



UNIVERSITY *of* DELAWARE



UNIVERSITY *of* DE

Language Models Approach

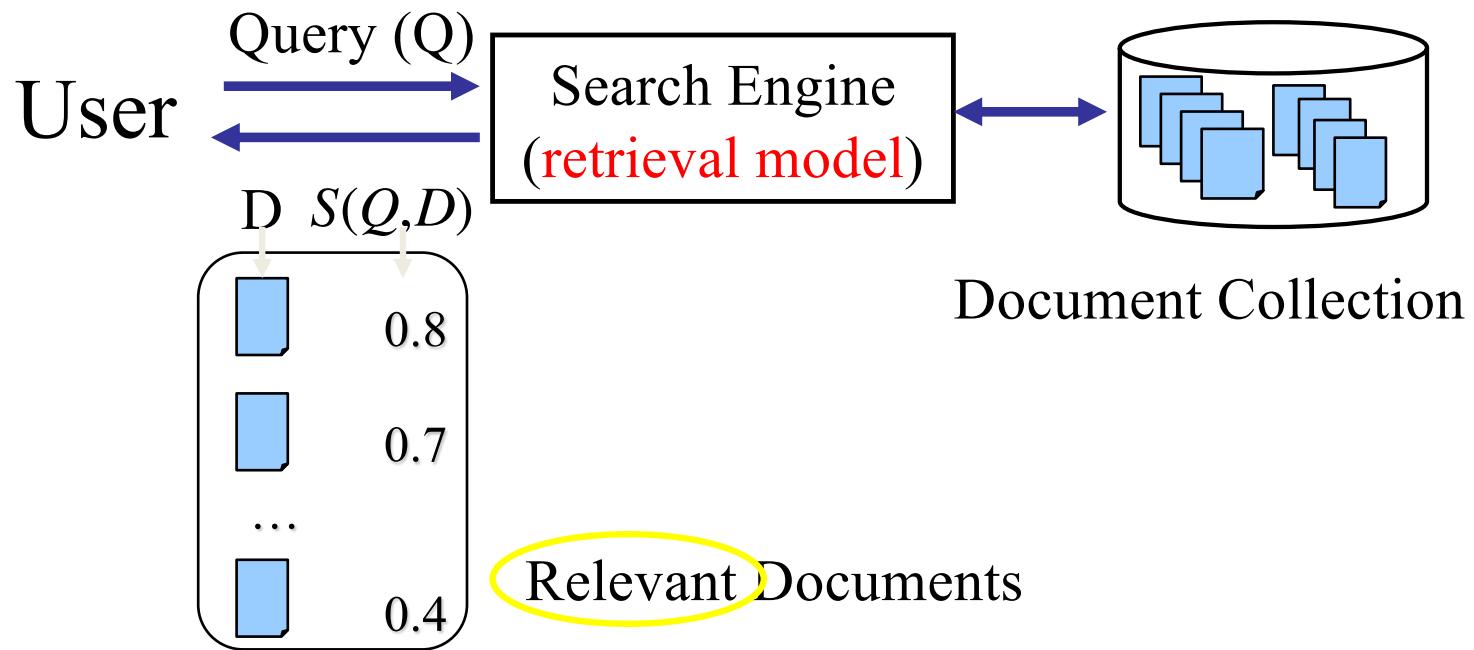
(Search and Data Mining)

Hui Fang

Department of Electrical and Computer Engineering
University of Delaware



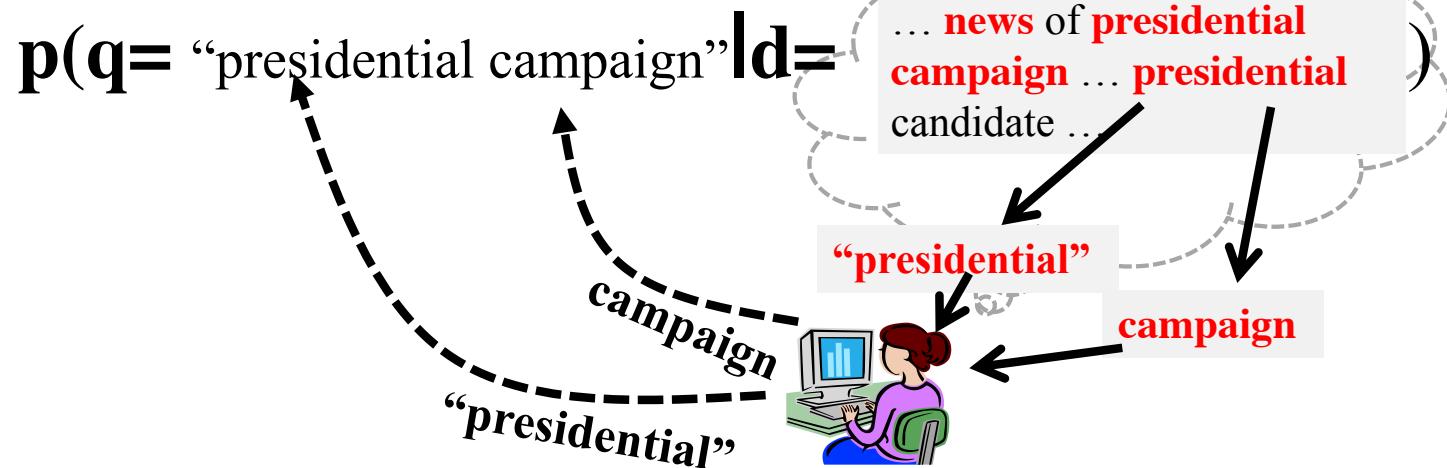
Key Problem: Retrieval Models



How to accurately model relevance, i.e., $S(Q,D)$?



Query Generation by Sampling Words from a Document



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



Unigram Query Likelihood

$$\begin{aligned} p(q = \text{"presidential campaign"} | d = & \text{ ... news of presidential} \\ & \text{campaign ... presidential} \\ & \text{candidate ...}) \\ = p(\text{"presidential"} | d) * p(\text{"campaign"} | d) \\ = \frac{c(\text{"presidential"}, d)}{|d|} * \frac{c(\text{"campaign"}, d)}{|d|} \end{aligned}$$

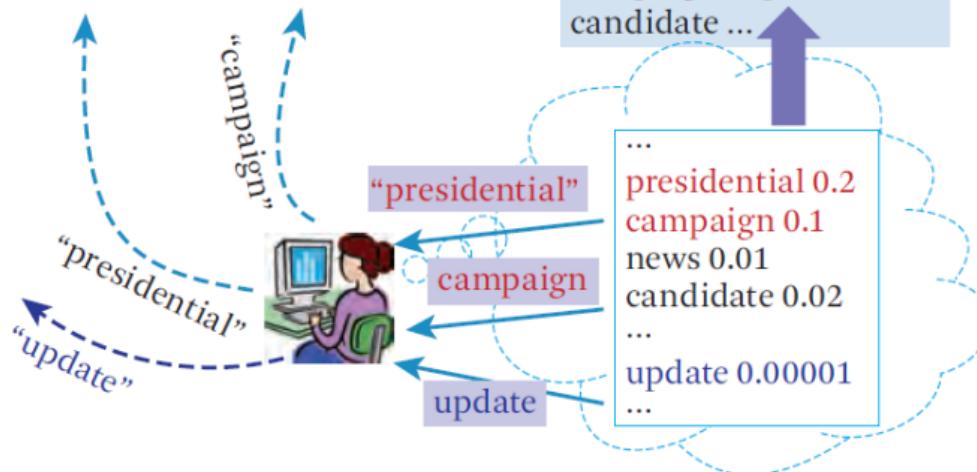
Assume each word is generated independently.

What if we have a query “presidential campaign update”?



Sampling Words from a Document Model

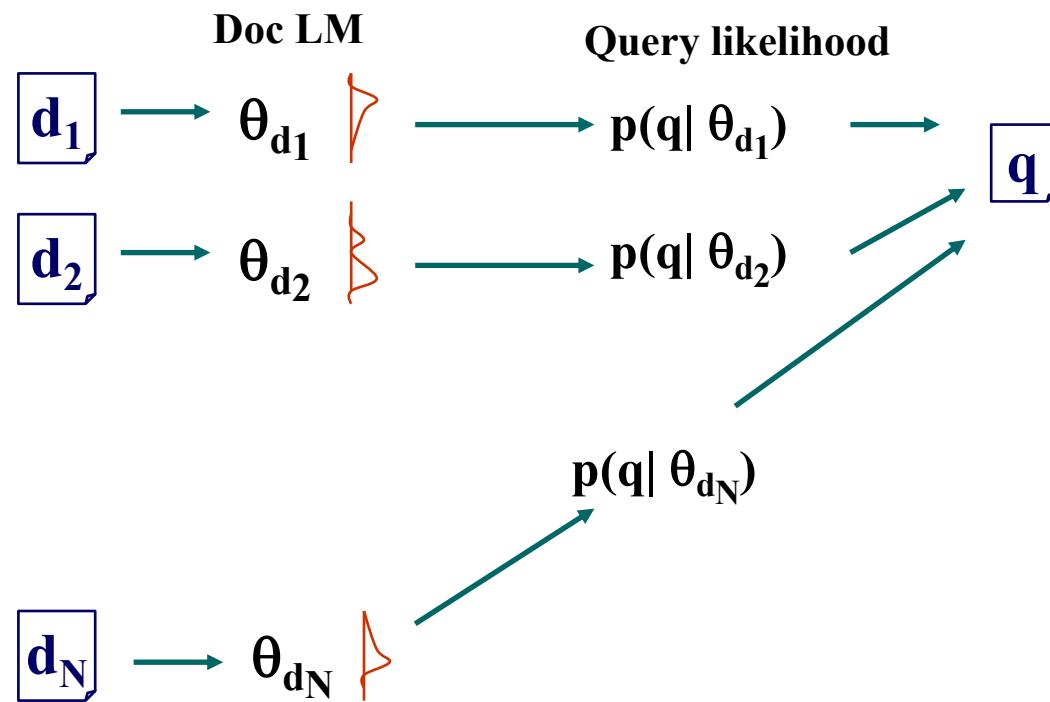
$$p(q = \text{"presidential campaign"} | d = \dots \text{news of presidential} \\ \text{campaign ... presidential candidate ...})$$



Computing the probability of a query given a document using a document language model



Ranking Documents based on Query Likelihood





Retrieval as Language Model Estimation

- Rank documents based on *query likelihood*

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

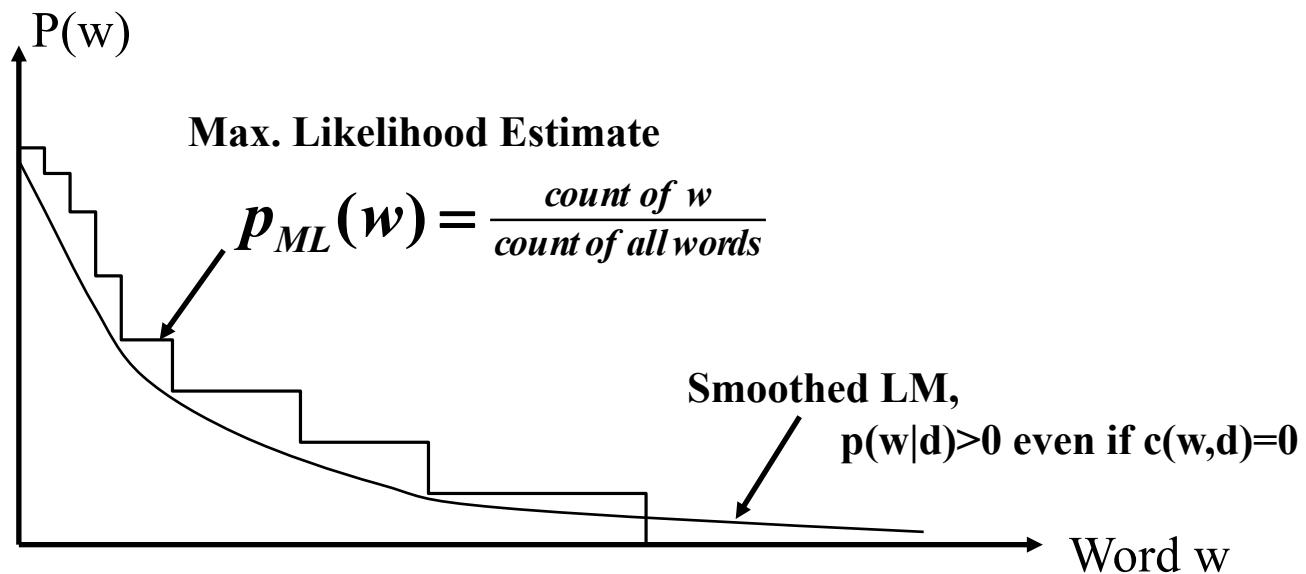
$$p(q | d) = p(w_1 | d) \times p(w_2 | d) \times \dots \times p(w_n | d).$$

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

- Retrieval problem \approx Estimation of $p(w|d)$



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





How to smooth a LM

- Key question of smoothing
 - What probability should be assigned to an unseen word?
- Let the probability of an unseen word be proportionally to its probability given by a reference LM.

$$p(w | d) = \begin{cases} p_{\text{seen}}(w | d) & \text{Discounted ML estimate} \\ \alpha_d p(w | C) & \text{if } w \text{ is seen in } d \\ & \text{otherwise} \end{cases}$$

Collection language model



Rewriting the Ranking Function with Smoothing

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

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

Query words matched in d

Query words not matched in d

$$\sum_{w \in V} c(w, q) \log \alpha_d p(w|C)$$

All query words

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

Query words matched in d

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



Benefit of Rewriting

- Efficient computation
- Better understanding of the ranking function
 - Smoothing with $p(w|C) \approx \text{TF-IDF} + \text{length norm.}$

$$\log p(q | d) = \sum_{\substack{w_i \in d \\ w_i \in q}} \left[\log \frac{p_{seen}(w_i | d)}{\alpha_d p(w_i | C)} \right] + |q| \log \alpha_d + \boxed{\sum_i \log p(w_i | C)}$$

↑
IDF weighting

TF weighting

Doc length normalization
(long doc is expected to have a smaller α_d)

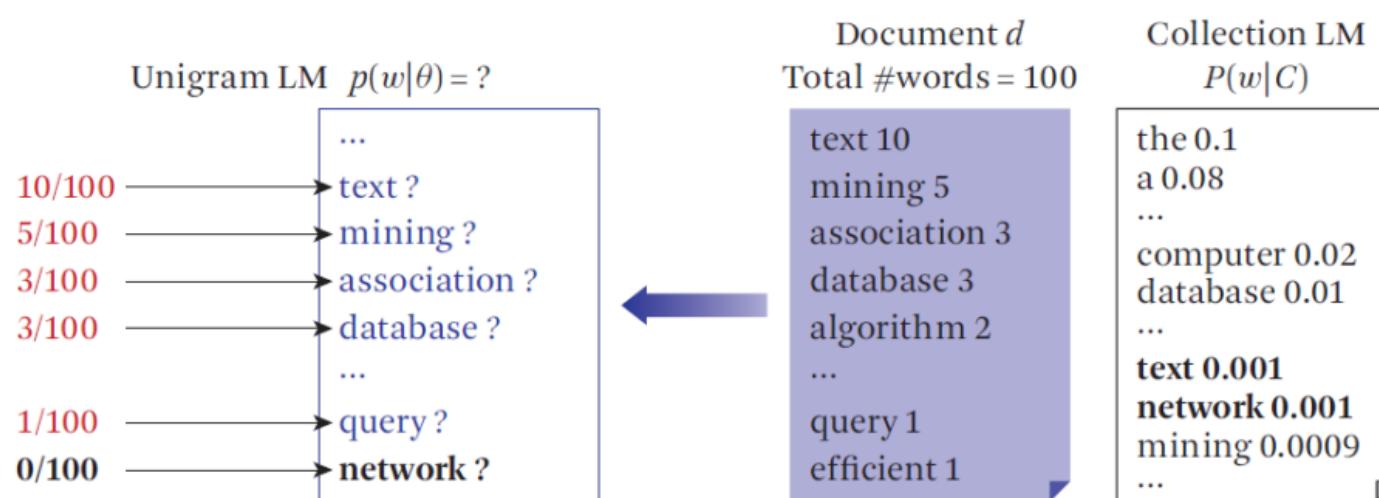
Ignore for ranking



Specific Smoothing Methods (1)

- Jelinek-Mercer (JM) smoothing:

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



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



Retrieval Function based on JM smoothing

$$\sum_{w \in d, q} c(w, q) \log \left(\frac{p_{\text{seen}}(w | d)}{\alpha_d \cdot p(w | C)} \right) + |q| \log \alpha_d.$$

- **Jelinek-Mercer (JM) smoothing:**

$$p(w | d) = (1 - \lambda)p_{ml}(w | d) + \lambda p(w | C)$$

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

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



Specific Smoothing Methods (2)

- Dirichlet Prior smoothing:

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

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

| | | |
|--------|---|---------------|
| 10/100 | → | text ? |
| 5/100 | → | mining ? |
| 3/100 | → | association ? |
| 3/100 | → | database ? |
| ... | | ... |
| 1/100 | → | query ? |
| 0/100 | → | network ? |

Document d
Total #words = 100

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

Collection LM
 $P(w|C)$

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

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

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



Dirichlet Prior Function

$$\sum_{w \in d, q} c(w, q) \log \left(\frac{p_{\text{seen}}(w \mid d)}{\alpha_d \cdot p(w \mid C)} \right) + |q| \log \alpha_d.$$

- **Dirichlet Prior Smoothing**

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

$$\alpha_d = \frac{\mu}{|d| + \mu}:$$

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



- More Info about Course Projects
 - March 13 and 15 (W, F)
 - No lectures
 - Project meeting between the instructor and each group
 - Working on the assignment 1
 - March 25 and 27 (M, W)
 - Oral presentation for the project proposal
 - March 29 (F)
 - Written proposals due