

Information Extraction and Integration

All 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
- Managing and utilizing text: information extraction and integration
- Scalability: a bottleneck for deployment
- Relevance to data mining community

Example: A Solution

Extracting Job Openings from the Web

Field	Value
Job Title	Ice Cream Guru
Employer	foodscience.com
Job Category	Travel/Hospitality
Job Function	Food Services
Job Location	Upper Midwest
Contact Phone	800-488-2611
Date Extracted	January 8, 2001
Source	www.foodscience.com/jobs_midwest.htm
Other Company Jobs	foodscience.com-Job1

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

Job Title	Location	Date Posted
Food Pantry Workers	Lutheran Social Services	October 11, 2002
Cooks	Lutheran Social Services	October 11, 2002
Bakers Assistants	Fine Catering by Russel Morn	October 11, 2002
Baker's Helper	Bird-in-Hand	October 11, 2002
Assistant Baker	Gourmet To Go	October 11, 2002
Host/Hostess	Sharis Restaurants	October 10, 2002
Cooks	Alta's Buncker Lodge	October 10, 2002
Line Attendant	Sun Valley Corporation	October 10, 2002
Food Service Worker II	Garden Grove Unified School District	October 10, 2002
Night Cook / Baker	SOHOCCO	October 10, 2002
Cooks/Prep Cooks	GrandView Lodge	October 10, 2002
Line Cook	Lone Mountain Ranch	October 10, 2002
Production Baker	Whole Foods Market	October 08, 2002
Cake Decorator/Baker	Mandalay Bay Hotel and Casino	October 08, 2002
Shift Supervisors	Bruegger's Bagels	October 08, 2002

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...



NAME	TITLE	ORGANIZATION
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	founder	Free Soft...

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

aka "named entity extraction"

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

NAME	TITLE	ORGANIZATION
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	Founder	Free Soft...

IE is different in different domains!

Example: on web there is less grammar, but more formatting & linking

Newsire

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.

Web

The screenshot shows the Apple website's directory structure. It includes a main navigation bar with links like 'Home', 'About Us', 'Products', 'Support', and 'Special Offers'. Below this, there are several columns of links and text, including a section for 'Apple Retail Store' and a 'Store Hours' section. The layout is clean and uses a grid system to organize the information.

Machine Learning Methods

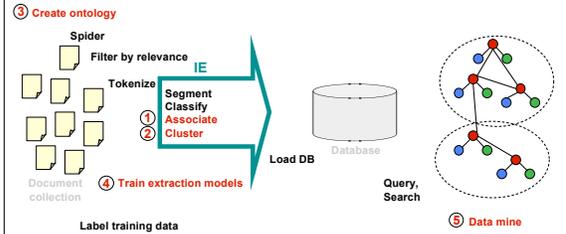
- Sequence models: HMMs, CMMs/MEMMs, CRFs
- Can work well when training data is easy to construct and is plentiful
- Can capture complex patterns that are hard to encode with hand-crafted rules
 - e.g., determine whether a review is positive or negative
 - extract long complex gene names

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 $\xrightarrow{\text{interact complex}}$ CBF-C

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

Relation Extraction

- Typically require 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
 - Learn rules/patterns from examples
 - Dan Roth 2005, Cardie 2006, Mooney 2005, ...
 - Partially-supervised: bootstrap from "seed" examples
 - Agichtein & Gravano 2000, Etzioni et al., 2004, ...
- Recently, hybrid models [Feldman2004, 2006]

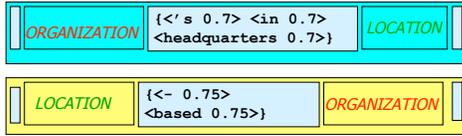
Example Extraction Rule [NYU Proteus]

```

::: 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.attributes,
  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 Extraction Patterns: Snowball [AG2000]



(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.

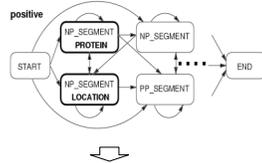
[Zelenko et al, 2002]

(1) Association with Finite State Machines

[Ray & Craven, 2001]

... This enzyme, UBC6, localizes to the endoplasmic reticulum, with the catalytic domain facing the cytosol. ...

DET this
N enzyme
N ubc6
V localizes
PREP to
ART the
ADJ endoplasmic
N reticulum
PREP with
ART the
ADJ catalytic
N domain
V facing
ART the
N cytosol

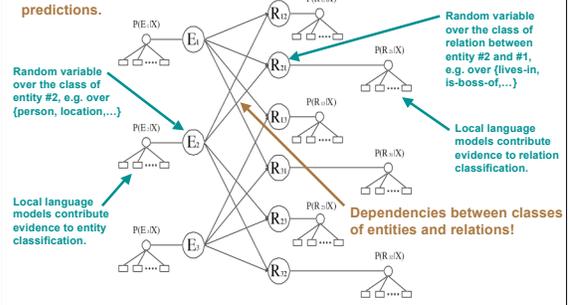


Subcellular-localization (UBC6, endoplasmic reticulum)

(1) Association with Graphical Models

[Roth & Yih 2002]

Capture arbitrary-distance dependencies among predictions.

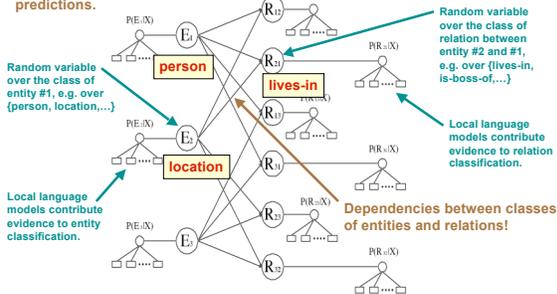


Inference with loopy belief propagation.

(1) Association with Graphical Models

[Roth & Yih 2002]

Also capture long-distance dependencies among predictions.



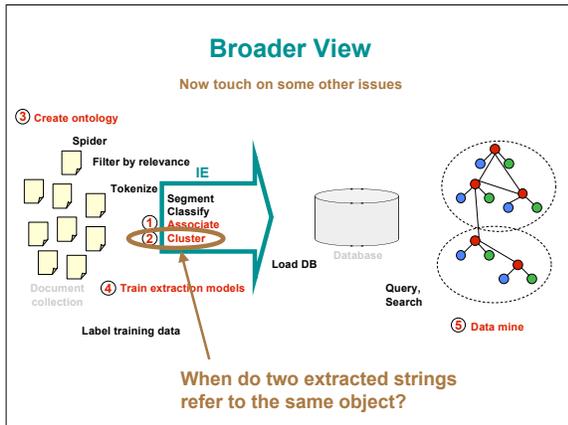
Inference with loopy belief propagation.

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 lower accuracy



Extracted Entities: Resolving Duplicates




Document 1: The Justice Department has officially ended its inquiry into the assassinations of **John F. Kennedy** and **Dr. 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