

# Discover Chinese Words from Unsegmented Text

*Xianping Ge, Wanda Pratt, Padhraic Smyth*  
Information and Computer Science  
University of California, Irvine  
{xge,pratt,smyth}@ics.uci.edu

# Problem

- No explicit word boundaries (i.e., spaces) in Chinese text.
- *But*, we need the words for
  - indexing in information retrieval (IR),
  - natural language understanding.

# Solution

- Use word probabilities to segment sentences into words.
- Discover the words and their probabilities from raw, unsegmented text using the Expectation-Maximization (EM) algorithm.

## Words are running together ...

- No explicit word boundaries (i.e., spaces) in a Chinese sentence. E.g.,  
今天，墨西哥与欧盟就签署自由贸易协定问题  
在此间正式开始谈判。
- The words in the above Chinese sentence:  
今天<sup>^</sup> 墨西哥<sup>^</sup> 与<sup>^</sup> 欧盟<sup>^</sup> 就<sup>^</sup> 签署<sup>^</sup> 自由<sup>^</sup> 贸易<sup>^</sup> 协定<sup>^</sup> 问题<sup>^</sup>  
<sup>^</sup> 在<sup>^</sup> 此间<sup>^</sup> 正式<sup>^</sup> 开始<sup>^</sup> 谈判
- Compare the English sentence WITHOUT / WITH the spaces  
between words:  
“*Today, the European Union (EU) and Mexico started negotiations  
to establish a zone of free trade between them.*”  
“Today, the European Union (EU) and Mexico started  
negotiations to establish a zone of free trade between them.”

## Segment the sentence into words: the maximum-likelihood approach

- Segment the sentence  $C_1 C_2 C_3 \dots C_n$  into words  $W_1 W_2 \dots W_m$  to maximize the likelihood  $P(W_1)P(W_2) \dots P(W_m)$ .

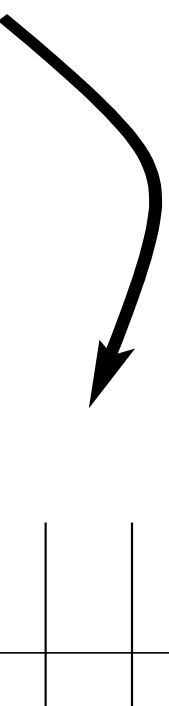
Segmentation	likelihood
$C_1 \wedge C_2 \wedge C_3$	0.01
$\vee C_1 C_2 \wedge C_3$	<u>0.09</u>
$C_1 \wedge C_2 C_3$	0.001
$C_1 C_2 C_3$	0.03

- This can be done by dynamic programming if we knew the words  $\{W_i\}$  and their probabilities  $P(W_i)$ .

## If we had a training corpus of segmented text ...

- then we can easily get the words and their probabilities!

*Count the words*

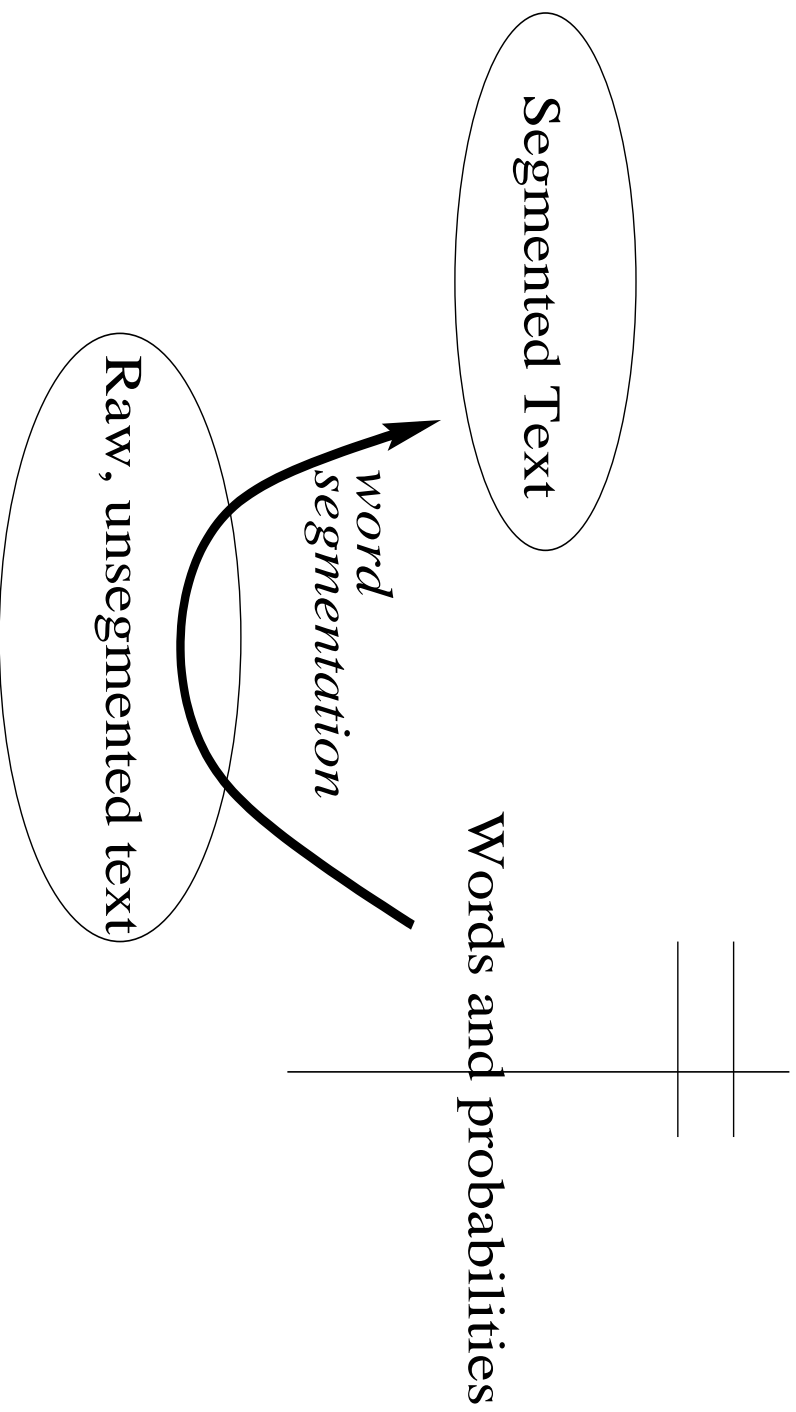


Segmented Text

Words and probabilities

## If we knew the words and their probabilities ...

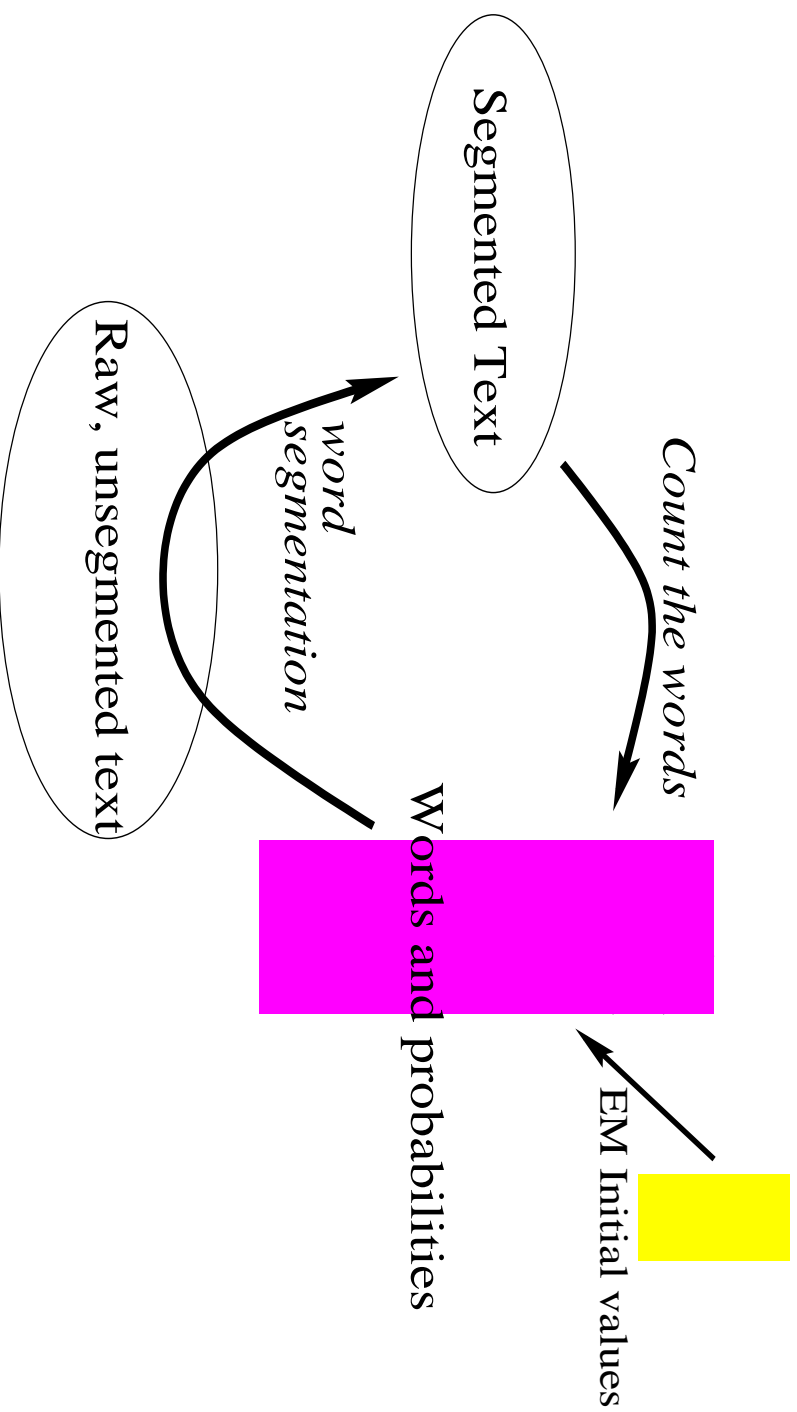
- then we can get the segmented text!



## Combine them together:

### Expectation-Maximization (EM) algorithm

- Solves this problem of “*Which came first, the chicken or the egg?*” by providing the first “egg”, i.e., the initial values for the word probabilities.

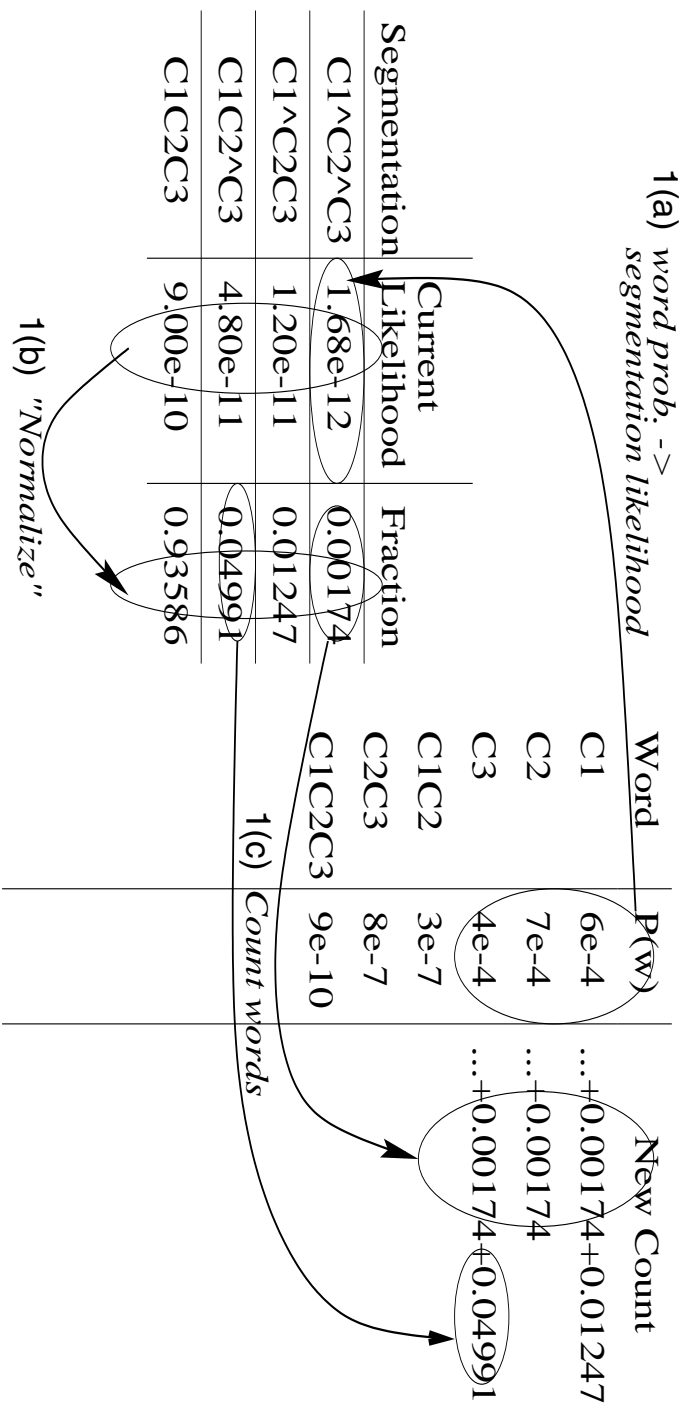




## **The main procedure**

1. For each sentence in the unsegmented text,
  - (a) Compute the likelihood of each possible segmentation using the current estimated values of the word probabilities.
  - (b) The segmentation likelihood is normalized as “fraction” that sums to 1.
  - (c) Count the words in each segmentation. I.e., add the “fraction” of the segmentation to the word count.
2. Update the word probabilities using the word counts.
3. Repeat until convergence.

# An example



## It works!

- The algorithm correctly discovers most words from the unsegmented text.
  - 但<sup>^</sup> 艺术<sup>^</sup>家<sup>^</sup> 要<sup>^</sup>尊重<sup>^</sup> 艺术<sup>^</sup> 规律
  - 城市<sup>^</sup> 是<sup>^</sup> 先进<sup>^</sup> 生产<sup>^</sup>力<sup>^</sup> 的<sup>^</sup> 集聚<sup>^</sup> 地<sup>^</sup> 和<sup>^</sup> 辐射<sup>^</sup>源
  - 香港<sup>^</sup> 特区<sup>^</sup> 是<sup>^</sup> 祖国<sup>^</sup> 的<sup>^</sup> 一个<sup>^</sup> 组成<sup>^</sup> 部分
  - 信用<sup>^</sup>社<sup>^</sup> 体制<sup>^</sup> 改革<sup>^</sup>和<sup>^</sup> 业务<sup>^</sup> 发展
- Recall/Precision=65.65%/71.91%
- After splitting single-character “stop” words from content words, Recall/Precision can be boosted up to 97.72%/91.05%

## Future work

- Some single-character stop words (like “and”, “of”, “should”, etc.) tend to cling to content words.

要尊重	要^尊重
shouldrespect	should respect

A more principled way of splitting these stop words from content words?

- Incorporate prior knowledge:
  - Distribution of word lengths
  - Existing word lists
  - Part of speech information
- Other applications

## References

- [1] Aitao Chen et al. Chinese text retrieval without using a dictionary. In *Proceedings of the 20th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 42–49, 1997.
- [2] R. Sproat and C. Shih. A statistical method for finding word boundaries in Chinese text. *Computer Processing of Chinese and Oriental Languages*, 4:336–351, March 1990.
- [3] Jay M. Ponte and W. Bruce Croft. Useg: A retargetable word segmentation procedure for information retrieval. In *Symposium on Document Analysis and Information Retrieval 96 (SDAIR)*, 1996.
- [4] Q. An and W. S. Wong. Automatic segmentation and tagging of Hanzi text using a hybrid algorithm. In *Proceedings of the 9th International Conference on Industrial & Engineering Applications of AI & Expert Systems*, June 4-7 1996.
- [5] M. R. Brent. An efficient, probabilistically sound algorithm for segmentation and word discovery. *Machine Learning* 34, 1-3, February 1999.