

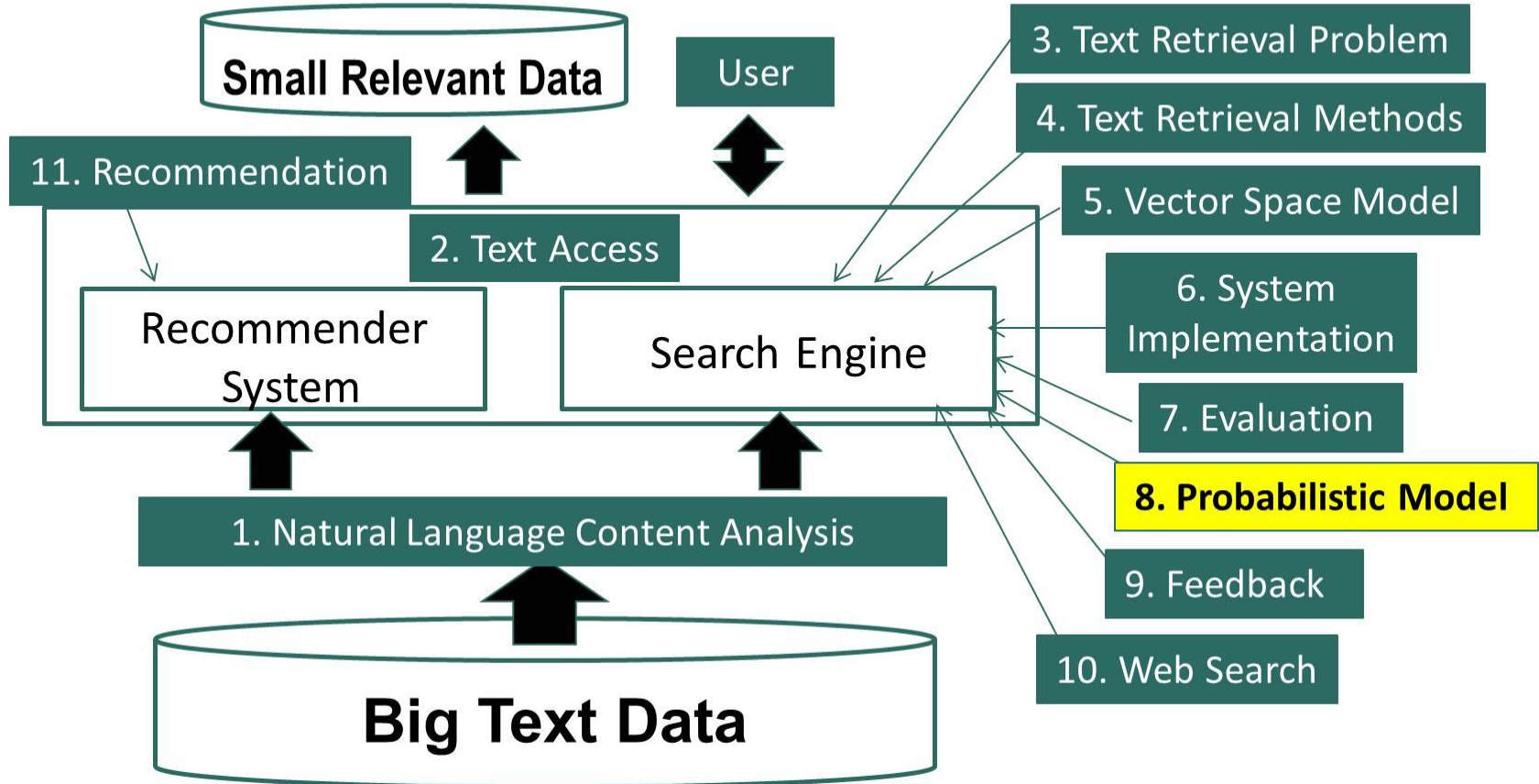


Text Retrieval and Search Engines

Probabilistic Retrieval Model: Statistical Language Model

ChengXiang “Cheng” Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

Probabilistic Retrieval Model: Statistical Language Model



Overview

- What is a Language Model?
- Unigram Language Model
- Uses of a Language Model

What is a Statistical Language Model (LM)?

- A probability distribution over word sequences
 - $p(\textit{“Today is Wednesday”}) \approx 0.001$
 - $p(\textit{“Today Wednesday is”}) \approx 0.0000000000000001$
 - $p(\textit{“The eigenvalue is positive”}) \approx 0.000001$
- Context-dependent!
- Can also be regarded as a probabilistic mechanism for “generating”



'so called a “generative” model

Today is Wednesday

Today Wednesday is

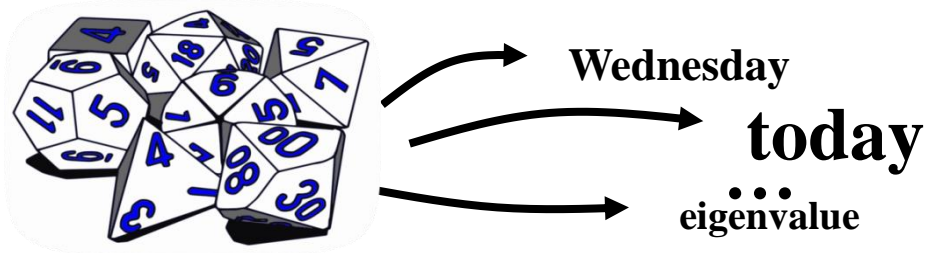
The eigenvalue is positive

Why is a LM Useful?

- Quantify the uncertainties in natural language
- Allows us to answer questions like:
 - Given that we see “*John*” and “*feels*”, how likely will we see “*happy*” as opposed to “*habit*” as the next word? (speech recognition)
 - Given that we observe “baseball” three times and “game” once in a news article, how likely is it about “sports”? (text categorization, information retrieval)
 - Given that a user is interested in sports news, how likely would the user use “baseball” in a query? (information retrieval)

The Simplest Language Model: Unigram LM

- Generate text by generating each word INDEPENDENTLY
- Thus, $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2)\dots p(w_n)$
- Parameters: $\{p(w_i)\}$ $p(w_1) + \dots + p(w_N) = 1$ (N is voc. size)
- Text = sample drawn according to this **word distribution**



$$\begin{aligned} p(\text{"today is Wed"}) \\ &= p(\text{"today"})p(\text{"is"})p(\text{"Wed"}) \\ &= 0.0002 \times 0.001 \times 0.000015 \end{aligned}$$

Text Generation with Unigram LM

Unigram LM $p(w|\theta)$ Sampling  Document =?

Topic 1:
Text mining

...
text 0.2
mining 0.1
association 0.01
clustering 0.02
food 0.00001
...



**Text mining
paper**

Topic 2:
Health

...
food 0.25
nutrition 0.1
healthy 0.05
diet 0.02
...



**Food nutrition
paper**



Estimation of Unigram LM

Unigram LM $p(w|\theta)=?$

Estimation

Text Mining Paper d

Total #words=**100**

10/100

5/100

3/100

3/100

1/100

text ?
mining ?
association
?
database ?
query ?

Maximum Likelihood (ML) Estimator:

$$p(w | \theta) = p(w | d) = \frac{c(w, d)}{|d|}$$



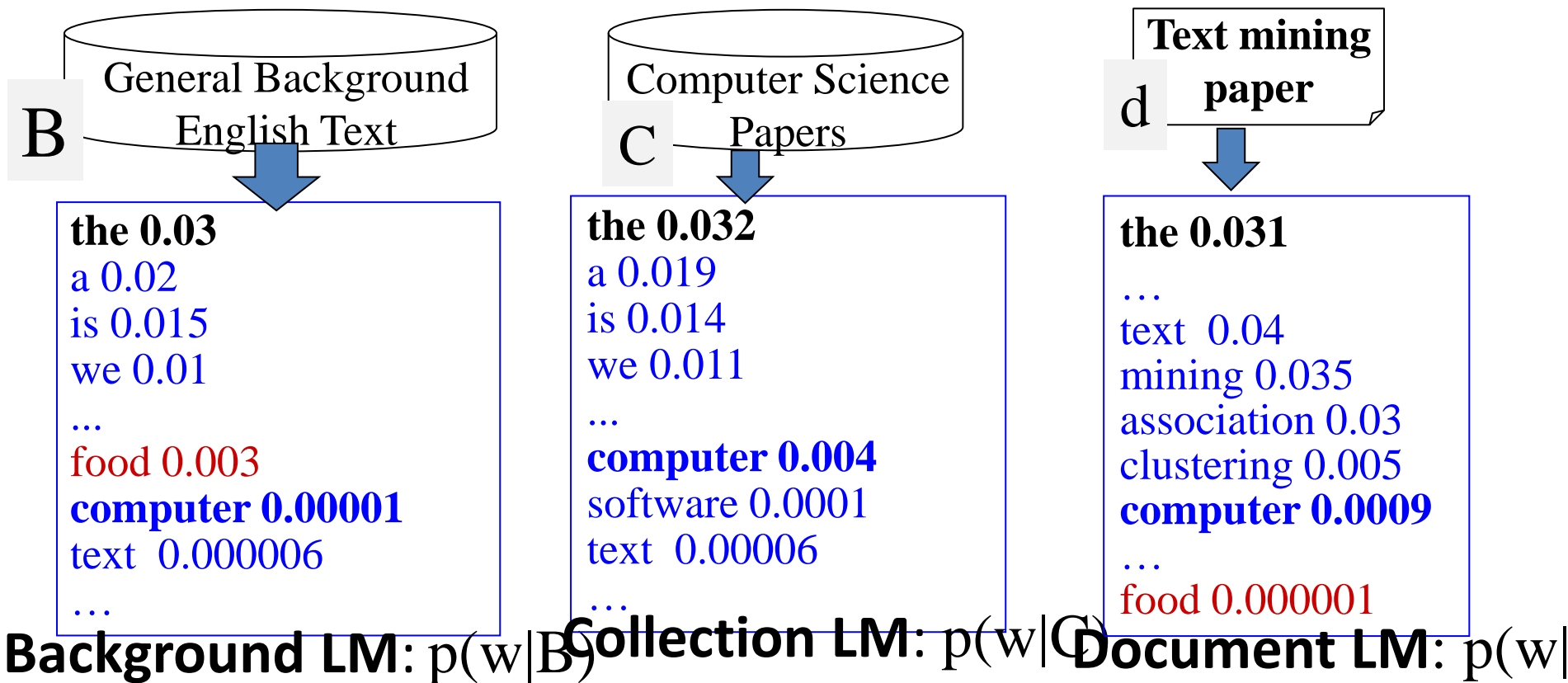
text 10
mining 5
association 3
database 3
algorithm 2

query 1
efficient 1

Is this the best estimate?

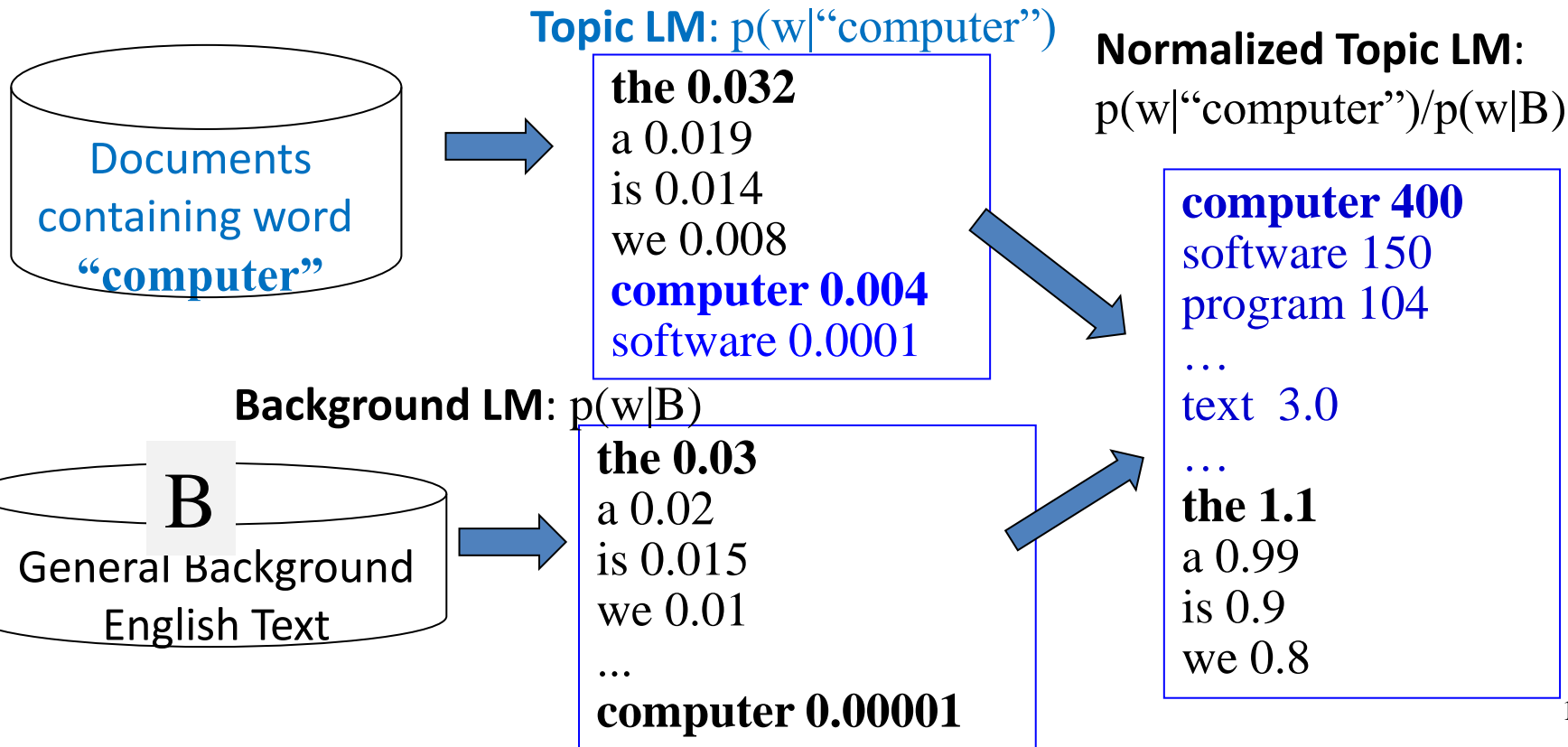


LMs for Topic Representation



LMs for Association Analysis

What words are semantically related to “computer”?



Summary

- Language Model = probability distribution over text
- Unigram Language Model = word distribution
- Uses of a Language Model
 - Representing topics
 - Discovering word associations

Additional Readings

- Chris Manning and Hinrich Schütze, Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999.
- Rosenfeld, R., "Two decades of statistical language modeling: where do we go from here?," *Proceedings of the IEEE* , vol.88, no.8, pp.1270,1278, Aug. 2000