

Chen Zhang czhang49

CS 598 CXZ Fall 2016 HW1

① a) Because different words have different weights to the meaning of query. e.g. in a search for "news about presidential campaign", about should be given less weights compared with "presidential".

b) $IDF(w) = \log\left(\frac{M+1}{K}\right)$, K being number of docs containing w , M being total number of docs.

If word w is contained in d , $M' = M+1$, $K' = K+1$, since $M > K \Rightarrow IDF(w)$ will decrease.

If word w is not contained in d , $M' = M+1$, $K' = K$, $\Rightarrow IDF(w)$ will increase

c) i) $precision = \frac{\text{relevant doc \#}}{\text{retrieved total doc \#}} = \frac{5}{10} = \frac{1}{2}$

ii) $Recall = \frac{\text{relevant doc \#}}{\text{total relevant \#}} = \frac{5}{16}$

iii) $F_1 = \frac{2PR}{P+R} = \frac{2 \times \frac{1}{2} \times \frac{5}{16}}{\frac{5}{16} + \frac{1}{2}} = \frac{\frac{5}{16}}{\frac{13}{16}} = \frac{5}{13}$

iv) $Average\ Precision = \frac{\frac{1}{1} + \frac{2}{2} + \frac{3}{4} + \frac{4}{5} + \frac{5}{8}}{10} = \frac{\frac{167}{40}}{10} = \frac{167}{400}$