

Topic Segmentation

David Kauchak
cs458
Fall 2012

Administrative

- Schedule
- Readings
- Lunch today!
- HW4 due tomorrow
- Attendance

Today's class

- Blend of introductory material and research talk
 - Problem of topic segmentation
 - Common data to work with
 - Approaches
 - Evaluation
 - Some initial results
- Represents a change in the direction of the course

[Data: Narrative texts]

Books

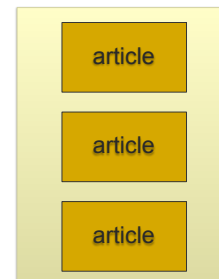
In the early nineteen seventies, a British photo retoucher named Robert Stevens arrived in south Florida to take a job at the National Enquirer, which is published in Palm Beach County. At the time, photo retouchers for supermarket tabloids used an airbrush (nowadays they use computers) to clarify news photographs of world leaders shaking hands with aliens or to give more punch to pictures of six-month-old babies who weigh three hundred pounds. Stevens was reputed to be one of the best photo retouchers in the business. The Enquirer was moving away from stories like "I Ate My Mother-in-Law's Head," and the editors recruited him to bring some class to the paper. Thy offered him much more than he made working for tabloids in Britain.

Stevens was in his early thirties when he moved to Florida. He brought a red Chevy pickup truck, and he put a CB radio in it and pasted an American-flag decal in the back window and installed a gun rack next to the flag.

Identify chapters or sections

[Data: Synthetic]

concatenate TDT articles



[How hard is this problem?]

Previous approaches have achieved error rates of 10%-20% on non-narrative data sets

(Hearst, 1994) examined the problem of paragraph identification

7 humans were asked to identify paragraphs

How well do you think people did?

Error rates were ~25%

[Data Sets]

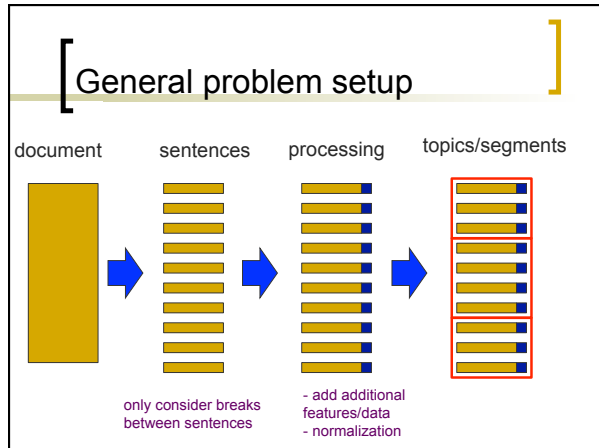
Broadcast news

Expository Texts

Narrative texts

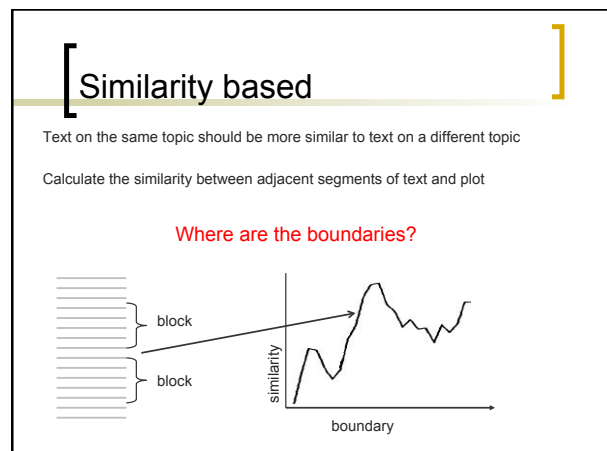
Synthetic Texts

How can we do this?



- ### [Data Set Cues]
- Broadcast news (Beeferman et al., 1999)
 - o Cues at boundaries, "Coming up..."
 - Synthetic Texts (Brants et al., 2002; Li and Yamashi, 2000)
 - o Strong topic shifts
 - Expository Texts (Hearst, 1994)
 - o Repetition of terms within segments
 - Narrative Texts (Kauchak, 2005)
 - o challenging

- ### [Previous Methods]
- Similarity based
 - Using Lexical Chains
 - Feature based



[Similarity Based]

Segment boundaries are identified by troughs in the similarities

[Similarity Measures]

Cosine Similarity of word frequency vectors

[Similarity Measures: PLSA]

Probabilistic latent semantic analysis (PLSA)

What is the probability of a word occurring in a particular document?
 $p(w,d)$

Rather than model directly, associate words with topics and documents as a blend of topics

[Similarity Measures: PLSA]

for each word, calculate the probability of occurring in each block

gives us a vector of word probabilities

[Lexical chains]

A lexical chain is a sequence of word occurrences where every word occurs within a predefined distance and each word is connected by a lexicographical relationship

Relationships: synonymy, part/whole, specialization/generalization

The dog and the cat are friends.

The cat likes the dog because they play together.

That furry feline loves to play so much.

It frolics around and around and around.

?

[Lexical chains]

A lexical chain is a sequence of word occurrences where every word occurs within a predefined distance and each word is connected by a lexicographical relationship

Relationships: synonymy, part/whole, specialization/generalization

The dog and the cat are friends.

The cat likes the dog because they play together.

That furry feline loves to play so much.

It frolics around and around and around.

[Lexical chains]

How might we use lexical chains to identify topic boundaries?

- Boundaries are located where there is a high density of chain beginnings and endings
- Boundaries have few lexical chains crossing them

[Feature based: Binary classifier setup]

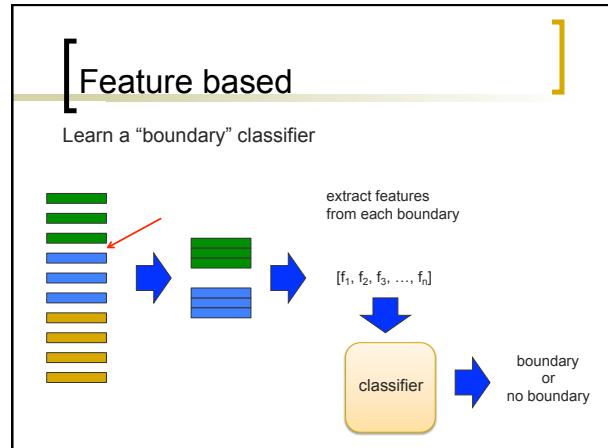
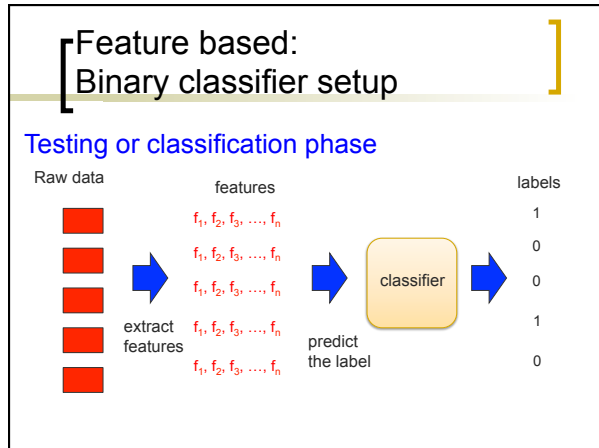
Training or learning phase

Raw data	Label	features	Label
■	0	$f_1, f_2, f_3, \dots, f_n$	0
■	0	$f_1, f_2, f_3, \dots, f_n$	0
■	1	$f_1, f_2, f_3, \dots, f_n$	1
■	1	$f_1, f_2, f_3, \dots, f_n$	1
■	0	$f_1, f_2, f_3, \dots, f_n$	0

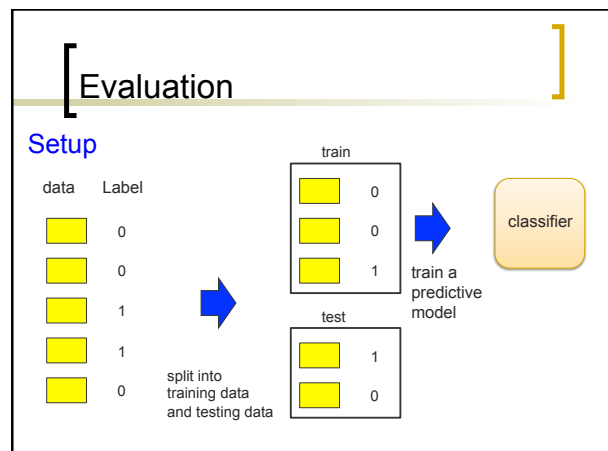
extract features →

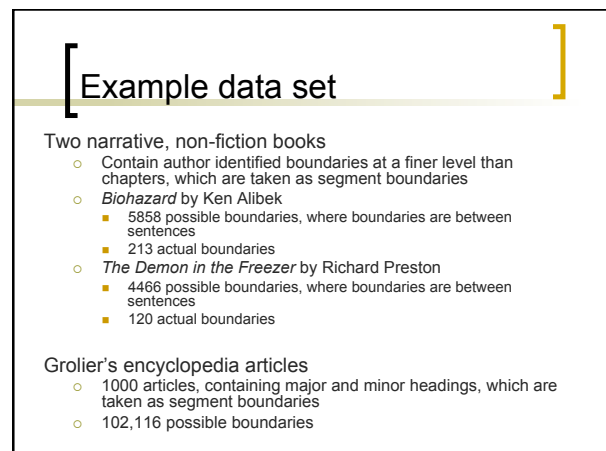
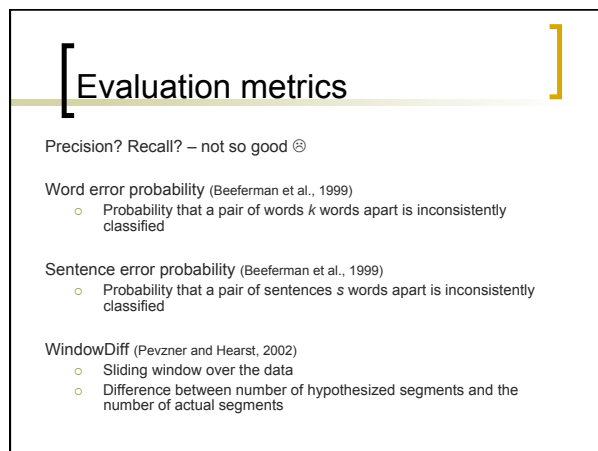
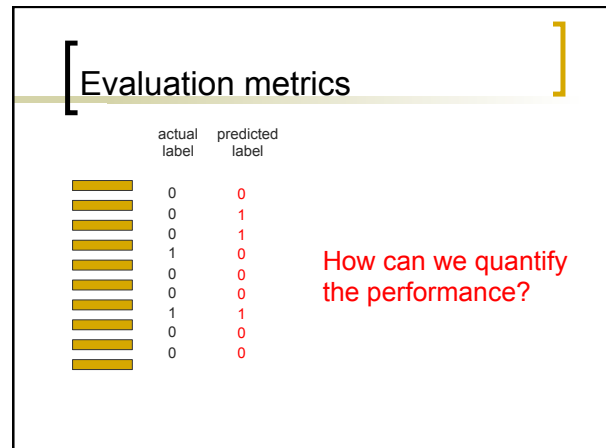
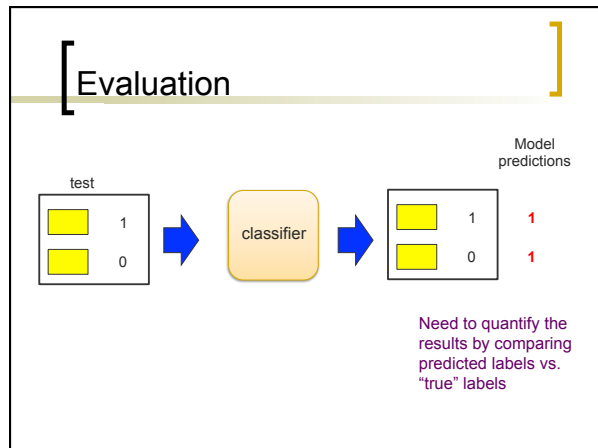
train a predictive model →

classifier



- ### Features?
- What types of things would be indicative of a shift in topic (or alternatively of staying within a topic)?
- Particular words
 - Word groups (e.g. synonyms)
 - Pronouns within 5 words from beginning
 - Lexical chain features
 - Part of a conversation?
 - Number of named entities
 - Number of synonyms to the right and left
 - Full name
 - numbers
 - ...





Performance of some approaches

Similarity

- Within segment similarity is similar to across boundary similarity
- PLSA and TextTiling (cosine similarity) perform similarly to random

Cue based

- No words occur significantly at both training and testing boundaries

Lexical chains

- Lexical chain occurrences are not correlated with boundaries

Narrative document properties

Segment Similarity

Examine adjacent block similarities and compare within segment similarities and similarities crossing segment boundaries

PLSA	Average	Standard Deviation
Within segment	0.903	0.074
Across boundary	0.914	0.041

Narrative document properties

Vocabulary

25% of the content words in the test set do not occur in the training set

33% of the content words in the test set occur two times or less

Narrative document properties

Boundary words

474 terms occur in first sentences of training boundaries


103 of these words occur at test boundaries

Only 9 *significant* words that occur in training occur in the test

No words occur significantly at both training and testing boundaries

[Narrative document properties]

Lexical chains



- Created chains using synonymy and repetition
- 219 chains
- 2 begin and 1 ends at a boundary
- 20% of the chains cross boundaries
- Average segment length is 185
- Average distance to closest beginning chain is 39 words
- Average distance to closest ending chain is 36 words
- About 4 chains per segment

[Biohazard and Demon in the Freezer (SVM)]

	Word Error	Sentence Error	Window Diff	Sent. Error improv.
<i>Biohazard</i>				
random (sent.)	0.488	0.485	0.539	
random (para.)	0.481	0.477	0.531	(baseline)
<i>Biohazard</i>				
exp1 → test	0.367	0.357	0.427	25%
exp2 → test	0.344	0.325	0.395	32%
3x cross validation	0.355	0.322	0.404	24%
Train <i>Biohazard</i>				
Test <i>Demon</i>	0.387	0.364	0.473	25%

[Grolier's results]

	Word Error	Sent. Error	Window Diff
random	0.482	0.483	0.532
Cosine Sim	0.407	0.412	0.479
PLSA	0.420	0.435	0.507
features (stumps)	0.387	0.400	0.495
features (SVM)	0.395	0.398	0.503

All methods given same number of segments (expected number of segments)

[Concluding comments]

Difficult problem: still room for improvement

Topic granularity level: **how do you predict the number of segments in the text?**

Hand crafting features is time consuming

Sequential data needs sequential models

References

- D. Beeferman, A. Berger and J. Lafferty. 1999. Statistical models for text segmentation. *Machine Learning*, 34:177-210
- T. Brants, F. Chen and I. Tsouchantaris. 2002. Topic-based document segmentation with probabilistic latent semantic analysis. In *Proceedings of CIKM*, 211-218.
- M.A. Hearst. 1994. Multi-paragraph segmentation of expository text. In *ACL*, pg. 9-16.
- H. Kozima and T. Furugori. 1994. Segmenting narrative text into coherent scenes. In *Literary and Linguistic Computing*, 9:13-19.
- H. Kozima. 1993. Text segmentation based on similarity between words. In *ACL*, pg. 286-288.
- H. Li and K. Yamanishi. 2000. Topic analysis using a finite mixture model. In *EMNLP*, pg. 35-44.
- J. Reynar. 1999. Statistical models for topic segmentation. In *ACL*, pg. 357-364.
- N. Stokes, J. Carthy, A. Smeaton. 2002. Segmenting broadcast news streams using lexical chains. In *STAIRS*, PG. 145-154/

Discussion

Google announces new search tools

<http://www.cnn.com/2012/08/08/tech/web/google-search-tools/index.html>

Good additions?

Siri?

Do you want to see your e-mail in search results?

An aside... 30 trillion unique URLs... crawls 20 billions sites per day (231,481 sites/second) ☺

Analysis of features

holdout: 74 actual boundaries and 2086 possible boundaries

	boundary	non-boundary
Paragraph	74	621
Entity groups	44	407
Word groups	39	505
Numbers	16	59
Full Name	2	109
Conversation	0	510
Pronoun	8	742
Pronoun ≤ 5	1	330

Analysis of features

Perfect recall

	boundary	non-boundary
Paragraph	74	621
Entity groups	44	407
Word groups	39	505
Numbers	16	59
Full Name	2	109
Conversation	0	510
Pronoun	8	742
Pronoun ≤ 5	1	330

[Example reweighting]

- Substantially more negative examples than positive
- Simply classifying all examples as negative results in reasonable classification performance
- Must reweight positive vs. negative examples
- Iteratively change weighting
 - Train
 - Test
 - Stop when expected number of segments based on the training data is approximately found in the test set