

LECTURE 24:

DOCUMENT AND WEB APPLICATIONS

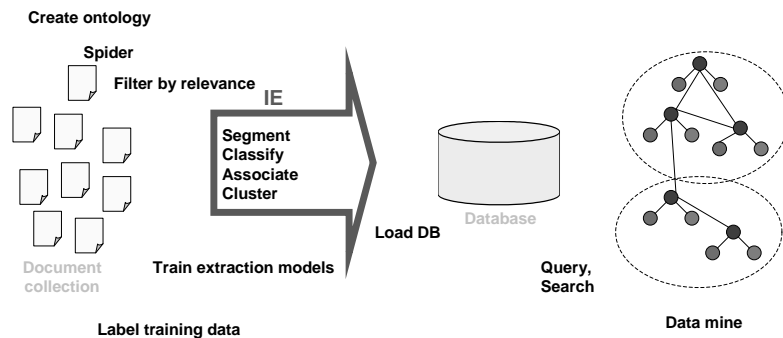
Sam Roweis

March 31, 2004

MACHINE LEARNING PROBLEMS FOR TEXT/WEB DATA

- Document / Web Page Classification or Detection
 1. Does this document/web page contain an example of thing X?
e.g. Job advertisements (FlipDog).
 2. Is this document/web page of type Y?
e.g. Course homepages (US government)?
 3. Should we block this page/document?
e.g. Spam detection, pornography content filtering.
 4. Directed crawling. Active crawling.

MACHINE LEARNING ON TEXT/WEB DATA



MACHINE LEARNING PROBLEMS FOR TEXT/WEB DATA

- Information Extraction / Entity Tagging / Disambiguation
 1. Get the author/title/company name/course title from this document or web page.
e.g. Citeseer, WebKB
 2. Find the contact info (phone number/fax/email/etc) for specific people or positions at an organization.
e.g. MarketIntelligence
 3. Find the salary/location/job title/description for a job posting on the web.
e.g. FlipDog
- Basic idea is documents → databases.

ENTITY EXTRACTION

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.

MACHINE LEARNING TECHNIQUES FOR TEXT/WEB DATA

- All of the above problems typically involve solving many separate problems in machine learning simultaneously.
 - segmentation
 - classification
 - association
 - clustering
- Why do Machine Learning on Text/Web Data?
 - Classic dream of AI: build a huge knowledge base and use it to reason about the world.
 - Cool machine learning problems.
 - Useful to companies/individuals in the real world.
- How are web documents different?
They have *links* and *rich formatting*.

MACHINE LEARNING PROBLEMS FOR TEXT/WEB DATA

- Searching / Indexing / Collaborative Filtering
 1. Sometimes called “information retrieval”.
 2. Find the most “relevant” document/product/movie/song given these search terms.
e.g. Google
 3. Find the documents/products most like the list given.
e.g. Amazon Recommendations
 4. Index collections of crosslinked items into an automatic hierarchy.
e.g. Citeseer

DOCUMENT CLASSIFICATION MODELS

- Basic models for classification:
 1. Naive Bayes
 2. Logistic Regression (MaxEnt)
 3. Support Vector Machines
 4. Decision Trees
 5. Winnow
- More sophisticated models:
 1. Mixtures of Naive Bayes
 2. Latent Probabilistic Semantic Indexing
 3. Latent Dirichlet Allocation
 4. Boosted Decision Trees

DOCUMENT FEATURES FOR CLASSIFICATION

- Binary word occurrence
- Word counts / log counts
- Binary presence on one or more lists of names, cities, companies, states, countries, products, television shows, etc.
- TF-IDF: Term Frequency * Inverse Document Frequency
- Binary indication of "Trigger Phrases".

SOME PRACTICAL CONSIDERATIONS

- Must use stop word lists and stemming.
- Naive Bayes is an excellent model, and often very hard to beat. Always add one to your counts.
- MaxEnt/Logistic Regression: don't use iterative scaling to update parameters, use conjugate gradient instead. Always use a quadratic prior on the weights.
- Feature selection is key. Some common approaches: max mutual info with class label; most frequent non-stopwords, non-stopwords appearing in most documents.

TF-IDF

- The TF-IDF measure counts how many times a word occurs (term-frequency), but normalizes that count by the proportion of documents containing a particular word (inverse-document-frequency).

- A typical measure is:

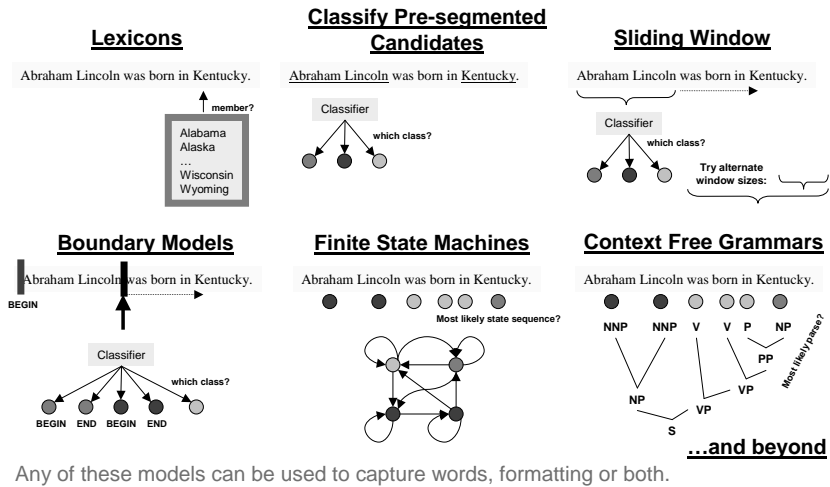
$$\text{TFIDF}(\text{word}, \text{document}) = n_{wd} \log \frac{N_{\text{documents}}}{\sum_d [n_{wd} > 0]}$$

- Very unusual words have their counts amplified, very common words have their counts multiplied by a very small number.
- Problem: hard to define this measure on a new test case or small test set. Can use IDF from training set, or pooled IDF from training plus testing sets.

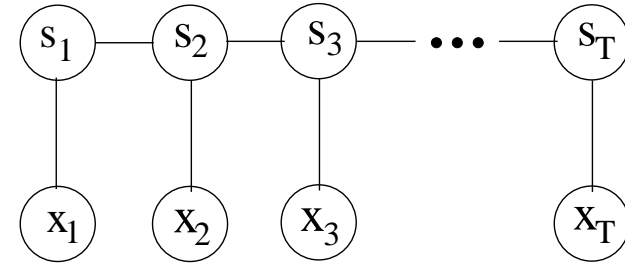
INFORMATION EXTRACTION MODELS

- Simple models: sliding windows, boundary finders.
- Use latent variable models where latent variable indicates entity groupings and observables are words or other document features.
- Hidden Markov Models
- Maximum Entropy Markov Models
- Conditional Random Fields
- Voted Perceptron
- Local-Global Models
- More sophisticated tree-based models...

INFORMATION EXTRACTION MODELS



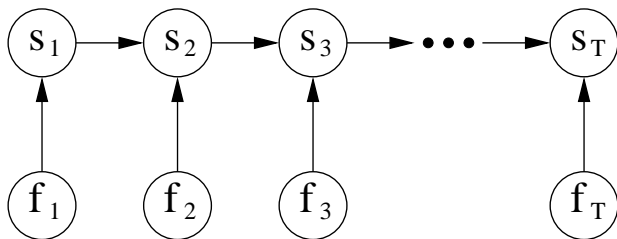
CONDITIONAL RANDOM FIELDS



How is this different than the MEMM?

Normalization is global and not local.

MAXIMUM ENTROPY MARKOV MODELS



Also known as "logistic regression through time".

What are the cliques here?

LABELLED VS. UNLABELLED DATA

- For information extraction / named-entity tagging, most models require labelled data, which can be very difficult to get in large quantities.
- Such data is often generated by a hand labelling pages and documents.
- Sometimes it is possible to "bootstrap" up from a small amount of labelled data to a larger amount.
e.g. word sense disambiguation

WEB SEARCHING MODELS

- Most web search engines work by using a combination of two technologies:
 1. Highly efficient index and search for retrieving a list of pages that contain the keywords you search for, or a set of words that mean roughly the same thing. (This is a databases problem.)
 2. A method of *ranking* the matching results so that the “best” or “most relevant” pages come earlier on the list. (This is a machine learning problem.)
- Virtually all ranking algorithms are eigenvector methods applied to the link matrix of the web.
e.g. Google, hubs and authorities
- “Bibliometrics” treats citations in documents like links on the web.

SUMMARIZATION

- Generate a small amount of text that summarizes a larger document.
e.g. Google news.
- Very hard problem because the computer has to generate some believable content.
- Easier problem: excerpt a small amount of original text or audio or video that best captures the entire document.

RECOMMENDATION/COLLABORATIVE FILTERING

- The most basic recommendation system is table-lookup into the past: given what you already like, recommend things that other people who liked what you do also liked.
- This only works if you have an enormous amount of data (e.g. weather prediction, Amazon).
- In general, we must group documents/products together based on co-occurrence and then extrapolated from our limited database to discover which items to recommend.
e.g. Aspect Model

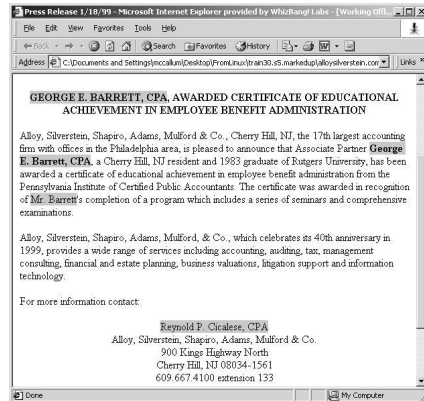
EVALUATION METRICS

- Perplexity: average number of plausible alternatives on test set.
- Precision vs. Recall curves.
- F-score: $(\beta^2 + 1)PR / (\beta^2 P + R)$
- N-best performance.
- Ranking performance.

EXAMPLE: PERSON NAME EXTRACTION

Person name Extraction

[McCallum 2001, unpublished]



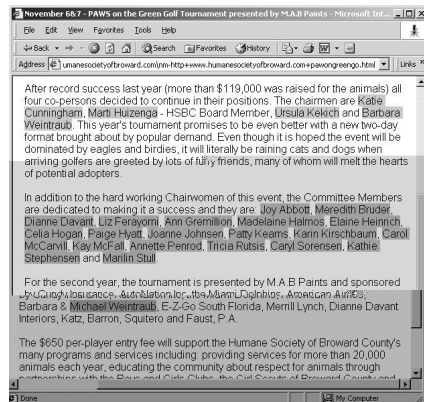
EXAMPLE: FEATURES USED

Capitalized	Xxxx	Character n-gram classifier says string is a person name (80% accurate)	Hand-built FSM person-name extractor says yes, (prec/recall ~ 30/95)
Mixed Caps	XxXxxx		Conjunctions of all previous feature pairs, evaluated at the current time step.
All Caps	XXXXX	In stopword list (the, of, their, etc)	Conjunctions of all previous feature pairs, evaluated at current step and one step ahead.
Initial Cap	X...	In honorific list (Mr, Mrs, Dr, Sen, etc)	All previous features, evaluated two steps ahead.
Contains Digit	xxx5	In person suffix list (Jr, Sr, PhD, etc)	All previous features, evaluated one step behind.
All lowercase	xxxx	In name particle list (de, la, van, der, etc)	
Initial	X	In Census lastname list; segmented by P(name)	
Punctuation	...!(), etc	In Census firstname list; segmented by P(name)	
Period	.	In locations lists (states, cities, countries)	
Comma	,	In company name list ("J. C. Penny")	
Apostrophe	'	In list of company suffixes (Inc, & Associates, Foundation)	
Dash	-		
Preceded by HTML tag			

Total number of features = ~200k

EXAMPLE: PERSON NAME EXTRACTION (2)

Person name Extraction



MORE RESOURCES

- Data
 - RISE, <http://www.isi.edu/~muslea/RISE/index.html>
 - Linguistic Data Consortium (LDC)
 - Penn Treebank, Named Entities, Relations, etc.
 - <http://www.biostat.wisc.edu/~craven/ie>
 - <http://www.cs.umass.edu/~mccallum/data>
- Code
 - TextPro, <http://www.ai.sri.com/~appelt/TextPro>
 - MALLET, <http://www.cs.umass.edu/~mccallum/mallet>
- Both
 - <http://www.cis.upenn.edu/~adwait/penntools.html>
 - <http://www.cs.umass.edu/~mccallum/ie>