

Assignment 5

1. a.) Raw term frequency gives too much reward to documents with high term frequency and document frequency. This is why it might be better to use $\log(TF)$ and also account for inverse document frequency. This ensures that rare terms get a higher reward than common words.

b.) Adding another instance of d to the corpus will reduce the IDF for all words in d and increase the IDF for all words not in d .

c.		precision	recall
1.	+	1/1	1/16
	+	2/3	2/16
	-	0	0
	+	3/4	3/16
	+	4/5	4/16
	-	0	0
	-	0	0
	+	5/8	5/16
	-	0	0
	-	0	0

$$\text{precision} = \frac{5}{10} = .5$$

$$\text{average precision} = \frac{1 + 1 + 3/4 + 4/5 + 5/8}{5} = .478$$

$$\text{recall} = \frac{5}{16}$$

$$F_1 \text{ score} = \frac{2 \left(\frac{5}{10} \right) \left(\frac{5}{16} \right)}{\left(\frac{5}{10} \right) + \left(\frac{5}{16} \right)} = .385$$

d.)		Cumulative Gain	Discounted Cumulative Gain
1	+	1	1
2	+	2	$1 + 1/\log(2)$
3	-	2	$1 + 1/\log(2) + 0$
4	+	3	$1 + 1/\log(2) + 1/\log(3)$
5	+	4	$1 + 1/\log(2) + 1/\log(3) + 1/\log(4)$
6	-	4	
7	-	4	$1 + 1/\log(2) + 1/\log(3) + 1/\log(4) = $ 2.93
	+		
	-		
	-		

Cumulative Gain at 7 = 4

$$\text{NDCG at } 7 = \frac{1}{1 + 2.93} = \frac{1}{3.93} = \text{.68}$$

2.) a.)
$$\sum_{w \in D} C(w, q) \log \frac{P_s(w|\theta_D)}{\alpha_D P(w|C)} + |Q| \cdot \log \alpha_D$$

Let $\alpha_D = \lambda$ and $P_s(w|\theta_D) = (1-\lambda)P_{ML}(w|\theta_D) + \lambda P(w|C)$

then

$$= \sum_{w \in D} \frac{C(w, q) \log ((1-\lambda)P_{ML}(w|\theta_D) + \lambda P(w|C))}{\lambda P(w|C)} + |Q| \cdot \log \lambda$$

b/c $|Q| \cdot \log \lambda$ does not depend on the document.
We can remove it

$$= \sum_{w \in D} \frac{c(w, q) \log((1-\lambda) P_{ML}(w|\Theta_D) + \lambda p(w|c))}{\lambda p(w|c)}$$

$$= \sum_{w \in D} \frac{c(w, q) \log((1-\lambda) P_{ML}(w|\Theta_D))}{\lambda p(w|c)} + \frac{\lambda p(w|c)}{\lambda p(w|c)}$$

$$= \sum_{w \in D} c(w, q) \log \left[\frac{(1-\lambda) P_{ML}(w|\Theta_D)}{\lambda p(w|c)} + 1 \right]$$

we only care about words in both doc and query \therefore

$$\text{Score}(Q, D) = \sum_{w \in Q \cap D} c(w, q) \cdot \log \left[1 + \frac{(1-\lambda) P_{ML}(w|\Theta_D)}{\lambda p(w|c)} \right]$$

$$P(w|c) = P(w|REF) \text{ and } P_{ML}(w|\Theta_D) = \frac{c(w, D)}{|D|}$$

\therefore

$$\text{Score}(Q, D) = \sum_{w \in Q \cap D} c(w, q) \cdot \log \left[1 + \frac{(1-\lambda) \cdot c(w, D)}{\lambda P(w|REF) |D|} \right]$$

b) $C(w, Q)$ is the query vector, D, D is the document vector. The similarity function is the entire function, TF is captured in $C(w, Q)$ and $C(w, D)$. IDF is captured in $\log \frac{1}{P(w|REF)/|D|}$

Document length normalization is $|D|$.

c) $\alpha_D = \frac{N}{|D|+N}$ $p_s(w|\theta_D) = \frac{C(w, D) + Np(w|C)}{|D|+N}$

Subbing these values in we get

$$= \sum_{w \in D} C(w, Q) \log \left[\frac{C(w, D) + Np(w|C)(|D|+N)}{|D| p(w|C)(|D|+N)} \right] +$$

$$|Q| \log \left[\frac{N}{|D|+N} \right]$$

$$= \sum_{w \in D} C(w, Q) \log \left[\frac{C(w, D) + Np(w|C)}{N p(w|C)} \right] + |Q| \log \left[\frac{N}{|D|+N} \right]$$

$$= \sum_{w \in D} C(w, Q) \log \left[\frac{C(w, D) + 1}{N p(w|C)} \right] + |Q| \log \left[\frac{N}{|D|+N} \right]$$

$$= \sum_{w \in D} C(w, Q) \log \left[\frac{C(w, D) + 1}{N p(w|C)} \right] + |Q| \log(N) - \log(|D|+N)$$

$\log(N)$ does not depend on d so we can remove it

$$= \sum_{w \in Q \cap D} C(w, Q) \log \left[1 + \frac{C(w, D)}{K \cdot p(w|C)} \right] - |Q| \log(|D| + N)$$

~~III~~

d. Q is the query vector D is the document vector,
 $C(w, Q)$ and $C(w, D)$ are term frequency. $-|Q| \log(|D| + N)$
 is document length normalization. IDF is $p(w|C)$.

e.) Jelinek-Mercer

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

$$\text{Score}(Q, D') = \sum_{w \in Q \cap D} C(w, Q) \log \left(\frac{1 + (1-\lambda) \cdot K C(w, D)}{\lambda \cdot p(w|REF) \cdot K |D|} \right)$$

K cancels out $\therefore \text{Score}(Q, D) = \text{Score}(Q, D')$

~~IV~~

Dirichlet prior:

$$\text{Score}(q, D) = \sum_{w \in q \cap D} c(w, q) \log \left[1 + \frac{c(w, D)}{N P(w|c)} \right] - |q| \log(|D|/N)$$

$$\text{Score}(q, D') = \sum_{w \in q \cap D} c(w, q) \log \left[1 + \frac{K c(w, D)}{N P(w|c)} \right] - |q| \log(K|D|/N)$$

$$\text{Score}(q, D) \neq \text{Score}(q, D')$$

Therefore Jelinek-Mercer does not overpenalize for a long document but Dirichlet prior does.