

19) To Compute their RSV for the doc 11, doc 12  
 Let analyze the relevance of each term in the query vector using document-term Matrix  
Summary of T1 to T6 in terms of their occurrence in relevant and no relevant documents.

Term	Relevance (R)	Non-Relevance (NR)
T1	3	3
T2	3	2
T3	2	3
T4	3	4
T5	1	5
T6	3	1

Given query vector  $(1, 1, 0, 1, 0, 1)$  and 0.5 to Smooth Count

$S$  = number of documents containing the term that are relevant.

$S$  = total number of relevant documents in the doc-term-matrix

$n-S$  = number of documents not containing the term that are relevant.

$N-S$  = total number of non-relevant documents minus

$N$  = total number of documents = 10

$x_i = 1$  if the term is in the document,  $x_i = 0$  if

For each term in the query, we compute  $P(R|x_i)$  where  $x_i = 1$  or 0  
 then multiply these probabilities together

RSV for doc11  $\Rightarrow (0, 1, 1, 0, 0, 1)$  and doc2  $\Rightarrow (1, 0, 1, 1, 0, 1)$

For Doc11

$$T_1 = 0 : S - s + 0.5 \text{ and } N - n - S + s + 0.5$$

$$T_2 = 1 : S + 0.5 \text{ and } n - s + 0.5$$

$$T_3 = 1 : S + 0.5 \text{ and } n - s + 0.5$$

$T_3 = 1 : S + 0.5$  and  $n - s + 0.5$

$T_4 = 0 : S - s + 0.5$  and  $N - n - S + s + 0.5$

$T_5 = 0 : S - s + 0.5$  and  $N - n - S + s + 0.5$

$T_6 = 0 : S + 0.5$  and  $n - s + 0.5$

For doc 12

$T_1 = 1 : S + 0.5$  and  $n - s + 0.5$

$T_2 = 0 : S - s + 0.5$  and  $N - n - S + s + 0.5$

$T_3 = 1 : S + 0.5$  and  $n - s + 0.5$

$T_4 = 1 : S + 0.5$  and  $n - s + 0.5$

$T_5 = 0 : S - s + 0.5$  and  $N - n - S + s + 0.5$

$T_6 = 1 : S + 0.5$  and  $n - s + 0.5$

For a term  $i$  where  $x_i = 1$

$$RSV = \log\left(\frac{S+0.5}{n-s+0.5}\right) + \log\left(\frac{N-n-S+s+0.5}{S-s+0.5}\right)$$

$x_i = 0$

$$RSV = \log\left(\frac{S-s+0.5}{N-n-S+s+0.5}\right) + \log\left(\frac{n-s+0.5}{s+0.5}\right)$$

For doc 11 ( $T_1 = 0$ )  $S = 3$  for  $T_1$   $S = 4$  in

$(NR+R) \Rightarrow n = \text{total number of documents where the term appears} = 6$

$N = \text{Total document (N)} = 10$

For doc 11 ( $T_1 = 0$ )

$$= \log\left(\frac{4-3+0.5}{10-6+4+3+0.5}\right) + \log\left(\frac{6-3+0.5}{3+0.5}\right)$$

$$= \log(1.5/3.5) + \log(1)$$

Using base 10  $\approx -0.85$

For doc 12 ( $T_1 = 1$ ) ~~without~~

$$= \log\left(\frac{3+0.5}{6-3+0.5}\right) + \log\left(\frac{10-6-4+3+0.5}{4-3+0.5}\right)$$

$\approx 0.37$

Doc 12 T1 has a higher RSV contribution

(19)

i) (T2) For doc 1:  $s=3, S=4, n=5, N=10$

$$RSV(T_2=1) = \log\left(\frac{3+0.5}{5-3+0.5}\right) + \log\left(\frac{10-5-4+3+0.5}{4-3+0.5}\right) \approx 0.62$$

For doc 12 ( $T_2=0$ )

$$= \log\left(\frac{4-3+0.5}{10-5-4+3+0.5}\right) + \log\left(\frac{5-3+0.5}{3+0.5}\right) \approx 0.3$$

(T3) Doc 11 T2 has a higher RSV contribution

For doc 11 for (T3)  $s=2, S=4, n=5, N=10$

Doc 11  $T_3=1$  and Doc 12  $T_3=1$  has equal RSV contribution

For doc 11 for (T4)  $s=3, S=4, n=7, N=10$

$$RSV(T_4=0) = \log\left(\frac{4-3+0.5}{10-7+4+0.5}\right) + \log\left(\frac{7-3+0.5}{3+0.5}\right) \approx -0.11$$

For doc 12,  $T_4=1$

$$= \log\left(\frac{3+0.5}{7-3+0.5}\right) + \log\left(\frac{10-7-4+3+0.5}{4-3+0.5}\right) \approx 0.11$$

Doc 12 T4 has a higher RSV contribution

(T5) For doc 11  $T_5=0$ , doc 12  $T_5=0$

$s=1, S=4, n=6, N=10$  they will have equal RSV contribution

(T6) For doc 11  $T_6=1$ , doc 12  $T_6=1$

$s=3, S=4, n=4, N=10$  they have equal RSV contribution

doc 11 RSV

$$= -0.85 + 0.62 + 0.11 = -0.34$$

doc 12 RSV

$$= 0.37 + 0.3 + 0.11 = 0.78$$

doc 12 RSV is higher than doc 11

$$BM25 = \sum_{t \in d} \left( IDF_{(t)} \times \frac{(K_1 + 1) \cdot tf_i}{tf_i + [(1 - b) + b \cdot \frac{dl}{avgdl}] \cdot K_1} \right)$$

$IDF_{(t)} \Rightarrow \log \frac{N}{df_i} \Rightarrow$  inverse document frequency of term t

$tf_i \Rightarrow$  term frequency of t in document d

$K_1 \approx 2$ ,  $b = 0.75$ ,  $avgdl = 1600$

$|dl| = 1000$ ,  $|d_1| = 1,500$ ,  $|d_2| = 1,200$

$IDF(\text{car}) = 1.65$ ,  $IDF(\text{Auto}) = 2.08$ ,  $IDF(\text{Insurance}) = 1.62$

$IDF(\text{Best}) = 1.5$

To calculate BM25 score for each documents

$$BM25(d_1) = \sum \left[ \begin{array}{l} 1.65 \times \frac{(2+1) \times 27}{27 + 2 \times [(1-0.75) + 0.75 \times \frac{1000}{1600}]} \\ + \\ 2.08 \times \frac{(2+1) \times 3}{3 + 2 \times [(1-0.75) + 0.75 \times \frac{1000}{1600}]} \\ + \\ 1.62 \times \frac{(2+1) \times 0}{0 + 2 \times [(1-0.75) + 0.75 \times \frac{1000}{1600}]} \\ + \\ 1.5 \times \frac{(2+1) \times 14}{14 + 2 \times [(1-0.75) + 0.75 \times \frac{1000}{1600}]} \end{array} \right]$$

~~BM25(d<sub>1</sub>)~~

$$BM25(d_2) = \sum \left[ \begin{array}{l} 1.65 \times \frac{(2+1) \times 4}{4 + 2 \times [(1-0.75) + 0.75 \times \frac{1500}{1600}]} \\ + \\ 2.08 \times \frac{(2+1) \times 33}{33 + 2 \times [(1-0.75) + 0.75 \times \frac{1500}{1600}]} \\ + \\ 1.62 \times \frac{(2+1) \times 33}{33 + 2 \times [(1-0.75) + 0.75 \times \frac{1500}{1600}]} \end{array} \right]$$

$$BM25(d_2) \left[ 0.5 \times \frac{(2+1) \times 0}{0 + 2 \times [(1 - 0.75) \times 0.75 \times \frac{1500}{1600}]} \right]$$

w

$$BM25(d_3) \left\{ \begin{array}{l} 1.65 \times \frac{(2+1) \times 24}{24 + 2 \times [(1 - 0.75) \times 0.75 \times \frac{1200}{1600}] } + \\ 2.08 \times \frac{(2+1) \times 0}{0 + 2 \times [(1 - 0.75) \times 0.75 \times \frac{1200}{1600}]} - \\ 1.62 \times \frac{(2+1) \times 29}{29 + 2 \times [(1 - 0.75) \times 0.75 \times \frac{1200}{1600}]} + \\ 1.5 \times \frac{(2+1) \times 17}{17 + 2 \times [(1 - 0.75) \times 0.75 \times \frac{1200}{1600}]} \end{array} \right\}$$

w

- 2)   
 $d_1$ : click go the shears boys click click click  
 $d_2$ : Click click  
 $d_3$ : metal here  
 $d_4$ : metal shears click here

$$\lambda = 0.5$$

The model will be MLE unigram model from the documents and collection, mixed.

Query 1 : click

$$P(q_1|d_1) = \lambda^{4/8} + (1-\lambda)^{7/16} = 0.5^{1/2} + 0.5^{7/16} \approx 0.47$$

$$P(q_1|d_2) = 0.5^1 + 0.5^{7/16} \approx 0.719$$

$$P(q_1|d_3) = 0.5^0 + 0.5^{7/16} \approx 0.22$$

$$P(q_1|d_4) = 0.5^{1/4} + 0.5^{7/16} \approx 0.345$$

So the ranking is  $d_2 > d_1 > d_3 > d_4$

Query 2 Shears

$$P(q_1|d_1) = 0.5^{1/8} + 0.5^{2/16} = 0.063 + 0.063 = 0.126$$

$$P(q_1|d_2) = 0.5^0 + 0.063 = 0.063$$

$$P(q_1|d_3) = 0.5^0 + 0.063 = 0.063$$

$$P(q_1|d_4) = 0.5^{1/8} + 0.5^{2/16} = 0.126$$

Ranking  $\Rightarrow d_1 \rightarrow d_2 \rightarrow d_3$

Query 3 Click Shears

$$P(q_1|d_1) = [0.5^{1/8} + 0.5^{9/16}] * [0.5^{1/8} + 0.5^{2/16}] = 0.47 + 0.126 = 0.593$$

$$P(q_1|d_2) = 0.719 + 0.063 = 0.782$$

$$P(q_1|d_3) = 0.22 + 0.063 = 0.283$$

$$P(q_1|d_4) = 0.345 + 0.126 = 0.471$$

Ranking  $d_2 > d_1 > d_4 > d_3$

### 3.) Evaluation of Unranked list

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\# \text{ retrieved items}} = P(\text{relevant} | \text{retrieved})$$

$$\text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{total relevant documents})} = P(\text{retrieved} | \text{relevant in coll})$$

These notion can be made clear by examining the following Contingency Table

	Relevant	Not Relevant
Retrieved	TP	FP
Not Retrieved	FN	TN

Total\_retrieved\_docs = 20, Total\_relevant\_items = 8

relevant\_items\_retrieved = 6

a. Precision of the system on the top 20 =  $6/20 = 0.3$

b. F1 Score on the top 20 =  $\frac{2 \times (\text{Precision}_{\text{at\_20}} \times \text{recall}_{\text{at\_20}})}{(\text{Precision}_{\text{at\_20}} + \text{recall}_{\text{at\_20}})}$

To calculate  $\text{Recall}_{\text{at\_20}} = 6/8 = 0.75$

$$= \frac{2 \times (0.3 \times 0.75)}{0.3 + 0.75} = 0.428$$

c. Uninterpolated precision at 25%, where required recall = 0.25

To calculate this we created a table to analyze precision and recall at each retrieved items.

i	item	Precision <sub>at i</sub>	Recall <sub>at i</sub>	Cumulative Cont. of Relevant
1	$\rightarrow R$	$\frac{1}{1} = 1$	$\rightarrow \frac{1}{8} \rightarrow 0.125$	1
2	$\rightarrow R$	$\frac{2}{2} = 1$	$\rightarrow \frac{2}{8} \rightarrow 0.25$	2
3	$\rightarrow N$	$\frac{2}{3} \rightarrow 0.67$	$\rightarrow 0.25$	
4	$\rightarrow N$	$\frac{2}{4} \rightarrow 0.5$	$\rightarrow 0.25$	
5	$\rightarrow N$	$\frac{2}{5} \rightarrow 0.4$	$\rightarrow 0.25$	
6	$\rightarrow N$	$\frac{2}{6} \rightarrow 0.33$	$\rightarrow 0.25$	
7	$\rightarrow N$	$\frac{2}{7} \rightarrow 0.286$	$\rightarrow 0.25$	
8	$\rightarrow N$	$\frac{2}{8} \rightarrow 0.25$	$\rightarrow 0.25$	
9	$\rightarrow R$	$\frac{3}{9} \rightarrow 0.33$	$\rightarrow 0.375$	3
10	$\rightarrow N$	$\frac{3}{10} \rightarrow 0.3$	$\rightarrow 0.375$	
11	$\rightarrow R$	$\frac{4}{11} \rightarrow 0.36$	$\rightarrow 0.5$	4
12	$\rightarrow N$	$\frac{4}{12} \rightarrow 0.33$	$\rightarrow 0.5$	
13	$\rightarrow N$	$\frac{4}{13} \rightarrow 0.31$	$\rightarrow 0.5$	
14	$\rightarrow N$	$\frac{4}{14} \rightarrow 0.29$	$\rightarrow 0.5$	
15	$\rightarrow R$	$\frac{5}{15} \rightarrow 0.33$	$\rightarrow 0.625$	5
16	$\rightarrow N$	$\frac{5}{16} \rightarrow 0.31$	$\rightarrow 0.625$	
17	$\rightarrow N$	$\frac{5}{17} \rightarrow 0.29$	$\rightarrow 0.625$	
18	$\rightarrow N$	$\frac{5}{18} \rightarrow 0.28$	$\rightarrow 0.625$	
19	$\rightarrow N$	$\frac{5}{19} \rightarrow 0.26$	$\rightarrow 0.625$	
20	$\rightarrow R$	$\frac{6}{20} \rightarrow 0.3$	$\rightarrow 0.75$	6

recall

first Value

c. Uninterpolated precision at 25% recall = 1 at which recall >= 25%

d. Interpolated precision at 33% recall : This will be

where taking the maximum of precision >= 33% of

recall  $\max[0.33, 0.3, 0.36, 0.33, 0.31, 0.29, 0.33, 0.31, 0.28, 0.26, 0.3]$

$$= 0.36$$

e. Mean Average Precision MAP for the query

We Isolate Precision based on relevant item

$$\frac{0.125 + 0.25 + 0.375 + 0.5 + 0.625 + 0.75}{6}$$

3e)  $\text{MAP}_{\text{for\_query}} = \frac{\text{Sum}(\text{average Precisions})}{\text{total\_relevant\_docs in coll}}$

$$= 1 + 1 + 0.33 + 0.36 + 0.33 + 0.3 = 3.32$$

$$= \cancel{3.32/20} \quad 3.32/8 \approx 0.415 \Rightarrow 41.5\%$$

3f.) To achieve the largest possible MAP that this system could have? - This could be achieved if all the relevant documents were ranked at the top of the retrieved list  
 $\Rightarrow 1 + 1 + 1 + 1 + 1 + 1 + 1 = 6/8 = 75\%$

3g) For Smallest  $\overset{\text{MAP}}{\underset{\text{could}}{\text{thus}}}$  achieved by ranking the relevant at the end of the list

$$\Rightarrow \frac{1/15 + 2/16 + 3/17 + 4/18 + 5/19 + 6/20}{8} \Rightarrow$$

$$= \frac{0.067 + 0.125 + 0.18 + 0.22 + 0.263 + 0.3}{8}$$

$$= 1.155/8 = 0.144 \Rightarrow 14.42\%$$

3h) The error from the estimated MAP to the largest possible MAP is:  $0.75 - 0.414 = 0.336$   
 The error from the estimated MAP to the smallest possible MAP:  $0.414 - 0.1442 = 0.2698$   
 $(0.336, 0.2698)$