

Discussion Section

APRIL 15

Syllabus (I)

- ▶ Search
 - ▶ Overview
 - ▶ Retrieval Models
 - ▶ Vector Space Models
 - ▶ Language Modeling
 - ▶ Axiomatic Approaches
 - ▶ Evaluation
 - ▶ Indexing + Query Processing
 - ▶ Feedback
 - ▶ Web Search
 - ▶ Basics
 - ▶ Link analysis
 - ▶ Machine learning for Web Search

Syllabus (II)

- ▶ Mining
 - ▶ Word association mining
 - ▶ Categorization
 - ▶ Clustering

Projects

- ▶ Comments due on April 13
- ▶ Many interesting ideas!!!
- ▶ Re-think about the scope of your project given the tighter timeline
- ▶ Schedule a meeting with me if you feel uncertain

Midterm

- ▶ Due on April 20
- ▶ Search related application
- ▶ Questions?

Inverted Index Compression

- ▶ Observations:
 - ▶ TF compression
 - ▶ Small numbers tend to occur more frequently.
 - ▶ Doc ID compression
 - ▶ “d-gap” (store differences): d1, d2-d1, d3-d2, ...
 - ▶ Feasible due to sequential access
- ▶ Implications
 - ▶ Exploit skewed frequency distribution and use variable-length encoding

Integer Compression Methods

- ▶ Binary: equal-length coding
 - ▶ 3=>00000011; 5=>00000101
- ▶ Unary: $x \geq 1$ is coded as $x-1$ one bits followed by 0
 - ▶ 3=>110; 5=>11110
- ▶ γ -code: $x=>$ unary code for $1 + \lfloor \log x \rfloor$ followed by binary code for $x - 2^{\lfloor \log x \rfloor}$ in $\lfloor \log x \rfloor$ bits
 - ▶ 3=>101, 5=>11001
- ▶ δ -code: same as γ -code, but replace the unary prefix with γ -code.
 - ▶ 3=>1001, 5=>10101

Examples

γ -code: $x=>$ unary code for $1 + \lfloor \log x \rfloor$ followed by binary code for $x - 2^{\lfloor \log x \rfloor}$ in $\lfloor \log x \rfloor$ bits

- ▶ Compute the compression results using gamma compression for the following values
 - ▶ 2 unary code for 2 [10]. binary code for 0 in 1 bit [0] → 100
 - ▶ 4 unary code for 3 [110]. binary code for 0 in 2 bits [00] → 11000
 - ▶ 14 unary code for 13[1110]. binary code for 6 in 3 bits [110] → 1110110
 - ▶ 63 unary code for 61[11110]. binary code 31 in 5 bits [11111] → 1111011111
 - ▶ 180

Examples

- ▶ Compute the compression results using gamma compression for the following values
 - ▶ 2 100
 - ▶ 4 11000
 - ▶ 14 1110110
 - ▶ 63 11111011111
 - ▶ 180 111111100110100

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 ≈ Estimation of $p(w|d)$

Dirichlet Prior Function

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

▶ Dirichlet Prior Smoothing

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

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

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

Dirichlet prior smoothing

Suppose we have a document collection with an extremely small vocabulary with only 6 words w_1, \dots, w_6 . Let $Q = w_1 w_2$ be a query.

1. Suppose we do not smooth the language model for d_1 and d_2 . Compute the likelihood of the query for both d_1 and d_2 , i.e., $P(Q|d_1)$ and $P(Q|d_2)$. (Do not compute the log-likelihood.) Show your calculations. Which document would be ranked higher?

$$\begin{aligned} s(Q, d_1) &= P(Q|d_1) = p(w_1|d_1) * p(w_2|d_1) = c(w_1, d_1)/|D| * c(w_2, d_1)/|D| = 2/10 * 3/10 \\ s(Q, d_2) &= 7/10 * 1/10 \quad d_2 \text{ will be ranked higher.} \end{aligned}$$

2. Suppose we now smooth the language model for d_1 and d_2 using Dirichlet prior smoothing method with $\mu = 10$. Recompute the likelihood of the query for both d_1 and d_2 , i.e., $P(Q|d_1)$ and $P(Q|d_2)$. Show your calculations. Which document would be ranked higher this time?

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

$$s(Q, d_1) = p(\text{dir}(w_1|d_1)) * p(\text{dir}(w_2|d_1)) = (10/(10+10)) * 2/10 + 10/(10+10) * 0.8 * \dots = 0.1$$

$$s(Q, d_2) = 0.07$$

3. Intuitively, which document do you think should be ranked higher? D1 or D2? Why?

Word	$P(w \text{REF})$	$C(w, d_1)$	$C(w, d_2)$
W1	0.8	2	0
W2	0.1	3	1
W3	0.025	2	1
W4	0.025	2	1
W5	0.025	1	0
W6	0.025	0	0
SUM	1.0	10	10