



Relation Extraction

Semi-supervised and unsupervised relation extraction



Seed-based or bootstrapping approaches to relation extraction

- No training set? Maybe you have:
 - A few seed tuples or
 - A few high-precision patterns
- Can you use those seeds to do something useful?
 - Bootstrapping: use the seeds to directly learn to populate a relation



Relation Bootstrapping (Hearst 1992)

- Gather a set of seed pairs that have relation R
- Iterate:
 1. Find sentences with these pairs
 2. Look at the context between or around the pair and generalize the context to create patterns
 3. Use the patterns for grep for more pairs



Bootstrapping

- <Mark Twain, Elmira> **Seed tuple**
 - Grep (google) for the environments of the seed tuple

"Mark Twain is buried in Elmira, NY."

X is buried in Y

"The grave of Mark Twain is in Elmira"

The grave of X is in Y

"Elmira is Mark Twain's final resting place"

Y is X's final resting place.
- Use those patterns to grep for new tuples
- Iterate



Dipre: Extract <author,book> pairs

Brin, Sergei. 1998. Extracting Patterns and Relations from the World Wide Web.

- Start with 5 seeds:

Author	Book
Isaac Asimov	The Robots of Dawn
David Brin	Startide Rising
James Gleick	Chaos: Making a New Science
Charles Dickens	Great Expectations
William Shakespeare	The Comedy of Errors

- Find Instances:

The Comedy of Errors, by William Shakespeare, was

The Comedy of Errors, by William Shakespeare, is

The Comedy of Errors, one of William Shakespeare's earliest attempts

The Comedy of Errors, one of William Shakespeare's most

- Extract patterns (group by middle, take longest common prefix/suffix)

?x , by ?y , ?x , one of ?y 's

- Now iterate, finding new seeds that match the pattern



Snowball

E. Agichtein and L. Gravano 2000. Snowball: Extracting Relations from Large Plain-Text Collections. ICDL

- Similar iterative algorithm

Organization	Location of Headquarters
Microsoft	Redmond
Exxon	Irving
IBM	Armonk

- Group instances w/similar prefix, middle, suffix, extract patterns
 - But require that X and Y be named entities
 - And compute a confidence for each pattern

.69 **ORGANIZATION** {'s, in, headquarters} **LOCATION**

.75 **LOCATION** {in, based} **ORGANIZATION**



Distant Supervision

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

Fei Wu and Daniel S. Weld. 2007. Autonomously Semantifying Wikipedia. CIKM 2007

Mintz, Bills, Snow, Jurafsky. 2009. Distant supervision for relation extraction without labeled data. ACL09

- Combine bootstrapping with supervised learning
 - Instead of 5 seeds,
 - Use a large database to get huge # of seed examples
 - Create lots of features from all these examples
 - Combine in a supervised classifier



Distant supervision paradigm

- Like supervised classification:
 - Uses a classifier with lots of features
 - Supervised by detailed hand-created knowledge
 - Doesn't require iteratively expanding patterns
- Like unsupervised classification:
 - Uses very large amounts of unlabeled data
 - Not sensitive to genre issues in training corpus



Distantly supervised learning of relation extraction patterns

- ① For each relation
- ② For each tuple in big database
- ③ Find sentences in large corpus with both entities
- ④ Extract frequent features (parse, words, etc)
- ⑤ Train supervised classifier using thousands of patterns

Born-In

<Edwin Hubble, Marshfield>

<Albert Einstein, Ulm>

Hubble was born in Marshfield

Einstein, born (1879), Ulm

Hubble's birthplace in Marshfield

PER was born in LOC

PER, born (XXXX), LOC

PER's birthplace in LOC

$P(\text{born-in} \mid f_1, f_2, f_3, \dots, f_{70000})$



Unsupervised relation extraction

M. Banko, M. Cararella, S. Soderland, M. Broadhead, and O. Etzioni.
2007. Open information extraction from the web. IJCAI

- Open Information Extraction:
 - extract relations from the web with no training data, no list of relations
- 1. Use parsed data to train a “trustworthy tuple” classifier
- 2. Single-pass extract all relations between NPs, keep if trustworthy
- 3. Assessor ranks relations based on text redundancy
 - (FCI, specializes in, software development)
 - (Tesla, invented, coil transformer)



Evaluation of Semi-supervised and Unsupervised Relation Extraction

- Since it extracts totally new relations from the web
 - There is no gold set of correct instances of relations!
 - Can't compute precision (don't know which ones are correct)
 - Can't compute recall (don't know which ones were missed)
 - Instead, we can approximate precision (only)
 - Draw a random sample of relations from output, check precision manually
- $$\hat{P} = \frac{\text{\# of correctly extracted relations in the sample}}{\text{Total \# of extracted relations in the sample}}$$
- Can also compute precision at different levels of recall.
 - Precision for top 1000 new relations, top 10,000 new relations, top 100,000
 - In each case taking a random sample of that set
- 49 But no way to evaluate recall



Relation Extraction

Semi-supervised and unsupervised relation extraction