

Probabilistic Latent Semantic Analysis for Broadcast News Story Segmentation

Mimi LU^{1,2}, Cheung-Chi LEUNG², Lei XIE¹, Bin MA² and
Haizhou LI²

¹Shaanxi Provincial Key Lab of Speech and Image Information
Processing, Northwestern Polytechnical University, China

²Institute for Infocomm Research, A*STAR, Singapore



Broadcast news story segmentation

- The task of dividing broadcast news (BN) programs into homogeneous units each addressing a main topic



- A key precursor to various tasks, such as spoken document retrieval and summarization
- Three categories of cues for story segmentation:
 - **Lexical**, acoustic and visual cues



- Lexical cohesion based methods
 - Words in a story hang together by semantic relations
 - Different stories deploy different set of words
 - Usually measured by rigid word counts
- Literal matching on individual terms is unreliable:
 - Synonym: “car”, “automobile”; ...
 - Polysemy: “china” can refer to a nation or porcelain; “apple” can refer to Apply Computer Inc or fruit; ...
- Conceptual matching is introduced:
 - E.g. Latent semantic analysis (LSA), Probabilistic latent semantic analysis (PLSA)



- Spoken document segmentation is different from text segmentation:
 - Task performs on LVCSR outputs where erroneous words exist, thus breaking the lexical cohesion
 - Many recognition errors are from Out-Of-Vocabulary (OOV) words, which are typically name entities that are key to topics
- Phoneme n -gram: partial matching
 - Incorrectly recognized words may contain subword units correctly recognized



- We use PLSA for story segmentation for broadcast news
- We use phoneme n -gram as the basic unit for lexical cohesion measure to handle erroneous LVCSR transcripts
- Cross entropy based approach is introduced for lexical cohesion measure and it is compared with cosine similarity
- We compare dynamic programming (DP) with TextTiling for story boundary identification



- Probabilistic latent semantic analysis

$$P(d, w) = P(d)P(w | d) \qquad P(w | d) = \sum_{z \in Z} P(w | z)P(z | d)$$

d: document, w: word, z: topic

- Maximum Likelihood Estimation

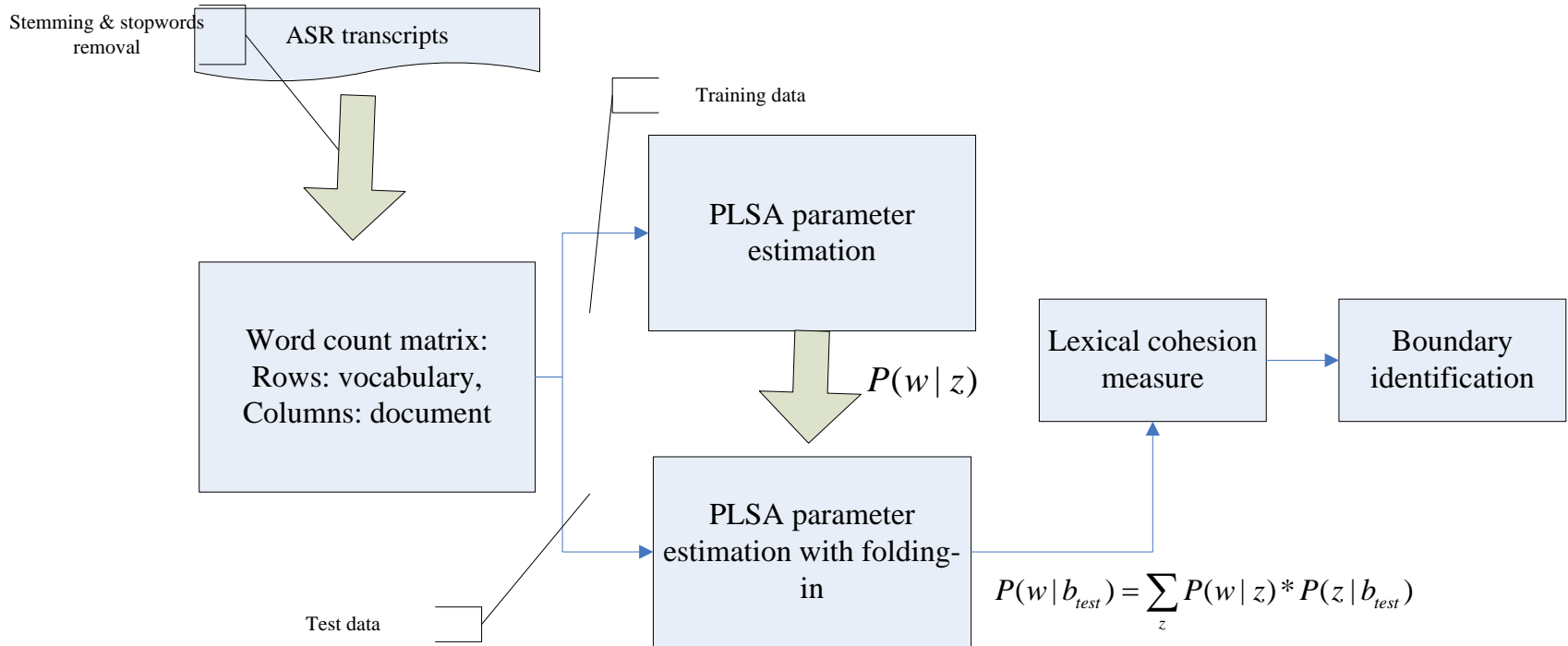
- Maximize log-likelihood of co-occurrence pairs $L = \sum_d \sum_w n(d, w) \log P(d, w)$

- E-step
$$P(z | d, w) = \frac{P(w | z)P(z | d)}{\sum_z P(w | z)P(z | d)}$$

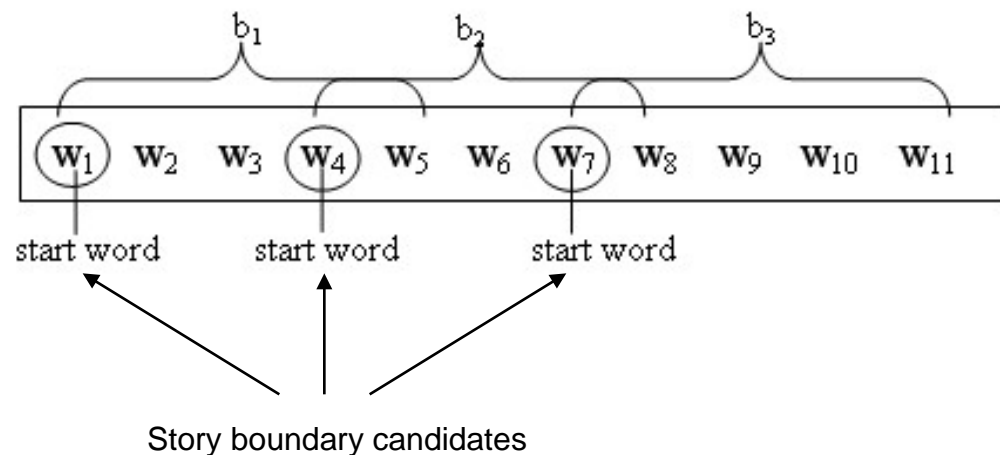
- M-step
$$P(w | z) = \frac{\sum_d n(d, w)P(z | d, w)}{\sum_w \sum_d n(d, w)P(z | d, w)} \qquad P(z | d) = \frac{\sum_w n(d, w)P(z | d, w)}{\sum_z \sum_w n(d, w)P(z | d, w)}$$

- Folding-in process for unseen test data: keep $P(w | z)$ fixed





- Sentence delimiters are not available in LVCSR transcripts
- Pseudo-sentence: each text block with a fixed number of consecutive words is formed



- Cosine similarity
 - Measure the closeness between two vectors, usually calculated on term frequencies
 - Apply with PLSA statistics:

$$Sim(i, j) = \frac{\sum_w P(w|b_i)P(w|b_j)}{\sqrt{\sum_w P(w|b_i)^2 \sum_w P(w|b_j)^2}}$$



- Cross entropy
 - A divergence measure to depict how different two distributions are

$$H(p, q) = -\sum_x p(x) \log q(x)$$

- Minimum obtained when $p = q$



- Cross entropy
 - Apply with PLSA statistics:

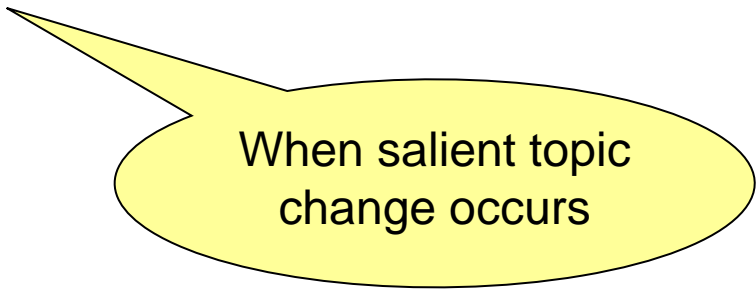
$$CrossEnt(i, j) = -\sum_w P(w|b_i) \log P(w|b_j)$$

- Normalization:

$$Dissim(i, j) = \frac{CrossEnt(i, j) - CrossEnt(i, i)}{CrossEnt(i, j)}$$



- Local comparison
 - Compute lexical scores between adjacent blocks
 - Locate valleys (similarity) or peaks (dissimilarity)
 - E.g. TextTiling

A yellow speech bubble with a black outline, pointing towards the top-left. It contains the text 'When salient topic change occurs' in black font.

When salient topic change occurs



- Global optimization

When topic transitions are smooth

- Minimize the cost of a specific segmentation

$$C(S) = \sum_{k=1}^K \text{Cost}(s_k)$$

$$\text{Cost}(s_k) = \frac{\sum_i \sum_j \text{Dissim}(i, j)}{N(\text{len}(s_k))}$$

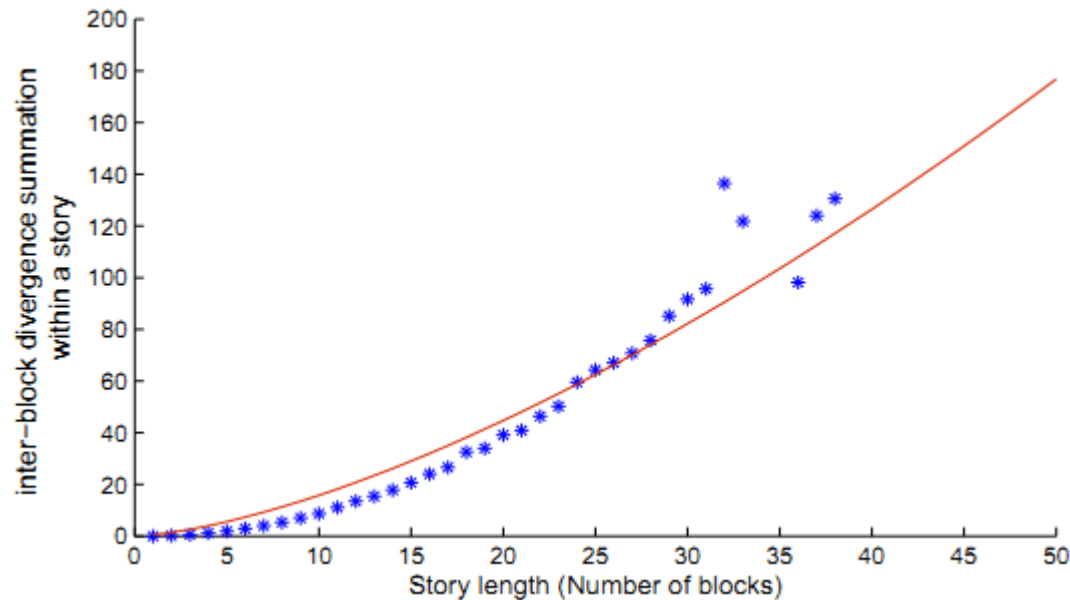
Normalization factor

$$\hat{S} = \arg \min_S C(S)$$

- $S = \{s_1, \dots, s_k, \dots, s_K\}$ is a segmentation of document D
- Implementation: dynamic programming (DP)



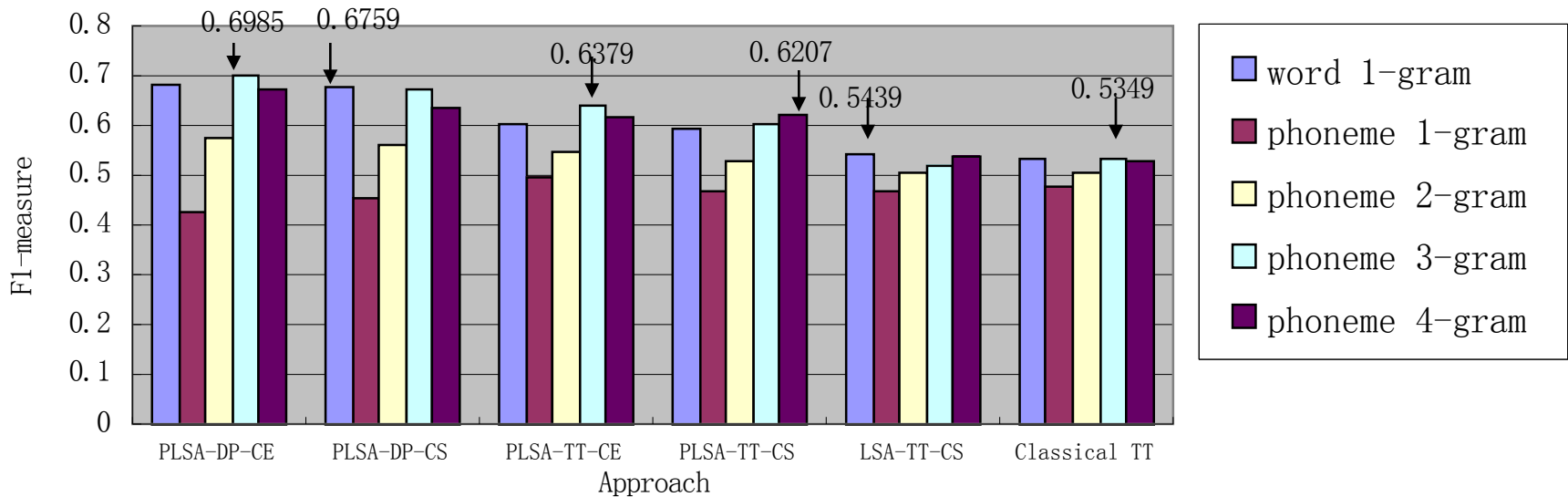
- Normalization factor $N(len(s_k))$
 - To make long and short segments comparable
 - Inter-block disparity distribution:



- **Corpus:**
 - LVCSR transcripts of TDT2 VOA English broadcast news
 - Data used (Number of programs):
training = 56, development = 27, test = 28
- **Tuning parameters:**
 - TextTiling: block length, sliding window shift, lexical score threshold
 - DP: block length, alpha in normalization factor
- **Phoneme n -gram sequences generated from word transcripts using the CMU dictionary**
- **Evaluation criterion: F1-measure**

$$F1\text{-measure} = \frac{2 * recall * precision}{recall + precision}$$





DP: dynamic programming; TT: TextTiling

CE: cross entropy;

CS: cosine similarity



- We investigate the use of PLSA for BN story segmentation
 - Phoneme subwords are adopted to address problems from LVCSR errors
 - Cross entropy and cosine similarity for lexical cohesion measure, and DP and TextTiling for story boundary identification are compared respectively
- Experimental results suggest:
 - PLSA can effectively boost story segmentation performance
 - Cross entropy shows advantages for describing distributional variation
 - DP provides better performance for story boundary identification
 - Performance gain using phoneme n -gram shows its ability to handle erroneous LVCSR transcripts



Thanks for your attention!

