

Example: The Problem

A screenshot of a Google search results page for the query "baker job opening". The search results are cluttered with various links, including job openings from FlipDog.com, Softpage Community Discussion Groups, and other websites. The page is annotated with three orange boxes and arrows pointing to specific search results:

- Martin Baker, a person**: Points to a search result for "Martin Baker" in a Softpage Community Discussion Group.
- Genomics job**: Points to a search result for "Genomics Job Opening" from CGI.
- Employers job posting form**: Points to a search result for "Post an Employee Benefits Job Opening" from Softpage.

Example: A Solution

A screenshot of the FlipDog job search website. The page features a search bar, navigation links, and a list of job openings. A prominent orange box highlights the text: "Job Openings: Category = Food Services Keyword = Baker Location = Continental U.S.". The website interface includes a search bar, navigation links, and a list of job openings.

Extracting Job Openings from the Web

A screenshot of a web browser displaying a job listing from foodscience.com. The listing is for a "Job Title: Ice Cream Guru" at "Employer: foodscience.com". The listing includes details such as "JobCategory: Travel/Hospitality", "JobFunction: Food Services", "JobLocation: Upper Midwest", and "Contact Phone: 800-488-2611". The listing was "Date/Extracted: January 8, 2001" and the source is "www.foodscience.com/jobs_midwest.htm". The page is annotated with several orange boxes and arrows pointing to specific elements of the job listing.

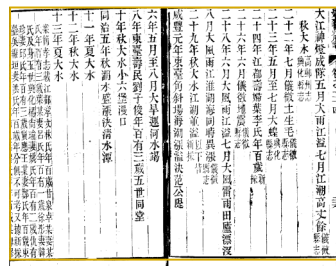
A screenshot of the FlipDog job search website showing a list of job openings. The list includes various roles such as "Food Pantry Workers at Lutheran Social Services", "Cooks at Lutheran Social Services", "Bakers Assistants at Fine Catering by Russell Moran", and "Baker's Helper at Bird-in-Hand". The list is annotated with an orange box containing the text: "Job Openings: Category = Food Services Keyword = Baker Location = Continental U.S.". The website interface includes a search bar, navigation links, and a list of job openings.

Data Mining the Extracted Job Information

A screenshot of the FlipDog Job Opportunity Index report for November 2001. The report features a map of the United States showing job supply by region. The text states: "The Job Opportunity Index™ (JOI) increased for the first time in three months in October — climbing on pain to risk and signaling a slight increase in U.S. job supply. However, numerous factors, including a dramatic half-point increase in the national unemployment rate, made October anything but normal." The report also includes a "Special Offer" for a limited-time JOI Subscriber Federal Program.

IE from Chinese Documents regarding WEATHER

Department of Terrestrial System, Chinese Academy of Sciences



200k+ documents
several millennia old

- Qing Dynasty Archives
- memos
- newspaper articles
- diaries

What is "Information Extraction"

As a family of techniques:

Information Extraction = segmentation + classification + association + clustering

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

What is "Information Extraction"

As a family of techniques:

Information Extraction = segmentation + classification + association + clustering

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

What is "Information Extraction"

As a family of techniques:

Information Extraction = segmentation + classification + association + clustering

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.

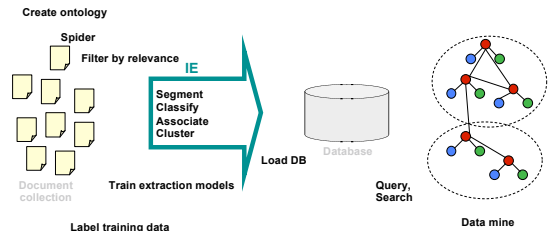
Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"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

IE in Context



IE History

Pre-Web

- Mostly news articles
 - De Jong's FRUMP [1982]
 - Hand-built system to fill Schank-style "scripts" from news wire
 - Message Understanding Conference (MUC) DARPA [87-'95], TIPSTER [92-'96]
- Most early work dominated by hand-built models
 - E.g. SRI's FASTUS, hand-built FSMs.
 - But by 1990's, some machine learning: Lehnert, Cardie, Grishman and then HMMs: Ekan [Leek '97], BBN [Bikel et al '98]

Web

- AAAI '94 Spring Symposium on "Software Agents"
 - Much discussion of ML applied to Web. Maes, Mitchell, Etzioni.
- Tom Mitchell's WebKB, '96
 - Build KB's from the Web.
- Wrapper Induction
 - Initially hand-build, then ML: [Soderland '96], [Kushneric '97],...

What makes IE from the Web Different?

Less grammar, but more formatting & linking

NewsWire

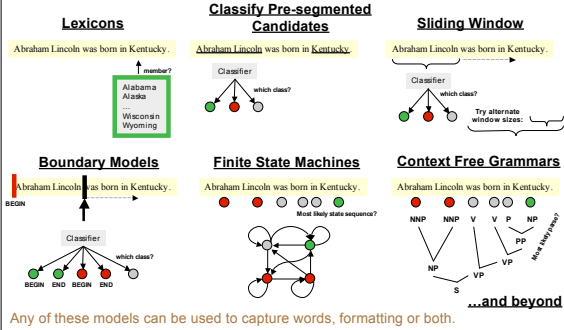
Apple to Open Its First Retail Store in New York City

MACWORLD EXPO, NEW YORK--July 17, 2002--Apple's first retail store in New York City will open in Manhattan's SoHo district on Thursday, July 18 at 8:00 a.m. EDT. The SoHo store will be Apple's largest retail store to date and is a stunning example of Apple's commitment to offering customers the world's best computer shopping experience.

"Fourteen months after opening our first retail store, our 31 stores are attracting over 100,000 visitors each week," said Steve Jobs, Apple's CEO. "We hope our SoHo store will surprise and delight both Mac and PC users who want to see everything the Mac can do to enhance their digital lifestyles."

The directory structure, link structure, formatting & layout of the Web is its own new grammar.

Landscape of IE Techniques (1/1): Models



Sliding Windows

Extraction by Sliding Window

E.g.
Looking for
seminar
location

GRAND CHALLENGES FOR MACHINE LEARNING

Jaime Carbonell
School of Computer Science
Carnegie Mellon University

3:30 pm
7500 Wean Hall

Machine learning has evolved from obscurity in the 1970s into a vibrant and popular discipline in artificial intelligence during the 1980s and 1990s. As a result of its success and growth, machine learning is evolving into a collection of related disciplines: inductive concept acquisition, analytic learning in problem solving (e.g. analogy, explanation-based learning), learning theory (e.g. PAC learning), genetic algorithms, connectionist learning, hybrid systems, and so on.

CMU UseNet Seminar Announcement

Extraction by Sliding Window

E.g.
Looking for
seminar
location

GRAND CHALLENGES FOR MACHINE LEARNING

Jaime Carbonell
School of Computer Science
Carnegie Mellon University

3:30 pm
7500 Wean Hall

Machine learning has evolved from obscurity in the 1970s into a vibrant and popular discipline in artificial intelligence during the 1980s and 1990s. As a result of its success and growth, machine learning is evolving into a collection of related disciplines: inductive concept acquisition, analytic learning in problem solving (e.g. analogy, explanation-based learning), learning theory (e.g. PAC learning), genetic algorithms, connectionist learning, hybrid systems, and so on.

CMU UseNet Seminar Announcement

Extraction by Sliding Window

E.g.
Looking for
seminar
location

GRAND CHALLENGES FOR MACHINE LEARNING

Jaime Carbonell
School of Computer Science
Carnegie Mellon University

3:30 pm
7500 Wean Hall

Machine learning has evolved from obscurity in the 1970s into a vibrant and popular discipline in artificial intelligence during the 1980s and 1990s. As a result of its success and growth, machine learning is evolving into a collection of related disciplines: inductive concept acquisition, analytic learning in problem solving (e.g. analogy, explanation-based learning), learning theory (e.g. PAC learning), genetic algorithms, connectionist learning, hybrid systems, and so on.

CMU UseNet Seminar Announcement

Extraction by Sliding Window

E.g.
Looking for
seminar
location

GRAND CHALLENGES FOR MACHINE LEARNING

Jaime Carbonell
School of Computer Science
Carnegie Mellon University

3:30 pm
7500 Wean Hall

Machine learning has evolved from obscurity in the 1970s into a vibrant and popular discipline in artificial intelligence during the 1980s and 1990s. As a result of its success and growth, machine learning is evolving into a collection of related disciplines: inductive concept acquisition, analytic learning in problem solving (e.g. analogy, explanation-based learning), learning theory (e.g. PAC learning), genetic algorithms, connectionist learning, hybrid systems, and so on.

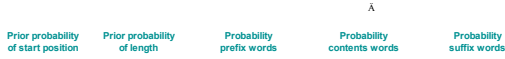
CMU UseNet Seminar Announcement

A "Naïve Bayes" Sliding Window Model

[Freitag 1997]



$P(\text{"Wean Hall Rm 5409"} = \text{LOCATION}) =$



Try all start positions and reasonable lengths Estimate these probabilities by (smoothed) counts from labeled training data.

If $P(\text{"Wean Hall Rm 5409"} = \text{LOCATION})$ is above some threshold, extract it.

Other examples of sliding window: [Baluja et al 2000] (decision tree over individual words & their context)

"Naïve Bayes" Sliding Window Results

Domain: CMU UseNet Seminar Announcements

GRAND CHALLENGES FOR MACHINE LEARNING

Jaime Carbonell
School of Computer Science
Carnegie Mellon University

3:30 pm
7500 Wean Hall

Machine learning has evolved from obscurity in the 1970s into a vibrant and popular discipline in artificial intelligence during the 1980s and 1990s. As a result of its success and growth, machine learning is evolving into a collection of related disciplines: inductive concept acquisition, analytic learning in problem solving (e.g. analogy, explanation-based learning), learning theory (e.g. PAC learning), genetic algorithms, connectionist learning, hybrid systems, and so on.

Field	F1
Person Name:	30%
Location:	61%
Start Time:	98%

Problems with Sliding Windows and Boundary Finders

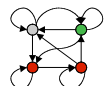
- Decisions in neighboring parts of the input are made independently from each other.
 - Naïve Bayes Sliding Window may predict a "seminar end time" before the "seminar start time".
 - It is possible for two overlapping windows to both be above threshold.
 - In a Boundary-Finding system, left boundaries are laid down independently from right boundaries, and their pairing happens as a separate step.

Finite State Machines

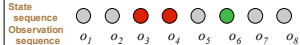
Hidden Markov Models

HMMs are the standard sequence modeling tool in genomics, music, speech, NLP, ...

Finite state model



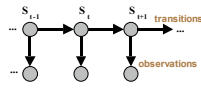
Generates:



State sequence

Observation sequence

Graphical model



$$P(\bar{s}, \bar{o}) \propto \prod_{t=1}^{|\bar{o}|} P(s_t | s_{t-1}) P(o_t | s_t)$$

Parameters: for all states $S = \{s_1, s_2, \dots\}$

Start state probabilities: $P(s_1)$

Transition probabilities: $P(s_t | s_{t-1})$

Observation (emission) probabilities: $P(o_t | s_t)$ Usually a multinomial over atomic, fixed alphabet

Training:

Maximize probability of training observations (w/ prior)

IE with Hidden Markov Models

Given a sequence of observations:

Yesterday Lawrence Saul spoke this example sentence.

and a trained HMM:



Find the most likely state sequence: (Viterbi)



Any words said to be generated by the designated "person name" state extract as a person name:

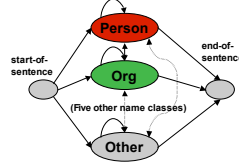
Person name: Lawrence Saul

HMMs for IE: A richer model, with backoff

HMM Example: "Nymble"

Task: Named Entity Extraction

[Bikel, et al 1998],
[BBN "IdentiFinder"]



Transition probabilities

$$P(s_t | s_{t-1}, o_{t-1})$$

Observation probabilities

$$P(o_t | s_t, s_{t-1})$$

or $P(o_t | s_t, o_{t-1})$

Back-off to:

$$P(s_t | s_{t-1})$$

Back-off to:

$$P(o_t | s_t)$$

Train on 450k words of news wire text.

$$P(s_t)$$

$$P(o_t)$$

Results:

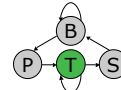
Case	Language	F1 _o
Mixed	English	93%
Upper	English	91%
Mixed	Spanish	90%

Other examples of shrinkage for HMMs in IE: [Freitag and McCallum '99]

HMMs for IE: Augmented finite-state structures with linear interpolation

Simple HMM structure for IE

- 4 state types:
 - Background (generates words not of interest),
 - Target (generates words to be extracted),
 - Prefix (generates typical words preceding target)
 - Suffix (words typically following target)



- Properties:
 - Extracts one type of target (e.g. target = person name), we will build one model for each extracted type.
 - Models different Markov-order n-grams for different predicted state contexts.
 - even though there are multiple states for "Background", state-path given labels is unambiguous. Therefore model parameters can all be computed using counts from labeled training data

More rich prefix and suffix structures

- In order to represent more context, add more state structure to prefix, target and suffix.
- But now overfitting becomes more of a problem.

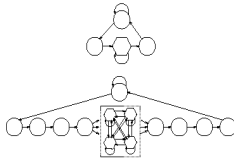
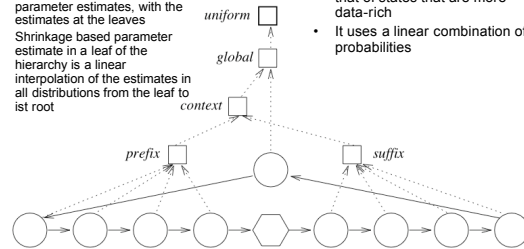


Figure 1: Two example HMM structures. Circle nodes represent non-target states; hexagon nodes represent target states.

Linear interpolation across states

- Is defined in terms of some hierarchy that represents the expected similarity between parameter estimates, with the estimates at the leaves
- Shrinkage based parameter estimate in a leaf of the hierarchy is a linear interpolation of the estimates in all distributions from the leaf to its root
- Shrinkage smoothes the distribution of a state towards that of states that are more data-rich
- It uses a linear combination of probabilities



Bayesian Model Merging

- Maximally Specific Model
- Neighbor-merging
- V-merging

Bayesian Model Merging

- Iterates merging states until an optimal tradeoff between fit to the data and model size has been reached

$P(M | D) \sim P(D | M) P(M)$

M = Model
D = Data

$P(D | M)$ can be calculated with the Forward algorithm
 $P(M)$ model prior can be formulated to reflect a preference for smaller models

HMM Emissions

2 million words of BibTeX data from the Web

HMM Information Extraction Results

Per-word error rate	Headers	References
One state/class Labeled data only	0.095	
Model Merging Labeled data only	0.087 (8% better)	
One state/class +BibTeX data	0.076 (20% better)	
Model Merging +BibTeX	0.071 (25% better)	0.066

Stochastic Optimization

- Start with a simple model
- Perform hill-climbing in the space of possible structures
- Make several runs and take the average to avoid local optima

State Operations

- Lengthen a prefix
- Split a prefix
- Lengthen a suffix
- Split a suffix
- Lengthen a target string
- Split a target string
- Add a background state

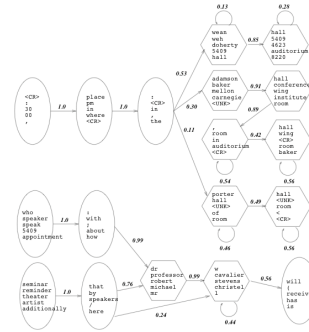
LearnStructure Algorithm

```

procedure LearnStructure(LabeledSet, Ops)
  ValidSet ← 1/3 of LabeledSet
  TrainSet ← LabeledSet - ValidSet
  CurModel ← the simple model
  Keepers ← {CurModel}
   $\bar{f} \leftarrow 0$ 
  while  $\bar{f} < 20$  and CurModel has fewer than 25 states
    Candidates ←  $\{M | M \in \text{op}(\text{CurModel}) \wedge \text{op} \in \text{Ops}\}$ 
    for  $M \in \text{Candidates}$ 
      score( $M$ ) ← average of 3 runs trained on
        TrainSet and scored for F1 on ValidSet
    CurModel ←  $M \in \text{Candidates}$  with highest score
    Keepers ← Keepers  $\cup$  {CurModel}
     $\bar{f} \leftarrow \bar{f} + 1$ 
  for  $M \in \text{Keepers}$ 
    score( $M$ ) ← average F1 from
      3-fold cross-validation on LabeledSet
  return  $M \in \text{Keepers}$  with highest score
  
```

Part of Example Learned Structure

Locations



Speakers

Accuracy of Automatically-Learned Structures

	speaker	location	acquired	dramt	title	company	conf	deadline	Average
Grown HMM	76.9	87.5	41.3	54.4	58.3	65.4	27.2	46.5	57.2
vs. SRV	+19.8	+16.0	+1.1	-1.6	—	—	—	—	+8.8
vs. Rapier	+23.9	+14.8	+12.5	+15.1	-11.7	+24.9	—	—	+13.3
vs. Simple HMM	+24.3	+5.6	+14.3	+5.6	+5.7	+11.1	+13.7	+6.7	+11.1
vs. Complex HMM	-2.1	+6.7	+7.3	-0.3	-0.3	+19.1	+0.0	-6.8	+3.0

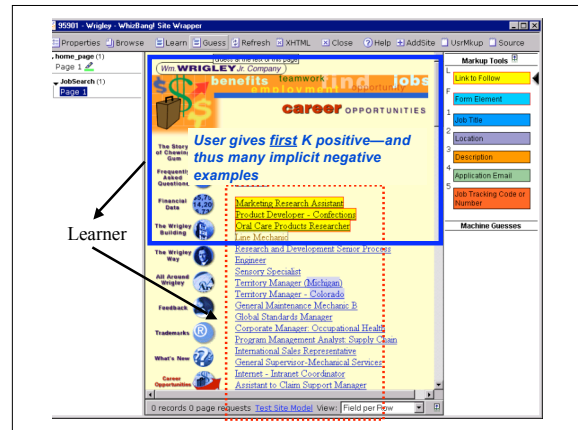
Table 2: Difference in F1 performance between the HMM using a learned structure and other methods. The + numbers indicate how much better our Grown HMM did than the alternative method.

Limitations of HMM/CRF models

- HMM/CRF models have a **linear** structure
- Web documents have a **hierarchical** structure
 - Are we suffering by not modeling this structure more explicitly?
- How can one learn a **hierarchical** extraction model?
 - Coming up: STALKER, a hierarchical **wrapper-learner**
 - But first: how do we train wrapper-learners?

Tree-based Models

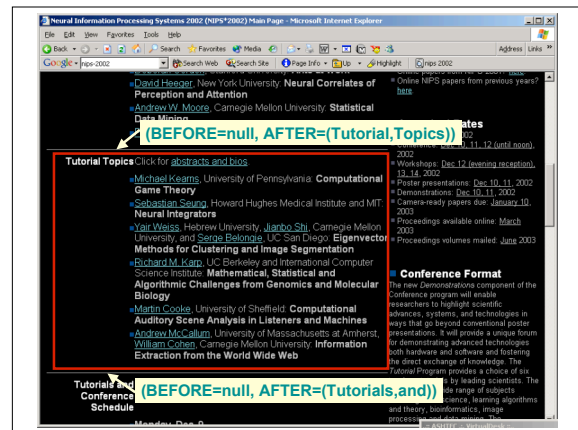
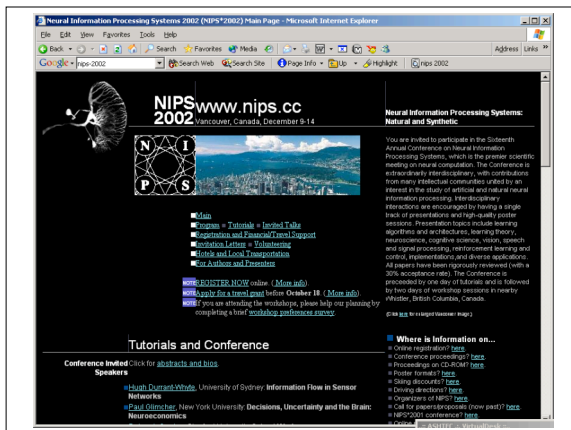
- Extracting from **one** web site
 - Use *site-specific* formatting information: e.g., "the JobTitle is a bold-faced paragraph in column 2"
 - For large well-structured sites, like parsing a **formal language**
- Extracting from **many** web sites:
 - Need general solutions to entity extraction, grouping into records, etc.
 - Primarily use *content* information
 - Must deal with a *wide range* of ways that users present data.
 - Analogous to parsing **natural language**
- Problems are **complementary**:
 - Site-dependent learning can **collect training data** for a site-independent learner
 - Site-dependent learning can **boost accuracy** of a site-independent learner on selected key sites

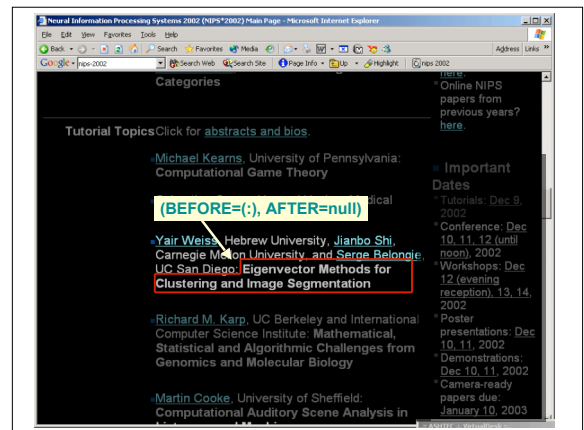
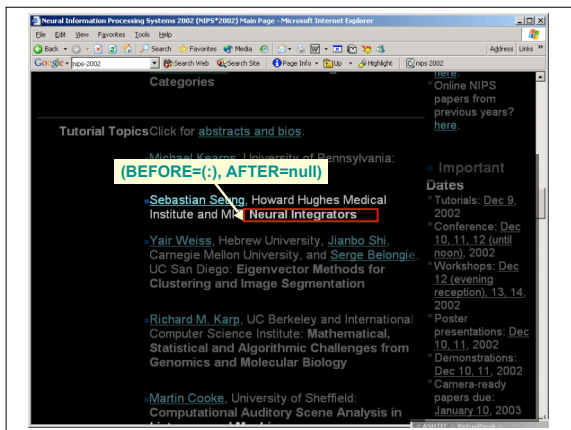
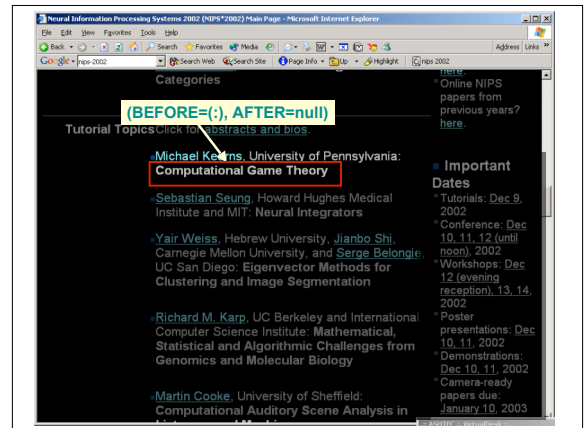
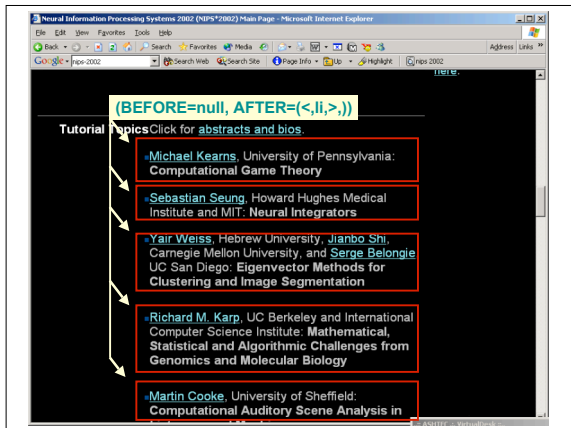


STALKER: Hierarchical boundary finding

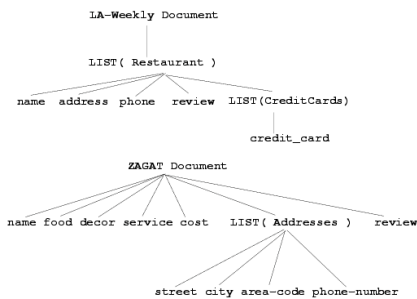
[Muslea, Minton & Knoblock 99]

- Main idea:
 - To train a hierarchical extractor, pose a **series** of learning problems, one for each node in the hierarchy
 - At each stage, extraction is simplified by knowing about the "context."





Stalker: hierarchical decomposition of two web sites



Stalker: summary and results

- Rule format:
 - “landmark automata” format for rules
 - E.g.: `<a>W. Cohen CMU: Web IE `
 - STALKER: `BEGIN = SkipTo(<, /, a, >), SkipTo(:)`
- Top-down rule learning algorithm
 - Carefully chosen ordering between types of rule specializations
- Very fast learning: e.g. 8 examples vs. 274
- **A lesson:** we often control the IE training data!

Learning Formatting Patterns "On the Fly": "Scoped Learning"

[Bagnell, Blei, McCallum, 2002]

Formatting is regular on each site, but there are too many different sites to wrap. Can we get the best of both worlds?

Scoped Learning Generative Model

- For each of the D documents:
 - Generate the multinomial formatting feature parameters f from $p(f|a)$
- For each of the N words in the document:
 - Generate the n th category c_n from $p(c_n)$.
 - Generate the n th word (global feature) from $p(w_n|c_n, \phi)$
 - Generate the n th formatting feature (local feature) from $p(f_n|c_n, f)$

$$p(\phi, \mathbf{c}, \mathbf{w}, \mathbf{f}) = p_\alpha(\phi) \prod_{n=1}^N p(c_n) p_\theta(w_n|c_n) p(f_n|c_n, \phi)$$

Inference

Given a new web page, we would like to classify each word resulting in $\mathbf{c} = \{c_1, c_2, \dots, c_n\}$

$$p(\mathbf{c}|\mathbf{w}, \mathbf{f}) = \frac{\int \prod_{n=1}^N p(w_n|c_n) p(f_n|c_n, \phi) p(c_n) p(\phi) d\phi}{\int \prod_{n=1}^N \sum_{c_n} p(w_n|c_n) p(f_n|c_n, \phi) p(c_n) p(\phi) d\phi}$$

This is not feasible to compute because of the integral and sum in the denominator. We experimented with two approximations:

- MAP point estimate of \mathbf{f}
- Variational inference

MAP Point Estimate

If we approximate \mathbf{f} with a point estimate, $\hat{\mathbf{f}}$, then the integral disappears and \mathbf{c} decouples. We can then label each word with:

$$\hat{c}_n = \arg \max_{c_n} p(w_n|c_n) p(f_n|c_n, \hat{\mathbf{f}}) p(c_n)$$

A natural point estimate is the posterior mode: a maximum likelihood estimate for the local parameters given the document in question:

$$\hat{\phi} = \arg \max_{\phi} p(\phi|\mathbf{f}, \mathbf{w})$$

E-step:

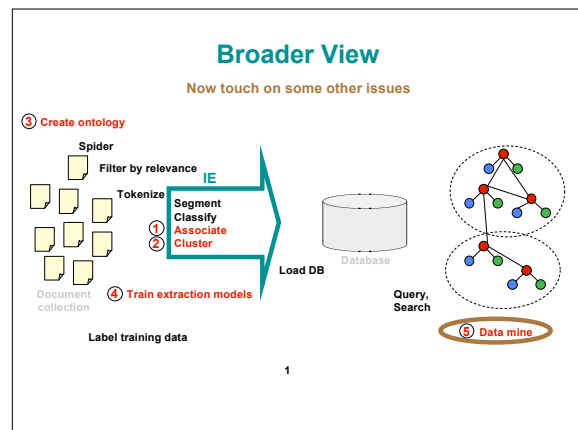
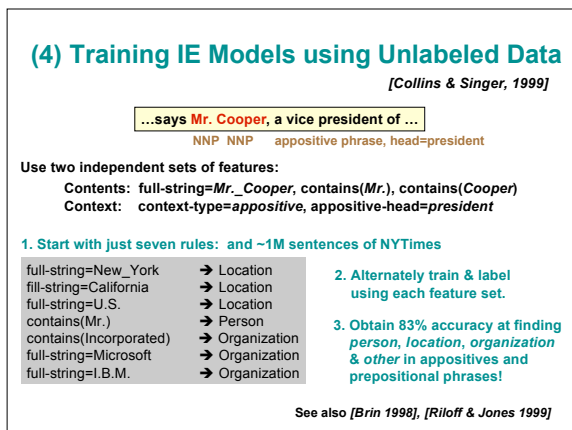
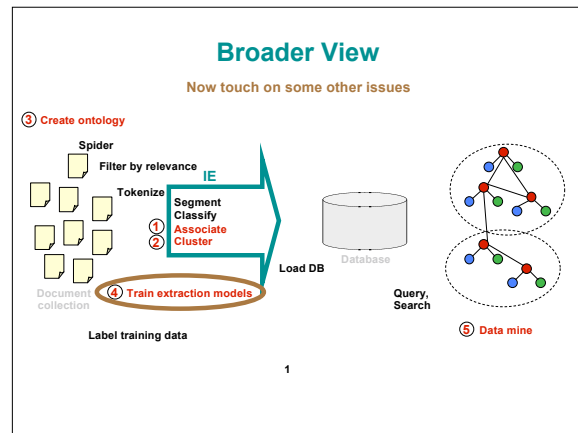
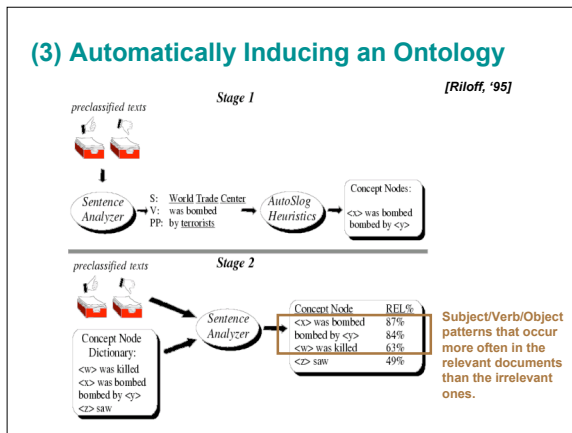
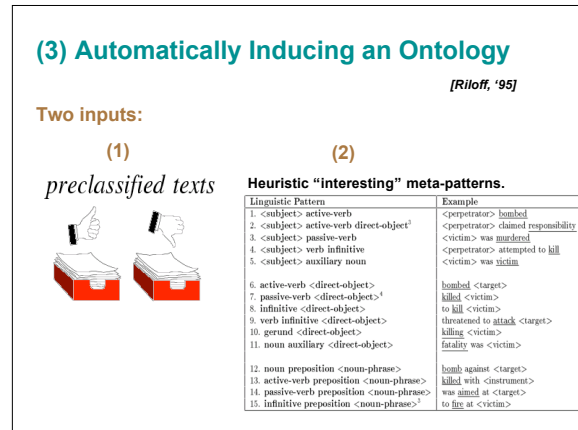
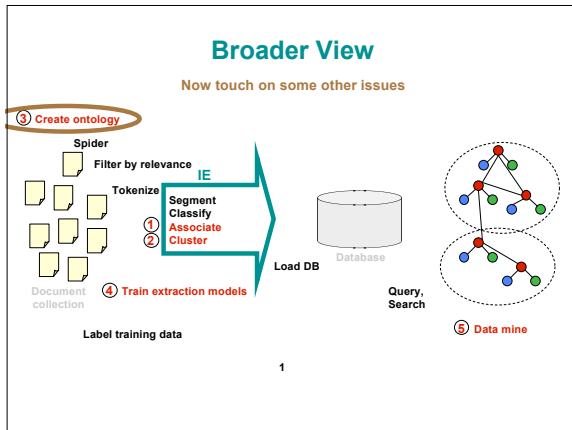
$$p^{(t+1)}(c_n|w_n, f_n; \phi) \propto p^{(t)}(f_n|c_n; \phi) p(w_n|c_n) p(c_n)$$

M-step:

$$\hat{\phi}_{c, f} = p^{(t+1)}(f|c; \phi) \propto \sum_{\{n: c_n=c, f_n=f\}} p^{(t)}(c_n|f_n, w_n)$$

Global Extractor: Precision = 46%, Recall = 75%

Scoped Learning Extractor: Precision = 58%, Recall = 75% DError = -22%



(5) Data Mining: Working with IE Data

- Some special properties of IE data:
 - It is based on extracted text
 - It is "dirty", (missing extraneous facts, improperly normalized entity names, etc.
 - May need cleaning before use
- What operations can be done on dirty, unnormalized databases?
 - Query it directly with a language that has "soft joins" across similar, but not identical keys. [Cohen 1998]
 - Construct features for learners [Cohen 2000]
 - Infer a "best" underlying clean database [Cohen, Kautz, MacAllister, KDD2000]

(5) Data Mining: Mutually supportive IE and Data Mining

[Nahm & Mooney, 2000]

Extract a large database
Learn rules to predict the value of each field from the other fields.
Use these rules to increase the accuracy of IE.

Example DB record

Filled Job Template

title: Senior DBMS Consultant
salary: Up to \$55K
state: TX
city: Dallas
country: US
language: Powerbuilder, Progress, C, C++, Visual Basic
platform: UNIX, NT
application: SQL Server, Oracle
area: Electronic Commerce, Customer Service
required years of experience: 3
desired years of experience: 5
required degree: BS

Sample Learned Rules

platform:AJX & application:Sybase &
application:DB2
→ application:Lotus Notes

language:C++ & language:C &
application:Corba &
title=SoftwareEngineer
→ platform:Windows

language:HTML & platform:WindowsNT &
application:ActiveServerPages
→ area:Database

Language:Java & area:ActiveX &
area:Graphics
→ area:Web