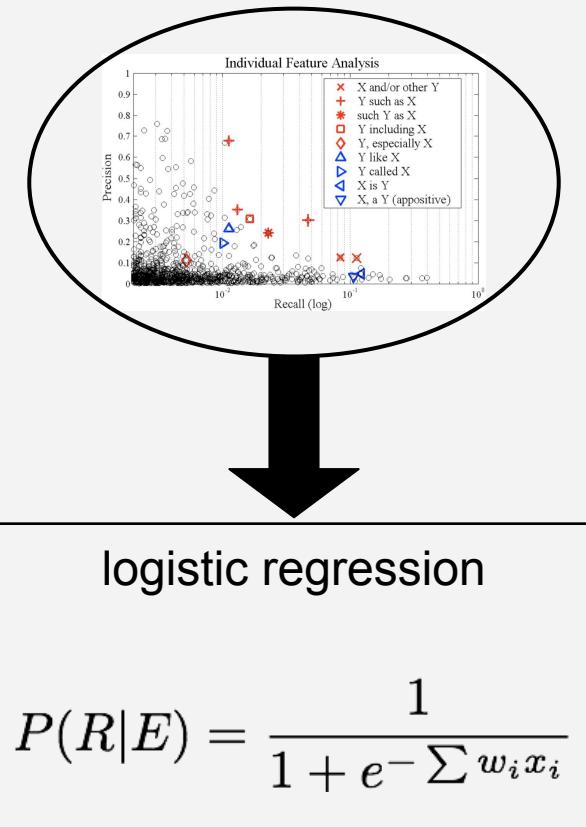
logistic regression

$$P(R|E) = \frac{1}{1 + e^{-\sum w_i x_i}}$$



10-fold Cross Validation on
14,000 WordNet-Labeled Pairs

P/R of hypernym classifier



10-fold Cross Validation on
14,000 WordNet-Labeled Pairs

What about other relations?

Mintz, Bills, Snow, Jurafsky (2009).

Distant supervision for relation extraction without labeled data.



Training set



102 relations
940,000 entities
1.8 million instances

Corpus



1.8 million articles
25.7 million sentences



Frequent Freebase relations

Relation name	Size	Example
/people/person/nationality	281,107	John Dugard, South Africa
/location/location/contains	253,223	Belgium, Nijlen
/people/person/profession	208,888	Dusa McDuff, Mathematician
/people/person/place_of_birth	105,799	Edwin Hubble, Marshfield
/dining/restaurant/cuisine	86,213	MacAyo's Mexican Kitchen, Mexican
/business/business_chain/location	66,529	Apple Inc., Apple Inc., South Park, NC
/biology/organism_classification_rank	42,806	Scorpaeniformes, Order
/film/film/genre	40,658	Where the Sidewalk Ends, Film noir
/film/film/language	31,103	Enter the Phoenix, Cantonese
/biology/organism_higher_classification	30,052	Calopteryx, Calopterygidae
/film/film/country	27,217	Turtle Diary, United States
/film/writer/film	23,856	Irving Shulman, Rebel Without a Cause
/film/director/film	23,539	Michael Mann, Collateral
/film/producer/film	22,079	Diane Eskenazi, Aladdin
/people/deceased_person/place_of_death	18,814	John W. Kern, Asheville
/music/artist/origin	18,619	The Octopus Project, Austin
/people/person/religion	17,582	Joseph Chartrand, Catholicism
/book/author/works_written	17,278	Paul Auster, Travels in the Scriptorium
/soccer/football_position/players	17,244	Midfielder, Chen Tao
/people/deceased_person/cause_of_death	16,709	Richard Daintree, Tuberculosis
/book/book/genre	16,431	Pony Soldiers, Science fiction
/film/film/music	14,070	Stavisky, Stephen Sondheim
/business/company/industry	13,805	ATS Medical, Health care



Collecting training data

Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

Training data

Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)



Collecting training data

Corpus text

Bill Gates founded Microsoft in 1975.

Bill Gates, founder of Microsoft, ...

Bill Gates attended Harvard from...

Google was founded by Larry Page ...

Training data

(Bill Gates, Microsoft)

Label: Founder

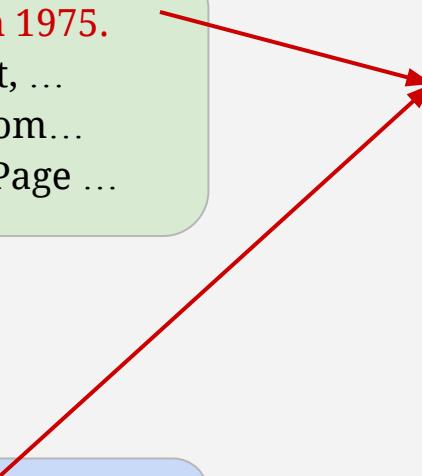
Feature: X founded Y

Freebase

Founder: (Bill Gates, Microsoft)

Founder: (Larry Page, Google)

CollegeAttended: (Bill Gates, Harvard)





Collecting training data

Corpus text

Bill Gates founded Microsoft in 1975.

Bill Gates, founder of Microsoft, ...

Bill Gates attended Harvard from...

Google was founded by Larry Page ...

Training data

(Bill Gates, Microsoft)

Label: Founder

Feature: X founded Y

Feature: X, founder of Y

Freebase

Founder: (Bill Gates, Microsoft)

Founder: (Larry Page, Google)

CollegeAttended: (Bill Gates, Harvard)



Collecting training data

Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

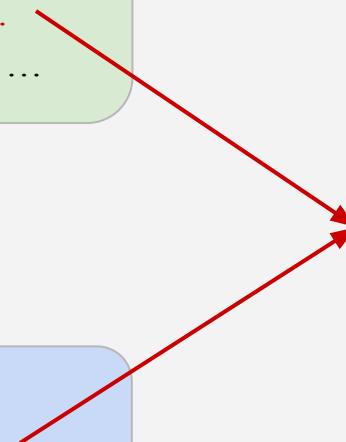
Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

Training data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

(Bill Gates, Harvard)
Label: CollegeAttended
Feature: X attended Y





Collecting training data

Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from...
Google was founded by Larry Page ...

Freebase

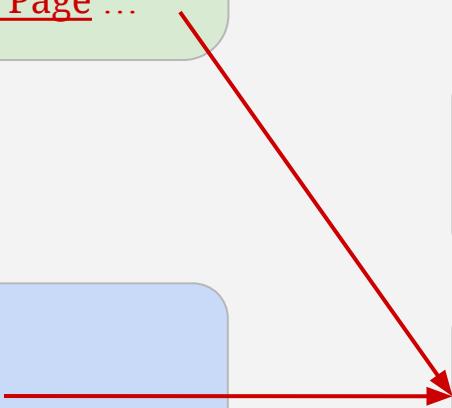
Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

Training data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

(Bill Gates, Harvard)
Label: CollegeAttended
Feature: X attended Y

(Larry Page, Google)
Label: Founder
Feature: Y was founded by X



Negative training data

Can't train a classifier with only positive data!

Need negative training data too!

Solution?

Sample 1% of unrelated pairs of entities.

Corpus text

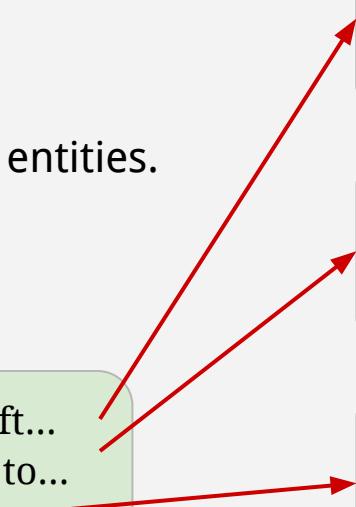
Larry Page took a swipe at Microsoft...
...after Harvard invited Larry Page to...
Google is Bill Gates' worst fear ...

Training data

(Larry Page, Microsoft)
Label: NO_RELATION
Feature: X took a swipe at Y

(Larry Page, Harvard)
Label: NO_RELATION
Feature: Y invited X

(Bill Gates, Google)
Label: NO_RELATION
Feature: Y is X's worst fear





Preparing test data

Corpus text

Henry Ford founded Ford Motor Co. in...
Ford Motor Co. was founded by Henry Ford...
Steve Jobs attended Reed College from...

Test data





Preparing test data

Corpus text

Henry Ford founded Ford Motor Co. in...

Ford Motor Co. was founded by Henry Ford...

Steve Jobs attended Reed College from...

Test data

(Henry Ford, Ford Motor Co.)
Label: ???
Feature: X founded Y





Preparing test data

Corpus text

Henry Ford founded Ford Motor Co. in...

Ford Motor Co. was founded by Henry Ford...

Steve Jobs attended Reed College from...

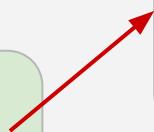
Test data

(Henry Ford, Ford Motor Co.)

Label: ???

Feature: X founded Y

Feature: Y was founded by X





Preparing test data

Corpus text

Henry Ford founded Ford Motor Co. in...

Ford Motor Co. was founded by Henry Ford...

Steve Jobs attended Reed College from...

Test data

(Henry Ford, Ford Motor Co.)

Label: ???

Feature: X founded Y

Feature: Y was founded by X

(Steve Jobs, Reed College)

Label: ???

Feature: X attended Y



The experiment

Positive training data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

(Bill Gates, Harvard)
Label: CollegeAttended
Feature: X attended Y

(Larry Page, Google)
Label: Founder
Feature: Y was founded by X

Negative training data

(Larry Page, Microsoft)
Label: NO_RELATION
Feature: X took a swipe at Y

(Larry Page, Harvard)
Label: NO_RELATION
Feature: Y invited X

(Bill Gates, Google)
Label: NO_RELATION
Feature: Y is X's worst fear

**Learning:
multiclass
logistic
regression**

Test data

(Henry Ford, Ford Motor Co.)
Label: ???
Feature: X founded Y
Feature: Y was founded by X

(Steve Jobs, Reed College)
Label: ???
Feature: X attended Y

**Trained
relation
classifier**

Predictions!

(Henry Ford, Ford Motor Co.)
Label: Founder

(Steve Jobs, Reed College)
Label: CollegeAttended

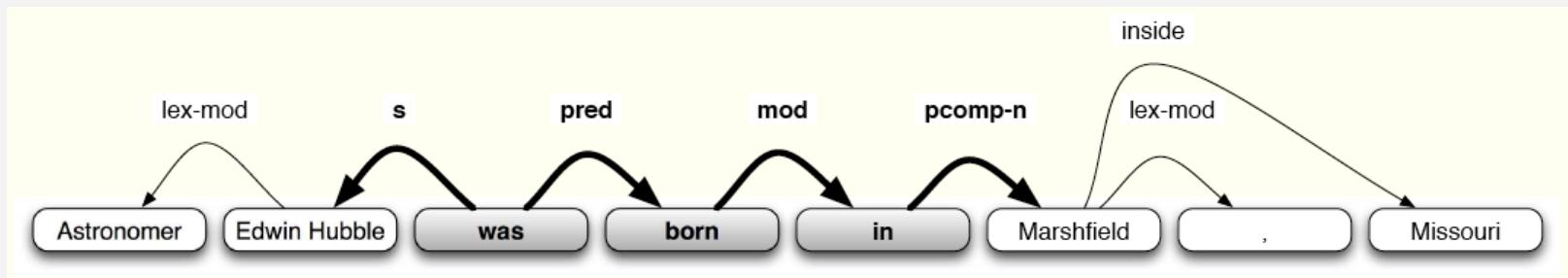


Advantages of the approach

- ACE paradigm: labeling sentences
- This paradigm: labeling entity pairs
- We make use of multiple appearances of entities
- If a pair of entities appears in 10 sentences, and each sentence has 5 features extracted from it, the entity pair will have 50 associated features

Lexical and syntactic features

Astronomer **Edwin Hubble** was born in **Marshfield**, Missouri.



Feature type	Left window	NE1	Middle	NE2	Right window
Lexical	[]	PER	[was/VERB born/VERB in/CLOSED]	LOC	[]
Lexical	[Astronomer]	PER	[was/VERB born/VERB in/CLOSED]	LOC	[,]
Lexical	[#PAD#, Astronomer]	PER	[was/VERB born/VERB in/CLOSED]	LOC	[, Missouri]
Syntactic	[]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[]
Syntactic	[Edwin Hubble ↓ _{lex-mod}]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[]
Syntactic	[Astronomer ↓ _{lex-mod}]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[]
Syntactic	[]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[↓ _{lex-mod} ,]
Syntactic	[Edwin Hubble ↓ _{lex-mod}]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[↓ _{lex-mod} ,]
Syntactic	[Astronomer ↓ _{lex-mod}]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[↓ _{lex-mod} ,]
Syntactic	[]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[↓ _{inside} Missouri]
Syntactic	[Edwin Hubble ↓ _{lex-mod}]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[↓ _{inside} Missouri]
Syntactic	[Astronomer ↓ _{lex-mod}]	PER	[↑ _s was ↓ _{pred} born ↓ _{mod} in ↓ _{pcomp-n}]	LOC	[↓ _{inside} Missouri]



High-weight features

Relation	Feature type	Left window	NE1	Middle	NE2	Right window
/architecture/structure/architect	LEX-~		ORG	, the designer of the	PER	
/book/author/works_written	SYN	designed \uparrow_s	ORG	\uparrow_s designed $\Downarrow_{by-subj}$ by \Downarrow_{pcn}	PER	\uparrow_s designed
/book/book_edition/author_editor	LEX		PER	s novel	ORG	
/business/company/founders	SYN		ORG	\uparrow_{pcn} by \uparrow_{mod} story \uparrow_{pred} is \Downarrow_s	ORG	
/business/company/place_founded	LEX-~		PER	s novel	PER	
/business/company/place_founded	SYN		ORG	\uparrow_{nn} series \Downarrow_{gen}	PER	
/film/film/country	LEX		ORG	co - founder	PER	
/film/film/country	SYN		ORG	\uparrow_{nn} owner \Downarrow_{person}	PER	
/geography/river/mouth	LEX		ORG	- based	LOC	
/government/political_party/country	SYN	opened \uparrow_s	PER	\uparrow_s founded \Downarrow_{mod} in \Downarrow_{pcn}	LOC	
/influence/influence_node/influenced	LEX-~		ORG	, released in	LOC	
/language/human_language/region	SYN		LOC	\uparrow_s opened \Downarrow_{mod} in \Downarrow_{pcn}	LOC	\uparrow_s opened
/language/human_language/region	LEX		LOC	, which flows into the	LOC	
/music/artist/origin	SYN	the \Downarrow_{det}	ORG	\uparrow_s is \Downarrow_{pred} tributary \Downarrow_{mod} of \Downarrow_{pcn}	LOC	\Downarrow_{det} the
/people/deceased_person/place_of_death	LEX-~		ORG	politician of the	LOC	
/people/person/nationality	SYN	candidate \uparrow_{nn}	PER	\uparrow_{nn} candidate \Downarrow_{mod} for \Downarrow_{pcn}	LOC	\uparrow_{nn} candidate
/people/person/nationality	LEX		PER	, a student of	PER	
/people/person/parents	SYN	of \uparrow_{pcn}	PER	\uparrow_{pcn} of \uparrow_{mod} student \uparrow_{appo}	PER	\uparrow_{pcn} of
/people/person/parents	LEX		LOC	- speaking areas of	LOC	
/people/person/place_of_birth	SYN		LOC	$\uparrow_{lex-mod}$ speaking areas \Downarrow_{mod} of \Downarrow_{pcn}	LOC	
/people/person/place_of_birth	LEX-~		ORG	based band	LOC	
/people/person/religion	SYN	is \uparrow_s	ORG	\uparrow_s is \Downarrow_{pred} band \Downarrow_{mod} from \Downarrow_{pcn}	LOC	\uparrow_s is
/people/person/religion	LEX		PER	died in	LOC	
/people/person/religion	SYN	hanged \uparrow_s	PER	\uparrow_s hanged \Downarrow_{mod} in \Downarrow_{pcn}	LOC	\uparrow_s hanged
/people/person/religion	LEX		PER	is a citizen of	LOC	
/people/person/religion	SYN		PER	\Downarrow_{mod} from \Downarrow_{pcn}	LOC	
/people/person/religion	LEX		PER	, son of	PER	
/people/person/religion	SYN	father \uparrow_{gen}	PER	\uparrow_{gen} father \Downarrow_{person}	PER	\uparrow_{gen} father
/people/person/religion	LEX-~		PER	is the birthplace of	PER	
/people/person/religion	SYN		PER	\uparrow_s born \Downarrow_{mod} in \Downarrow_{pcn}	LOC	
/people/person/religion	LEX		PER	embraced	LOC	
/people/person/religion	SYN	convert \Downarrow_{appo}	PER	\Downarrow_{appo} convert \Downarrow_{mod} to \Downarrow_{pcn}	LOC	\Downarrow_{appo} convert



Implementation

- Classifier: multi-class logistic regression optimized using L-BFGS with Gaussian regularization (Manning & Klein 2003)
- Parser: MINIPAR (Lin 1998)
- POS tagger: MaxEnt tagger trained on the Penn Treebank (Toutanova et al. 2003)
- NER tagger: Stanford four-class tagger {PER, LOC, ORG, MISC, NONE} (Finkel et al. 2005)
- 3 configurations: lexical features, syntax features, both



Experiment

- 1.8 million relation instances used for training
 - Compared to 17,000 relation instances in ACE
- 800,000 Wikipedia articles used for training,
400,000 different articles used for testing
- Only extract relation instances not already in Freebase



Newly discovered instances

Ten relation instances extracted by the system that weren't in Freebase

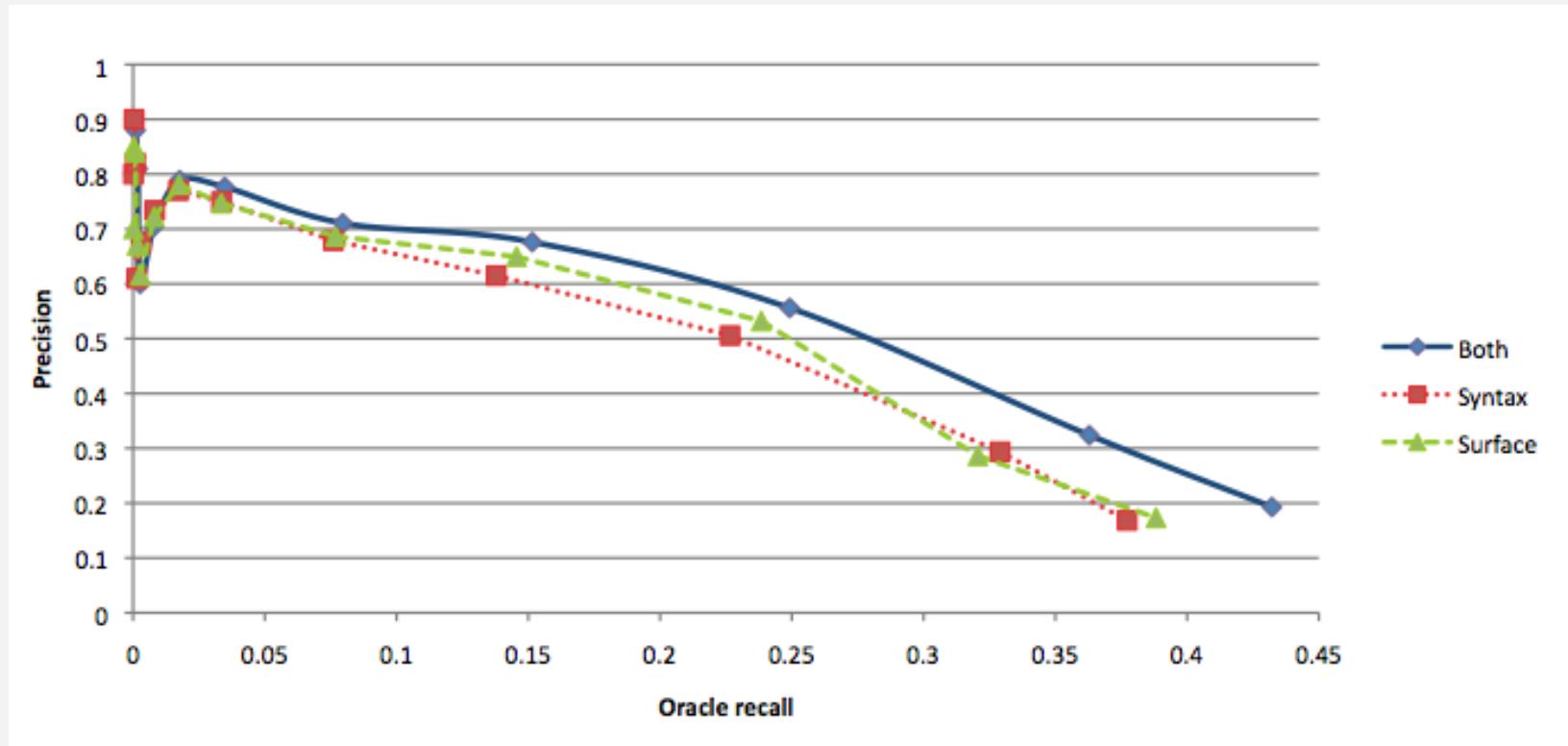
Relation name	New instance
/location/location/contains	Paris, Montmartre
/location/location/contains	Ontario, Fort Erie
/music/artist/origin	Mighty Wagon, Cincinnati
/people/deceased_person/place_of_death	Fyodor Kamensky, Clearwater
/people/person/nationality	Marianne Yvonne Heemskerk, Netherlands
/people/person/place_of_birth	Wavell Wayne Hinds, Kingston
/book/author/works_written	Upton Sinclair, Lanny Budd
/business/company/founders	WWE, Vince McMahon
/people/person/profession	Thomas Mellon, judge



Evaluation

- Held-out evaluation
 - Train on 50% of gold-standard Freebase relation instances, test on other 50%
 - Used to tune parameters quickly without having to wait for human evaluation
- Human evaluation
 - Performed by evaluators on Amazon Mechanical Turk
 - Calculated precision at 100 and 1000 recall levels for the ten most common relations

Held-out evaluation



Automatic evaluation on 900K instances of 102 Freebase relations.
Precision for three different feature sets is reported at various recall levels.

Human evaluation

Precision, using Mechanical Turk labelers:

Relation name	100 instances			1000 instances		
	Syn	Lex	Both	Syn	Lex	Both
/film/director/film	0.49	0.43	0.44	0.49	0.41	0.46
/film/writer/film	0.70	0.60	0.65	0.71	0.61	0.69
/geography/river/basin_countries	0.65	0.64	0.67	0.73	0.71	0.64
/location/country/administrative_divisions	0.68	0.59	0.70	0.72	0.68	0.72
/location/location/contains	0.81	0.89	0.84	0.85	0.83	0.84
/location/us_county/county_seat	0.51	0.51	0.53	0.47	0.57	0.42
/music/artist/origin	0.64	0.66	0.71	0.61	0.63	0.60
/people/deceased_person/place_of_death	0.80	0.79	0.81	0.80	0.81	0.78
/people/person/nationality	0.61	0.70	0.72	0.56	0.61	0.63
/people/person/place_of_birth	0.78	0.77	0.78	0.88	0.85	0.91
Average	0.67	0.66	0.69	0.68	0.67	0.67

- At recall of 100 instances, using both feature sets (lexical and syntax) offers the best performance for a majority of the relations
- At recall of 1000 instances, using syntax features improves performance for a majority of the relations



Where syntax helps

Back Street is a 1932 film made by Universal Pictures, directed by John M. Stahl, and produced by Carl Laemmle Jr.

Back Street and *John M. Stahl* are far apart in surface string, but close together in dependency parse



Where syntax doesn't help

Beaverton is a city in *Washington County*, Oregon ...

Beaverton and *Washington County* are close together in the surface string.



Distant supervision: conclusions

- Distant supervision extracts high-precision patterns for a variety of relations
- Can make use of 1000x more data than simple supervised algorithms
- Syntax features almost always help
- The combination of syntax and lexical features is sometimes even better
- Syntax features are probably most useful when entities are far apart, often when there are modifiers in between



Distant supervision: discussion

- Relation extraction → learning by reading
- Suppose we could do relation extraction perfectly?
- What would we still be missing?
- What knowledge could we still not gather from the web?



Relation extraction: 5 easy methods

1. Hand-built patterns
2. Bootstrapping methods
3. Supervised methods
4. Distant supervision
5. Unsupervised methods

KnowItAll (Etzioni et al. 2005)

- Input: target class labels and relation labels

Predicate	class label	relation label
City	“city”, “town”	—
Country	“country”, “nation”	—
capitalOf(City,Country)	—	“capital of”



- Use Hearst patterns to find instances of classes
 - ENTITY and/or other CLASS; such CLASS as ENTITY; etc.
- Now use new pattern templates to find relations
 - CLASS1 is the RELATION CLASS2
 - CLASS1, RELATION CLASS2
- So once you learn Paris and Berlin are cities:
 - “Paris is the capital of France” → capitalOf(Paris, France)



KnowItAll PMI-based assessor

- Validate candidate instances using “discriminators”
- Compare # search engine hits for

- instance alone
- instance + discriminator

$$\text{PMI}(I, D) = \frac{|\text{Hits}(D + I)|}{|\text{Hits}(I)|}$$

- (This is *not* the conventional definition of PMI!)
- Use “PMI” scores as features for Naïve Bayes classifier
- Example: linguists such as
 - $\text{PMI}(\text{linguists such as, Chomsky}) = 275K / 16.9M = 1.63 \text{ E-02}$
 - $\text{PMI}(\text{linguists such as, Potts}) = 3120 / 27.7M = 1.13 \text{ E-04}$

TextRunner (Banko et al. 2007)



1. **Self-supervised learner**: automatically labels +/- examples & learns a crude relation extractor
2. **Single-pass extractor**: makes one pass over corpus, extracting candidate relations in each sentence
3. **Redundancy-based assessor**: assigns a probability to each extraction, based on frequency counts



Step 1: Self-supervised learner

- Run a parser over 2000 sentences
 - Parsing is relatively expensive, so can't run on whole web
 - For each pair of base noun phrases NP_i and NP_j
 - Extract all tuples $t = (NP_i, \text{relation}_{i,j}, NP_j)$
- Label each tuple based on features of parse:
 - Positive iff the dependency path between the NPs is short, and doesn't cross a clause boundary, and neither NP is a pronoun
- Now train a Naïve Bayes classifier on the labeled tuples
 - Using *lightweight* features like POS tags nearby, stop words, etc.



Step 2: Single-pass extractor

- Over a huge (web-sized) corpus:
 - Run a dumb POS tagger
 - Run a dumb Base Noun Phrase chunker
 - Extract all text strings between base NPs
 - Run heuristic rules to simplify text strings
Scientists from many universities are intently studying stars
→ ⟨**scientists**, **are studying**, **stars**⟩
- Pass candidate tuples to Naïve Bayes classifier
- Save only those predicted to be “trustworthy”



Step 3: Redundancy-based assessor

- Collect counts for each simplified tuple
 $\langle \text{scientists}, \text{are studying}, \text{stars} \rangle \rightarrow 17$
- Compute likelihood of each tuple
 - given the counts for each relation
 - and the number of sentences
 - and a combinatoric balls-and-urns model [Downey et al. 05]

$$P(x \in C | x \text{ appears } k \text{ times in } n \text{ draws}) \approx \frac{1}{1 + \frac{|E|}{|C|} \left(\frac{p_E}{p_C} \right)^k e^{n(p_C - p_E)}}$$



TextRunner demo

<http://www.cs.washington.edu/research/textrunner/>

(Note that they've re-branded TextRunner as ReVerb,
but it's largely the same system.)



TextRunner examples

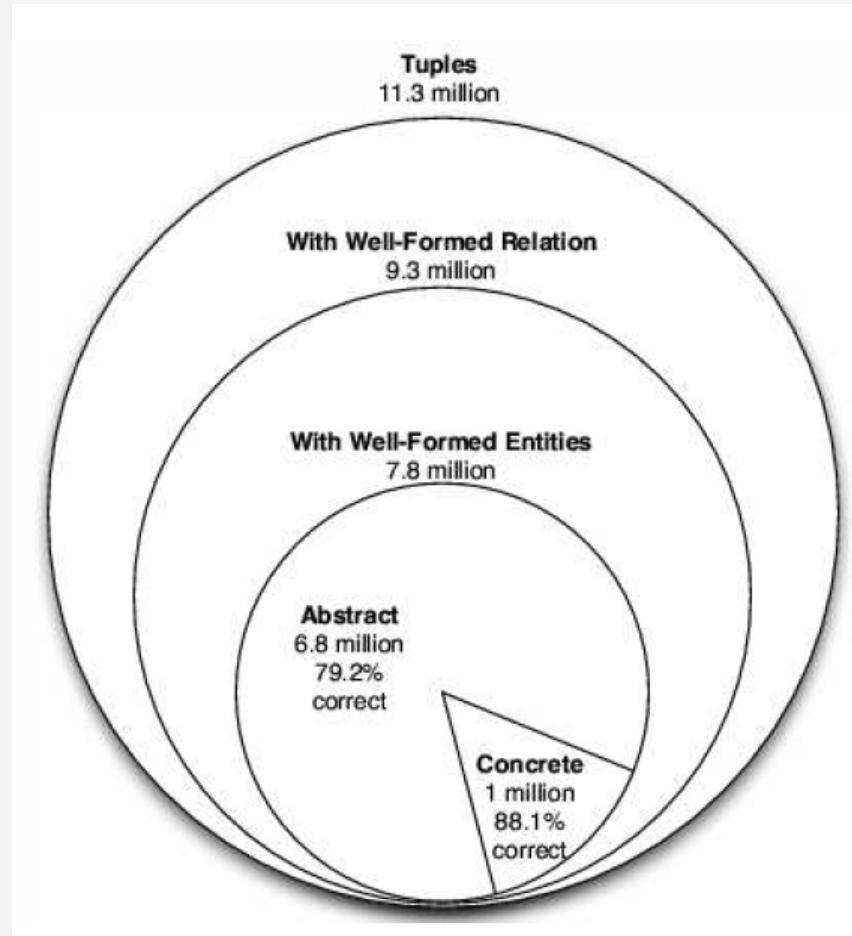
Probability	Count	Arg1	Predicate	Arg2
0.98	59	Smith	invented	the margherita
0.97	49	Al Gore	invented	the Internet
0.97	44	manufacturing plant	first invented	the automatic revolver
0.97	41	Alexander Graham Bell	invented	the telephone
0.97	36	Thomas Edison	invented	light bulbs
0.97	29	Eli Whitney	invented	the cotton gin
0.96	23	C. Smith	invented	the margherita
0.96	19	the Digital Equipment Corporation manufacturing plant	first invented	the automatic revolver
0.96	18	Edison	invented	the phonograph



TextRunner results

- From corpus of 9M web pages, containing 133M sentences
- Extracted 60.5 million tuples
 - $\langle FCI, \text{specializes in}, \text{software development} \rangle$
- Evaluation
 - Not well formed:
 - $\langle \text{demands, of securing, border} \rangle \langle 29, \text{dropped, instruments} \rangle$
 - Abstract:
 - $\langle Einstein, \text{derived, theory} \rangle \langle \text{executive, hired by, company} \rangle$
 - True, concrete:
 - $\langle Tesla, \text{invented, coil transformer} \rangle$

Evaluating TextRunner





DIRT (Lin & Pantel 2003)

- DIRT = Discovery of Inference Rules from Text
- Looks at MINIPAR dependency paths between noun pairs
 - N:subj:V←find→V:obj:N→solution→N:to:N
 - i.e., X finds solution to Y
- Applies "extended distributional hypothesis"
 - If two paths tend to occur in similar contexts, the meanings of the paths tend to be similar.
- So, defines path similarity in terms of cooccurrence counts with various slot fillers
- Thus, extends ideas of (Lin 1998) from *words* to *paths*



DIRT examples

The top-20 most similar paths to “X solves Y”:

Y is solved by X
X resolves Y
X finds a solution to Y
X tries to solve Y
X deals with Y
Y is resolved by X
X addresses Y
X seeks a solution to Y
X do something about Y
X solution to Y

Y is resolved in X
Y is solved through X
X rectifies Y
X copes with Y
X overcomes Y
X eases Y
X tackles Y
X alleviates Y
X corrects Y
X is a solution to Y



Ambiguous paths in DIRT

- X **addresses** Y
 - I **addressed** my letter to him personally.
 - She **addressed** an audience of Shawnee chiefs.
 - Will Congress finally **address** the immigration issue?
- X **tackles** Y
 - Foley **tackled** the quarterback in the endzone.
 - Police are beginning to **tackle** rising crime.
- X **is a solution to** Y
 - (5, 1) **is a solution to** the equation $2x - 3y = 7$
 - Nuclear energy **is a solution to** the energy crisis.



Yao et al. 2012: motivation

- Goal: induce clusters of dependency paths which express the same semantic relation, like DIRT
- But, improve upon DIRT by properly handling semantic ambiguity of individual paths



Yao et al. 2012: approach

1. Extract tuples (entity, path, entity) from corpus
2. Construct feature representations of every tuple
3. Group the tuples for each path into sense clusters
4. Cluster the sense clusters into semantic relations



Extracting tuples

- Start with NYT corpus
- Apply lemmatization, NER tagging, dependency parsing
- For each pair of entities in a sentence:
 - Extract dependency path between them, as in Lin
 - Form a tuple consisting of the two entities and the path
- Filter rare tuples, tuples with two direct objects, etc.
- Result: 1M tuples, 500K entities, 1300 patterns



Feature representation

- Entity names, as bags of words, prefixed with "l:" or "r:"
 - ex: ("LA Lakers", "NY Knicks") => {l:LA, l:Lakers, r:NY, r:Knicks}
 - Using bag-of-words encourages overlap, i.e., combats sparsity
- Words between and around the two entities
 - Exclude stop words, words with capital letters
 - Include two words to the left and right
- Document theme (e.g. sports, politics, finance)
 - Assigned by an LDA topic model which treats NYTimes topic descriptors as words in a synthetic document
- Sentence theme
 - Assigned by a standard LDA topic model

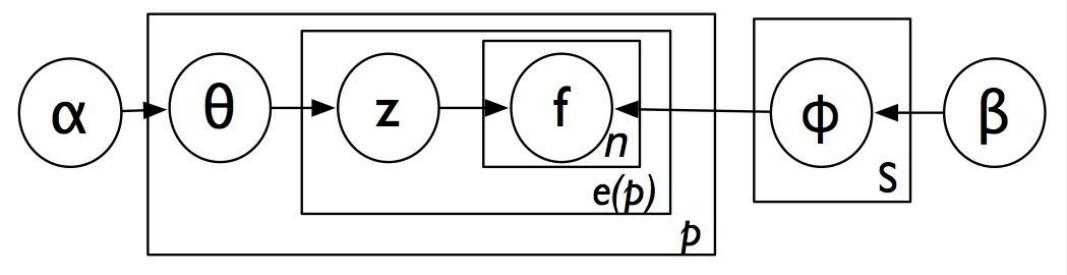


Clustering tuples into senses

- Goal: group tuples for each path into coherent sense clusters
- To do this, we apply yet another LDA topic model
 - Not vanilla LDA this time — rather, a slight variant
 - Details on next slide
- Use Gibbs sampling for inference
- Result: each tuple is assigned one topic/sense
- Tuples with the same topic/sense constitute a cluster

The Sense-LDA model

$$\begin{aligned}
 \theta_{p_i} &\sim \text{Dirichlet}(\alpha) \\
 \phi_z &\sim \text{Dirichlet}(\beta) \\
 z_e &\sim \text{Multinomial}(\theta_{p_i}) \\
 f_k &\sim \text{Multinomial}(\phi_{z_e})
 \end{aligned}$$



- A slight variation on standard LDA (Blei et al. 2003)
- For each path, form “document” of all its tuples, with features
- For each path/document, sample a multinomial distribution θ over topics/senses from a Dirichlet prior
- For each tuple, sample a topic/sense from θ
- Features are sampled from a topic/sense-specific multinomial
- Features are conditionally independent, given topic/sense



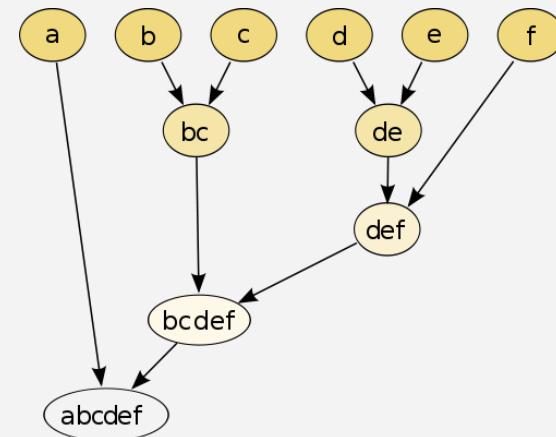
Sense cluster examples

Path	20:sports	30:entertainment	25:music/art
A play B	Americans, Ireland Yankees, Angels Ecuador, England Redskins, Detroit Red Bulls, F.C. Barcelona	Jean-Pierre Bacri, Jacques Rita Benton, Gay Head Dance Jeanie, Scrabble Meryl Streep, Leilah Kevin Kline, Douglas Fairbanks	Daniel Barenboim, recital of Mozart Mr. Rose, Ballade Gil Shaham, Violin Romance Ms. Golabek, Steinways Bruce Springsteen, Saints
doc theme	sports	music books television	music theater
sen theme	game yankees	theater production book film show	music reviews opera
lexical words	beat victory num-num won	played plays directed artistic r:theater	director conducted production r:theater r:hall r:york l:opera
entity names	-		

Sense clusters for path "A play B",
 along with sample entity pairs and top features.

Clustering the clusters!

- Now cluster sense clusters from different paths into semantic relations — this is the part most similar to Lin & Pantel 2003
- Uses Hierarchical Agglomerative Clustering (HAC)
- Start with minimal clustering, then merge progressively
- Uses cosine similarity between sense-cluster feature vectors
- Uses complete-linkage strategy





Semantic relation results

relation	paths
entertainment	A, who play B:30; A play B 30; star A as B:30
sports	lead A to victory over B:20; A play to B:20; A play B 20; A's loss to B:20; A beat B:20; A trail B:20; A face B:26; A hold B:26; A play B:26; A acquire (X) from B:26; A send (X) to B:26;
politics	A nominate B:39; A name B:39; A select B:39; A name B:42; A select B:42; A ask B:42; A choose B:42; A nominate B:42; A turn to B:42;
law	A charge B:39; A file against B:39; A accuse B:39; A sue B:39

Just like DIRT, each semantic relation has multiple paths.

But, one path can now appear in multiple semantic relations.

DIRT can't do that!



Evaluation against Freebase

System	Pairwise				B^3		
	Prec.	Rec.	F-0.5	MCC	Prec.	Rec.	F-0.5
Rel-LDA/300	0.593	0.077	0.254	0.191	0.558	0.183	0.396
Rel-LDA/1000	0.638	0.061	0.220	0.177	0.626	0.160	0.396
HAC	0.567	0.152	0.367	0.261	0.523	0.248	0.428
Local	0.625	0.136	0.364	0.264	0.626	0.225	0.462
Local+Type	0.718	0.115	0.350	0.265	0.704	0.201	0.469
Our Approach	0.736	0.156	0.422	0.314	0.677	0.233	0.490
Our Approach+Type	0.682	0.110	0.334	0.250	0.687	0.199	0.460

Automatic evaluation against Freebase

HAC = hierarchical agglomerative clustering alone

(i.e. no sense disambiguation — most similar to DIRT)

Sense clustering adds 17% to precision!