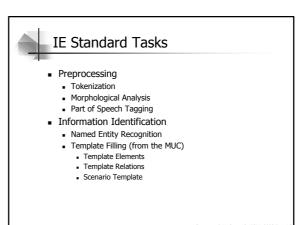
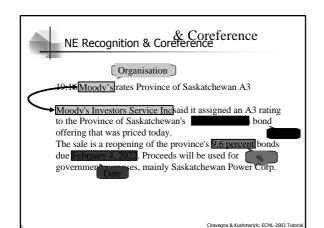
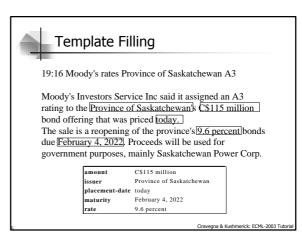
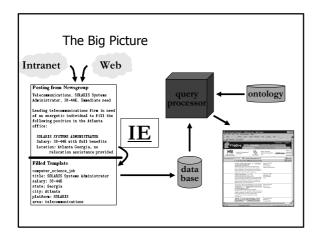


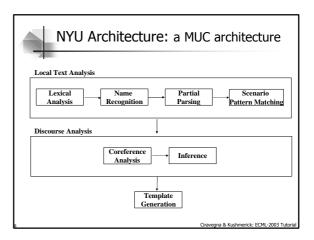
- - Document
 - newspaper article, Web page, email message, ...
 - Pre-defined "information need"
 - frame slots, template fillers, database tuples, ..
- Output
 - The specific substrings/fragments of the document or labels that satisfy the stated information need, possibly organised in a template
- DARPA's 'Message Understanding Conferences/Competitions' since late1980's; most recent: MUC-7, 1998.
- Recent interest in the machine learning and Web communities













Semantic Web

- A brain for Human Kind
- From Information-based to Knowledge-Based
- Processable Knowledge means:
 - Better Retrieval
 - Reasoning
- Where can IE contribute?



Building the SW

- Document annotation
 - Manually associate documents (or parts) to ontological descriptions
 - Document classification for retrieval
 Where can I buy an Hamster?

 - Pet shop web page -> pet shop concept -> hamster
 - Knowledge annotation
 - Where can I find a hotel in Berlin where single rooms cost less than 400€?
 - The Hotel is located in central <u>Berlin</u> and the cost for a single room is <u>300€</u>
 - Editors are currently available for manual annotation of texts



IE for Annotating Documents

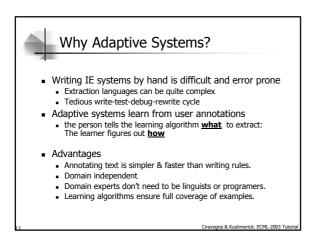
- Manual annotation is
 - Expensive
 - Error prone
- IE can be used for annotating documents
 - Automatically
 - Semi-Automatically As user support
- Advantages
 - Speed

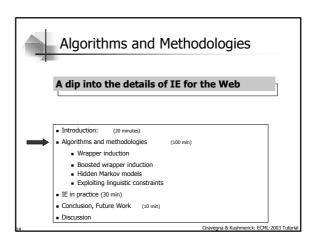
 - Low cost Consistency
 - Can provide automatic annotation different from the one provided by the author(!)

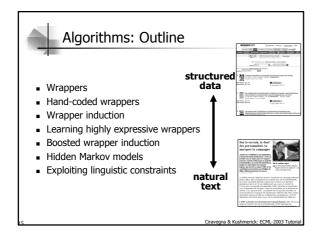


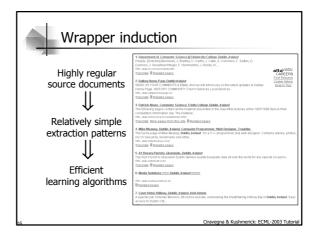
SW for Knowledge Management

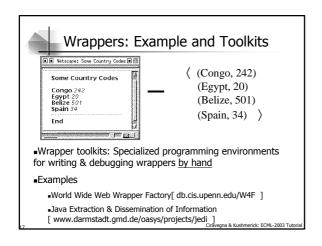
- SW is important for everyday Internet users
- SW is <u>necessary</u> for large companies
 - Millions of documents where knowledge is interspersed
 - Most documents are now
 - web-based
 - Available over an Intranet
 - Companies are valued for their
 - Tangible assets (e.g. plants)
 - Intangible assets (e.g. knowledge) Knowledge is stored in
 - mind of employees
 - Documentation
 - Companies spend 7-10% of revenues for KM

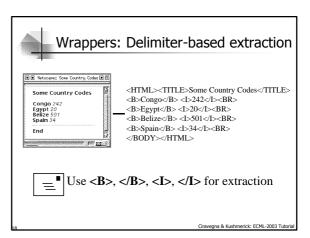


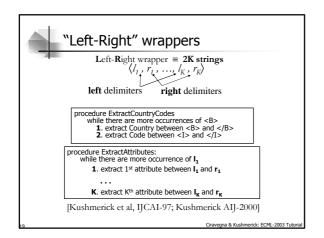


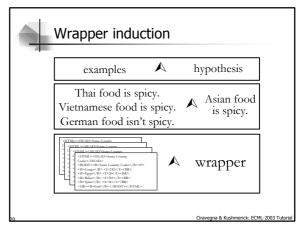


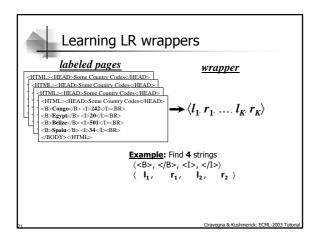


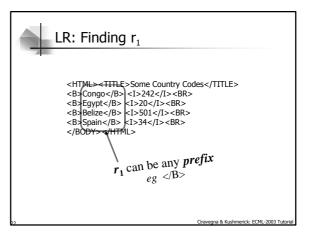


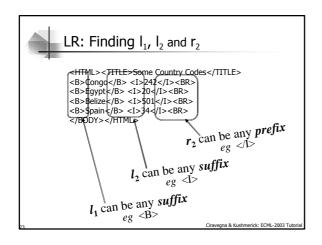


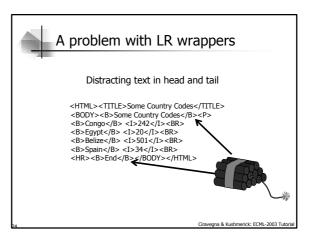


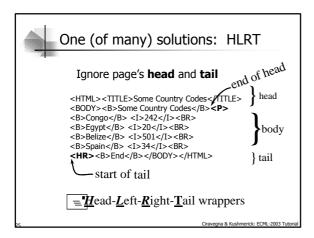


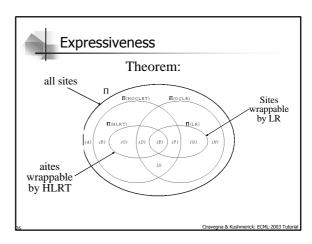


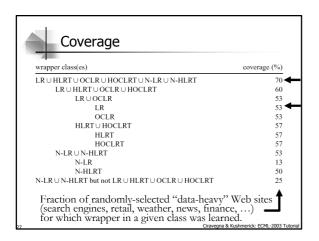


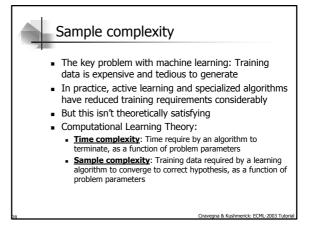


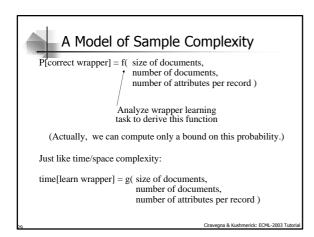


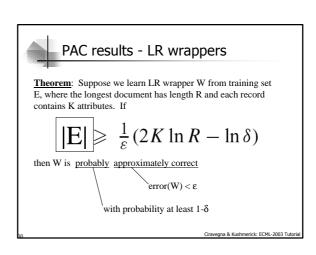














More sophisticated wrappers

- LR & HLRT wrappers are extremely simple
- Recent wrapper induction research has explored...
 - more expressive wrapper classes
 - [Muslea et al, Agents-98; Hsu et al, JIS-98; Thomas et al, JIS-00, \ldots]
 - Disjunctive delimiters
 - Sequential/landmark-based delimiters
 - Multiple attribute orderings
 - Missing attributes
 - Multiple-valued attributes
 - Hierarchically nested data
 - Wrapper verification/maintenance

[Kushmerick, AAAI-1999; Kushmerick WWWJ-00; Cohen, AAAI-1999; Minton et al, AAAI-00]

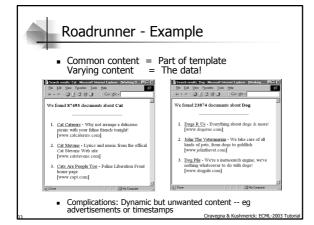


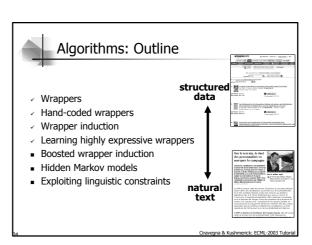
One of my favorites

Roadrunner

[Valter Crescenzi et al; Univ Roma 3]

- <u>Unsupervised</u> wrapper induction
 - They research databases, not machine learning, so they didn't realize training data was needed :-)
- Intuition:
 - Pose two different queries
 - The common bits of the documents come from the template and can be ignored
 - The bits that are different are the data that we're looking for







Boosted wrapper induction

[Freitag & Kushmerick, AAAI-00]

- Wrapper induction is suitable only for rigidly-structured machine-generated HŤMĹ...
- ... or is it?!
- Can we use simple patterns to extract from natural language documents?

... Name: Dr. Jeffrey D. Hermes Who: Professor Manfred Paul will be given by Dr. R. J. Pangborn ...

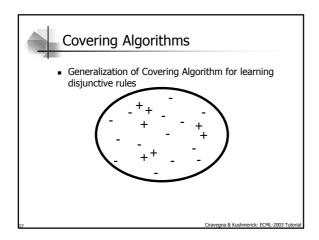
.. Ms. Scott will be speaking ...

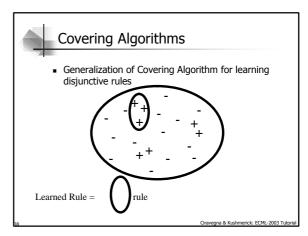
... Karen Shriver, **Dept. of ...**... Maria Klawe, **University of ...**

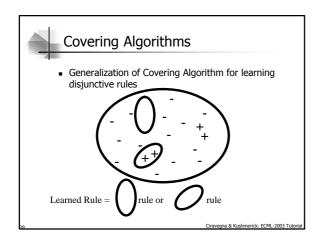


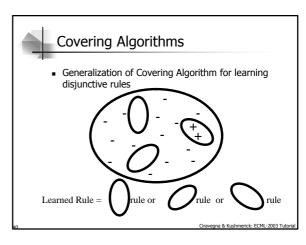
BWI: The basic idea

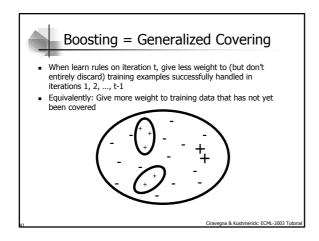
- Learn "wrapper-like" patterns for natural texts pattern = exact token sequence
- Learn many such "weak" patterns
- Combine with boosting to build "strong" ensemble pattern
- Of course, not all natural text is sufficiently regular!
- Demo: www.smi.ucd.ie/bwi

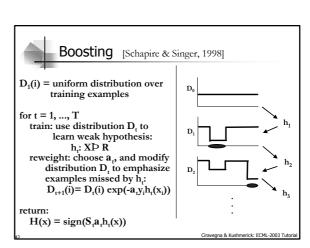














Weak hypotheses: Boundary Detectors

Boundary Detector: [who:][dr. < Capitalized>]

prefix

suffix

matches (e.g.) "... Who: Dr. Richard Nixon ..."

Weak Learning Algorithm

Greedy growth from null detector Pick best prefix/suffix extension at each step Stop when no further extension improves accuracy

 $\frac{\text{Weighting}}{a_t} = \frac{1}{2} \ln[(W^+ + e) / (W^- + e)]$

[Cohen & Singer, 1999]



Boosted Wrapper Induction

Training

input: labeled documents

Fore = Adaboost fore detectors Aft = Adaboost aft detectors Lengths = length histogram

output: Extractor = <Fore, Aft, Lengths>

Execution

input: Document, Extractor, t

 $F = \{\langle i, c_i \rangle \mid \text{token } i \text{ matches } Fore$ with confidence c_i } $A = \{\langle i, c \rangle \mid \text{ token } j \text{ matches } Aft \}$ with confidence c_i }

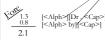
 $\begin{cases}
\langle i,j \rangle \mid \langle i, c_i \rangle \in F, \langle j, c_j \rangle \in A, \\
c_i c_j L(j-i) > t
\end{cases}$



BWI execution example

<0.26.4.95.09.31.03.bg02+@andrew.cmu.edu.0>
Type: cmu.andrew.official.cmu-news
Topic: Chem. Eng. Seminar
Dates: 2-May-95
Time: 10:45 AM
PostedBy: Bruce Gerson on 26-Apr-95 at 09:31 from andrew.cmu.edu

The Chemical Engineering Department will offer a seminar entitled "Creating Value in the Chemical Industry," at 10:45 a.m., Tuesday, May 2 in Doherty Hall 1112.
The seminar will be given by br. R. J. (Bob) Pangborn, Director, Central Research and Development, John Dow Chemical Company.





Confidence of "Dr. R. J. (Bob) Pangborn" = 2.1 0.7 0.05 = 0.074



Samples of learned patterns

Speaker: Reid Simmons, School of ... [speaker :][<Alph>] [speaker <Any>][<FName>]

Presentation Abstract Joe Cascio, IBM Set Constraints Alex Aiken (IBM, Almaden) [<Cap>][<FName> <Any> <Punc> ibm]

John C. Akbari is a Masters student at Michael A. Cusumano is an Associate Professor of Lawrence C. Stewart is a Consultant Engineer at [. <Any>][is <ANum> <Cap>]



Evaluation

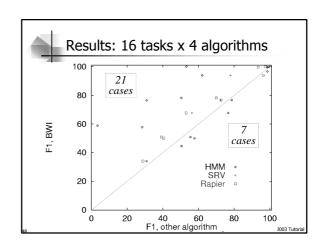
- Wrappers are usually 100% accurate, but perfection is generally impossible with natural text
- ML/IE community has a well developed evaluation
 - methodology

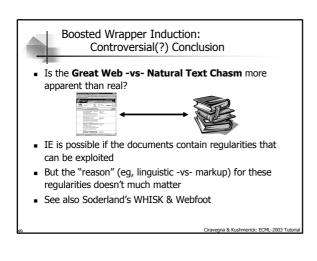
 Cross-validation: Repeat many times randomly select 2/3 of the data for training, test on remaining 1/3. **Precision**: fraction of extracted items that are correct
 - Recall: fraction of actual items extracted
 - $F_1 = 2 / (1/P + 1/R)$
- 16 IE tasks from 8 document collections

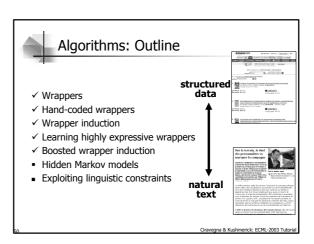
seminar announcements job listings Reuters corporate acquisitions CS department faculty lists

Zagats restaurant reviews LA Times restaurant reviews Internet Address Finder Stock quote server

· Competitors: SRV, Rapier, HMM







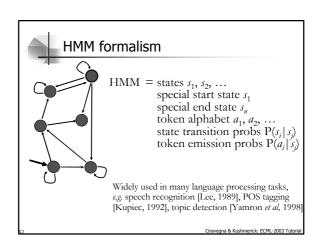


Hidden Markov models

- Previous discussion examine systems that use explicit extraction patterns/rules
- ■HMMs are a powerful alternative based on statistical token models rather than explicit extraction patterns.

[Leek, UC San Diego, 1997; Bikel et al, ANLP-97, MLJ 99; Freitag & McCallum, AAAI-99 MLIE Workshop; Seymore, McCallum & Rosenfeld, AAAI-99 MLIE Workshop; Freitag & McCallum, AAAI-2000]

Ciravegna & Kushmerick: ECML-2003 Tuto

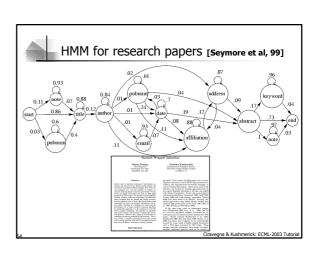


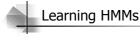


Applying HMMs to IE

- Document ⇒ generated by a stochastic process modelled by an HMM
- Token ⇒ word
- **State** ⇒ "reason/explanation" for a given token
 - 'Background' state emits tokens like 'the', 'said', ...
 - 'Money' state emits tokens like 'million', 'euro', ...
 - 'Organization' state emits tokens like 'university', 'company',
- Extraction: The Viterbi algorithm is a dynamic programming technique for efficiently computing the most likely sequence of states that generated a document.

Ciravegna & Kushmerick: ECML-2003 Tutor





Good news:

 If training data tokens are tagged with their generating states, then simple frequency ratios are a maximum-likelihood estimate of transition/emission probabilities. (Use smoothing to avoid zero probs for emissions/transitions absent in the training data.)

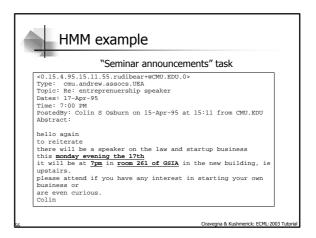
Great news:

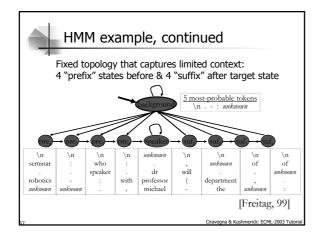
■ Baum-Welch algorithm trains HMM using <u>unlabelled</u> training data!

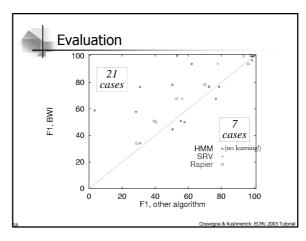
Bad news:

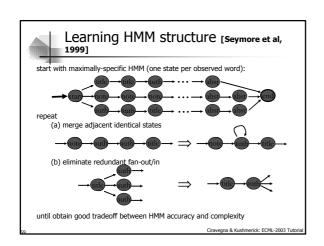
- How many states should the HMM contain?
- How are transitions constrained?
- $\, \blacksquare \,$ Insufficiently expressive \Rightarrow Unable to model important distinctions
- $\, \bullet \,$ Overly-expressive \Rightarrow sparse training data, overfitting

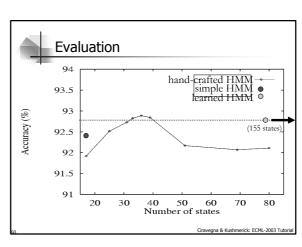
Ciravegna & Kushmerick: ECML-2003 Tutori

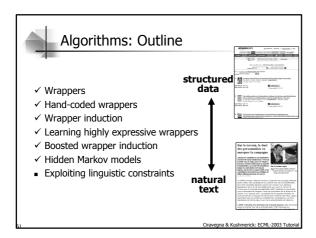


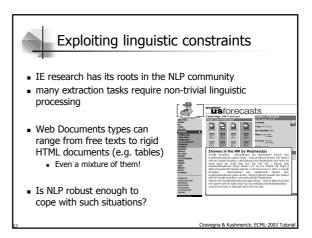














Current Approaches

- NLP Approaches (MUC-like Approaches)
 - Ineffective on most Web-related texts:
 - web pages/emails
 - stereotypical but ungrammatical texts
 - Extra-linguistic structures convey information
 - HTML tags, Document formatting, Regular stereotypical language
- Wrapper induction systems
 - Designed for rigidly structured HTML texts
 - Ineffective on unstructured texts
 - Approaches avoid generalization over flat word sequence
 - Data Sparseness on free texts

Ciravegna & Kushmerick: ECML-2003 Tutori



Lazy NLP based Algorithm

- Learns the best level of language analysis for a specific IE task mixing deep linguistic and shallow strategies
 - 1. Initial rules: shallow wrapper-like rules
 - 2. Linguistic Information (LI) progressively added to rules
 - Addition stopped when LI becomes
 - unreliable
 - ineffective
- Lazy NLP learns best strategy for each information/context separately
 - Example:
 - Using parsing for recognising the speaker in seminar announcements,
 - Using shallow approaches to spot the seminar location

Ciravegna & Kushmerick: ECML-2003 Tutori

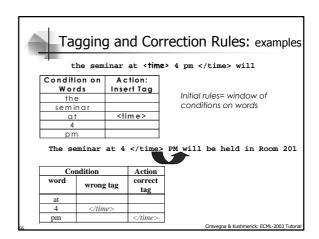


[Ciravegna 2001 – IJCAI 01- ATEM01]

- Covering algorithm based on LazyNlp
- Single tag learning (e.g. </speaker>)
- Tagging Rules
 - Insert annotation in texts
- Correction Rules
 - Correct imprecision in information identification by shifting tags to the correct position

TBL-like, with some fundamental differences

Ciravegna & Kushmerick: ECML-2003 Tuto





Rule Generalisation

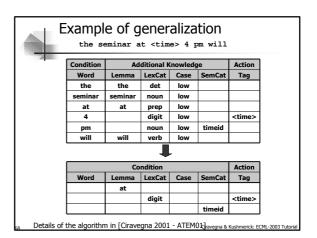
- Each instance is generalised by reducing its pattern in length
- Generalizations are tested on training corpus
- Best k rules generated from each instance reporting:
 - Smallest error rate (wrong/matches)
 - Greatest number of matches
 - Cover different examples
- Conditions on words are replaced by information from NLP modules
 - Capitalisation
 - Morphological analysis
 - Generalizes over gender/number
 - POS tagging
 - Generalizes over lexical categories
 - User-defined dictionary or gazetteer
 - Named Entity Recognizer

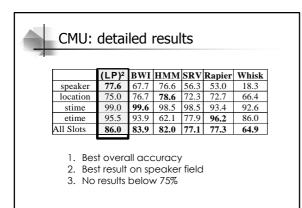
Ciravegna & Kushmerick: ECML-2003 Tuto

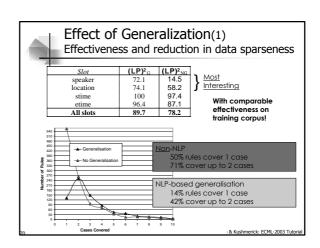
Implemented as a

general to specific beam search with

pruning (AQ-like)









Best level of Generalization

- ITC seminar announcements (mixed Italian/English)
 - Date, time, location generally in Italian
 - Speaker, title and abstract generally in English
 - English POS also for the Italian part
 - NLP-based outperforms other version

	Words	POS	NE
speaker	74.1	75.4	84.3
title	62.8	62.4	62.8
date	90.8	93.4	93.9
time	100	100	100
location	95.0	95.0	95.5

Ciravegna & Kushmerick: ECML-2003 Tuto

Linguisti

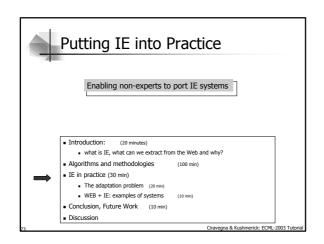
Linguistic constraints: Conclusions

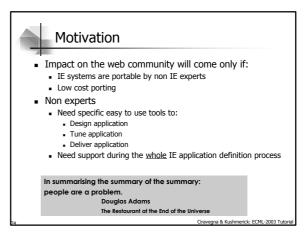
- Linguistic phenomena can't be handled by simple wrapperlike extraction patterns
- Even shallow linguistic processing (eg POS tagging) can improve performance dramatically.
 - NOTE: linguistic processing must be regular, not necessarily correct!
 - Example

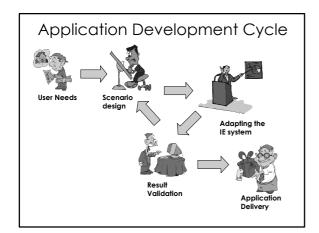
 $\label{lem:cat:NNP} $$ (LexCat:NNP +
 +
) <SPEAKER>(NER:<person>) $$ none of the covered 32 examples starts actually with an NNP $$$

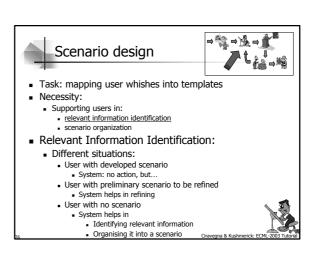
- What about more sophisticated NLP techniques?
 - Extension to parsing and corefernce resolution?

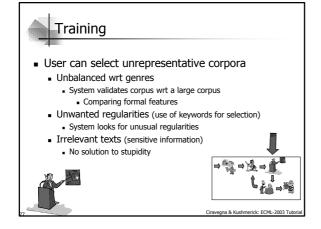
Ciravegna & Kushmerick: ECML-2003 Tuto

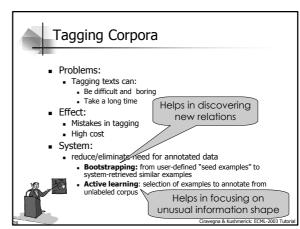














Result Validation

accurate!

threshold

- How well does the system perform?
 - Solution:
 - Facilities for:
 - Inspecting tagged corpus
 - Showing details on correctness
 - Statistics on corpus
 - Details on errors (highlight correct/incorrect/missing)
 (e.g. MUC scorer is an excellent tool)
- Influencing system behavior
 - Solution
 - Interface for bridging the user's qualitative vision and the system's numerical vision



Application Delivery



- Problem:
 - Incoming texts deviate from training data
 - Training corpus non representative
 - Document features change in time
- Solution:
 - Monitoring application.
 - Warn user if incoming texts' features are statistically different from training corpus:
 - Formal features: texts length, distribution of nouns
 - Semantic features: distribution of template fillers



Putting IE into Practice (2)

Some examples of Adaptive User-driven IE for real world applications



Learning Pinocchio

- · Commercial tool for adaptive IE
 - Based on the (LP)2 algorithm
 - Adaptable to new scenarios/applications by:
 - Corpus tagging via SGML
 - A user with analyst's knowledge
- Applications
 - "Tombstone" data from Resumees (Canadian company) (E)
- IE from financial news (Kataweb) (I) IE from classified ads (Kataweb) (I)
- Information highlighting (intelligence) (Many others I have lost track of...)
- A number of licenses released around the world for

application development

[Ciravegna 2001 - IJCAI] http://tcc.itc.it/research/textec/tools-resources/learningpinocchio/



Application development time

Resumees:

- Scenario definition: 10 person hours Tagging 250 texts: 14 person hours
- Rule induction: 72 hours on 450MHz computer
- Result validation: 4 hours

Contact:

Alberto Lavelli ITC-Irst lavelli@itc.it

http://tcc.itc.it/research/textec/tools-resources/learningpinocchio/



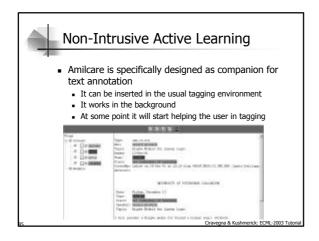
Amilcare

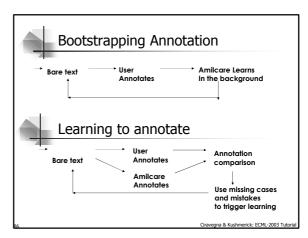
active annotation for the Semantic Web

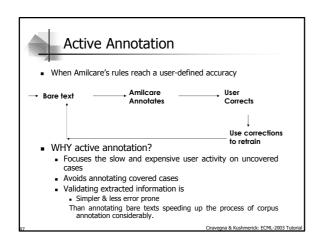


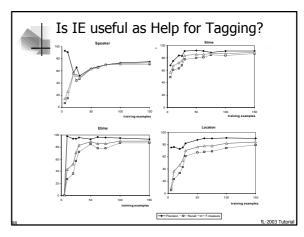
- Tool for adaptive IE from Web-related texts
 - Based on (LP)²
 - Uses Gate and Annie for preprocessing
 - Effective on different text types
 - From free texts to rigid docs (XML,HTML, etc.) Integrated with
 MnM (Open University) Ontomat (University of Karlsruhe)
 - - Gate (U Sheffield)
- Adapting Amilcare:
- Define a scenario (ontology)
 - Define a Corpus of documents
 - Annotate texts
 - Via MnM, Gate, Ontomat ■ Train the system
 - Tune the application (*) Deliver the application
- Ciravegna & Kushmerick: ECML-2003 Tu

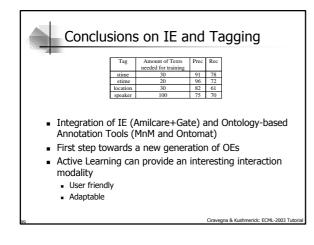
[Ciravegna 2002 -SIGIR] www.dcs.shef.ac.uk/~fabio/A

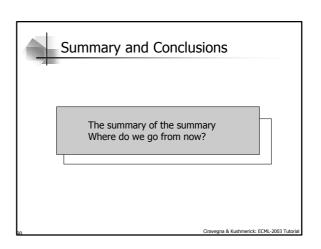














Summary

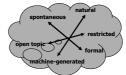
- Information extraction:
 - core enable technology for variety of next-generation information
 - Data integration agents

 - Semantic WebKnowledge Management
- Scalable IE systems must be adaptive
 - automatically learn extraction rules from examples
- Dozens of algorithms to choose from
- State of the art is 70-100% extraction accuracy (after hand-tuning!) across numerous domains.
 - Is this good enough? Depends your application.
- Yeah, but does it really work?!
 - Several companies sell IE products.
 - SW ontology editors start including IE



Open issues, Future directions

- Knob-tuning will continue to deliver substantial incremental performance increments
- Grand Unified Theory of text "structuredness", to automatically select optimal IE algorithm for a given task





■ Cross-Document Extraction





Open issues, Future directions

- Adaptive only?
- Mentioned systems are designed for non experts
 - E.g. do not require users to revise or contribute rules.
 - Is this a limitation? What about experts or even the whole spectrum of skills?
 - Future direction: making the best use of user's knowledge
- Expressive enough?
 - What about filling templates?

 - Coreferences
 (ACME is producing part for YMB Inc. The company will deliver...)

 - Reasoning (if X retires then X leaves his/her company)