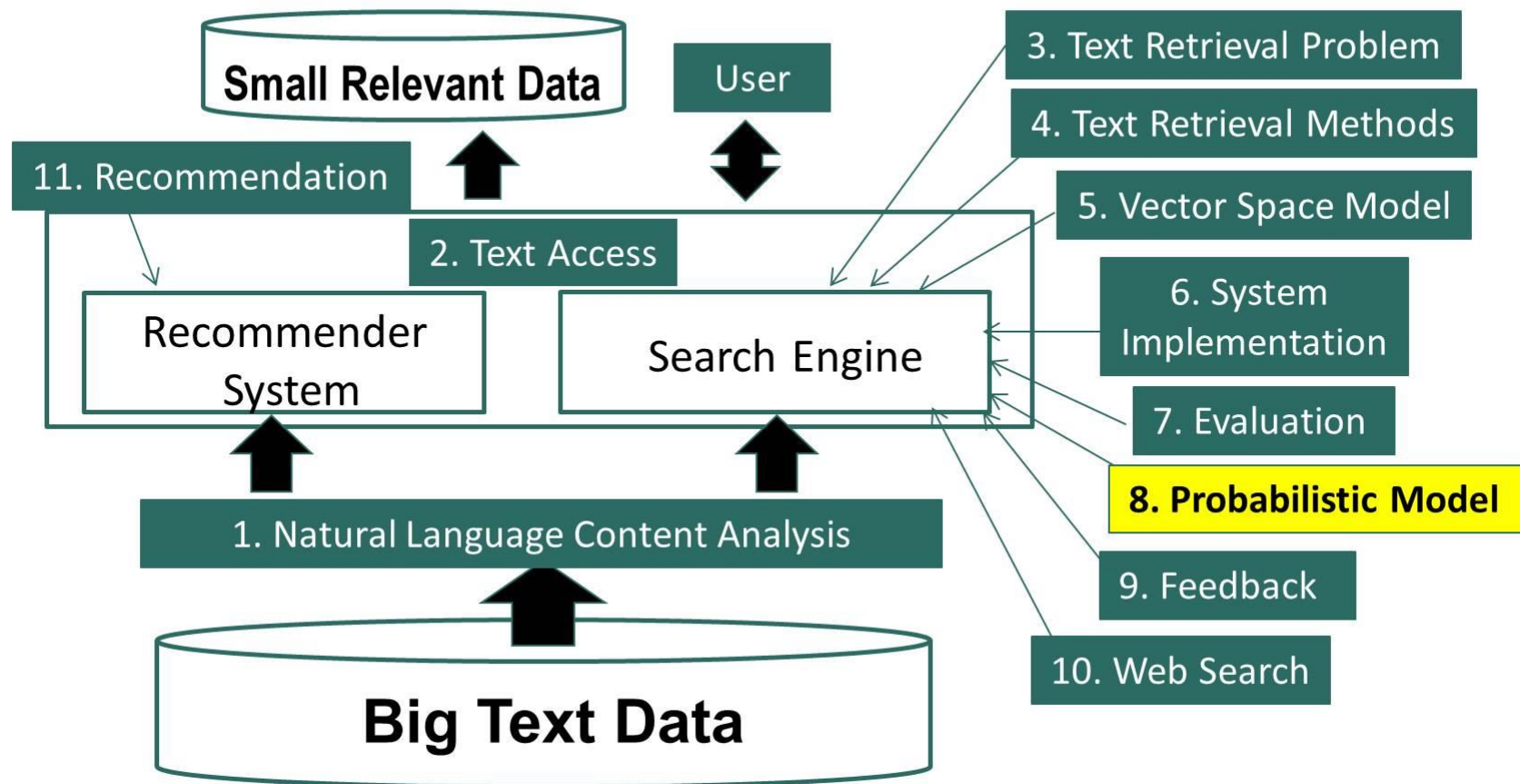# Probabilistic Retrieval Model: Basic Idea

# Probabilistic Retrieval Model: Basic Idea

# Many Different Retrieval Models

- **Probabilistic models**: f(d,q) = p(R=1|d,q),     R $\in$ {0,1}
  - Classic probabilistic model ➜ BM25
  - **Language model ➜ Query Likelihood**
  - Divergence-from-randomness model ➜ PL2

$$p(R=1|d,q)\approx p(q|d,R=1)$$

If a user likes document d, how likely would the user enter query q (in order to retrieve d)?

# Probabilistic Retrieval Models: Basic Idea

| Query | Doc | Rel |
|-------|-----|-----|
| **q** | **d** | **R** |
| q1 | d1 | 1 |
| q1 | d2 | 1 |
| q1 | d3 | 0 |
| q1 | d4 | 0 |
| q1 | d5 | 1 |
| … | | |
| q1 | d1 | 0 |
| q1 | d2 | 1 |
| q1 | d3 | 0 |
| q2 | d3 | 1 |
| q3 | d1 | 1 |
| q4 | d2 | 1 |
| q4 | d3 | 0 |

$$f(q,d)=p(R=1|d,q)=?\quad \frac{count(q,d,R=1)}{count(q,d)}$$

P(R=1|q1,d1) = ?  1/2

P(R=1|q1,d2) = ?  2/2

P(R=1|q1,d3) = ?  0/2

What about unseen documents?
Unseen queries?

# Query Likelihood Retrieval Model

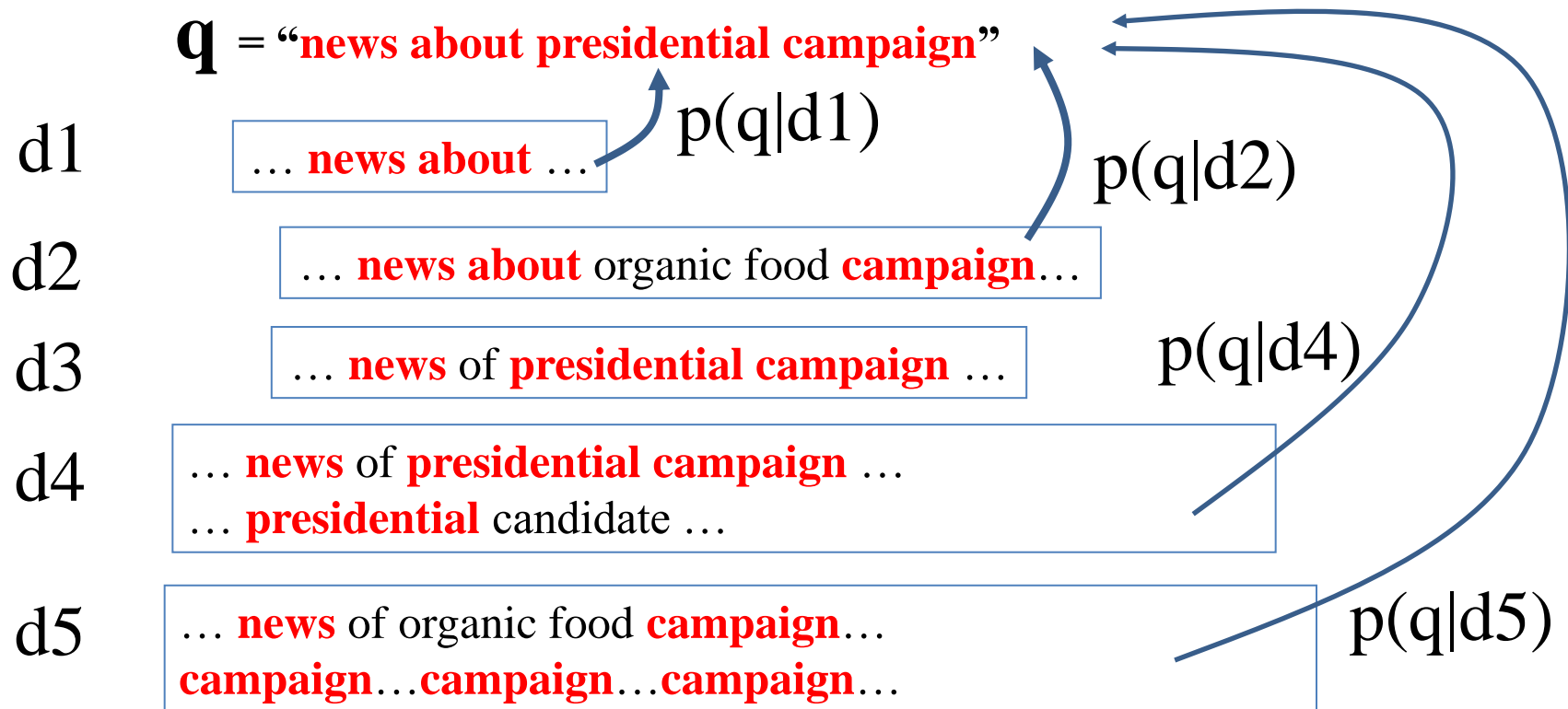| Query | Doc | Rel |
|-------|-----|-----|
| **q** | **d** | **R** |
| q1 | d1 | 1 |
| q1 | d2 | 1 |
| q1 | d3 | 0 |
| q1 | d4 | 0 |
| q1 | d5 | 1 |
| … | | |
| q1 | d1 | 0 |
| q1 | d2 | 1 |
| q1 | d3 | 0 |
| q2 | d3 | 1 |
| q3 | d1 | 1 |
| q4 | d2 | 1 |
| q4 | d3 | 0 |

User likes d

$$f(q,d)=p(R=1|d,q)\approx \mathbf{p(q|d,R=1)}$$

How likely the user enters q

Assumption:
A user formulates a query based on an "**imaginary relevant document**"

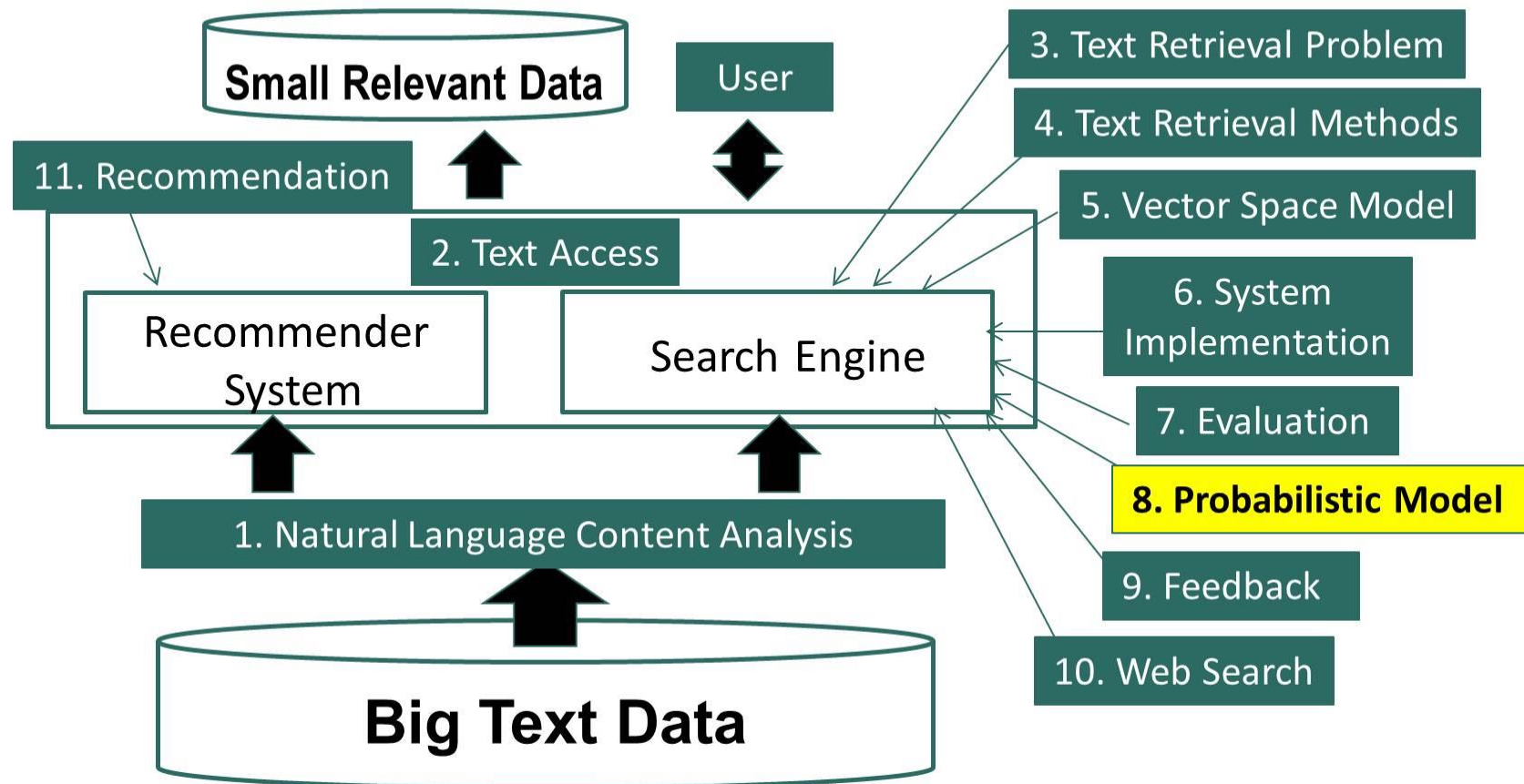# Which doc is Most Likely the "Imaginary Relevant Doc"?

**q** = "**news about presidential campaign**"

p(q|d1)

d1 | … **news about** …

p(q|d2)

d2 | … **news about** organic food **campaign**…

p(q|d4)

d3 | … **news** of **presidential campaign** …

d4 | … **news** of **presidential campaign** …
… **presidential** candidate …

p(q|d5)

d5 | … **news** of organic food **campaign**…
**campaign**…**campaign**…**campaign**…

# Summary

- Relevance(q,d) = p(R=1|q,d) ➔ p(q|d,R=1)

- **Query likelihood** ranking function: f(q,d)=p(q|d)
  - Probability that a user who likes d would pose query q

- How to compute p(q|d)? How to compute probability of text in general? ➔ Language Model

$$\mathbf{p(q}= \text{"presidential campaign"}\mathbf{|d}=\boxed{\ldots \textbf{ news } \text{of} \textbf{ presidential campaign} \ldots \textbf{ presidential} \text{ candidate} \ldots}\mathbf{)}$$

# Probabilistic Retrieval Model: Statistical Language Model

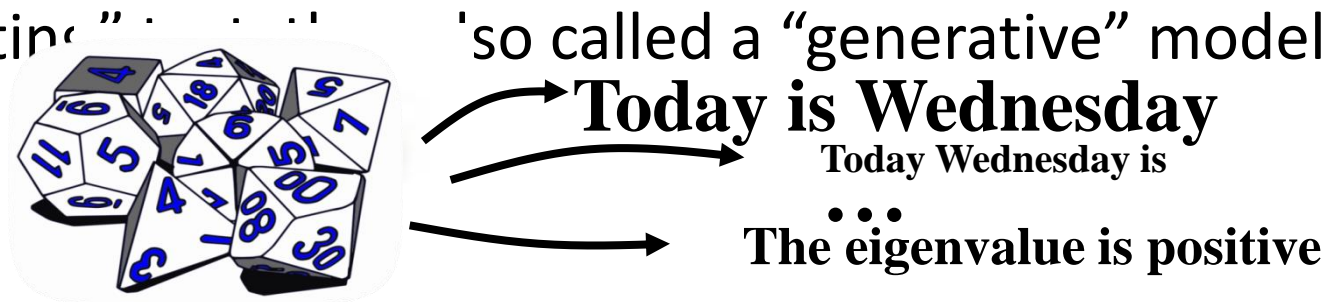# Probabilistic Retrieval Model: Statistical Language Model



Small Relevant Data

User

3. Text Retrieval Problem

4. Text Retrieval Methods

11. Recommendation

5. Vector Space Model

2. Text Access

Recommender System

Search Engine

6. System Implementation

7. Evaluation

8. Probabilistic Model

1. Natural Language Content Analysis

9. Feedback

10. Web Search

Big Text Data

# 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("*Today is Wednesday*") ≈ 0.001
  - p("*Today Wednesday is*") ≈ 0.0000000000001
  - p("*The eigenvalue is positive"*) ≈ 0.00001
- Context-dependent!
- Can also be regarded as a probabilistic mechanism for "generating" text, thus also called a "generative" model

**Today is Wednesday**

**Today Wednesday is**
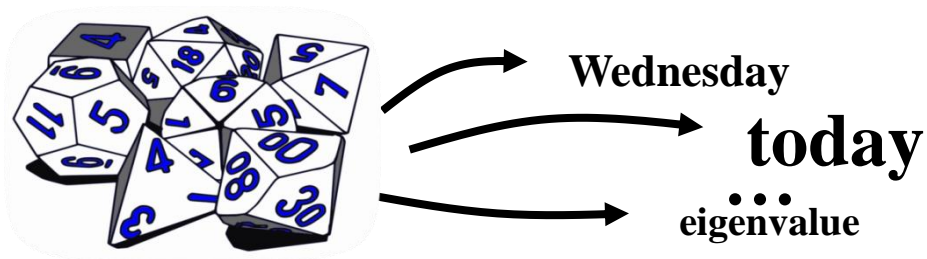
• • •

**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, $p(w_1\ w_2\ ...\ w_n) = p(w_1)p(w_2)...p(w_n)$

- Parameters: $\{p(w_i)\}$  $p(w_1)+...+p(w_N)=1$ (N is voc. size)

- Text = sample drawn according to this **word distribution**

Wednesday

**today**
...
eigenvalue

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

# Text Generation with Unigram LM

**Unigram LM  p(w|θ)**                   **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|θ)=?**

**Estimation**

**Text Mining Paper  d**

Total #words=**100**

10/100 → ...
text  ?

5/100 → mining ?

3/100 → association

3/100 → ?

database ?

1/100 → ...
query ?

text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1

**Maximum Likelihood (ML) Estimator:**

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

Is this the best estimate?

# LMs for Topic Representation



**B** General Background English Text

**the 0.03**
a 0.02
is 0.015
we 0.01
...
food 0.003
**computer 0.00001**
text  0.000006
…

**C** Computer Science Papers

**the 0.032**
a 0.019
is 0.014
we 0.011
...
**computer 0.004**
software 0.0001
text  0.00006
…

**d** Text mining paper

**the 0.031**
…
text  0.04
mining 0.035
association 0.03
clustering 0.005
**computer 0.0009**
…
food 0.000001

**Background LM**: $p(w|B)$  **Collection LM**: $p(w|C)$  **Document LM**: $p(w|d)$

# LMs for Association Analysis

## What words are semantically related to "computer"?

**Topic LM**: $p(w|\text{"computer"})$

the **0.032**
a 0.019
is 0.014
we 0.008
**computer 0.004**
software 0.0001

**Background LM**: $p(w|B)$

**B**
General Background English Text

the **0.03**
a 0.02
is 0.015
we 0.01
…
**computer 0.00001**

Documents containing word **"computer"**

**Normalized Topic LM**:

$p(w|\text{"computer"})/p(w|B)$

**computer 400**
software 150
program 104
…
text 3.0
…
**the 1.1**
a 0.99
is 0.9
we 0.8

# Summary

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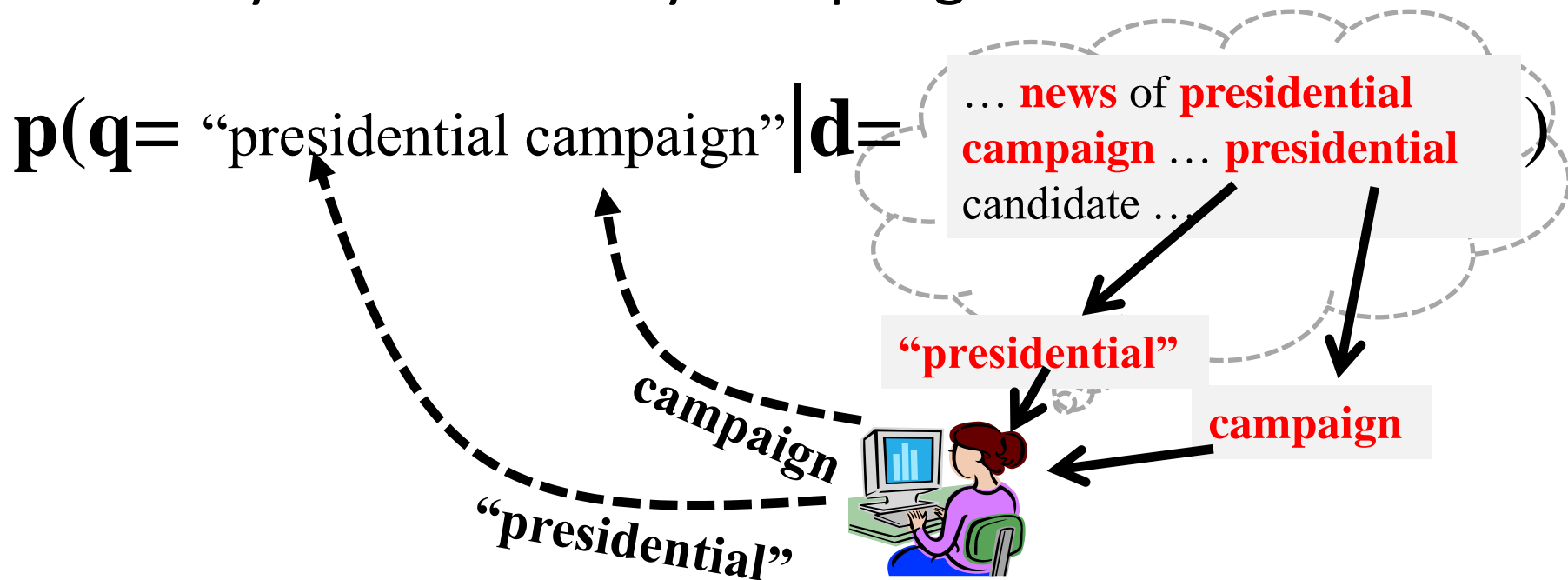
# Additional Readings

- 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

# Probabilistic Retrieval Model: Query Lilkelihood

Probabilistic Retrieval Model: Query Likelihood

# Query Generation by Sampling Words from Doc



$$p(q= \text{"presidential campaign"} | d= \ldots)$$

… **news** of **presidential campaign** … **presidential** candidate …

"presidential"

campaign

"presidential"

campaign

If the user is **thinking of this doc** ,
how likely would she **pose this query**?

# Unigram Query Likelihood

$$\mathbf{p}(\mathbf{q}=\text{ “presidential campaign”}|\mathbf{d}= \boxed{\text{… \textbf{news} of \textbf{presidential campaign} … \textbf{presidential} candidate …}} )$$

$$= \mathbf{p}(\text{“presidential ”}|\mathbf{d})*\mathbf{p}(\text{“campaign ”}|\mathbf{d})$$

$$= \frac{c("presidential", d)}{|d|} * \frac{c("campaign", d)}{|d|}$$

**Assumption:**
Each query word is generated independently

# Does Query Likelihood Make Sense?

$$p(q =" \textit{presidential campaign}"|d) = \frac{c("\textit{presidential}",d)}{|d|} * \frac{c("\textit{campaign}",d)}{|d|}$$

**p(q|d4=** [ … news of **presidential campaign** … **presidential** candidate … ] **)** $= \frac{2}{|d4|} * \frac{1}{|d4|}$

**p(q|d3=** [ … news of **presidential campaign** … ] **)** $= \frac{1}{|d3|} * \frac{1}{|d3|}$

**p(q|d2=** [ … news about organic food **campaign**… ] **)** $= \frac{0}{|d2|} * \frac{1}{|d2|} = 0$

**d4> d3 > d2** as we expected

# Try a Different Query?

$\mathbf{q}$ = "**presidential campaign update**"

$\mathbf{p(q|d4}$ = … news of **presidential campaign** … **presidential** candidate … ) $= \dfrac{2}{|d4|} * \dfrac{1}{|d4|} * \dfrac{0}{|d4|} = 0!$

$\mathbf{p(q|d3}$ = … news of **presidential campaign** … ) $= \dfrac{1}{|d3|} * \dfrac{1}{|d3|} * \dfrac{0}{|d3|} = 0!$

$\mathbf{p(q|d2}$ = … news about organic food **campaign**… ) $= \dfrac{0}{|d2|} * \dfrac{1}{|d2|} * \dfrac{0}{|d2|} = 0$

What assumption has caused this problem? How do we fix it?

# Improved Model: Sampling Words from a **Doc Model**

How likely would we observe **this query** from **this doc model**?

$$p(q = \text{"presidential campaign"} | d = )$$

… **news** of **presidential campaign** … **presidential** candidate …

"campaign"

"presidential"

"update"

**"presidential"**

**campaign**

**update**

**…**
**presidential  0.2**
**campaign 0.1**
**news 0.01**
**candidate 0.02**
**…**
**update 0.00001**
**…**

# Computation of Query Likelihood

Document
d1

**Text mining paper**

Document LM
**p(w|d1)**

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

d2

**Food nutrition paper**

**p(w|d2)**

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

**Query q =**
**"data mining algorithms"**

p("data mining alg"|d1)
=   p("data"|d1)
  × p("mining"|d1)
  × p("alg"|d1)

p("data mining alg"|d2)
=   p("data"|d2)
  × p("mining"|d2)
  × p("alg"|d2)

# **Summary**: Ranking based on Query Likelihood

$$q = w_1 w_2 \ldots w_n \qquad p(q \mid d) = p(w_1 \mid d) \times \ldots \times p(w_n \mid d)$$

$$f(q,d) = \log p(q \mid d) = \sum_{i=1}^{n} \log p(w_i \mid d) = \sum_{w \in V} c(w,q) \log \boxed{p(w \mid d)}$$
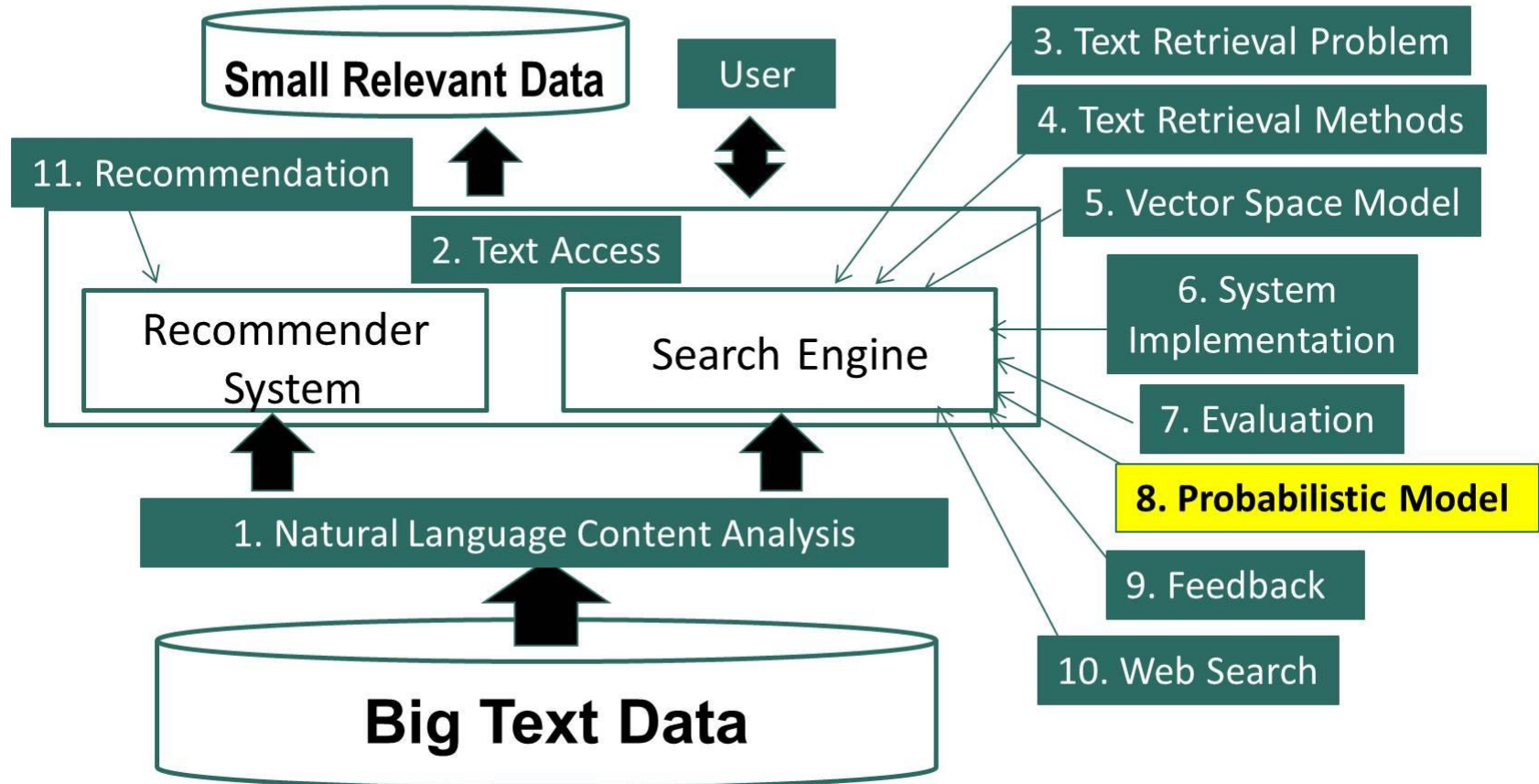
**Document language model**

Retrieval  problem  ➜ Estimation of $p(w_i/d)$

Different estimation methods ➜ different ranking functions

# Probabilistic Retrieval Model: Smoothing

# Probabilistic Retrieval Model: Smoothing
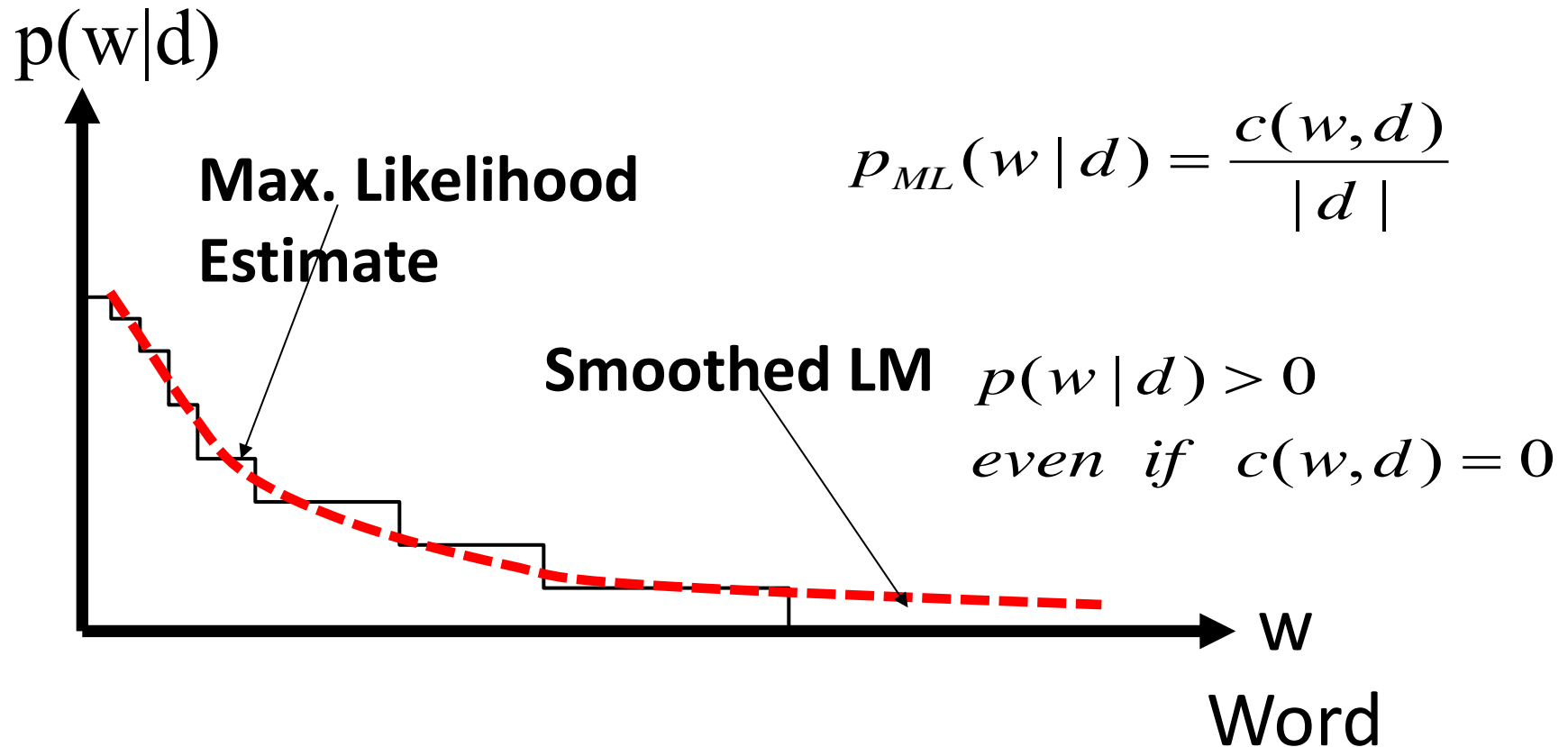
# Ranking Function based on Query Likelihood

$$q = w_1 w_2 \ldots w_n \qquad p(q \mid d) = p(w_1 \mid d) \times \ldots \times p(w_n \mid d)$$

$$f(q,d) = \log p(q \mid d) = \sum_{i=1}^{n} \log p(w_i \mid d) = \sum_{w \in V} c(w,q) \log \boxed{p(w \mid d)}$$

How should we estimate *p(w|d)?*

# How to Estimate p(w|d)



p(w|d)

**Max. Likelihood Estimate**

$$p_{ML}(w \mid d) = \frac{c(w,d)}{|d|}$$

**Smoothed LM**

$$p(w \mid d) > 0$$
$$even \ \ if \ \ c(w,d) = 0$$

w

Word

# How to smooth a LM

- Key Question: what probability should be assigned to an unseen word?

- Let the probability of an unseen word be proportional to its probability given by a reference LM

- One possibility: Reference LM = Collection LM

$$p(w \mid d) = \begin{cases} p_{Seen}(w \mid d) & if \ w \ is \ seen \ in \ d \\ \alpha_d \, p(w \mid C) & otherwise \end{cases}$$

Discounted ML estimate

Collection language model

# Rewriting the Ranking Function with Smoothing

$$\log p(q \mid d) = \sum_{w \in V} c(w,q) \log p(w \mid d)$$

$$= \sum_{w \in V, c(w,d)>0} c(w,q) \log p_{Seen}(w \mid d) + \boxed{\sum_{w \in V, c(w,d)=0} c(w,q) \log \alpha_d p(w \mid C)}$$

Query words **matched** in d          Query words **not matched** in d

$$\sum_{w \in V} c(w,q) \log \alpha_d p(w \mid C) - \sum_{w \in V, c(w,d)>0} c(w,q) \log \alpha_d p(w \mid C)$$

**All** query words          Query words **matched** in d

$$= \sum_{w \in V, c(w,d)>0} c(w,q) \log \frac{p_{Seen}(w \mid d)}{\alpha_d p(w \mid C)} + \mid q \mid \log \alpha_d + \sum_{w \in V} c(w,q) \log p(w \mid C)$$
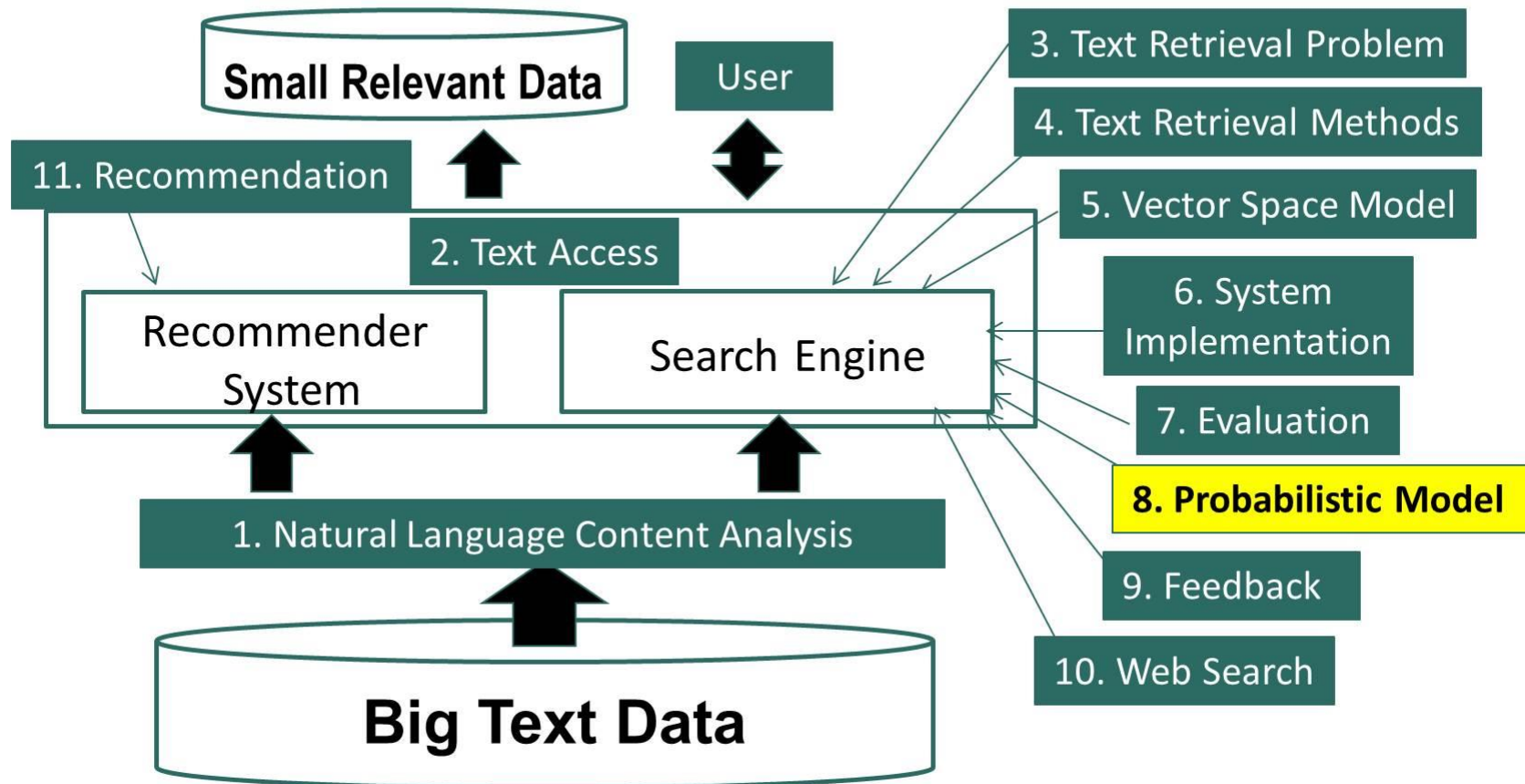
# Benefit of Rewriting

- Better understanding of the ranking function
  - Smoothing with p(w|C) ➡ TF-IDF weighting + length norm.

$$\log p(q \mid d) = \sum_{\substack{w_i \in d \\ w_i \in q}} [\log \frac{p_{Seen}(w_i \mid d)}{\alpha_d \, p(w_i \mid C)}] + n \log \alpha_d + \boxed{\sum_{i=1}^{n} \log p(w_i \mid C)}$$

- Enable efficient computation

# Probabilistic Retrieval Model: Smoothing

# Probabilistic Retrieval Model: Smoothing

# Benefit of Rewriting

- Better understanding of the ranking function
  - Smoothing with p(w|C) ➜ TF-IDF weighting + length norm.

**TF weighting**

**Doc length normalization**

$$\log p(q \mid d) = \sum_{\substack{w_i \in d \\ w_i \in q}} [\log \frac{p_{Seen}(w_i \mid d)}{\alpha_d\, p(w_i \mid C)}] + n \log \alpha_d + \boxed{\sum_{i=1}^{n} \log p(w_i \mid C)}$$

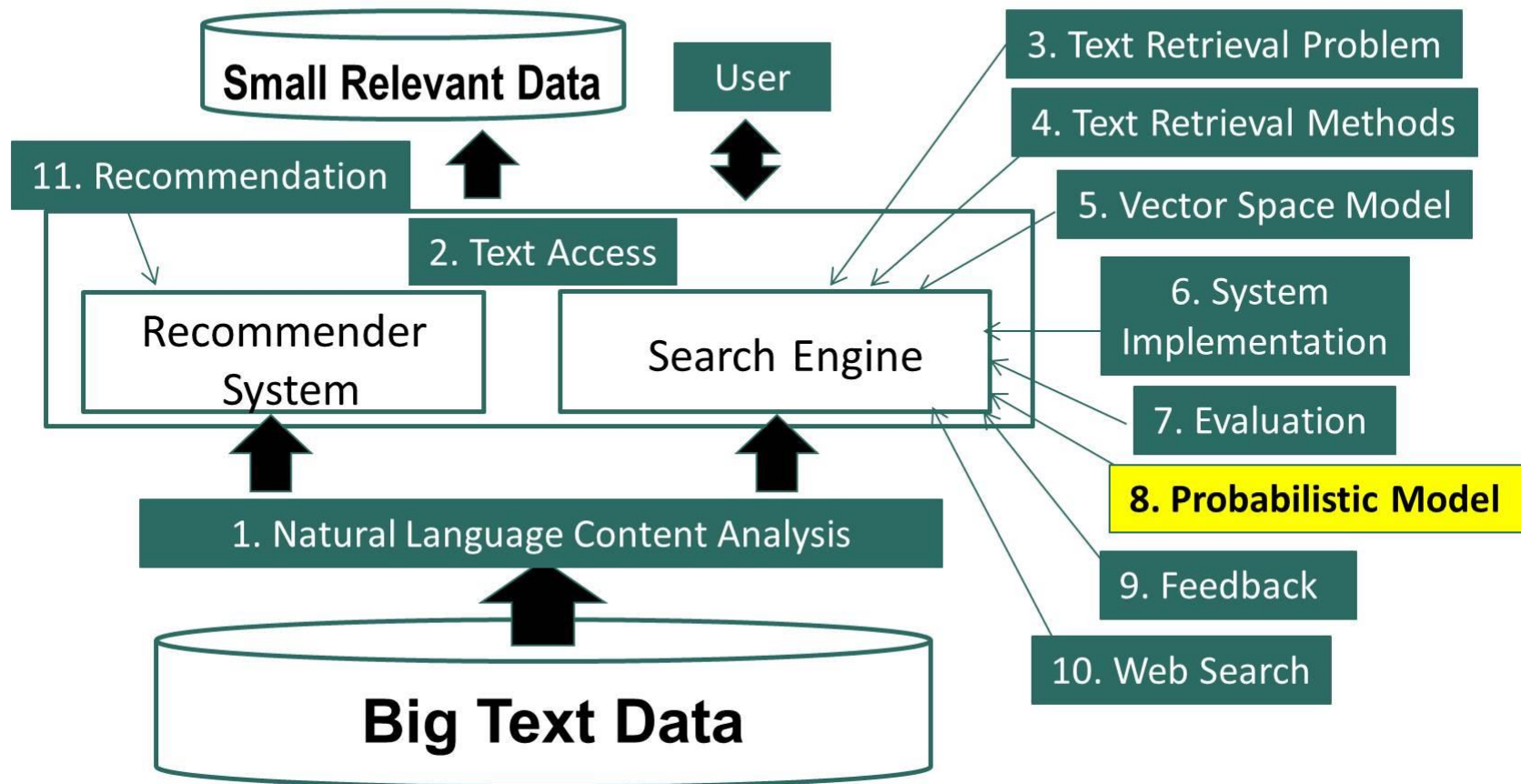**matched query terms**

**IDF weighting**

**Ignore for ranking**

- Enable efficient computation

# Summary

- Smoothing of p(w|d) is necessary for query likelihood
- General idea: smoothing with p(w|C)
  - The probability of an unseen word in d is assumed to be proportional to p(w|C)
  - Leads to a general ranking formula for query likelihood with TF-IDF weighting and document length normalization
  - Scoring is primarily based on sum of weights on matched query terms
- However, how exactly should we smooth?

# Probabilistic Retrieval Model: Smoothing Methods

# Probabilistic Retrieval Model: Smoothing Methods

# Query Likelihood + Smoothing with p(w|C)

$$\log p(q \mid d) = \sum_{\substack{w_i \in d \\ w_i \in q}} c(w,q)[\log \frac{p_{Seen}(w_i \mid d)}{\alpha_d \, p(w_i \mid C)}] + n \log \alpha_d + \boxed{\sum_{i=1}^{n} \log p(w_i \mid C)}$$

$$f(q,d) = \sum_{\substack{w_i \in d \\ w_i \in q}} c(w,q)[\log \frac{p_{Seen}(w_i \mid d)}{\alpha_d \, p(w_i \mid C)}] + n \log \alpha_d$$
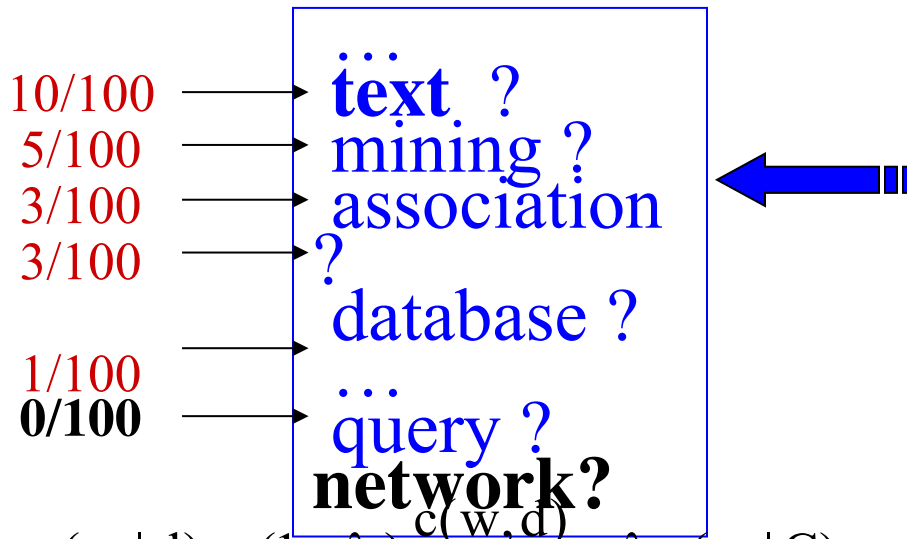
$$\boxed{\begin{aligned} &p_{Seen}(w_i \mid d) = ? \\ &\alpha_d = ? \end{aligned}}$$

How to smooth p(w|d)?

# Linear Interpolation (Jelinek-Mercer) Smoothing

**Document  d**

**Unigram LM  p(w|θ)=?**

Total #words=**100**

Collection LM

**P(w|C)**

10/100 → … **text** ?

5/100 → mining ?

3/100 → association ?

3/100 →

database ?

1/100 →

0/100 → query ?

**network?**

text 10
mining 5
association 3
database 3
algorithm 2
…
query  1
efficient  1

the 0.1
a   0.08
..
computer 0.02
database 0.01
……
**text 0.001**
**network 0.001**
mining 0.0009
…

$$p(w \mid d) = (1-\lambda)\frac{c(w,d)}{|d|} + \lambda\, p(w \mid C)$$

$$\lambda \in [0,1]$$

$$p("text"|d) = (1-\lambda)\frac{10}{100} + \lambda * 0.001$$

$$p("network"|d) = \lambda * 0.001$$

# Dirichlet Prior (Bayesian) Smoothing

**Unigram LM  p(w|θ)=?**

**Document  d**
Total #words=**100**

Collection LM
**P(w|C)**

| | | |
|---|---|---|
| 10/100 | text ? | text 10 |
| 5/100 | mining ? | mining 5 |
| 3/100 | association ? | association 3 |
| 3/100 | database ? | database 3 |
| 1/100 | ... | algorithm 2 |
| **0/100** | query ? | ... query 1 |
| | **network?** | efficient 1 |

the 0.1
a   0.08
..
computer 0.02
database 0.01
......
**text 0.001**
**network 0.001**
mining 0.0009
...

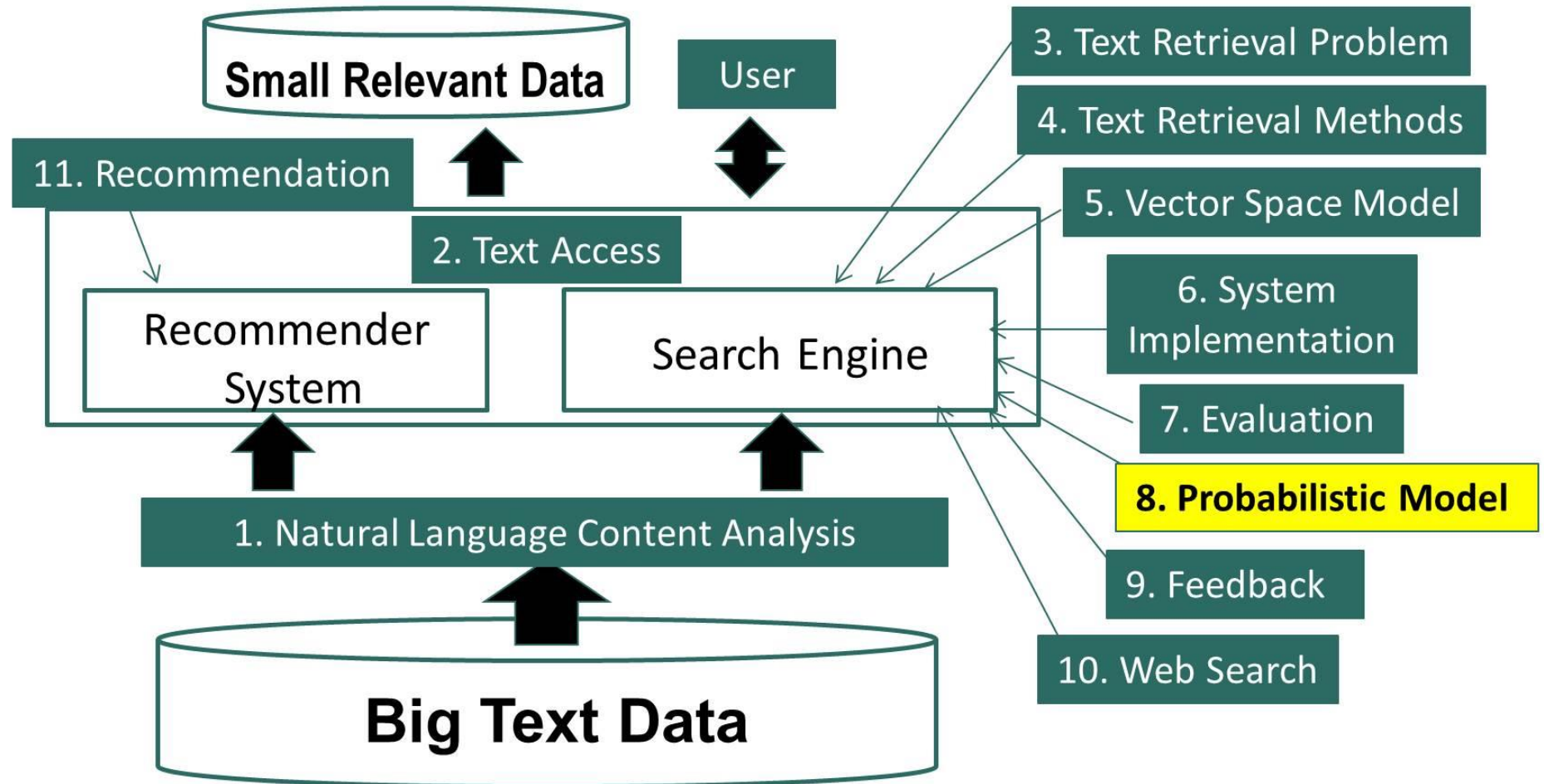$$p(w|d) = \frac{c(w,d) + \mu p(w|C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w,d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|C)$$

$$\mu \in [0, +\infty)$$

$$p("text"|d) = \frac{10 + \mu * 0.001}{100 + \mu}$$

$$p("network"|d) = \frac{\mu}{100 + \mu} * 0.001$$

# Probabilistic Retrieval Model: Smoothing Methods

# Probabilistic Retrieval Model: Smoothing Methods

# Ranking Function for JM Smoothing

$$f(q,d) = \sum_{\substack{w_i \in d \\ w_i \in q}} c(w,q)[\log \frac{p_{\text{Seen}}(w_i \mid d)}{\alpha_d \, p(w_i \mid C)}] + n \log \alpha_d$$

$$p(w \mid d) = (1-\lambda)\frac{c(w,d)}{|d|} + \lambda \, p(w \mid C) \qquad \lambda \in [0,1]$$

$$\frac{p_{\text{seen}}(w \mid d)}{\alpha_d p(w \mid C)} = \frac{(1-\lambda)p_{\text{ML}}(w \mid d) + \lambda p(w \mid C)}{\lambda p(w \mid C)} = 1 + \frac{1-\lambda}{\lambda}\frac{c(w,d)}{|d| \, p(w \mid C)}$$

$$f_{\text{JM}}(q,d) = \sum_{\substack{w \in d \\ w \in q}} c(w,q) \log[1 + \frac{1-\lambda}{\lambda}\frac{c(w,d)}{|d| p(w \mid C)}]$$

# Ranking Function for Dirichlet Prior Smoothing

$$f(q,d) = \sum_{\substack{w_i \in d \\ w_i \in q}} c(w,q)[\log \frac{p_{Seen}(w_i \mid d)}{\alpha_d \, p(w_i \mid C)}] + n \log \alpha_d$$

$$p(w \mid d) = \frac{c(w;d) + \mu \, p(w \mid C)}{d \mid + \mu} = \frac{\mid d \mid}{\mid d \mid + \mu} \frac{c(w,d)}{\mid d \mid} + \frac{\mu}{\mid d \mid + \mu} p(w \mid C)$$

$$\mu \in [0, +\infty)$$

$$\frac{p_{seen}(w \mid d)}{\alpha_d p(w \mid C)} = \frac{\dfrac{c(w,d) + \mu p(w \mid C)}{\mid d \mid + \mu}}{\dfrac{\mu p(w \mid C)}{\mid d \mid + \mu}} = 1 + \frac{c(w,d)}{\mu p(w \mid C)} \qquad \alpha_d = \frac{\mu}{\mid d \mid + \mu}$$

$$f_{DIR}(q,d) = [\sum_{\substack{w \in d \\ w \in q}} c(w,q) \log[1 + \frac{c(w,d)}{\mu p(w \mid C)}]] + n \log \frac{\mu}{\mu + \mid d \mid}$$

# Summary

- Two smoothing methods
  - Jelinek-Mercer: Fixed coefficient linear interpolation
  - Dirichlet Prior: Adding pseudo counts; adaptive interpolation
- Both lead to state of the art retrieval functions with assumptions clearly articulated (less heuristic)
  - Also implementing TF-IDF weighting and doc length normalization
  - Has precisely one (smoothing) parameter

# Summary of Query Likelihood Probabilistic Model

- Effective ranking functions obtained using pure probabilistic modeling
  - Assumption 1: Relevance(q,d) = $p(R=1|q,d) \approx p(q|d,R=1)$ **$\approx p(q|d)$**
  - Assumption 2: Query words are generated independently
  - Assumption 3: Smoothing with $p(w|C)$
  - Assumption 4: JM **or** Dirichlet prior smoothing
- Less heuristic compared with VSM
- Many extensions have been made [Zhai 08]

# Additional Readings

- ChengXiang Zhai, *Statistical Language Models for Information Retrieval* (Synthesis Lectures Series on Human Language Technologies), Morgan & Claypool Publishers, 2008.

http://www.morganclaypool.com/doi/abs/10.2200/S00158ED1V01Y200811HLT001