

Information Extraction & Named Entity Recognition



Christopher Manning
CS224N



NLP for IR/web search?

- It's a no-brainer that NLP should be useful and used for web search (and IR in general):
 - Search for 'Jaguar'
 - the computer should know or ask whether you're interested in big cats [scarce on the web], cars, or, perhaps a molecule geometry and solvation energy package, or a package for fast network I/O in Java
 - Search for 'Michael Jordan'
 - The basketballer or the machine learning guy?
 - Search for laptop, don't find notebook
 - [Google used to not even *stem*:
 - Searching *probabilistic model* didn't even match pages with *probabilistic models* - but it does now.]



NLP for IR/web search?

- Word sense disambiguation technology generally works well (like text categorization)
- Synonyms can be found or listed
- Lots of people were into fixing this
 - Especially around 1999-2000
 - Lots of (ex-)startups:
 - LingoMotors
 - iPhrase "Traditional keyword search technology is hopelessly outdated"



NLP for IR/web search?

- But in practice it's an idea that hasn't gotten much traction
 - Correctly finding linguistic base forms is straightforward, but produces little advantage over crude stemming which just slightly over equivalence classes words
 - Word sense disambiguation only helps on average in IR if over 90% accurate (Sanderson 1994), and that's about/above where we are
 - Syntactic phrases should help, but people have been able to get most of the mileage with "statistical phrases" - which have been aggressively integrated into systems recently (covert phrases; proximity weighting)



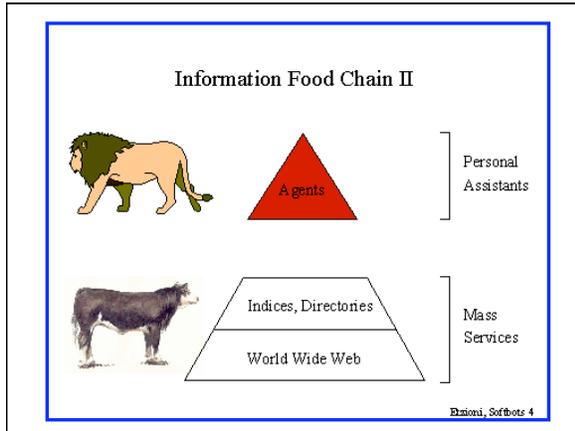
NLP for IR/web search?

- Much more progress has been made in link analysis, the use of anchor text, etc.
- Anchor text gives human-provided synonyms
- Using human intelligence always beats artificial intelligence
- People can easily scan among results (on their 24" monitor) ... if you're above the fold
- Link or click stream analysis gives a form of pragmatics: what do people find correct or important (in a default context)
- Focus on short, popular queries, news, etc.



NLP for IR/web search?

- Methods which use rich ontologies, etc., can work very well for intranet search within a customer's site (where anchor-text, link, and click patterns are much less relevant)
- But don't really scale to the whole web
- *Moral: it's hard to beat keyword search for the task of general ad hoc document retrieval*
- *Conclusion: one should move up the food chain to tasks where finer-grained understanding of meaning is needed*
- One possibility: information extraction



Product information/ Comparison shopping, etc.

- Need to learn to extract info from online vendors
- Can exploit uniformity of layout, and (partial) knowledge of domain by querying with known products
- Early e.g., Jango Shopbot (Etzioni and Weld)
 - Gives convenient aggregation of online content
- Bug: originally not popular with vendors
 - Make personal agents rather than web services?
- This seems to have changed (e.g., Froogle)

Results 1 - 10 of about 3 confirmed 16,830 total results for **lego fire engine**. (0.78 seconds)

Gold Planet Micro Urban Rescue HW Set Hot Wheels
\$11.78 - Yahoo! Auctions - Buy it now & shipping

LEGO RESCUE FIRE ENGINE HOT WHEELS HW PLANET MICRO LIMITED EDITION GOLD - RELINQUISHED EXCOVA, PS2 PLAYSTATION 2, MICROGLOBE, MICRO LEGO, MICRO MICRO...

LEGO Community Transport Set - 50 Pieces - SmartSet4ids
\$24.99 - Buy.com - Buy it now & shipping

LEGO Vehicles Set
\$76.31 - Toysrus.com - Buy it now & shipping

Passion Fire engine 0674-2
\$125.00 - www.pastich.com

Radio Flyer Red Fire Engine
\$24.99 - www.sagehill.com

Commercial information...

A book, Not a toy (green arrow pointing to the book title)

Title (green arrow pointing to the book title)

Need this price (yellow arrow pointing to the price)

Product Code: 0960781200
 USA/Canada: US\$ 43.40
 Australia/NZ: A\$ 124.50
 Other Countries: US\$ 50.00

Information Extraction

- Information extraction systems
 - Find and understand the limited relevant parts of texts
 - Clear, factual information (*who did what to whom when?*)
 - Produce a structured representation of the relevant information: *relations* (in the DB sense)
 - Combine knowledge about language and a domain
 - Automatically extract the desired information
- E.g.
 - Gathering earnings, profits, board members, etc. from company reports
 - Learn drug-gene product interactions from medical research literature
 - "Smart Tags" (Microsoft) inside documents

Classified Advertisements (Real Estate)

- Background:
 - Advertisements are plain text
 - Lowest common denominator: only thing that 70+ newspapers with 20+ publishing systems can all handle

```

<ADNUM>2067206v1</ADNUM>
<DATE>March 02, 1998</DATE>
<ADTITLE>MADDINGTON $89,000</ADTITLE>
<ADTEXT>
OPEN 1.00 - 1.45<BR>
U 11 / 10 BERTRAM ST<BR>
NEW TO MARKET Beautiful<BR>
3 brm freestanding<BR>
villa, close to shops & bus<BR>
Owner moved to Melbourne<BR>
ideally suit 1st home buyer,<BR>
investor & 55 and over.<BR>
Brian Hazelden 0418 958 996<BR>
R WHITE LEEING 9332 3477
  
```

The screenshot shows a real estate website interface. On the left, there's a navigation menu with options like 'Home', 'Real Estate', 'Property', and 'Business'. The main content area features a map of a residential area with a blue circle highlighting a specific location. Below the map, there's a text box with the address '10 BERTRAM ST, Maddington, WA'. The website header includes the logo for 'S NLP' and navigation links like 'subscribe', 'search', 'feedback', and 'help'.

Why doesn't text search (IR) work?

What you search for in real estate advertisements:

- **Town/suburb.** You might think easy, but:
 - **Real estate agents:** Coldwell Banker, Mosman
 - **Phrases:** Only 45 minutes from Parramatta
 - **Multiple property ads have different suburbs**
- **Money:** want a range not a textual match
 - **Multiple amounts:** was \$155K, now \$145K
 - **Variations:** offers in the high 700s [but not rents for \$270]
- **Bedrooms:** similar issues (br, bdr, beds, B/R)

Canonicalization: Product information

The screenshot shows a search results page on CNET Reviews for 'IBM ThinkPad X31'. The page displays a list of product listings with columns for 'Product Name', 'Price', 'Average Value', and 'Product Life'. The first listing is for the 'IBM ThinkPad X31' with a price of \$2004-\$2235 and an average value of 4.7. The page also includes a search bar and navigation links.

Canonicalization: Product information

This screenshot shows a detailed view of the 'IBM ThinkPad X31' product page. It includes a table of specifications and a list of related products. The table has columns for 'Product Name', 'Price', 'Average Value', and 'Product Life'. The first product listed is the 'IBM ThinkPad X31' with a price of \$2004-\$2235 and an average value of 4.7. The page also includes a search bar and navigation links.

Inconsistency: digital cameras

- **Image Capture Device:** 1.68 million pixel 1/2-inch CCD sensor
- **Image Capture Device** Total Pixels Approx. 3.34 million Effective Pixels Approx. 3.24 million
- **Image sensor** Total Pixels: Approx. 2.11 million-pixel
- **Imaging sensor** Total Pixels: Approx. 2.11 million 1,688 (H) x 1,248 (V)
- **CCD** Total Pixels: Approx. 3,340,000 (2,140[H] x 1,560 [V])
 - **Effective Pixels:** Approx. 3,240,000 (2,088 [H] x 1,550 [V])
 - **Recording Pixels:** Approx. 3,145,000 (2,048 [H] x 1,536 [V])
- *These all came off the same manufacturer's website!!*
- *And this is a very technical domain. Try sofa beds.*

Using information extraction to populate knowledge bases

The screenshot shows the Protege software interface, which is used for creating and editing knowledge bases. The interface displays a hierarchical structure of concepts and instances, along with a list of extracted information from a source document. The extracted information includes details about a person's education and research interests.

<http://protege.stanford.edu/>



Named Entity Extraction

- The task: **find** and **classify** names in text, for example:

The **European Commission** [ORG] said on Thursday it disagreed with **German** [MISC] advice.

Only **France** [LOC] and **Britain** [LOC] backed **Fischler** [PER] 's proposal .

"What we have to be extremely careful of is how other countries are going to take **Germany** 's lead", **Welsh National Farmers ' Union** [ORG] (**NFU** [ORG]) chairman **John Lloyd Jones** [PER] said on **BBC** [ORG] radio .

- The purpose:
 - ... a lot of information is really associations between named entities.
 - ... for question answering, answers are usually named entities.
 - ... the same techniques apply to other slot-filling classifications.



CoNLL (2003) Named Entity Recognition task

Task: Predict semantic label of each word in text

Foreign	NNP	I-NP	ORG	} Standard evaluation is per entity, not per token
Ministry	NNP	I-NP	ORG	
spokesman	NN	I-NP	O	
Shen	NNP	I-NP	PER	
Guofang	NNP	I-NP	PER	
told	VBD	I-VP	O	
Reuters	NNP	I-NP	ORG	
:	:	:	:	



Precision and recall

- Precision**: fraction of retrieved items that are relevant = P(correct|selected)
- Recall**: fraction of relevant docs that are retrieved = P(selected|correct)

	Correct	Not Correct
Selected	tp	fp
Not Selected	fn	tn

- Precision $P = tp / (tp + fp)$
- Recall $R = tp / (tp + fn)$



A combined measure: F

- Combined measure that assesses this tradeoff is F measure (weighted harmonic mean):

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- People usually use balanced F_1 measure
 - i.e., with $\beta = 1$ or $\alpha = 1/2$: $F = 2PR / (P+R)$
- Harmonic mean is conservative average
 - See CJ van Rijsbergen, *Information Retrieval*



Precision/Recall/F1 for IE

- Recall and precision are straightforward for tasks like IR and text categorization, where there is only one grain size (documents)
- The measure behaves a bit funny for IE/NER when there are *boundary errors* (which are common):
 - First **Bank of Chicago** announced earnings ...
- This counts as both a fp and a fn
- Selecting *nothing* would have been better
- Some other systems (e.g., MUC scorer) give partial credit (according to complex rules)



Natural Language Processing-based Hand-written Information Extraction

- If extracting from automatically generated web pages, simple regex patterns usually work.
- If extracting from more natural, unstructured, human-written text, some NLP may help.
 - Part-of-speech (POS) tagging
 - Mark each word as a noun, verb, preposition, etc.
 - Syntactic parsing
 - Identify phrases: NP, VP, PP
 - Semantic word categories (e.g. from WordNet)
 - KILL: kill, murder, assassinate, strangle, suffocate
- Extraction patterns can use POS or phrase tags.
 - Crime victim:
 - Prefiller: [POS: V, Hypernym: KILL]
 - Filler: [Phrase: NP]



MUC: the NLP genesis of IE

- DARPA funded significant efforts in IE in the early to mid 1990's.
- Message Understanding Conference (MUC) was an annual event/competition where results were presented.
- Focused on extracting information from news articles:
 - Terrorist events
 - Industrial joint ventures
 - Company management changes
- Information extraction is of particular interest to the intelligence community ...
 - Though also to all other "information professionals"

Example of IE from FASTUS (1993)

Bridgestone Sports Co. said Friday it had set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be supplied to Japan.

The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990 with production of 20,000 iron and "metal wood" clubs a month.

TIE-UP-1 Relationship: TIE-UP Entities: "Bridgestone Sport Co." "a local concern" "a Japanese trading house" Joint Venture Company: "Bridgestone Sports Taiwan Co." Activity: ACTIVITY-1 Amount: NT\$200000000	ACTIVITY-1 Activity: PRODUCTION Company: "Bridgestone Sports Taiwan Co." Product: "iron and 'metal wood' clubs" Start Date: DURING: January 1990
---	---

Example of IE: FASTUS(1993)

Bridgestone Sports Co. said Friday it had set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be supplied to Japan.

The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990 with production of 20,000 iron and "metal wood" clubs a month.

TIE-UP-1 Relationship: TIE-UP Entities: "Bridgestone Sport Co." "a local concern" "a Japanese trading house" Joint Venture Company: "Bridgestone Sports Taiwan Co." Activity: ACTIVITY-1 Amount: NT\$200000000	ACTIVITY-1 Activity: PRODUCTION Company: "Bridgestone Sports Taiwan Co." Product: "iron and 'metal wood' clubs" Start Date: DURING: January 1990
---	---

FASTUS

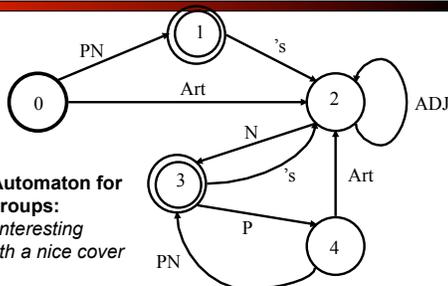


Based on finite state automata (FSA) transductions

1. Complex Words: Recognition of multi-words and proper names
2. Basic Phrases: Simple noun groups, verb groups and particles
3. Complex phrases: Complex noun groups and verb groups
4. Domain Events: Patterns for events of interest to the application. Basic templates are to be built.
5. Merging Structures: Templates from different parts of the texts are merged if they provide information about the same entity or event.



Grep++ = Cascaded grepping



Rule-based Extraction Examples

Determining which person holds what office in what organization

- [person] , [office] of [org]
 - Vuk Draskovic, leader of the Serbian Renewal Movement
- [org] (named, appointed, etc.) [person] P [office]
 - NATO appointed Wesley Clark as Commander in Chief

Determining where an organization is located

- [org] in [loc]
 - NATO headquarters in Brussels
- [org] [loc] (division, branch, headquarters, etc.)
 - KFOR Kosovo headquarters



Naive Bayes Classifiers

Task: Classify a new instance based on a tuple of attribute values

$$\langle x_1, x_2, \dots, x_n \rangle$$

$$c_{MAP} = \operatorname{argmax}_{c_j \in C} P(c_j | x_1, x_2, \dots, x_n)$$

$$c_{MAP} = \operatorname{argmax}_{c_j \in C} \frac{P(x_1, x_2, \dots, x_n | c_j) P(c_j)}{P(c_1, c_2, \dots, c_n)}$$

$$c_{MAP} = \operatorname{argmax}_{c_j \in C} P(x_1, x_2, \dots, x_n | c_j) P(c_j)$$



Naïve Bayes Classifier: Assumptions

- $P(c_j)$
 - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_n | c_j)$
 - $O(|X|^n \cdot |C|)$
 - Could only be estimated if a very, very large number of training examples was available.

Conditional Independence Assumption:

⇒ Assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities.

$$P(X_1, \dots, X_5 | C) = P(X_1 | C) \cdot P(X_2 | C) \cdot \dots \cdot P(X_5 | C)$$



Naïve Bayes in NLP

- For us, the x_i are usually bags of occurring words
- A class-conditional unigram language model!
 - Different from having a variable for each word type
- As usual, we need to smooth $P(x_i | c_j)$

- Zero probabilities cannot be conditioned away, no matter what other evidence there is

$$\hat{P}(x_{i,k} | c_j) = \frac{N(X_i = x_{i,k}, C = c_j) + mp_{i,k}}{N(C = c_j) + m}$$

- As before, multiplying lots of small numbers can cause floating-point underflow.
 - As $\log(xy) = \log(x) + \log(y)$ and \log is monotonic, it is faster and better to work by summing logs probabilities



'Change of Address' email

```
From: Robert Kubinsky <robert@lousycorp.com>
Subject: Email update

Hi all - I'm moving jobs and wanted to stay in touch
with everyone so...
My new email address is robert@cubeimedia.com
Hope all is well :)
>>R
```

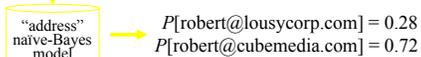


CoA: Details

1. Classification

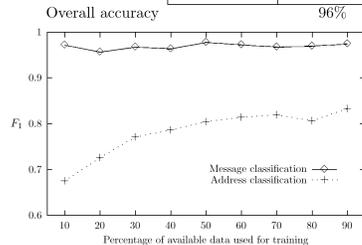


2. Extraction



Kushmerick et al. 2001 ATEM: Change of Address Results

	Words			Phrases		
	P	R	F ₁	P	R	F ₁
Message classification	.96	.66	.78	.98	.97	.98
Address classification	.96	.62	.76	.98	.68	.80



36 CoA messages
86 addresses
55 old, 31 new
5720 non-CoA