

# Probabilistic Retrieval Models

Quach Dinh Hoang

Slides are obtained from [Zhai and Massung, 2016]

# Probabilistic Retrieval Models

- Framework

$$f(d,q) = p(R = 1 \mid d,q), R \in \{0,1\}$$

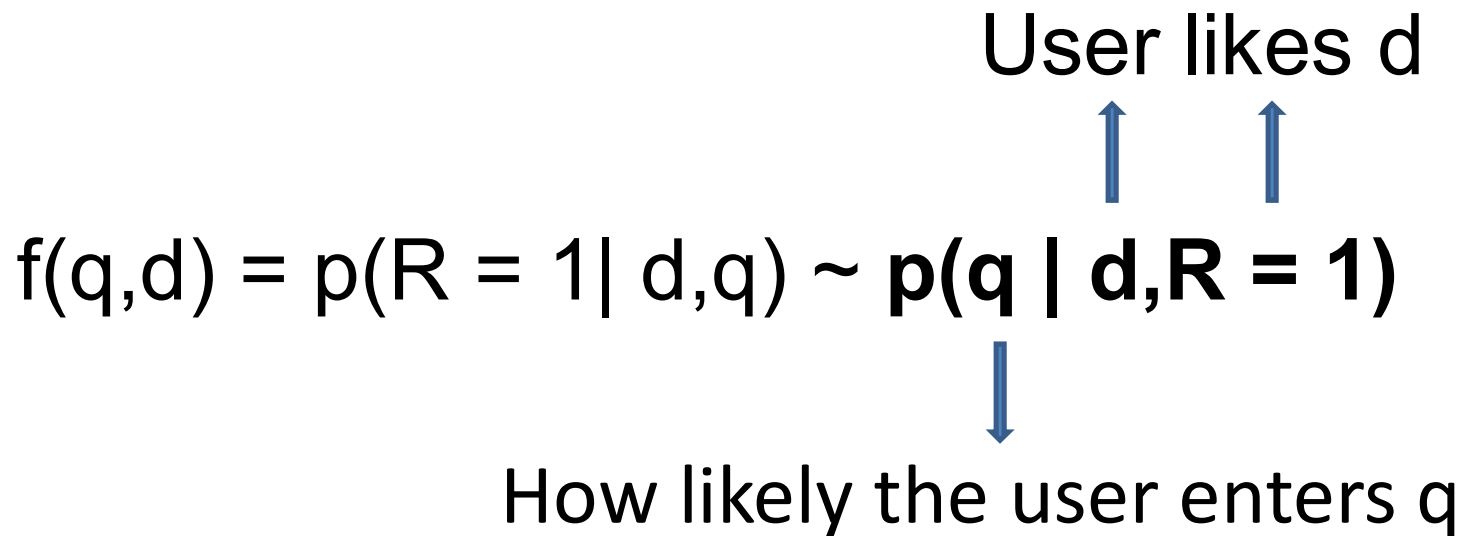
- Classic probabilistic model: BM25
- **Language model: Query Likelihood**
- Divergence-from-randomness model: PL2

- Language model

$$p(R = 1 \mid d,q) \sim p(q \mid d, R = 1)$$

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

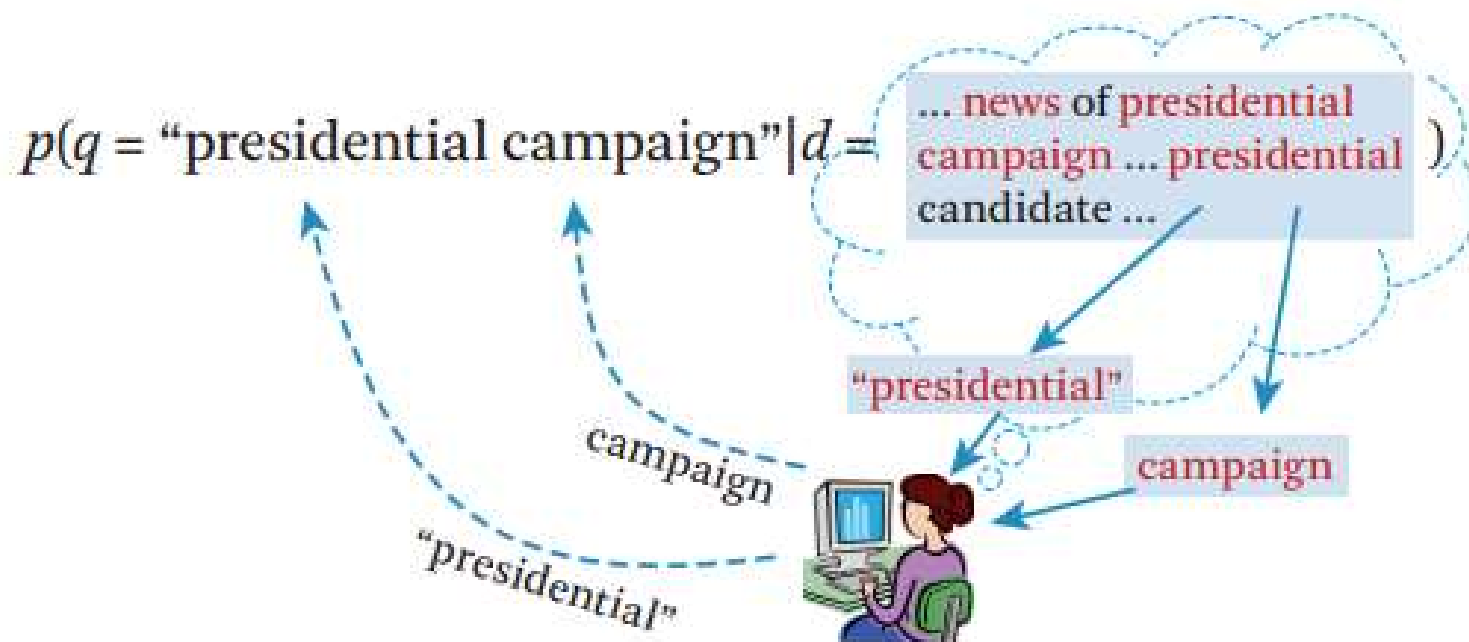
# Query Likelihood Retrieval Model



## Assumption:

A user formulates a query based on an  
**“imaginary relevant document”**

# Query Likelihood Retrieval Model



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

# Language Model

- It is a probability distribution over word sequences
- A document is considered as a sample have drawn from an underline language model
- Each document is ranked based on the capability of its language model generate the words in query

# The Simplest Language Model: Unigram LM

- Generate text by generating each word  
INDEPENDENTLY

$$q = w_1, w_2, \dots, w_n$$

$$p(q \mid d) = p(w_1 \mid d) \times p(w_2 \mid d) \times \dots \times p(w_n \mid d).$$

- Parameters:  $\{p(w_i)\}$

$$p(w_1) + \dots + p(w_N) = 1$$

– N is vocabulary size

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

# Unigram Query Likelihood Model Example

**Assumption:** each query word is independent

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

$$p(q | d_4 = \begin{array}{l} \text{... news of presidential campaign} \\ \text{... presidential candidate ...} \end{array}) = \frac{2}{|d_4|} * \frac{1}{|d_4|}$$

$$p(q | d_3 = \text{... news of presidential campaign ...}) = \frac{1}{|d_3|} * \frac{1}{|d_3|}$$

$$p(q | d_2 = \begin{array}{l} \text{... news about organic food} \\ \text{campaign ...} \end{array}) = \frac{0}{|d_2|} * \frac{1}{|d_2|} = 0$$

# Try a different query

**q** = “**presidential campaign** update”

$$p(q|d4 = \text{... news of } \textbf{presidential campaign} \text{ ... } \textbf{presidential} \text{ candidate ...}) = \frac{2}{|d4|} * \frac{1}{|d4|} * \frac{0}{|d4|} = 0!$$

$$p(q|d3 = \text{... news of } \textbf{presidential campaign} \text{ ...}) = \frac{1}{|d3|} * \frac{1}{|d3|} * \frac{0}{|d3|} = 0!$$

$$p(q|d2 = \text{... news about organic food } \textbf{campaign} \text{ ...}) = \frac{0}{|d2|} * \frac{1}{|d2|} * \frac{0}{|d2|} = 0$$

- All of documents have zero probability of generating this query.

Clearly, this's not desirable. How to fix it?

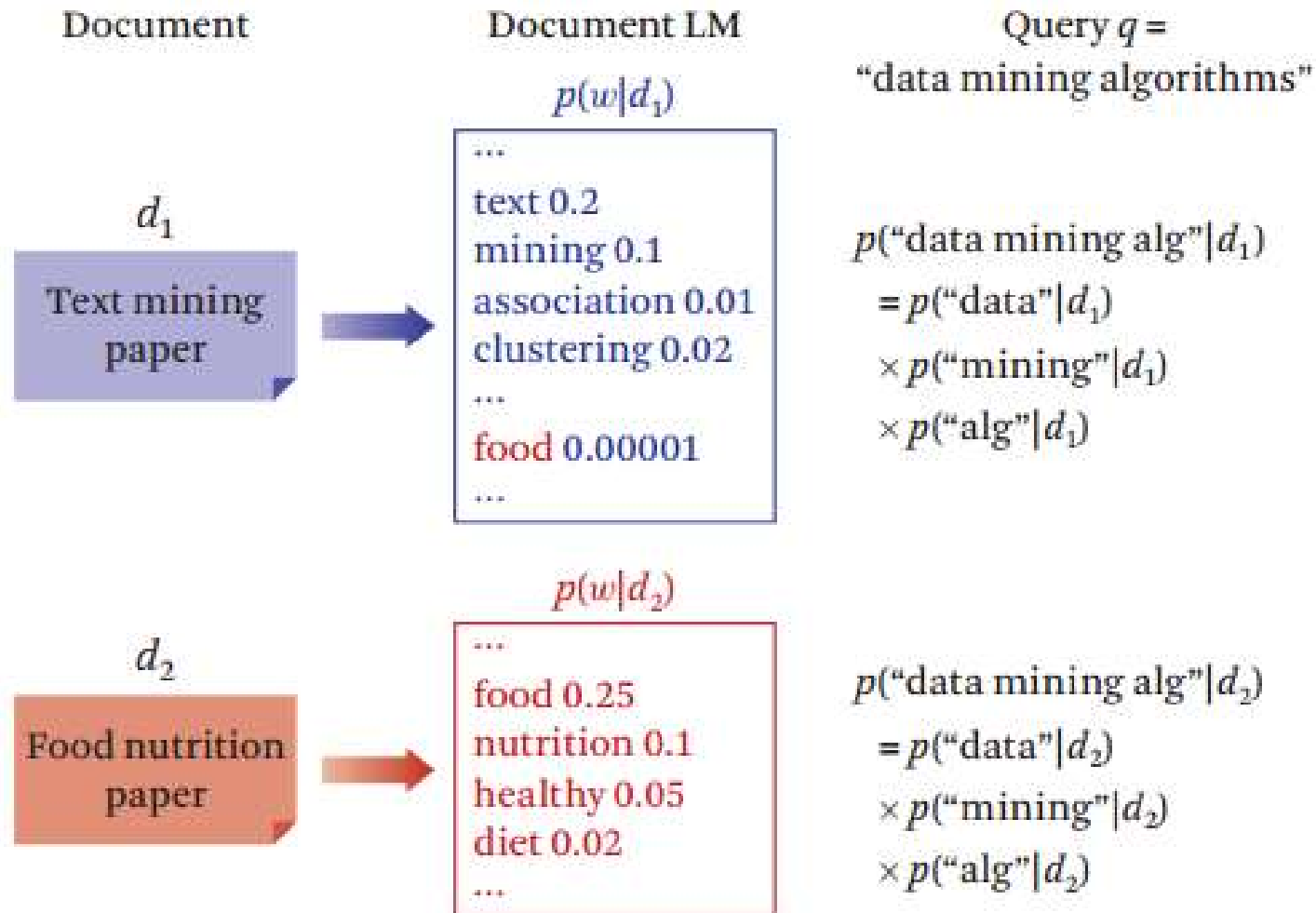


# Document Language Model

- Assume that the user could have drawn a query from a document language model



# Document Language Model Example



# Summary of Language Model

- Generate text by generating each word  
INDEPENDENTLY

$$q = w_1, w_2, \dots, w_n \quad p(q | d) = p(w_1 | d) \times p(w_2 | d) \times \dots \times p(w_n | d).$$

- Parameters:  $\{p(w_i)\}$

$$p(w_1) + \dots + p(w_N) = 1$$

– N is vocabulary size

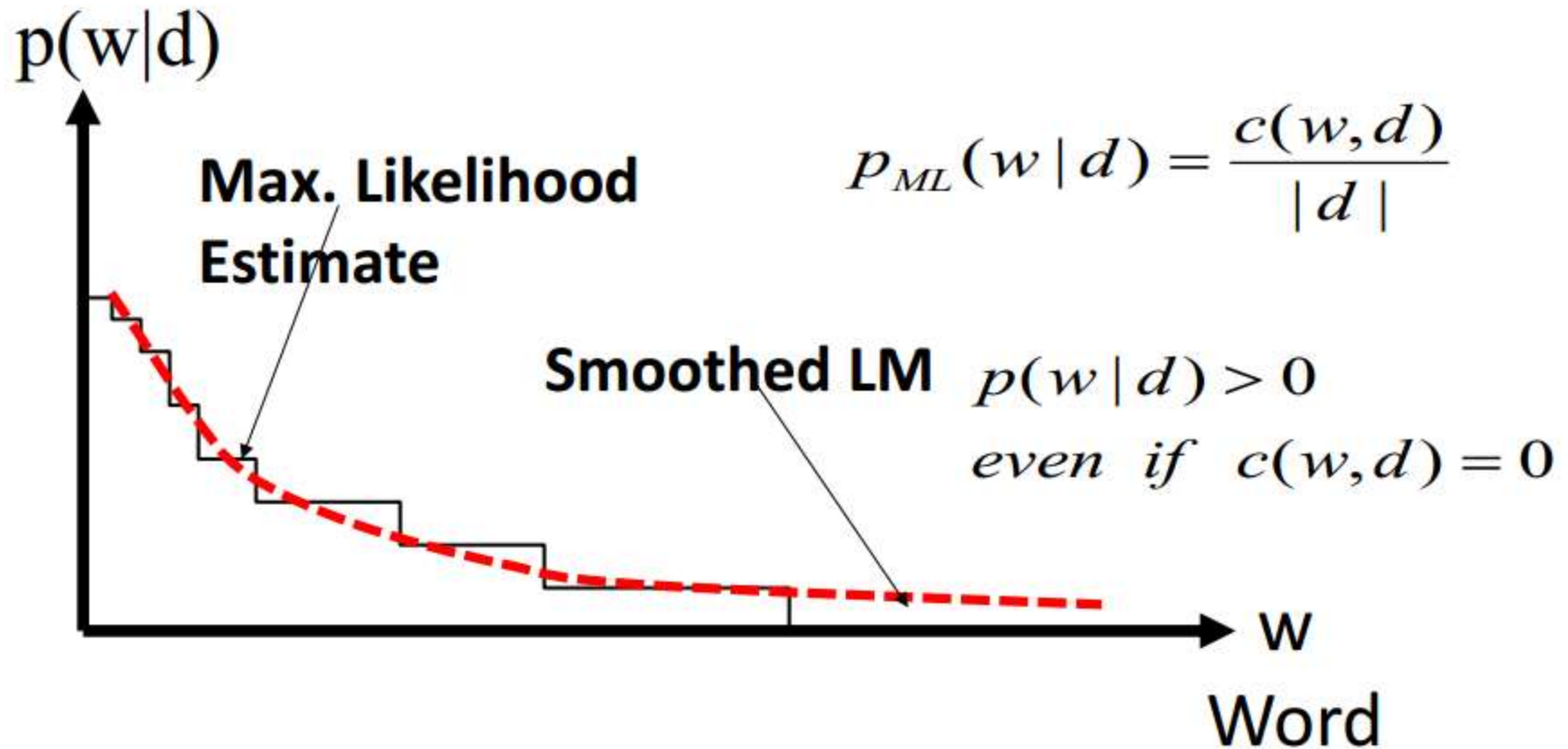
- In practice

$$\text{score}(q, d) = \log p(q | d) = \sum_{i=1}^n \log p(w_i | d) = \sum_{w \in V} c(w, q) \log p(w | d).$$

Document language model



# How to Estimate $p(w|d)$



# Smoothing a Language Model

- 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 language model (LM)
- One possibility: Reference LM = Collection LM

$$p(w | d) = \begin{cases} p_{\text{seen}}(w | d) & \text{if } w \text{ seen in } d \\ \alpha_d \cdot p(w | C) & \text{otherwise.} \end{cases}$$

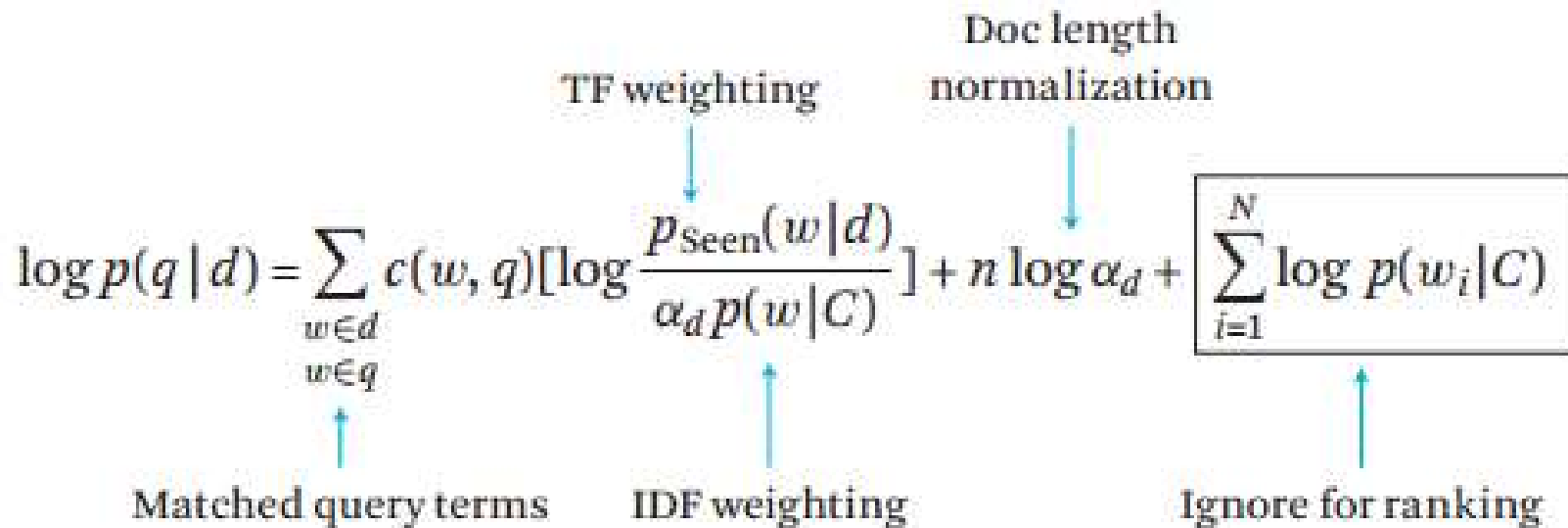
Document language model

Collection language model

# Rewriting the Ranking Function with Smoothing

$$\begin{aligned}
 \log p(q|d) &= \sum_{w \in V} c(w, q) \log p(w|d) \\
 &= \sum_{w \in V, c(w, d) > 0} c(w, q) \log p_{\text{Seen}}(w|d) + \boxed{\sum_{w \in V, c(w, d) = 0} c(w, q) \log \alpha_d p(w|C)} \\
 &\quad \begin{array}{cc} \uparrow & \uparrow \\ \text{Query words matched in } d & \text{Query words not matched in } d \end{array} \\
 &\quad \begin{array}{c} \downarrow \\ \boxed{\sum_{w \in V} c(w, q) \log \alpha_d p(w|C)} - \sum_{w \in V, c(w, d) > 0} c(w, q) \log \alpha_d p(w|C) \\ \begin{array}{cc} \uparrow & \uparrow \\ \text{All query words} & \text{Query words matched in } d \end{array} \end{array} \\
 &= \sum_{w \in V, c(w, d) > 0} c(w, q) \log \frac{p_{\text{Seen}}(w|d)}{\alpha_d p(w|C)} + |q| \log \alpha_d + \sum_{w \in V} c(w, q) \log p(w|C)
 \end{aligned}$$

# Rewriting the Ranking Function with Smoothing



The diagram illustrates the components of the ranking function  $\log p(q | d)$ . The equation is: 
$$\log p(q | d) = \sum_{\substack{w \in d \\ w \in q}} c(w, q) \left[ \log \frac{p_{\text{Seen}}(w | d)}{\alpha_d p(w | C)} \right] + n \log \alpha_d + \boxed{\sum_{i=1}^N \log p(w_i | C)}$$
 Annotations with arrows: 

- Matched query terms**: points to the intersection of  $w \in d$  and  $w \in q$  in the summation index.
- TF weighting**: points to the  $c(w, q)$  term.
- IDF weighting**: points to the  $\frac{1}{\alpha_d}$  term in the log ratio.
- Doc length normalization**: points to the  $\log \alpha_d$  term.
- Ignore for ranking**: points to the boxed summation term  $\sum_{i=1}^N \log p(w_i | C)$ .

# Summary of Smoothing

- 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 TFIDF weighting and document length normalization
  - Scoring is primarily based on sum of weights on matched query terms



# Smoothing Methods

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

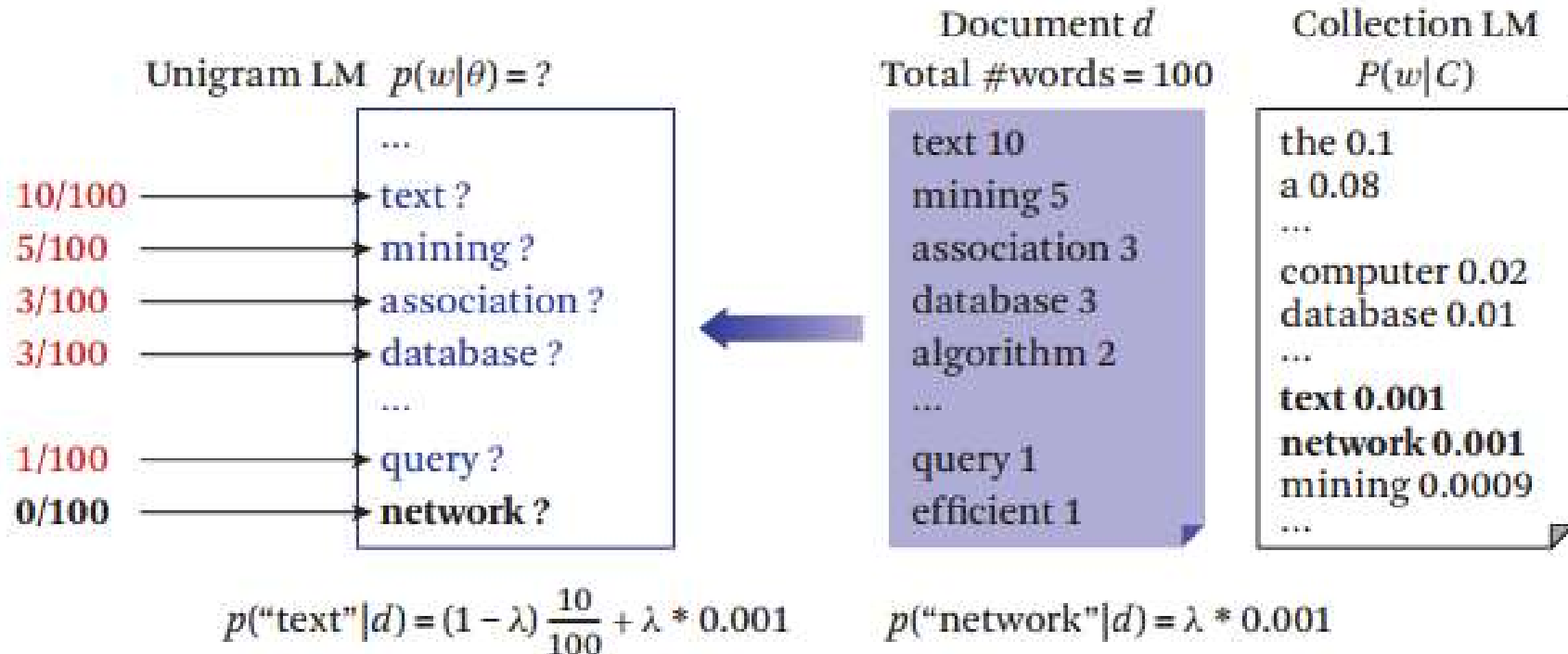
$$f(d, q) = \sum_{\substack{w \in d \\ w \in q}} c(w, q) \left[ \log \frac{p_{\text{Seen}}(w | d)}{\alpha_d p(w | C)} \right] + n \log \alpha_d$$

$$\begin{cases} p_{\text{seen}}(w_i | d) = ? \\ \alpha_d = ? \end{cases}$$

How to smooth  $p(w|d)$ ?

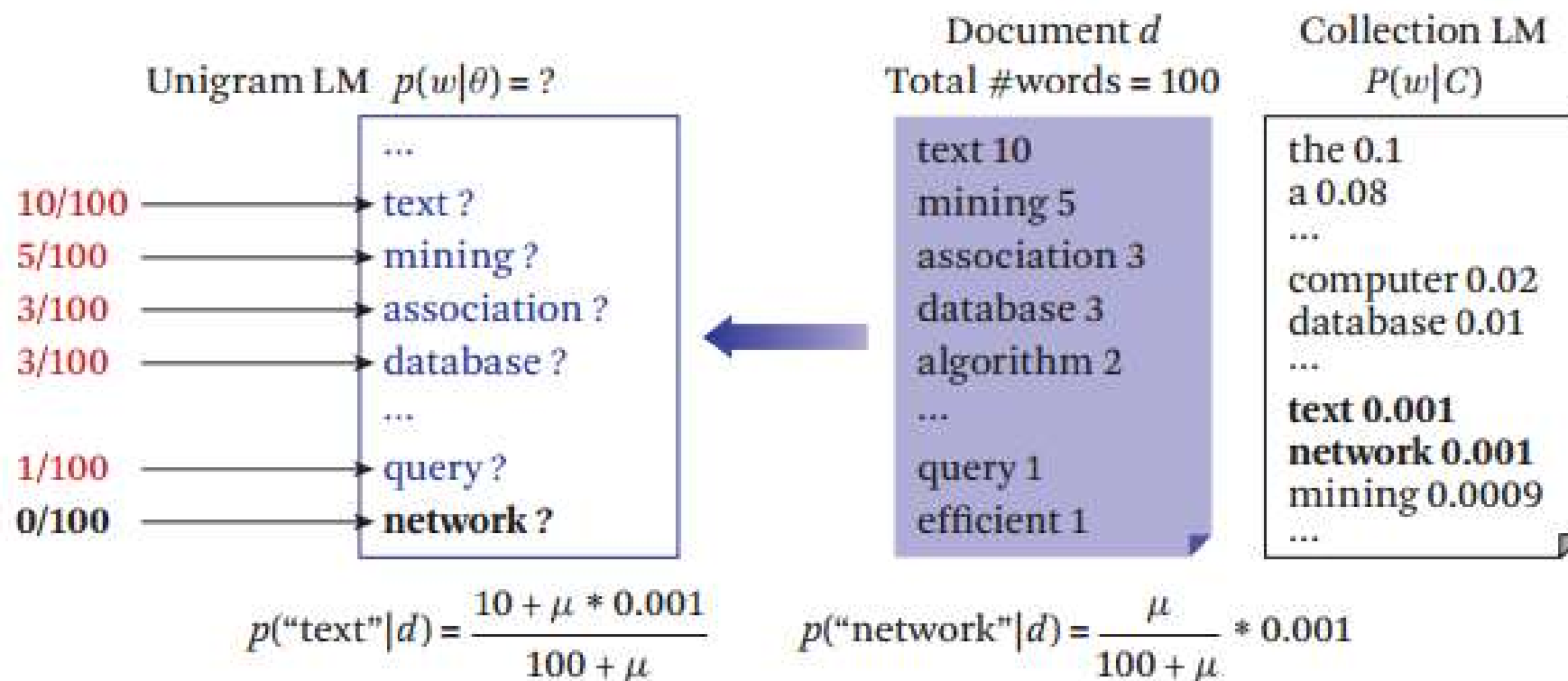
# Jelinek-Mercer Smoothing

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



# Dirichlet Prior Smoothing

$$p(w|d) = \frac{c(w, d) + \mu p(w|C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|C) \quad \mu \in [0, +\infty)$$



# Ranking Function for Jelinek-Mercer Smoothing

$$f(d, q) = \sum_{\substack{w \in d \\ w \in q}} c(w, q) \left[ \log \frac{p_{\text{seen}}(w | d)}{\alpha_d p(w | C)} \right] + n \log \alpha_d$$

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

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

$$\text{score}_{JM}(q, d) = \sum_{w \in q, d} c(w, q) \log \left( 1 + \frac{1 - \lambda}{\lambda} \cdot \frac{c(w, d)}{|d| \cdot p(w | C)} \right)$$

# Ranking Function for Dirichlet Prior Smoothing

$$f(d, q) = \sum_{\substack{w \in d \\ w \in q}} c(w, q) \left[ \log \frac{p_{\text{seen}}(w | d)}{\alpha_d p(w | C)} \right] + n \log \alpha_d$$

$$p(w | d) = \frac{c(w, d) + \mu p(w | C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} p(w | C) \quad \mu \in [0, +\infty)$$

$$\frac{p_{\text{seen}}(w | d)}{\alpha_d \cdot p(w | C)} = \frac{\frac{c(w, d) + \mu \cdot p(w | C)}{|d| + \mu}}{\frac{\mu \cdot p(w | C)}{|d| + \mu}} = 1 + \frac{c(w, d)}{\mu \cdot p(w | C)}$$

$$\text{score}_{DIR}(q, d) = \sum_{w \in q, d} c(w, q) \log \left( 1 + \frac{c(w, d)}{\mu \cdot p(w | C)} \right) + |q| \log \frac{\mu}{\mu + |d|}$$

# Summary of Smoothing Methods

- 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 Language Model

- Effective ranking functions obtained using pure probabilistic modeling
  - Assumption 1:  $\text{Relevance}(q,d) = p(R=1|q,d) \sim p(q|d, R=1) \sim \mathbf{p(q|d)}$
  - 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