

1) a) Raw term frequency doesn't account for how big the document is in terms of the number of words.

In addition, stop words will occur more frequently but they should be equally likely among most documents in the corpus. Therefore, a better measure would be the term-frequency-inverse document frequency which reflects how important a word is with respect to other documents in the corpus. Usually, the cosine similarity is used rather than the dot product because the magnitude doesn't matter.

So cosine similarity normalizes the document length of the two vectors that we are comparing.

b) Suppose the corpus that we're interested in has  $X$  instances of the term. Originally if  $IDF = \log \frac{N}{|\{d \in D : t \in d\}|} = \log \frac{N}{B}$

then now  $IDF = \log \frac{N+1}{B+X}$ , so the IDF would decrease since the term is now more common in the corpus since we made a copy of  $d$ .

c)

	Retrieved	Not Retrieved
Relevant	5	6
Not Relevant	5	0

$N = 16$   
 $\leftarrow$  all 16 docs are relevant

i)  $P = \frac{5}{10} = 0.5$

ii)  $R = \frac{5}{11} = 0.4545$

iii)  $F_1 = \frac{2 \cdot \frac{5}{10} \cdot \frac{5}{11}}{\frac{5}{10} + \frac{5}{11}} = 0.476$

Rank as follows

iv)  $\begin{matrix} + & + & - & + & + & - & - & + & - & - \\ 1 & 1 & 0 & \frac{3}{4} & \frac{4}{5} & 0 & 0 & \frac{5}{8} & 0 & 0 \end{matrix} \Rightarrow \text{sum} = 4.175 \xrightarrow{\div 16} \text{avg} = 0.26$

d) i)  $CG = 5$   $IDCG = 1 + 1 + \frac{1}{\log_2 3} + \frac{1}{2} + \frac{1}{\log_2 5} = 3.562$

ii)  $DCG = 1 + 1 + \frac{1}{\log_2 4} + \frac{1}{\log_2 5} + \frac{1}{\log_2 8} = 3.264$

$NDCG = \frac{3.264}{3.562} = 0.916$

2)	V	$P(A=1 V)$	$P(K=1 V)$	$P(L=1 V)$	$P(V)$
a)	0	2/6	3/6	2/6	1/2
	1	4/6	5/6	2/6	1/2

$$b) P(V=1|A=0, K=1, L=0) = P(V=1|A=0) \cdot P(V=1|K=1) \cdot P(V=1|L=0)$$

$$= \frac{P(A=0|V=1)P(V)}{P(A)} \cdot \frac{P(K=1|V=1)P(V)}{P(K=1)} \cdot \frac{P(L=0|V=1)P(V)}{P(L=0)}$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$= \frac{2/6 \cdot 1/2}{1/2} \cdot \frac{5/6 \cdot 1/2}{8/12} \cdot \frac{4/6 \cdot 1/2}{8/12}$$

$$= \frac{5}{48} = 0.104$$

$$P(V=0|A=0, K=1, L=0) = P(V=0|A=0) \cdot P(V=0|K=1) \cdot P(V=0|L=0)$$

$$P(V=0|A=0, K=1, L=0) + P(V=1|A=0, K=1, L=0) \stackrel{\text{should be 1}}{=} \frac{P(A=0|V=0)P(V)}{P(A)} \cdot \frac{P(K=1|V=0)P(V)}{P(K=1)} \cdot \frac{P(L=0|V=0)P(V)}{P(L=0)}$$

$$\because P(V=0|A=0, K=1, L=0) > P(V=1|A=0, K=1, L=0)$$

Therefore the message most likely does not contain a virus

$$= \frac{4/6 \cdot 1/2}{1/2} \cdot \frac{3/6 \cdot 1/2}{8/12} \cdot \frac{4/6 \cdot 1/2}{8/12}$$

$$= \frac{1}{8} = 0.125$$

c) If computed in the direct way as in a) then

$$P(V=1|A=0, K=1, L=0) = 0.5, P(V=0|A=0, K=1, L=0) = 0.5$$

this is because the pair  $(A, K, L) = (0, 1, 0)$  is treated as one event but rather in the b) calculation they are treated as separate events.

d) For the prior  $P(V=0)$  and  $P(V=1)$  must sum to 1

so we know that  $P(V=1) = 1/2$ . Other than that

the posteriors should be arbitrary

2) e) Changing the  $A$  value in sample #0 would mean that  $p(V=1 | A=0)$  would change so the  $p(V=1 | A=0, K=1, L=0)$  result would change.

$$f) p(A, L, K | V) = p(A | V) p(L | V) p(K | V) \\ = \frac{p(V | A) p(A)}{p(V)} \frac{p(V | L) p(L)}{p(V)} \frac{p(V | K) p(K)}{p(V)}$$

We need to specify only two posteriors  $\because p(V | A) + p(V | L) + p(V | K) = 1$   
then we need to specify 4 priors  $p(A), p(L), p(K), p(V)$

So we need 6 probability values as a minimum.

g) The events  $A, L, K$  may not be completely independent.  
For example, maybe messages that are shorter than 10 words are more likely to contain an attachment because the content is in the attachment so that the messages simply refer to the documents. So if event  $A$  &  $L$  are dependent, then the independence assumption  $p(A, L, K | V) = p(A | V) p(L | V) p(K | V)$  is not completely valid

$$3) \ a) \quad l = p(X = \{x_1, \dots, x_n\}) = \frac{u^{x_1} e^{-u}}{x_1!} \cdot \frac{u^{x_2} e^{-u}}{x_2!} \cdot \dots \cdot \frac{u^{x_n} e^{-u}}{x_n!}$$

$$l = \frac{u^{\sum x_i} e^{-nu}}{\prod x_i!} \Rightarrow L = \log l = \sum x_i \log u - nu \log e - \log \prod x_i!$$

$$L = \log u \sum x_i - nu - \log \prod x_i!$$

$$\frac{\partial L}{\partial u} = \frac{\sum x_i}{u} - n = 0 \Rightarrow \hat{u} = \frac{\sum x_i}{n},$$

Best parameter is sample average

$$b) \quad l = p(U = \{u_1, \dots, u_n\}) = \lambda e^{-\lambda u_1} \cdot \lambda e^{-\lambda u_2} \cdot \dots \cdot \lambda e^{-\lambda u_n}$$

$$l = \lambda^n e^{-\lambda \sum u_i}$$

$$L = \log l = n \log \lambda - \lambda \sum u_i \log e$$

$$\frac{\partial L}{\partial \lambda} = \frac{n}{\lambda} - \sum u_i = 0 \Rightarrow \left| \hat{\lambda} = \frac{n}{\sum u_i} \right|$$