Mining Meaning From Wikipedia – part 2

PD Dr. Günter Neumann LT-lab, DFKI, Saarbrücken

Outline

- 1. Namend Entity Disambiguation
- 2. Relation Extraction
- 3. Ontology Building and the Semantic Web

Disambiguation of named entities

- Wikipedia is recognized as the largest available resource of named entities!
- Approaches considered here:
 - Linking named entities appearing in text or in search queries to corresponding Wikipedia articles for disambiguation.

Bunescu & Pasca, 2006

- Disambiguate NE in search queries in order to group search results by the corresponding senses.
- Core idea: Use dictionary for identifying possible NE meanings, that can also be used for disambiguation.
- Use Wikipedia as a resource for creating such a dictionary:
 - A dictionary of 500.000 NE that appear in Wikipedia
 - Note that NE are mainly mentioned as titles in articles, and hence some preprocessing of articles are required, e.g., capitalization, multi words
 - add redirects and disambiguated names to each entry
 - Capture the many-to-many relationship between names and entities

Examples

TITLE	REDIRECT	DISAMBIG	CATEGORIES
			Star Wars music,
John Williams (composer)	John Towner Williams	John Williams	Film score composers,
			20th century classical composers
John Williams (wrestler)	Ian Rotten	John Williams	Professional wrestlers,
			People living in Baltimore
John Williams (VC)	none	John Williams	British Army soldiers,
			British Victoria Cross recipients
Boston Pops Orchestra	Boston Pops,	Pops	American orchestras,
	The Boston Pops Orchestra		Massachusetts musicians
United States	US, USA,	US, USA,	North American countries,
	United States of America	United States	Republics, United States
		Venus,	Venus
Venus (planet)	Planet Venus	Morning Star,	Planets of the Solar System,
_		Evening Star	Planets, Solar System,

Table 1: Examples of Wikipedia titles, aliases and categories

Approach

- If a query contains a term that corresponds to 2 or more entries, chose the one whose Wikipedia article has greatest cosine similarity with the query.
- If similarity value is too low, use category to which the article belongs.
- Result: 55%-85% accuracies for Wikipedia's People by occupation category.

Cucerzan, 2007

- Disambiguate NE in free text
- Also uses a compiled dictionary from Wikipedia with two parts
 - Surface forms, s.a. Article titles, redirects, ...
 - Associated entities together with contextual information about them
 - ~1.4M entries, average of 2.4 surface forms each

Further extracted data structures

- «NE,tag» entries from Wikipedia list articles
 - Texas (band) \rightarrow LIST_band name etymologies
 - Because it appears in a list with this title
 - 540.000 entries
- Categories assigned to Wikipedia articles decribing NE used as tag too
 - 2.65M entries
- A context for each NE is collected from its article
 - 38M entries <NE, context>

Identification of NE in text

- Capitalization rules indicate which phrases are surface forms of NE
- Co-occurrence statistics from web by a search engine are used to identify boundaries (Google n-grams)
 - "Whitney Museum of American Art" vs. "Whitney Museum in New York" (about 2.450.000 vs. 405.000 hits)
- Lexical analysis to collect identical NE (Mr. Brown & Brown)
- Disambiguation on basis of similarity of
 - the document in which the surface form appears
 - with Wikipedia articles that represent all NE that have been identified in it, and their context terms
- Result:
 - 88% accuracy/5000 Wikipedia entities
 - 91% accuracy on 750 news article entities

Disambiguating thesaurus & ontology terms

cardiovascular cardiovascular s. = circulatory s. system Then add blood circulation to Agrovoc USE BT NT 02 circulatory system of information as heart domain-specific circulatory thesaurus. redirects system Thus it can also be used blood belongs-to circulation to extend and improve cardiovascular

resources.

blood Category Article Article Redirect blood system cross-link heart

Descriptor

Non-Descriptor

RT

Figure 4. Comparison of organization structure in Agrovoc and Wikipedia.

Establishing Mappings between Wikipedia and other Resources

- Mapping Wikipedia articles to WordNet
- Idea, cf. Ruiz-Casado et al. 2005
 - If a Wikipedia article matches several WordNet synsets, the appropriate one is chosen by computing the similarity between the Wikipedia entry word-bag and the WordNet synset gloss.
- Gets 84% accuracy using dot product on stemmed texts.
- Problem: if (Simple) Wikipedia is growing
 - Ambiguities increase as well
 - Mapping gaps
 - We have observed similar problem when using WordNet for crosslingual QA!

Intermediate summary

y 1st 2011...

http://mostpopularwebsites.net/

Today's Most Popular Websites on the Internet:

		Ads by Google	<u>2010</u>	WWW Google com	<u>31 12 2010</u>	Best Bar 2010
The fo	ollow	ving list of the Most P	opular V	Vebsites was upda	ated on Satu	rday, Januar
		Carala and				
1.		Google.com		www.go	ogie.com	
2.	Tev	Pacebook.com		www.rad	ebook.com	
з. 4	labe	Yohan ees		www.yo	utube.com	
4.		tanoo.com		www.ya	noo.com	
5.		Live.com		www.iiv	e.com	
ь. ¬	20	Baldu.com		www.ba	iau.com	
<i>/</i> .	W	Wikipedia.org		www.wi	kipedia.org	
o. 0	0	Diogspot.com		www.bic	gspot.com	
9.		Qq.com		www.qq	.com	
10.	~	Twitter.com		www.tw	itter.com	
11.	Θ	Blogger.com		www.bic	ogger.com	
12.	- 27	Yahoo.co.jp		www.ya	noo.co.jp	
13.		Amazon.com		www.an	nazon.com	
14.	193	laobao.com		www.tao	obao.com	
15.	1	Google.co.in		www.go	ogle.co.in	
16.	0	Sina.com.cn		www.sin	a.com.cn	
17.	×.	Google.de		www.go	ogle.de	
18.	×.	Google.com.hk		www.go	ogle.com.hk	
19.	av.	Wordpress.com		www.wo	ordpress.com	
20.	in	Linkedin.com		www.lin	kedin.com	
21.	e \$X	Ebay.com		www.eb	ay.com	
22.	×	Google.co.uk		www.go	ogle.co.uk	
23.		Microsoft.com		www.mi	crosoft.com	
24.		Yandex.ru		www.ya	ndex.ru	
25.	×	Google.fr		www.go	ogle.fr	
26.	Þ	Bing.com		www.bir	ig.com	
27.	×	Google.co.jp		www.go	ogle.co.jp	
28.	8	163.com		www.16	3.com	
29.	8	Google.com.br		www.go	ogle.com.br	
30.	2	Fc2.com		www.fc2	2.com	
31.	\odot	Craigslist.com		www.cra	igslist.com	
32.	0	Conduit.com		www.co	nduit.com	
33.	8	Google.it		www.go	ogle.it	
34.	••	Flickr.com		www.flic	kr.com	
35.	В	Vkontakte.ru		www.vk	ontakte.ru	
36.	Ð	Charlotte.craigslist.or	rg	www.ch	arlotte.craigs	list.org
37.	3	Apple.com		www.ap	ple.com	
38.		Googleusercontent.co	om	www.go	ogleusercont	ent.com
39.	Ð	Craigslist.org		www.cra	igslist.org	
40.	8	Google.es		www.go	ogle.es	
		-				

Benefit to Wikipedia: Tools

- Internal link maintenance
- Infobox Creation
- Schema Management
- Reference suggestion & fact checking
- Disambiguation page maintenance
- Translation across languages
- Vandalism Alerts

Motivating Vision Next-Generation Search = Information Extraction + Ontology + Inference

. . .

Which German Scientists Taught at US Universities?



Next-Generation Search

Information Extraction

- <Einstein, Born-In, Germany>
- <Einstein, ISA, Physicist>
- <Einstein, Lectured-At, IAS>
- <IAS, In, New-Jersey>
- <New-Jersey, In, United-States>
- Ontology
 - Physicist (x) →
 Scientist(x)

••

...

- Inference
 - Einstein = Einstein

. . .

Set Means to the was born phy

Albert Einstein was a Ge<u>Mozilla Download - Mozilla Firefox</u> born the physicis New Jersey is a state in the Northeastern region of the United States

...

Mozilla Download - Mozilla Firefox

Edit

<u>V</u>iew Hi<u>s</u>tory <u>B</u>ookmarks <u>T</u>ools <u>H</u>elp

Extracting more Structure

Microsoft

From Wikipedia, the free encyclopedia

Microsoft Corporation is an American public multinational corporation headquartered in Redmond, Washington, USA that develops, manufactures, licenses, and supports a wide range of products and services predominantly related to computing through its various product divisions. Established on April 4, 1975 to develop and sell BASIC interpreters for the Altair 1880. Microsoft rese to dominate the home computer operating system (QS) market with MS-DOS in the mid-1980s, followed by the Microsoft Windows line of OSs. Microsoft would also come to dominate the office suite market with Microsoft Office. The company has diversified in recent years into the video game industry with the Xbox and its successor, the Xbox 360 as well as into the consumer electronics market with Zune and the Windows Phone OS. The ensuing rise of stock in the company's 1986 initial public offering (IPO) made an estimated four billionaires and 12,000 millionaires form Microsoft empress.

Primarily in the 1990s, critics contend Microsoft used monopolistic business practices and anti-competitive strategies including refusal to deal and tying, put unreasonable restrictions in the use of its software, and used misrepresentative marketing tactics; both the U.S. Department of Justice and European Commission found the company in violation of antitrust laws. Known for its interviewing process with obscure questions, various studies and ratings were generally favorable to Microsoft's diversity within the company as well as its overall environmental impact with the exception of the electronics protion of the business.

C	ontents [hide]	
1 History		
1.1 1984-1994: Windows and Ol	ffice	
1.2 1995-2005: Internet and the	32-bit era	
1.3 2006 on: Vista and Cloud cor	mputing	
2 Product divisions		
2.1 Windows & Windows Live Di	vision, Server and Tools, Online Services Division	
2.2 Business Division, Entertainn	nent and Devices Division	
3 Culture		
4 Corporate affairs		
4.1 Environment		
4.2 Marketing		
5 See also		
6 References		
7 External links		
History		
Main articles: History of Microso	ft and History of Microsoft Windows	
6	Paul Allen and Bill Gates, childhood frien programming, were seeking to make a su	ds with a passion in computer accessful business utilizing the
	programming, were seeking to make a su	ccessful business u

Instrumentation and Telemetry Systems's (MITS) Altair 8800 microcomputer.

Dow Jones Industrial Average Component S&P 500 Component Industry Computer software Consumer electronics Digital distribution Computer hardware Video games IT consulting Online advertising Retail stores Automotive software Founded Albuquerque, New Mexico April 4, 1975 **Bill Gates** Founder(s) Paul Allen Headquarters One Microsoft Way Redmond, Washington, United States Worldwide Area served Key people Steve Ballmer (CEO) Brian Kevin Turner (COO) Bill Gates (Chairman) Ray Ozzie (CSA) Craig Mundie (CRSO) Products See products listing Services See services listing \$62,484 billion (2010) Revenue \$24.098 billion (2010) Operating income

\$18.760 billion (2010)

\$86.113 billion (2010)
 \$46.175 billion (2010)

89,000 (2010)

Profit

Total assets

Total equity

Employees



Coordinates: 🙆 47°38'22.55"N 122°7'42.42"W

Microsoft Corporation



Relation Extraction

- Basically means:
 - extract semantic relations from unstructured text
- Example:
 - Apple Inc.'s world corporate headquarters are located in the middle of Silicon Valley, at 1 Infinite Loop, Cupertino, California.
 - hasHeadquarters(Apple Inc., 1 Infinite Loop-Cupertino-California)
- Challenges:
 - Extract this relation from sentences expressing the same information about Apple Inc., regardless of the actual wording.
 - Same relation should be determined with different arguments from similar sentences
 - hasHeadquarters(Google Inc., Google Campus-Mountain View-California)

Wikipedia: Relation Extraction

- Semantic relations in Wikipedia's raw text
- Semantic relations in structured parts of Wikipedia
- Typing Wikipedia's named entities

Semantic relations in Wikipedia's raw text

- Standard approach:
 - Take known binary relations as seeds
 - Extract patterns from their textual representation, e.g.,
 - X's * headquarters are located in * at Y
 - Patterns are then applied to identify new relations as values from X and Y
- Known difficulties of this approach:
 - Enumerating over all pairs of entities yields a low density of correct relations even when restricted to a single sentence
 - Errors in the entity recognition stage create inaccuracies in relation classification.

Anchestors of Wikipedia-based relation extraction from text

- Brin, S.: Extracting patterns and relations from the World Wide Web, 1998
- Agichtein, E., Gravano, L.: Snowball: Extracting relations from large plain-text collections, 2002.
- Ravichandran, D., Hovy, E.: Learning surface text patterns for a question answering system, 2002.
- There is also close relationship to Machine Reading

Snowball : Extracting Relations from Large Plain-Text Collections

Eugene Agichtein Luis Gravano

Department of Computer Science Columbia University, ACM paper 2000

Example Task: Organization/Location

Redundancy	Organization	Location	
<i>Microsoft</i> 's central headquarters in <i>Redmond</i> is home to almost every product group and division.	Microsoft	Redmond	
Brent Barlow, 27, a software analyst and beta-tester at <i>Apple Computer</i> headquarters	Apple Computer	Cupertino	
in <i>Cupertino</i> , was fired Monday for "thinking a little too different."	Nike	Portland	
		AND SHI NESSON NAMESAN	
Apple's programmers "think different" on a "campus" in			
Cupertino, Cal. Nike emp	<i>Cupertino, Cal. Nike</i> employees "just do it" at what the		
company refers to as its "W	company refers to as its "World Campus," near <i>Portland</i> ,		
Ore.			

Related Work

- Bootstrapping
 - Riloff et al. ('99), Collins & Singer ('99)
 - (Named-entity recognition see previous lectures of this course,)
 - Brin (DIPRE) ('98)

Initial Seed Tuples:	ORGANIZATION	LOCATION
	MICROSOFT	REDMOND
	IBM	ARMONK
	BOEING	SEATTLE
	INTEL	SANTA CLARA







•<*STRING1*>'s headquarters in <*STRING2*>

DIPRE Patterns:

•<STRING2> -based <STRING1>

•<STRING1> , <STRING2>



	ORGANIZATION	LOCATION
Generate new seed	AG EDWARDS	ST LUIS
	157TH STREET	MANHATTAN
tuples;	7TH LEVEL	RICHARDSON
start new iteration	3COM CORP	SANTA CLARA
	3DO	REDWOOD CITY
	JELLIES	APPLE
	MACWEEK	SAN FRANCISCO



Extracting Relations from Text: Potential Pitfalls

- Invalid tuples generated
 - Degrade quality of tuples on subsequent iterations
 - Must have automatic way to select high quality tuples to use as new seed
- Pattern representation
 - Patterns must generalize

Extracting Relations from Text Collections

- Related WorkDIPRE
- The Snowball System:
 - Pattern representation and generation
 - Tuple generation
 - Automatic pattern and tuple evaluation
- Evaluation Metrics
- Experimental Results

Extracting Relations from Text: *Snowball*

Initial Seed Tuples:	ORGANIZATION	LOCATION
	MICROSOFT	REDMOND
	IBM	ARMONK
	BOEING	SEATTLE
	INTEL	SANTA CLARA





Problem: Patterns Excessively General

Pattern: <*STRING2*>-based <*STRING1*>

Today's merger with McDonnell Douglas positions **Seattle** <u>-based</u> **Boeing** to make major money in space.



Extracting Relations from Text: Snowhall

Tag Entities

Use MITRE's Alembic Named Entity tagger

Augment Table



Extracting Relations from Text



•<LOCATION> -based <ORGANIZATION>

•<*ORGANIZATION*> , <*LOCATION*>

PROBLEM: PATTERNS TOO SPECIFIC:

HAVE TO MATCH TEXT *exactly*.



Snowball: Pattern Representation

A Snowball pattern vector is a 5-tuple <*left, tag1, middle, tag2, right*

- *tag1*, *tag2* are named-entity tags
- *left*, *middle*, and *right* are vectors of

The weight of a term in each vector is a function of the frequency of the term in the corresponding context. These vectors are scaled so their norm is one. FInally, they are mulitplied by a scaling factor to indicate each vector's relative importance. From experiments, terms in the middle are more important, and hence get higher weight

ORGANIZATION 's central headquarter hence get higher weight.



Snowball: Pattern Generation

Tagged Occurrences of seed tuples:



In mid-afternoon trading, share of **Redmond**-based **Microsoft** fell...





Snowball Pattern Generation: **Cluster Similar Occurrences**

Occurrences of seed tuples converted to *Snowball* representation:



Similarity Metric

- $P = \langle Lp, tag |, Mp, tag |, Rp \rangle$
- $S = \langle Ls, tag \rangle, Ms, tag \rangle, Rs \rangle$

$Match(P, S) = \int_{0}^{P} \cdot [LS] + [Mp] \cdot [MS] + [Rp] \cdot [RS]$ if the tags match otherwise

Snowball Pattern Generation: Clustering

Cluster 1				
<pre>{<servers 0.75=""> <at 0.75="">} </at></servers></pre> <pre> ORGANIZATION { 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0</pre>		<pre>{<'s 0.5> <central 0.5=""> <headquarters 0.5=""> <in 0.5="">}</in></headquarters></central></pre>	LOCATION	
{ <operate 0.75=""> <from 0.75="">}</from></operate>	ORGANIZATION	<pre>{<'s 0.7> <headquarters 0.7=""> <in 0.7="">}</in></headquarters></pre>	LOCATION	



Snowball: Pattern Generation

Patterns are formed as *centroids* of the clusters. Filtered by minimum number of supporting tuples.



Snowball: Tuple Extraction

Using the patterns, scan the collection to generate new seed tuples:



Snowball: Tuple Extraction

Represent each new text segment in the collection as the context 5-tuple:



Find most similar pattern (if any)

	ORGANIZATION	{<'s 0.7>, <headquarters 0.7="">, < in 0.7>}</headquarters>	LOCATION	
--	--------------	--	----------	--

Snowball: Automatic Pattern Evaluation

Seed tuples



Snowball: Automatic Tuple Evaluation

Brent Barlow, 27, a software analyst and beta-tester at *Apple Computer* headquarters in *Cupertino*, was fired Monday for "thinking a little too different."

Apple's programmers "think different" on a "campus" in Cupertino, Cal.

<Apple Computer, Cupertino>

$Conf(Tuple) = 1 - (1 - Conf(P_i))$

- Estimation of Probability (Correct (Tuple))
- A tuple will have high confidence if generated by multiple high-confidence patterns (P_i).

Similar to Yangarber et al (cf. Session number 6).

Snowball: Filtering Seed Tuples

Generate new seed tuples:

ORGANIZATION	LOCATION	CONF
AG EDWARDS	ST LUIS	0.93
AIR CANADA	MONTREAL	0.89
7TH LEVEL	RICHARDSON	0.88
3COM CORP	SANTA CLARA	0.8
3DO	REDWOOD CITY	0.8
3M	MINNEAPOLIS	0.8
MACWORLD	SAN FRANCISCO	0.7
157TH STREET	MANHATTAN	0.52
15TH CENTURY EUROPE	NAPOLEON	0.3
15TH PARTY CONGRESS	CHINA	0.3
MAD	SMITH	0.3



Extracting Relations from Text Collections

- Related Work
- The Snowball System:
 - Pattern representation and generation
 - Tuple generation
 - Automatic pattern and tuple evaluation
- Evaluation Metrics
- Experimental Results

Task Evaluation Methodology

- Data: Large collection, extracted tables contain many tuples (> 80,000)
- Need scalable methodology:
 - Ideal set of tuples
 - Automatic recall/precision estimation
- Estimated precision using sampling

Collections used in Experiments

More than 300,000 real newspaper articles

Collection	Source	Year
	The New York Times	1996
Training	The Wall Street Journal	1996
	The Los Angeles Times	1996
	The New York Times	1995
Test	The Wall Street Journal	1995
	The Los Angeles Times	1995,'97

The *Ideal* Metric (1)

Creating the Ideal set of tuples



* A perfect, (*ideal*) system would be able to extract all these tuples



Precision:

| Correct (*Extracted* \cap *Ideal*) | Extracted ∩ Ideal

Recall:

Correct (*Extracted* \cap *Ideal*) Ideal

Estimate Precision by Sampling

Sample extracted table
 Random samples, each 100 tuples
 Manually check validity of tuples in each sample

Extracting Relations from Text Collections

- Related Work
- The Snowball System:
 - Pattern representation and generation
 - Tuple generation
 - Automatic pattern and tuple validation
- Evaluation Metrics
- Experimental Results

Experimental results: Test Collection



(a)

(b)

Recall (a) and precision (a) using the *Ideal* metric, plotted against the minimal number of occurrences of test tuples in the collection

Experimental results: Sample and Check



Recall (a) and precision (b) for varying minimum confidence threshold T_t . **NOTE**: Recall is estimated using the *Ideal* metric, precision is estimated by **manually checking random samples** of result table.

Conclusions

- Snowball system:
 - Requires minimal training (handful of seed tuples)
 - Uses a flexible pattern representation
 - Achieves high recall/precision

> 80% of test tuples extracted

Scalable evaluation methodology