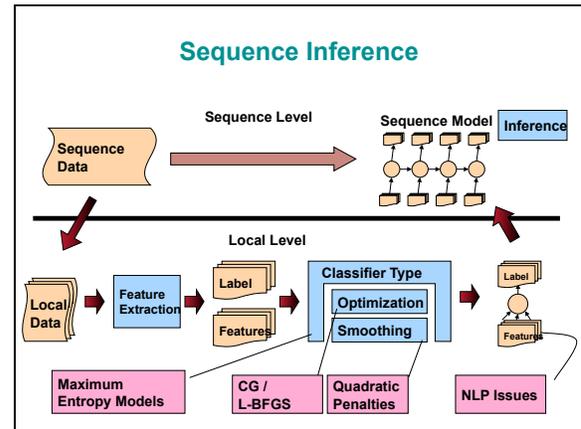


## Information Extraction: Sequence Models, Information Extraction Tasks and Information Integration

CS 224N  
2009



### MEMM inference in systems

- For a Conditional Markov Model (CMM) a.k.a. a Maximum Entropy Markov Model (MEMM), the classifier makes a single decision at a time, conditioned on evidence from observations and previous decisions.
- A larger space of sequences is explored via search

| Local Context |     |      |      |     | Decision Point | Features         |         |
|---------------|-----|------|------|-----|----------------|------------------|---------|
| -3            | -2  | -1   | 0    | +1  |                |                  |         |
| DT            | NNP | VBD  | ???  | ??? |                | $W_0$            | 22.6    |
| The           | Dow | fell | 22.6 | %   |                | $W_{+1}$         | %       |
|               |     |      |      |     |                | $W_{-1}$         | fell    |
|               |     |      |      |     |                | $T_{-1}$         | VBD     |
|               |     |      |      |     |                | $T_{-1}, T_{-2}$ | NNP-VBD |
|               |     |      |      |     |                | hasDigit?        | true    |
|               |     |      |      |     |                | ...              | ...     |

(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)

### Two ways to Search: Beam Inference



- Beam inference:**
  - At each position keep the top  $k$  complete sequences.
  - Extend each sequence in each local way.
  - The extensions compete for the  $k$  slots at the next position.
- Advantages:**
  - Fast; and beam sizes of 3–5 are as good or almost as good as exact inference in many cases.
  - Easy to implement (no dynamic programming required).
- Disadvantage:**
  - Inexact: the globally best sequence can fall off the beam.

### Two ways to Search: Viterbi Inference



- Viterbi inference:**
  - Dynamic programming or memoization.
  - Requires small window of state influence (e.g., past two states are relevant).
- Advantage:**
  - Exact: the global best sequence is returned.
- Disadvantage:**
  - Harder to implement long-distance state-state interactions (but beam inference tends not to successfully capture long-distance resurrection of sequences anyway).

### Viterbi Inference: J&M Ch. 6

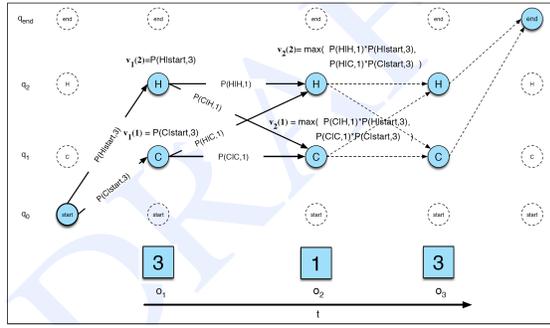
- I'm basically punting on this ... read Ch. 6.
  - I'll do dynamic programming for parsing
- It's a small change from HMM Viterbi
  - From:

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) P(s_j | s_i) P(o_t | s_j) \quad 1 \leq j \leq N, 1 < t \leq T$$

– To:

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) P(s_j | s_i, o_t) \quad 1 \leq j \leq N, 1 < t \leq T$$

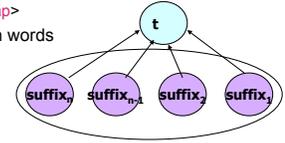
## Viterbi Inference: J&M Ch. 6



## Part-of-speech tagging: HMM Tagging Models of Brants 2000

- Highly competitive with other state-of-the-art models
- Trigram HMM with smoothed transition probabilities
- Capitalization feature becomes part of the state – each tag state is split into two e.g.  
 $NN \rightarrow \langle NN, cap \rangle, \langle NN, not\ cap \rangle$
- Suffix features for unknown words

$$P(w | tag) = P(suffix | tag)(w | suffix) \\ \approx \hat{P}(suffix) \tilde{P}(tag | suffix) / \hat{P}(tag)$$



$$\tilde{P}(tag | suffix_n) = \lambda_1 \hat{P}(tag | suffix_n) + \lambda_2 \hat{P}(tag | suffix_{n-1}) + \dots + \lambda_n \hat{P}(tag)$$

## MEMM Tagging Models -II

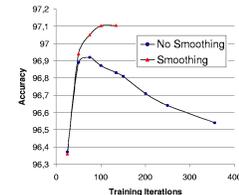
- Ratnaparkhi (1996): local distributions are estimated using maximum entropy models
  - Previous two tags, current word, previous two words, next two words, suffix, prefix, hyphenation, and capitalization features for unknown words
- Toutanova et al. (2003)
  - Richer features, bidirectional inference, better smoothing, better unknown word handling

| Model                | Overall Accuracy | Unknown Words |
|----------------------|------------------|---------------|
| MEMM (Ratn. 1996)    | 96.63            | 85.56         |
| HMM (Brants 2000)    | 96.7             | 85.5          |
| MEMM (T. et al 2003) | 97.24            | 89.04         |

## Smoothing: POS Tagging

- From (Toutanova et al., 2003):

|                   | Overall Accuracy | Unknown Word Acc |
|-------------------|------------------|------------------|
| Without Smoothing | 96.54            | 85.20            |
| With Smoothing    | 97.10            | 88.20            |



- Smoothing helps:
  - Softens distributions.
  - Pushes weight onto more explanatory features.
  - Allows many features to be dumped safely into the mix.
  - Speeds up convergence (if both are allowed to converge)!

## Summary of POS Tagging

For POS tagging, the change from generative to discriminative model **does not by itself** result in great improvement

One profits from discriminative models for specifying dependence on **overlapping features of the observation** such as spelling, suffix analysis, etc.

This additional power (of the CMM, CRF, Perceptron models) has been shown to result in improvements in accuracy

A CMM allows integration of rich features of the observations, but can suffer from assuming independence from following observations; this effect can be relieved by moving to a CRF, but also by adding dependence on following words

The **higher accuracy** of discriminative models comes at the price of **much slower training**

## CRFs [Lafferty, Pereira, and McCallum 2001]

- Another sequence model: Conditional Random Fields (CRFs)
- A whole-sequence conditional model rather than a chaining of local models.

$$P(c | d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_c \exp \sum_i \lambda_i f_i(c', d)}$$

- The space of  $c$ 's is now the space of sequences
  - But if the features  $f_i$  remain local, the conditional sequence likelihood can still be calculated exactly using dynamic programming
- Training is slower, but CRFs avoid causal-competition biases
- These (or a variant using a max margin criterion) are seen as the state-of-the-art these days, and fairly standardly used

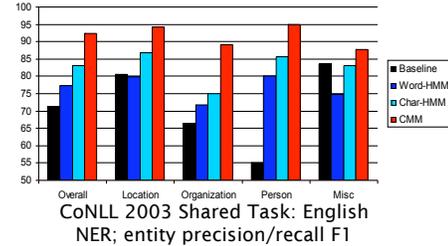
## NER Results: CoNLL (2003) Named Entity Recognition task

Task: Predict semantic label of each word in text

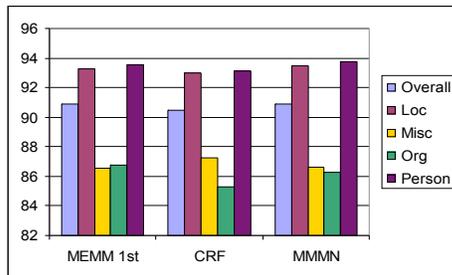
|           |     |      |     |   |
|-----------|-----|------|-----|---|
| Foreign   | NNP | I-NP | ORG | } Standard<br>evaluation<br>is per<br>entity, <i>not</i><br>per token |
| Ministry  | NNP | I-NP | ORG |   |
| spokesman | NN  | I-NP | O   |   |
| Shen      | NNP | I-NP | PER |   |
| Guofang 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

## NER Results: Discriminative Model

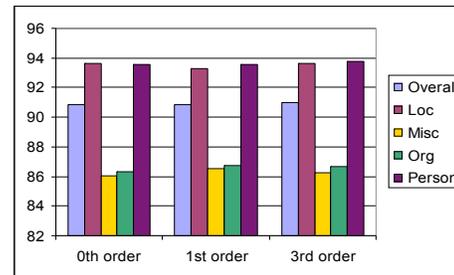
- Increases from better features, a better classification model.



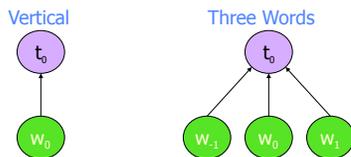
## Sequence models? CoNLL 2003 NER shared task Results on English Devset



## CoNLL NER Results: CMM Order



## Sequence Tagging Without Sequence Information: POS tagging

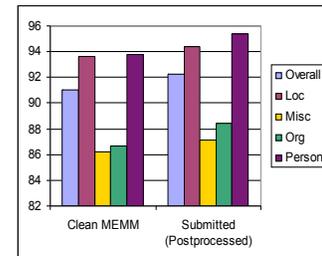


| Model    | Features | Token         | Unknown | Sentence |
|----------|----------|---------------|---------|----------|
| Vertical | 56,805   | <b>93.69%</b> | 82.61%  | 26.74%   |
| 3Words   | 239,767  | <b>96.57%</b> | 86.78%  | 48.27%   |

Using 3 words works significantly better than using only the current word and the previous two or three tags instead! (Toutanova et al. 2003)

## CoNLL NER: A real difference

- A difference of about 0.7% gives significance among good CoNLL results
- Here we get one!
- It was done with some Perl regular expressions



## Biomedical NER Motivation

- The biomedical world has a huge body of information, which is growing rapidly.
  - MEDLINE, the primary research database serving the biomedical community, currently contains over 12 million abstracts, with 60,000 new abstracts appearing each month.
  - There is also an impressive number of biological databases containing information on genes, proteins, nucleotide and amino acid sequences, including *GenBank*, *Swiss-Prot*, and *Fly-Base*; each contains entries numbering from the thousands to the millions and are multiplying rapidly.

## Motivation

- Currently, all of these resources are curated by hand by expert annotators at enormous expense.
- The information overload from the massive growth in the scientific literature has shown the necessity to automatically locate, organize and manage facts relating to experimental results
- Natural Language Processing can aid researchers and curators of biomedical databases by automating these tasks.

## Named Entity Recognition

- General NER vs. Biomedical NER

<PER> Christopher Manning </PER> is a professor at <ORG> Stanford University </ORG>, in <LOC> Palo Alto </LOC>.

<RNA> TAR </RNA> independent transactivation by <PROTEIN> Tat </PROTEIN> in cells derived from the <CELL> CNS </CELL> - a novel mechanism of <DNA> HIV-1 gene </DNA> regulation.

## Why is this difficult?

- The list of biomedical entities is growing.
  - New genes and proteins are constantly being discovered, so explicitly enumerating and searching against a list of known entities is not scalable.
  - Part of the difficulty lies in identifying previously unseen entities based on contextual, orthographic, and other clues.
- Biomedical entities don't have strict naming conventions.
  - Common English words such as *period*, *curved*, and *for* are used for gene names.
  - Entity names can be ambiguous. For example, in FlyBase, "clk" is the gene symbol for the "Clock" gene but it also is used as a synonym of the "period" gene.
- Biomedical entity names are ambiguous
  - Experts only agree on whether a word is even a gene or protein 69% of the time. (Krauthammer *et al.*, 2000)
  - Often systematic polysemies between gene, RNA, DNA, etc.

## Interesting Features

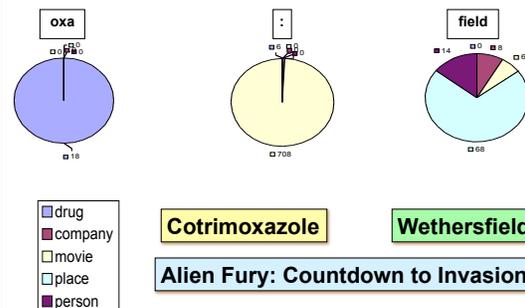
- Word, and surrounding context
- Word Shapes
  - Map words to simplified representation that encodes attributes such as length, capitalization, numerals, Greek letters, internal punctuation, etc.

|                  |        |
|------------------|--------|
| Varicella-zoster | Xx-xxx |
| mRNA             | xXXX   |
| CPA1             | XXXd   |

- Character substrings

#Dpp# → #Dpp#, #Dpp, Dpp#, #Dp, Dpp, pp#, #D, Dp, pp, p#, D, p

## Features: What's in a Name?



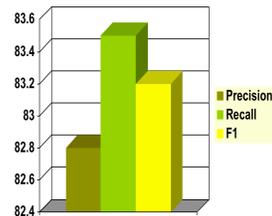
## Interesting Features

- Part-of-Speech tags
- Parsing information
- Searching the web for the word in a given context
  - *X gene, X mutation, X antagonist*
- Gazetteer
  - list words whose classification is known
- Abbreviation extraction (Schwartz and Hearst, 2003)
  - Identify short and long forms when occurring together in text

... Zn finger homeodomain 2 (Zfh 2) ...

## Finkel et al. (2004) Results

- BioCreative task – Identify genes and proteins



| Precision | Recall | F1    |
|-----------|--------|-------|
| 81.3%     | 86.1%  | 83.6% |

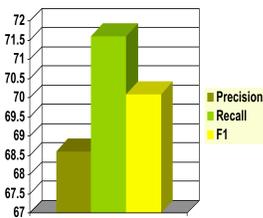
$$\text{precision} = \text{tp} / (\text{tp} + \text{fp})$$

$$\text{recall} = \text{tp} / (\text{tp} + \text{fn})$$

$$\text{F1} = 2(\text{precision})(\text{recall}) / (\text{precision} + \text{recall})$$

## Finkel et al. (2004) Results

- BioNLP task – Identify genes, proteins, DNA, RNA, and cell types



| Precision | Recall | F1    |
|-----------|--------|-------|
| 68.6%     | 71.6%  | 70.1% |

$$\text{precision} = \text{tp} / (\text{tp} + \text{fp})$$

$$\text{recall} = \text{tp} / (\text{tp} + \text{fn})$$

$$\text{F1} = 2(\text{precision})(\text{recall}) / (\text{precision} + \text{recall})$$

## Quiz question!

- Answer in one sentence:

In discriminative sequence labeling tasks like NER, why do sequence models (that condition on other labels) often offer little value over using straight classifiers (which don't condition on other labels)?

## Information Extraction and Integration

Following slides from:  
 William Cohen  
 Andrew McCallum  
 Eugene Agichtein  
 Sunita Sarawagi

## The Value of Text Data

- “Unstructured” text data is the primary source of human-generated information
  - Citeseer, comparison shopping, PIM systems, web search, data warehousing, scientific literature
- Managing and utilizing text: information extraction and integration
- Scalability: a bottleneck for deployment
- Relevance to data mining community

## Example: A Solution

The screenshot shows the FlipDog website interface. At the top, it displays '647,514' jobs from '63,641' employers. Below this, there are navigation links like 'Home', 'Find Jobs', 'Your Account', and 'Resource Center'. A search bar is visible with the text 'Find a Job!'. The main content area shows a list of job categories and a 'Find a Job!' button. On the right, there are links for 'About FlipDog', 'Privacy Policy', 'Terms of Service', and 'Contact Us'.

## Extracting Job Openings from the Web

The screenshot shows a job opening on the foodscience.com website. The job title is 'Ice Cream Guru'. The employer is 'foodscience.com'. The job category is 'Travel/Hospitality'. The job function is 'Food Services'. The job location is 'Upper Midwest'. The contact phone is '800-488-2611'. The date extracted is 'January 8, 2001'. The source is 'www.foodscience.com/jobs\_midwest.htm'. The other company jobs are 'foodscience.com-Job1'. The job description includes 'Major food manufacturer... Chicago area... food professional...'. The job requirements include 'Requires a BS in Food Science or Dairy...'. The contact information is 'E:800-488-2611'.

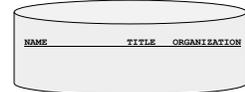
Job Openings:  
Category = Food Services  
Keyword = Baker  
Location = Continental U.S.

The screenshot shows the FlipDog website search results for '1-26 of 47 jobs'. The search criteria are 'Category = Food Services', 'Keyword = Baker', and 'Location = Continental U.S.'. The results are displayed in a table with columns for job title, location, date, and employer. The jobs listed include 'Food Pantry Workers at Lutheran Social Services', 'Cooks at Lutheran Social Services', 'Bakers Assistants at Fine Catering by Russell Momm', 'Baker's helper at Berlin-land', 'Assistant Baker at Gourmet To Go', 'Host/Hostess at Shania Restaurants', 'Cooks at Alta's Rustler Lodge', 'Line Attendant at Sun Valley Corporation', 'Food Service Worker II at Garden Grove Unified School District', 'Night Cook / Baker at SOMOCO', 'Cooks/Brew Cooks at Grandview Lodge', 'Line Cook at Lone Mountain Ranch', 'Production Baker at Whole Foods Market', 'Cafe Decorator/Baker at Mandalay Bay Hotel and Casino', and 'Shift Supervisors at Bueggers Bagels'.

## What is "Information Extraction"

As a task: **Filling slots in a database from sub-segments of text.**

October 14, 2002, 4:00 a.m. PT  
For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.  
Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels—the coveted code behind the Windows operating system—to select customers.  
"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."  
Richard Stallman, founder of the Free Software Foundation, countered saying...



## What is "Information Extraction"

As a task: **Filling slots in a database from sub-segments of text.**

October 14, 2002, 4:00 a.m. PT  
For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.  
Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels—the coveted code behind the Windows operating system—to select customers.  
"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."  
Richard Stallman, founder of the Free Software Foundation, countered saying...



## What is "Information Extraction"

As a family of techniques: **Information Extraction = segmentation + classification + clustering + association**

October 14, 2002, 4:00 a.m. PT  
For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.  
Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels—the coveted code behind the Windows operating system—to select customers.  
"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."  
Richard Stallman, founder of the Free Software Foundation, countered saying...

Microsoft Corporation  
CEO  
Bill Gates  
Microsoft  
Gates  
Microsoft  
Bill Veghte  
Microsoft  
VP  
Richard Stallman  
founder  
Free Software Foundation



## Landscape of IE Tasks (3/4): Complexity

E.g. word patterns:

### Closed set

U.S. states

He was born in **Alabama**...

The big **Wyoming** sky...

### Complex pattern

U.S. postal addresses

University of Arkansas  
P.O. Box 140  
**Hope, AR 71802**

Headquarters:  
1128 Main Street, 4th Floor  
Cincinnati, Ohio 45210

### Regular set

U.S. phone numbers

Phone: **(413) 545-1323**

The CALD main office can be reached at **412-268-1299**

### Ambiguous patterns, needing context and many sources of evidence

Person names

...was among the six houses sold by **Hope Feldman** that year.

**Pawel Opalinski**, Software Engineer at WhizBang Labs.

## Landscape of IE Tasks (4/4): Single Field/Record

**Jack Welch will retire as CEO of General Electric tomorrow. The top role at the Connecticut company will be filled by Jeffrey Immelt.**

### Single entity

Person: Jack Welch

Person: Jeffrey Immelt

Location: Connecticut

### Binary relationship

Relation: Person-Title  
Person: Jack Welch  
Title: CEO

Relation: Company-Location  
Company: General Electric  
Location: Connecticut

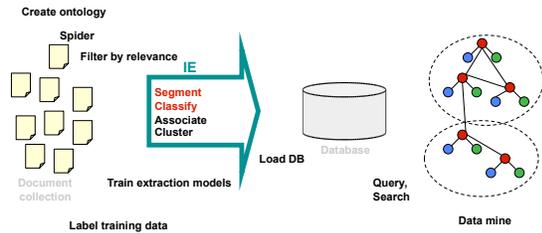
### N-ary record

Relation: Succession  
Company: General Electric  
CEO  
Out: Jack Welch  
In: Jeffrey Immelt

"Named entity" extraction

## Broader View

Up to now we have been focused on ML methods for segmentation and classification



## Steps 1 & 2: Hand Coded Rule Example: Conference Name

```
# These are subordinate patterns
SwordOrdinals="(?:first|second|third|fourth|fifth|sixth|seventh|eighth|ninth|tenth|eleventh|twelfth|thirteenth|fourteenth|fifteenth)";
my $numberOrdinals="(?:\d(?:?:1st|2nd|3rd|1th|2th|3th|4th|5th|6th|7th|8th|9th|0th))";
my $ordinals="(?:SwordOrdinals|$numberOrdinals)";
my $confTypes="(?:Conference|Workshop|Symposium)";
my $words="(?:[A-Z][w+vs]*)"; # A word starting with a capital letter and ending with 0 or more spaces
my $confDescriptors="(?:international|s+[A-Z]+|s+)"; # e.g "International Conference ..." or the conference name for workshops (e.g. "VLDB Workshop ...")
my $connectors="(?:,|and)";
my $abbreviations="(?:\s+([A-Z][w+|Nw|s]?(?:\d+)?))"; # Conference abbreviations like "(SIGMOD'06)"
# The actual pattern we search for. A typical conference name this pattern will find is
# "3rd International Conference on Blah Blah Blah (ICBBB-05)"
my $fullNamePattern="(?:SwordOrdinals|$words|$confDescriptors)?$confTypes(?:\s+$connectors|s+)*\s+";
$abbreviations?(?:\s+|s+)$";
#####
# Given a <dbworldMessage> look for the conference pattern
#####
lookForPattern($dbworldMessage, $fullNamePattern);
#####
# In a given <file>, look for occurrences of <pattern>
# <pattern> is a regular expression
#####
sub lookForPattern {
    my ($file,$pattern) = @_;
```

## Machine Learning Methods

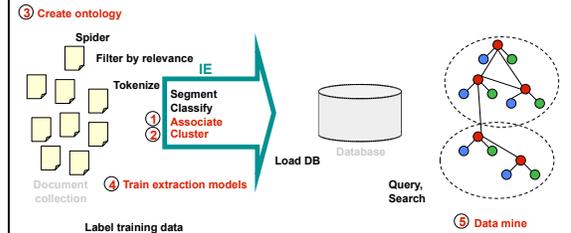
- Sequence models: HMMs, CMMs/MEMMs, CRFs
- Can work well when training data is easy to construct and is plentiful
- Can capture complex patterns that are hard to encode with hand-crafted rules
  - e.g., determine whether a review is positive or negative
  - extract long complex gene names

*The human T cell leukemia lymphotropic virus type 1 Tax protein represses MyoD-dependent transcription by inhibiting MyoD-binding to the KIX domain of p300.*

- Can be labor intensive to construct training data
  - Question: how much training data is sufficient?

## Broader View

Now touch on some other issues



## Relation Extraction: Disease Outbreaks

- Extract structured relations from text

**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...

Disease Outbreaks in *The New York Times*

| Date      | Disease Name    | Location |
|-----------|-----------------|----------|
| Jan. 1995 | Malaria         | Ethiopia |
| July 1995 | Mad Cow Disease | U.K.     |
| Feb. 1995 | Pneumonia       | U.S.     |

Information Extraction System (e.g., NYU's Proteus)

## Example: 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**.”

CBF-A  $\xrightarrow[\text{complex}]{\text{interact}}$  CBF-C

CBF-B  $\xrightarrow{\text{associates}}$  CBF-A-CBF-C complex

## Relation Extraction

- Typically requires Entity Tagging as preprocessing
- Knowledge Engineering
  - Rules defined over lexical items
    - \* <company> located in <location>
  - Rules defined over parsed text
    - \* \*(Obj <company>) (Verb located) (\*) (Subj <location>))
  - Proteus, GATE, ...
- Machine Learning-based
  - Supervised: Learn rules/patterns from examples
    - Roth 2005, Cardie 2006, Mooney 2005, Bunescu 2007, ...
  - Partially-supervised: bootstrap from "seed" examples
    - Agichtein & Gravano 2000, Etzioni et al., 2004, ...

## Example Extraction Rule [NYU Proteus]

```

;; For <company> appoints <person> <position>

(defpattern appoint
  "np-sem(C-company)? rn? sa? vg(C-appoint) np-sem(C-person) ',?'
  to-be? np(C-position) to-succeed?:"
  company-at=1.attributes, sa=3.span, lv=4.span, person-at=5.attribute:
  position-at=8.attributes |
  ...
  (defun when-appoint (phrase-type)
    (let ((person-at (binding 'person-at))
          (company-entity (entity-bound 'company-at))
          (person-entity (essential-entity-bound 'person-at 'C-person))
          (position-entity (entity-bound 'position-at))
          (predecessor-entity (entity-bound 'predecessor-at))
          new-event)
      (not-an-antecedent position-entity)
      ;; if no company is specified for position, use agent
      ...
    )
  )
  )
  
```

## Example of Learned Extraction Patterns: Snowball [AG2000]

**ORGANIZATION** {<'s 0.7> <in 0.7> <headquarters 0.7> **LOCATION**

**LOCATION** {<- 0.75> <based 0.75> **ORGANIZATION**

## (1) Association as Binary Classification

**Christos Faloutsos conferred with Ted Senator, the KDD 2003 General Chair.**

Person Person Role

Person-Role (Christos Faloutsos, KDD 2003 General Chair) → NO

Person-Role ( Ted Senator, KDD 2003 General Chair) → YES

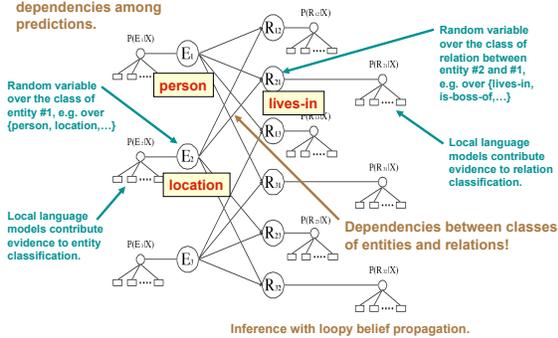
Do this with SVMs and tree kernels over parse trees.

[Zelenko et al, 2002]

## (2) Association with Graphical Models

Capture arbitrary-distance dependencies among predictions.

[Roth & Yih 2002]



## Accuracy of Information Extraction

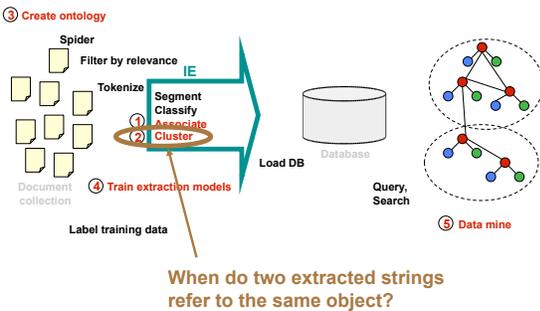
| Information Type | Accuracy |
|------------------|----------|
| Entities         | 90-98%   |
| Attributes       | 80%      |
| Facts            | 60-70%   |
| Events           | 50-60%   |

[Feldman, ICML 2006 tutorial]

- Errors cascade (error in entity tag → error in relation extraction)
- This estimate is optimistic:
  - Holds for well-established tasks
  - Many specific/novel IE tasks exhibit much lower accuracy

## Broader View

Now touch on some other issues



## Extracted Entities: Resolving Duplicates



Document 1: The Justice Department has officially ended its inquiry into the assassinations of John F. Kennedy and Martin Luther King Jr., finding "no persuasive evidence" to support conspiracy theories, according to department documents. The House Assassinations Committee concluded in 1978 that Kennedy was "probably" assassinated as the result of a conspiracy involving a second gunman, a finding that broke from the Warren Commission's belief that Lee Harvey Oswald acted alone in Dallas on Nov. 22, 1963.

Document 2: In 1953, Massachusetts Sen. John F. Kennedy married Jacqueline Lee Bouvier in Newport, R.I. In 1960, Democratic presidential candidate John F. Kennedy confronted the issue of his Roman Catholic faith by telling a Protestant group in Houston, "I do not speak for my church on public matters, and the church does not speak for me."

Document 3: David Kennedy was born in Leicester, England in 1959. ...Kennedy co-edited The New Poetry (Bloodaxe Books 1993), and is the author of New Relations: The Refashioning Of British Poetry 1980-1994 (Seren 1996).

[From Li, Morie, & Roth, AI Magazine, 2005]

## Important Problem

- Appears in numerous real-world contexts
- Plagues many applications
  - Citeseer, DBLife, AliBaba, Rexa, etc.

## (2) Information Integration

[Minton, Knoblock, et al 2001], [Doan, Domingos, Halevy 2001], [Richardson & Domingos 2003]

Goal might be to merge results of two IE systems:

|          |                                  |   |             |                      |
|----------|----------------------------------|---|-------------|----------------------|
| Name:    | Introduction to Computer Science | → | Title:      | Intro. to Comp. Sci. |
| Number:  | CS 101                           | → | Num:        | 101                  |
| Teacher: | M. A. Kludge                     | → | Dept:       | Computer Science     |
| Time:    | 9-11am                           | → | Teacher:    | Dr. Kludge           |
| Name:    | Data Structures in Java          | → | TA:         | John Smith           |
| Room:    | 5032 Wean Hall                   | → | Topic:      | Java Programming     |
|          |                                  |   | Start time: | 9:10 AM              |