

Text Mining 3 - Relation Extraction

March 2018

Plan

- Relation Extraction
 - Introduction
 - Relation types
- Relation Extraction Methods
 - Hand-built patterns
 - Bootstrapping methods
 - Supervised learning methods
 - Distant supervision

Relation Extraction - Example

CHICAGO (AP) — Citing high fuel prices, United Airlines said Friday it has increased fares by \$6 per round trip on flights to some cities also served by lower-cost carriers. American Airlines, a unit of AMR, immediately matched the move, spokesman Tim Wagner said. United, a unit of UAL, said the increase took effect Thursday night and applies to most routes where it competes against discount carriers, such as Chicago to Dallas and Atlanta and Denver to San Francisco, Los Angeles and New York.

Question: What relations should we extract?

Relation Extraction - Example

CHICAGO (AP) — Citing high fuel prices, United Airlines said Friday it has increased fares by \$6 per round trip on flights to some cities also served by lower-cost carriers. **American Airlines**, a **unit of AMR**, immediately matched the move, spokesman **Tim Wagner** said. **United**, a **unit of UAL**, said the increase took effect Thursday night and applies to most routes where it competes against discount carriers, such as Chicago to Dallas and Atlanta and Denver to San Francisco, Los Angeles and New York.

Subject	Relation	Object
American Airlines	subsidiary	AMR
Tim Wagner	employee	American Airlines
United Airlines	subsidiary	UAL

slide adapted from Luke Zettlemoyer

Relation Types

For generic news texts ...

Relations	Examples	Types
Affiliations	Personal	<i>married to, mother of</i>
	Organizational	<i>spokesman for, president of</i>
	Artifactual	<i>owns, invented, produces</i>
Geospatial	Proximity	<i>near, on outskirts</i>
	Directional	<i>southeast of</i>
Part-Of	Organizational	<i>a unit of, parent of</i>
	Political	<i>annexed, acquired</i>

slide adapted from Jim Martin

Relation Types - ACE 2003

ROLE: relates a person to an organization or a geopolitical entity

subtypes: member, owner, affiliate, client, citizen

PART: generalized containment

subtypes: subsidiary, physical part-of, set membership

AT: permanent and transient locations

subtypes: located, based-in, residence

SOCIAL: social relations among persons

subtypes: parent, sibling, spouse, grandparent, associate

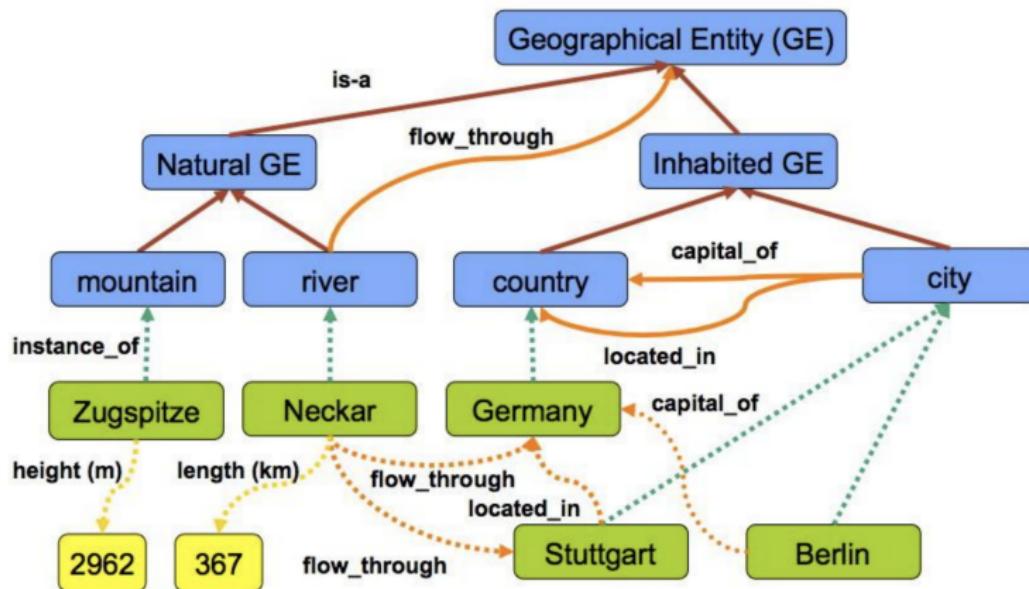
slide adapted from Doug Appelt

Relation Types - Freebase

23 Million Entities, thousands of 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

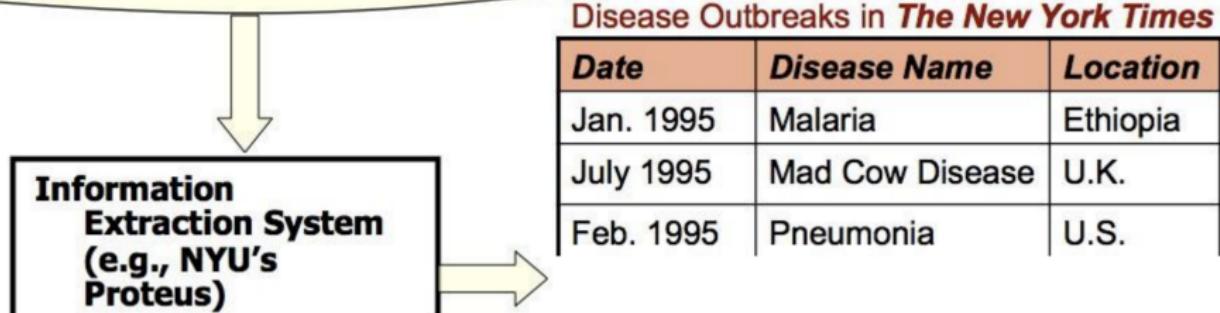
Relation Types - Geospatial



slide adapted from Paul Buitelaar

Relation Types - Disease Outbreaks

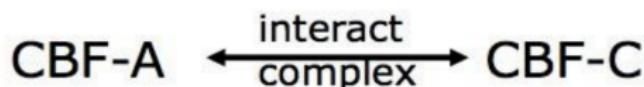
May 19 1995, Atlanta -- The Centers for Disease Control and Prevention, which is in the front line of the world's response to the deadly **Ebola** epidemic in **Zaire**, is finding itself hard pressed to cope with the crisis...



slide adapted from Eugene Agichtein

Relation Types - Protein Interactions

„We show that CBF-A and CBF-C interact with each other to form a CBF-A-CBF-C complex and that CBF-B does not interact with CBF-A or CBF-C individually but that it associates with the CBF-A-CBF-C complex.“



slide adapted from Rosario & Hearst

Relation Extraction - Need

- NLP applications need word meaning!
 - Question answering
 - Conversational agents
 - Summarization
- One key meaning component: word relations
 - Hyponymy: San Francisco is an instance of a city
 - Antonymy: acidic is the opposite of basic
 - Meronymy: an alternator is a part of a car

slide adapted from Luke Zettlemoyer

Relation Extraction - Need

WordNet is incomplete

Ontological relations are missing for many words:

In WordNet 3.1	Not in WordNet 3.1
insulin progesterone	leptin pregnenolone
combustibility navigability	affordability reusability
HTML	XML
Google, Yahoo	Microsoft, IBM

Esp. for specific domains: restaurants, auto parts, finance

slide adapted from Luke Zettlemoyer

Relation Extraction - Methods

- Hand-built patterns
- Bootstrapping methods
- Supervised learning methods
- Distant supervision method

Hand-built Patterns

What does *Gelidium* mean?

Hand-built Patterns

What does *Gelidium* mean?

What does *Agar is a substance prepared from a mixture of red algae, such as Gelidium, for laboratory or industrial use*

Hearst's Lexico-Syntactic Patterns

Y such as X ((, X)* (, and/or) X)

such Y as X...

X... or other Y

X... and other Y

Y including X...

Y, especially X...

Examples of the Hearst's Patterns

Hearst pattern	Example occurrences
X and other Y	...temples, treasuries, and other important civic buildings.
X or other Y	bruises, wounds, broken bones or other injuries...
Y such as X	The bow lute, such as the Bambara ndang...
such Y as X	...such authors as Herrick, Goldsmith, and Shakespeare.
Y including X	...common-law countries, including Canada and England...
Y, especially X	European countries, especially France, England, and Spain...

slide adapted from Luke Zettlemoyer

Problems with Hand-built Patterns

- Requires hand-building patterns for each relation!
 - hard to write; hard to maintain
 - there are zillions of them
 - domain-dependent
- Don't want to do this for all possible relations!
- Plus, we'd like better accuracy
 - Hearst: 66% accuracy on hyponym extraction

slide adapted from Luke Zettlemoyer

Relation Extraction - Methods

- Hand-built patterns
- Bootstrapping methods
- Supervised learning methods
- Distant supervision method

Bootstrapping Approach

- If you don't have enough annotated text to train on ...
- But you do have:
 - some **seed instances** of the relation
 - (or some patterns that work pretty well)
 - and lots & lots of **unannotated text** (e.g., the web)
- ... can you use those seeds to do something useful?
- Bootstrapping can be considered *semi-supervised*

Bootstrapping Example

Seed: (Arthur Conan Doyle, The Adventures of Sherlock Holmes)

A Web crawler finds all documents contain the pair.

Bootstrapping - Matched Document 1

...

Read **The Adventures of Sherlock Holmes** by **Arthur Conan Doyle**
online or in you email

...



Extract **tuple**:

[0, Arthur Conan Doyle, The Adventures of Sherlock Holmes,
Read, online or, by]

A tuple of 6 elements: [order, author, book, prefix, suffix, middle]

order = 1 if the author string occurs before the book string, = 0 otherwise

prefix and suffix are strings contain the 10 characters occurring to the left/right of the match

middle is the string occurring between the author and book

slide adapted from Nguyen Bach and Sameer Badaskar

Bootstrapping - Matched Document 2

...

know that Sir Arthur Conan Doyle wrote The Adventures of
Sherlock Holmes, in 1892

...



Extract tuple:

[1, Arthur Conan Doyle, The Adventures of Sherlock Holmes,
now that Sir, in 1892, wrote]

slide adapted from Nguyen Bach and Sameer Badaskar

Bootstrapping - Developing Patterns

Extracted list of tuples:

[0, Arthur Conan Doyle, The Adventures of Sherlock Holmes, Read, online or, by]

[1, Arthur Conan Doyle, The Adventures of Sherlock Holmes, now that Sir, in 1892, wrote]

[1, Arthur Conan Doyle, The Adventures of Sherlock Holmes, When Sir, in 1892 he, wrote]

...

Group tuples by matching *order* and *middle* and induce *patterns*

Induce patterns from group of tuples:

[longest-common-suffix of prefix strings, author, middle, book, longest-common-prefix of suffix strings]

Pattern:

[Sir, Arthur Conan Doyle, wrote, The Adventures of Sherlock Holmes, in 1892]

Pattern with wild card expression:

[Sir, .*?, wrote, .*?, in 1892]

slide adapted from Nguyen Bach and Sameer Badaskar

Bootstrapping - Search for more patterns and new tuples

Use the wild card patterns [Sir, .*?, wrote, .*?, in 1892]

search the Web to find more documents

...

Sir Arthur Conan Doyle wrote Speckled Band in 1892, that is around 62 years apart which would make the stories

...

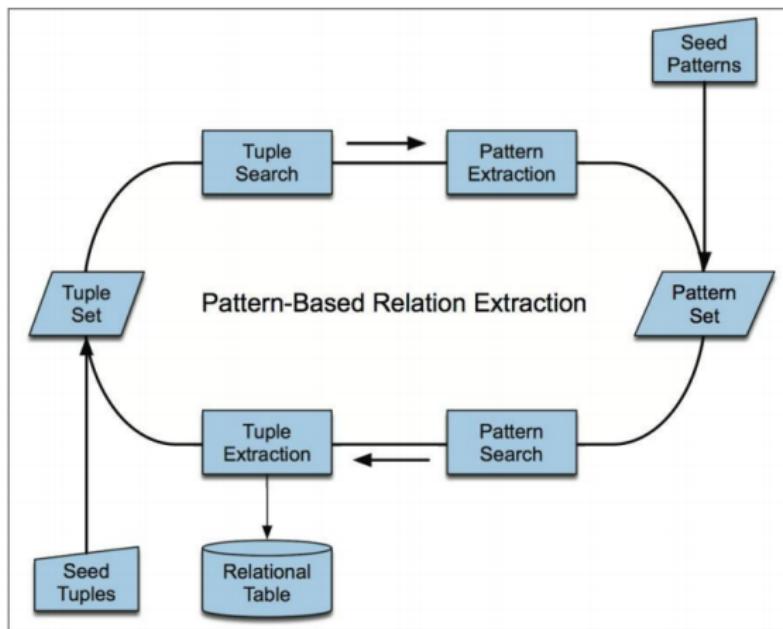
Extract new relations:

(Arthur Conan Doyle, Speckled Band)

Repeat the algorithm with the new relation.

slide adapted from Nguyen Bach and Sameer Badaskar

Bootstrapping System



slide adapted from Jim Martin

Bootstrapping - Problems

- Requires that we have seeds for each relation
 - Sensitive to original set of seeds
- Big problem of semantic drift at each iteration
- Precision tends to be not that high
- Generally have lots of parameters to be tuned
- No probabilistic interpretation
 - Hard to know how confident to be in each result

Relation Extraction - Methods

- Hand-built patterns
- Bootstrapping methods
- Supervised learning methods
- Distant supervision method

Supervised Relation Extraction - Approach

The supervised approach requires:

- Defining an inventory of output labels
 - Relation detection: true/false
 - Relation classification: located-in, employee-of, inventor-of, ...
- Collecting labeled training data: MUC, ACE, ...
- Defining a feature representation: words, entity types, ...
- Choosing a classifier: Naïve Bayes, MaxEnt, SVM,
...
• Evaluating the results

Word Features for Relation Extraction

American Airlines, a unit of AMR, immediately matched the move, spokesman **Tim Wagner** said
Mention 1
Mention 2

- Named-entity types
 - M1: ORG
 - M2: PERSON
 - Concatenation of the two named-entity types
 - ORG-PERSON
 - Entity Level of M1 and M2 (NAME, NOMINAL, PRONOUN)
 - M1: NAME [it or he would be PRONOUN]
 - M2: NAME [the company would be NOMINAL]

Parse Features for Relation Extraction

American Airlines, a unit of AMR, immediately matched the move, spokesman **Tim Wagner** said
Mention 1
Mention 2

- Base syntactic chunk sequence from one to the other
NP NP PP VP NP NP
 - Constituent path through the tree from one to the other
NP ↑ NP ↑ S ↑ S ↓ NP
 - Dependency path

Airlines matched Wagner said

Classifiers for Supervised Learning Methods

- Now you can use any classifier you like
 - MaxEnt
 - Naïve Bayes
 - SVM
 - ...
- Train it on the training set, tune on the dev set, test on the test set

Supervised Learning Methods - Evaluation

Compute P/R/F₁ for each relation

$$P = \frac{\text{\# of correctly extracted relations}}{\text{Total \# of extracted relations}}$$

$$R = \frac{\text{\# of correctly extracted relations}}{\text{Total \# of gold relations}}$$

$$F_1 = \frac{2PR}{P+R}$$

Supervised Learning Methods - Summary

- + Can get high accuracies with enough hand-labeled training data, if test similar enough to training
 - Labeling a large training set is expensive
 - Supervised models are brittle, don't generalize well to different genres

Relation Extraction - Methods

- Hand-built patterns
- Bootstrapping methods
- Supervised learning methods
- Distant supervision method

Distant supervision method

- 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

Distant supervision - approach

- 1 For each relation Born-In
- 2 For each tuple in big database
`<Edwin Hubble, Marshfield>`
`<Albert Einstein, Ulm>`
- 3 Find sentences in large corpus
with both entities
`Hubble was born in Marshfield`
`Einstein, born (1879), Ulm`
`Hubble's birthplace in Marshfield`
- 4 Extract frequent features
(parse, words, etc)
`PER was born in LOC`
`PER, born (XXXX), LOC`
`PER's birthplace in LOC`
- 5 Train supervised classifier using
thousands of patterns
 $P(\text{born-in} \mid f_1, f_2, f_3, \dots, f_{70000})$

slide adapted from Dan Jurafsky

Distant supervision - approach

- 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