

Problem Set 3

chen zhang, NetId: czhang49

Handed In: September 24, 2016

1 Classic Probabilistic Retrieval Model

1.1 (a)

For document generation, using multinomial model, the score can be written as:

$$\begin{aligned} \text{score}(Q, D) &= \frac{P(D|Q, R = 1)}{P(D|Q, R = 0)} \\ &= \frac{\prod_{j=1}^{|V|} P(w_j|Q, R = 1)^{c(w_j|Q, R=1)}}{\prod_{j=1}^{|V|} P(w_j|Q, R = 0)^{c(w_j|Q, R=0)}} \\ &= - \sum_{x \in \Omega} P(x) \log_2 P(x) \end{aligned}$$

And then we have

$$\begin{aligned} \text{score}(Q, D) &\propto \log \frac{P(D|Q, R = 1)}{P(D|Q, R = 0)} \\ &\propto \log \frac{\prod_{j=1}^{|V|} P(w_j|Q, R = 1)^{c(w_j|Q, R=1)}}{\prod_{j=1}^{|V|} P(w_j|Q, R = 0)^{c(w_j|Q, R=0)}} \end{aligned}$$

Since the occurrence of a word in the document is independent of the query, we have

$$c(w_j|Q, R = 0) = c(w_j|Q, R = 1) = c(w_j, D)$$

Then we have

$$\begin{aligned} \text{score}(Q, D) &\propto c(w_j, D) \log \frac{\prod_{j=1}^{|V|} P(w_j|Q, R = 1)}{\prod_{j=1}^{|V|} P(w_j|Q, R = 0)} \\ &\propto c(w_j, D) \log \prod_{j=1}^{|V|} \frac{P(w_j|Q, R = 1)}{P(w_j|Q, R = 0)} \\ &\propto \sum_{w \in V} c(w, D) \log \frac{P(w|Q, R = 1)}{P(w|Q, R = 0)} \end{aligned}$$

There's too many equations for this HW, I will write the answers very clearly starting from 1.(b). I apologize for that.

$$(b) P(w|Q, R=0) = \frac{c(w|c)}{|c|}$$

with $|c|$ representing the total number of words in $c = D_1, \dots, D_n$

$$(c) P(w|Q, R=1) = \frac{c(w|q)}{|q|}$$

with $|q|$ representing the total number of words in q

$$(d) P(w|Q, R=1) = (1-\lambda) \frac{c(w|q)}{|q|} + \lambda P(w|\text{REF})$$

(e)

$$\text{score}(Q, D) \propto \sum_{w \in V} c(w, D) \log \frac{P(w|Q, R=1)}{P(w|Q, R=0)}$$

$$= \sum_{\substack{w \in V \\ w \in Q}} c(w, D) \log P(w|Q, R=1) + \sum_{\substack{w \in V \\ w \notin q}} c(w, D) P(w|Q, R=1)$$

$$- \sum_{w \in V} c(w, D) \log P(w|Q, R=0)$$

$$\begin{aligned}
 &= \sum_{\substack{w \in V \\ w \in q}} c(w, D) \log \left[(1-\lambda) \frac{c(w|q)}{|q|} + \lambda p(w|\text{REF}) \right] \\
 &\quad - \sum_{\substack{w \in V \\ w \notin q}} c(w, D) \log [\lambda p(w|\text{REF})] \\
 &\quad + \sum_{\substack{w \in V \\ w \in q}} c(w, D) \log [\lambda p(w|\text{REF})] - \sum_{w \in V} c(w, D) \log p(w|Q, R=0) \\
 &= \sum_{\substack{w \in V \\ w \in q}} c(w, D) \log \left[1 + \frac{(1-\lambda)c(w|q)}{|q| \cdot \lambda p(w|\text{REF})} \right]^{\textcircled{A}} \\
 &\quad + \sum_{\substack{w \in V \\ w \notin q}} c(w, D) \log [\lambda p(w|\text{REF})] \\
 &\quad - \sum_{w \in V} c(w, D) \log p(w|Q, R=0).
 \end{aligned}$$

\textcircled{A} captures TF, \textcircled{B} captures IDF,

Document length normalization not captured

2. (a). For Jelleh-Mercer,

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

$$\begin{aligned} \log p(q|d) &= \sum_{i=1}^m \log p(q_i|d) = \sum_{i=1}^{|U|} c(w_i, q) \log p(w|d) \\ &= \sum_{w \in q} c(w, q) \log p(w|d). \end{aligned}$$

plug in $p(w|d)$ we have

$$\begin{aligned} \log p(q|d) &= \sum_{\substack{w \in q \\ w \notin d}} c(w, q) \log \left[(1-\lambda) \frac{c(w,d)}{|d|} + \lambda p(w|\text{REF}) \right] \\ &\quad + \sum_{\substack{w \in q \\ w \in d}} c(w, q) \log [\lambda p(w|\text{REF})] \\ &= \sum_{\substack{w \in q \\ w \notin d}} c(w, q) \log \left[(1-\lambda) \frac{c(w,d)}{|d|} + \lambda p(w|\text{REF}) \right] \\ &\quad - \sum_{\substack{w \in q \\ w \in d}} c(w, q) \log [\lambda p(w|\text{REF})] \\ &\quad + \sum_{w \in q} c(w, q) \log [\lambda p(w|\text{REF})] \end{aligned}$$

$$= \sum_{w \in Q \cap D} c(w, Q) \log \left[1 + \frac{(1-\lambda) c(w, D)}{\lambda p(w | \text{REF}) |D|} \right] \\ + \underbrace{\sum_{w \in Q} c(w, Q) \log [\lambda p(w | \text{REF})]}_{\text{independent of document}}$$

$$\propto \sum_{w \in Q \cap D} c(w, Q) \log \left[1 + \frac{(1-\lambda) c(w, D)}{\lambda p(w | \text{REF}) |D|} \right]$$

(b) $q = (w_1, q, w_2, q \dots w_n, q)$

if word i in query then $w_i, q = 1$, else $w_i, q = 0$

$$d = (w_1, d, w_2, d \dots w_t, d)$$

if word j in doc then $w_j, d = 1$, else $w_j, d = 0$

Weight $w_t, d = \log \left[1 + \frac{(1-\lambda) c(w, D)}{\lambda p(w | \text{REF}) \times |D|} \right]$ This captures

TF & IDF with $c(w, D)$ and $p(w | \text{REF})$

(c) Since for Jelinek-Mercer,

$$\text{Score} \propto \sum_{w \in Q \cap D} C(w, Q) \log \left(1 + \frac{(1-\lambda) C(w, D)}{\lambda p(w|RDF) \times |D|} \right)$$

$\frac{C(w, D)}{|D|} = \text{const}$ when doubling doc

& $C(w, Q)$ will increase \Rightarrow

Score should increase when doubling doc \Rightarrow

\Rightarrow no over constraint

For Dirichlet prior Ⓐ

Ⓑ

$$\text{Score} \propto \sum_{\substack{w \in D \\ w \notin Q}} \log \left[1 + \frac{C(w, D)}{N p(w|C)} \right] + \alpha \log \frac{p}{|D| + p}$$

when double doc, the increase in Ⓐ is
more than the decrease in Ⓑ \Rightarrow

\Rightarrow score increases \Rightarrow no over constraints.