1    **logical view of document**

Documents in a collection are represented through set of index terms or key words. Such keywords might be extracted directly from the text of the document. No matter if representative keywords are derived automatically or generated by a specialist they provide a logical view of the document

The user of retrieval system has to translate his information need into a query in the language provided by the system. We make a clear distinction between the different task the user of the retrieval system might be engaged in the task might be of 2 different type Information or data retrieval and browsing. Data retrieval system a query expression is used to convey the constraints that must be satisfied by objects in the answer set. Browsing in language of world wide web pulling action where the user request the information in an interactive manner, or by using software agent which push information towards the user.

4   **Precision :** It is defined as the ratio or fraction of the retrieved documents which is relevant

$$\text{Precision} = \frac{|Ra|}{|A|}$$

$|A|$ no of doc in this set

$|Ra|$ No. of doc intersection of used $R \cap A$

**Recall :** It is defined as the fraction of the relevant documents which has been retrieved

$|R|$ No of doc in this set

$$\text{Recall} = \frac{|Ra|}{|R|}$$

$$A = \{ d_{28}, d_{39}, d_1, d_{32}, d_{98}, d_{29}, d_{81}, d_{93} \}$$

$$R = \{ d_9, d_{18}, d_{32}, d_{40}, d_{29}, d_{39}, d_{92}, d_{28} \}$$

$$|A| = 8 \qquad |R| = 8$$

$$\{Ra\} = \{ d_{28}, d_{29}, d_{32}, d_{39} \}$$

$$|Ra| = 4.$$

1 of 3

Nikitha M S
ANII7J8061
ISE 8 'A'
J80588 - IRS

$$Precision = \frac{4}{8} = \frac{1}{2}$$

$$Recall = \frac{|Ra|}{|R|} = \frac{4}{8} = \frac{1}{2}$$

5   R - precision : To generate a single value summary of the ranking by computing precision of $p^{th}$ position in the ranking, where R is the total number of relevant documents for a current query. The R-precision measure is useful parameter for observing the behavior of an algorithm for each individual query in an experiment.

Answered Alg A → $\{d_9, d_{39}, d_{38}, d_{26}, d_1\}$

Answered alg B → $\{d_{28}, d_{18}, d_{40}, d_{27}, d_9, d_{20}\}$

R → $\{d_9, d_{18}, d_{32}, d_{40}, d_{29}, d_{39}, d_{92}, d_{28}\}$

$$|R| = 8$$

$$RP_A = \frac{2}{8} = \frac{2}{8} = \frac{1}{4} = 0.25$$

$$RP_B = \frac{4}{8} = \frac{1}{2} = 0.5$$

2   $D_1$ = cat and rat are animals

$D_2$ = Animals chased their Pray

$Q$ = cat chased rat

freq $(i,j)$

$N$ - Total no of doc

$n_i$ - no of doc which contain the keyword $k_i$

$$idf = log\left(\frac{N}{n_i}\right)$$

$$w_{i,j} = freq(i,j) * idf$$

$$Sim(d_j, q) = \frac{\vec{d_j} \cdot \vec{q}}{|\vec{d_j}| \cdot |\vec{q}|}$$

$$= \frac{\sum_{i=1}^{t} w_{i,j} \times w_{i,q}}{\sqrt{\sum_{i=1}^{t} w_{i,j}^2} \times \sqrt{\sum_{i=1}^{t} w_{i,q}^2}}$$

2 of 3

3    document similarity    $Sim(D, Q)$

$$Sim(d_j, q) = \frac{\vec{d_j} \cdot \vec{q}}{|\vec{d_j}| \cdot |\vec{q}|}$$

$$= \frac{\sum_{i=1}^{t} w_{i,j} \times w_{i,q}}{\sqrt{\sum_{i=1}^{t} w_{i,j}^2} \times \sqrt{\sum_{i=1}^{t} w_{i,q}^2}}$$

$$Sim(Q, D) = \Pi \; \frac{P(t_i/R)}{P(t_i/\bar{R})} \; \frac{P(R)}{P(\bar{R})}$$

| Documents | Revelence | $t_1$ | $t_2$ |
|-----------|-----------|-------|-------|
| $D_1$ | R | 0 | 1 |
| $D_2$ | NR | 1 | 0 |
| $D_3$ | NR | 1 | 1 |
| $D_4$ | R | 1 | 0 |
| $D_5$ | NR | 1 | 1 |

Relevant $\rightarrow$   $\overset{t_2}{D_1}$   $\overset{t_1}{D_4}$

non relevant $\rightarrow$   $\underset{t_1}{D_2}$   $\underset{t_1, t_2}{D_3}$   $\underset{t_1, t_2}{D_5}$

$$Sim(Q, D_2) = \frac{1/2}{2/3} \cdot \frac{2/5}{3/5} = 0.5$$

$$Sim(Q, D_3) = \frac{1/2}{2/3} \cdot \frac{2/5}{3/5} \cdot \frac{2/5}{1/3} \cdot \frac{2/5}{3/5}$$

$$= \frac{5}{4} \times \frac{2}{3} \times 3 \times \frac{2}{3}$$

$$= 0.6 \quad 1$$