

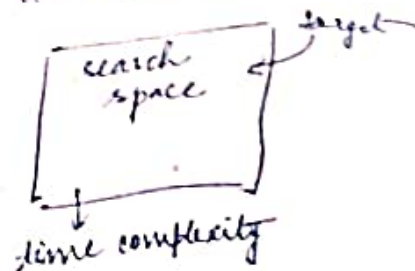
So. 09.21

Search Technique

What is searching?

Eg: - linear
- binary.

— target in a search space



Storage

i) Alphabetical order (Dictionary arrangement)

- The data are arranged in alphabetical

ii) Numerical arrangement.

- ascending / descending order

iii) Classified arrangement
↳ using categories

40-41
Steps that involve in search process.

- i) Recognize & state the need.
- ii) Development of search strategy.
- iii) Execution of search strategy.
- iv) Review search results.
- v) Edit search results.
- vi) Evaluation & feedback.

← efficient method eg: for many types of search

Search string → structure

- i) Syntactic value
- ii) Semantic value
- iii) Boolean operators

meaning

OR
AND
NOT

Different types of search

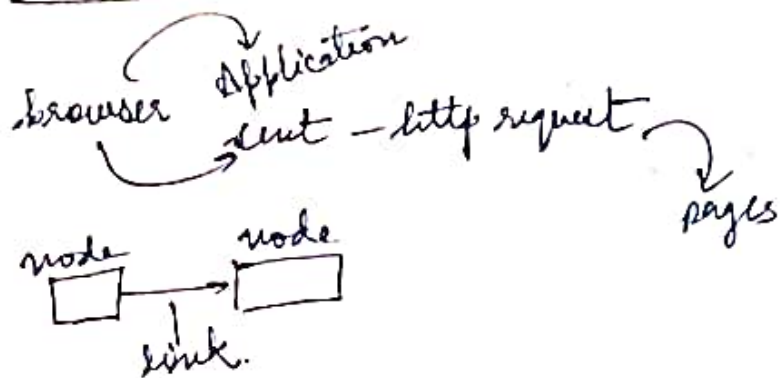
- i) Field based search
- ii) Full text search
- iii) Truncation search
- iv) Proximity search
- v) Limiting search
- vi) Range search
- vii) Simple search
- viii) Advance search

DB (SQL) query

→ search is conducted diff. form of word having the same root.

01/10/24

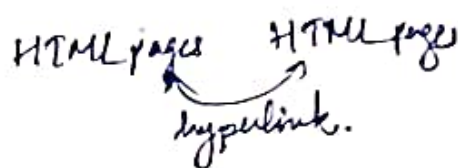
Web Search basics



Unit 4 -

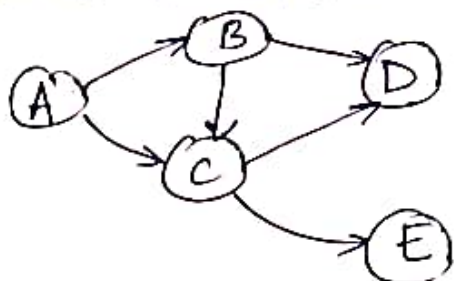
Advanced topic

Ranking
Relevance feedback
Ranking & Query Expansion



WEB GRAPH

It consists of HTML pages together with the hyperlink represented as directed graph where each web page is a node & each hyperlink as a directed edge.



Indegree

$$C = 2$$

$$B = 1$$

$$D = 2$$

POWER LAW

Total no. of web pages with indegree is proportional to $\frac{1}{i^\alpha}$, α is Zipf's law of distribution of word.

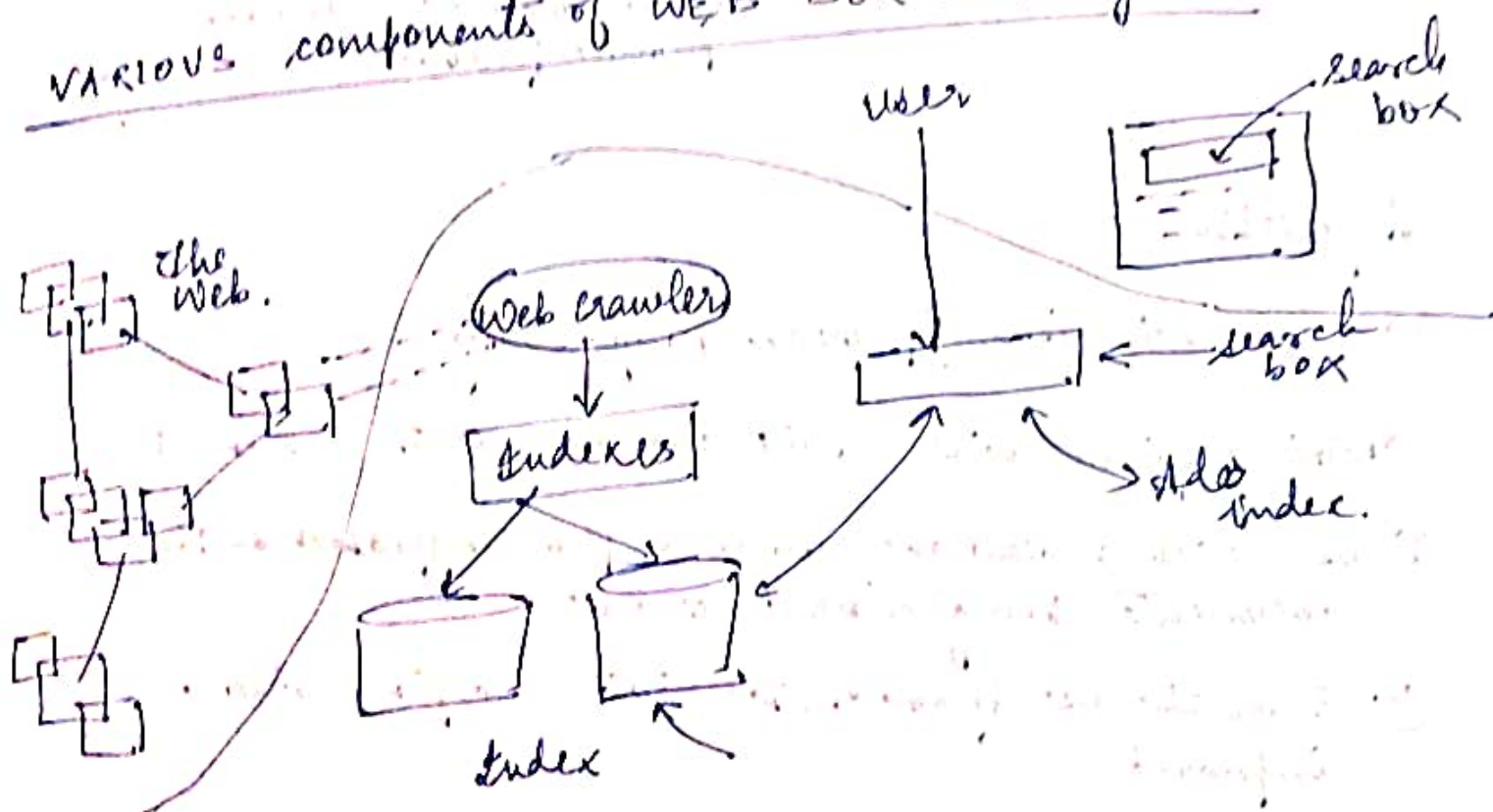
01.10.21

BOWTIE STRUCTURE OF WEB



Spam — advertisements
unwanted/junk
sent in bulk.

VARIOUS components of WEB SEARCH engine

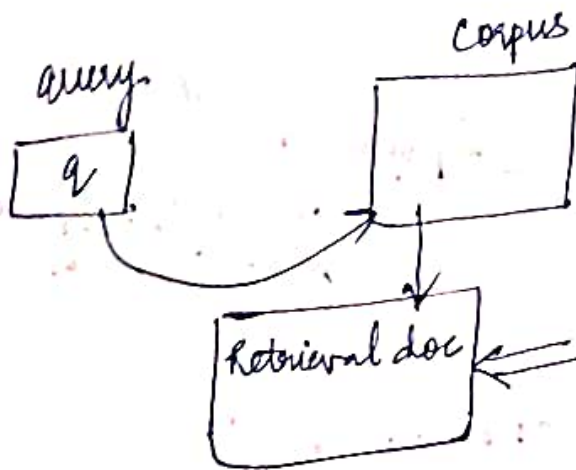


07/10/24

~~Relevance Feedback~~

RELEVANCE FEEDBACK

① Issue a query string



Algorithm:

- ① user will issue / provide query string
- ② Retrieval system will result the documents as relevant
- ③ user will review the results & provide feedback on the documents being relevant or not.
- ④ From the user feedback, the retrieval system will be improved.

07/10/24

Near duplicates & Shingling.

(prev. chap. we
- topic was missed)

Multiple web pages

any 2
duplicates

diff. web pages
same content

Shingling technique

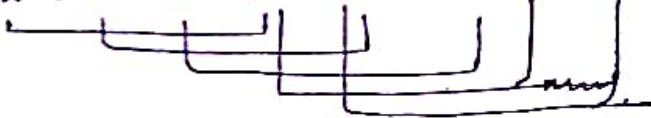
mechanism / tech.

identify the near duplicates web pages.

Given: an integer k & sequence of terms in document d
shingling of document $d \Rightarrow$ set of all consecutive sequence
of k terms in document.

$k=4$

docd - a rose is a rose is a rose.

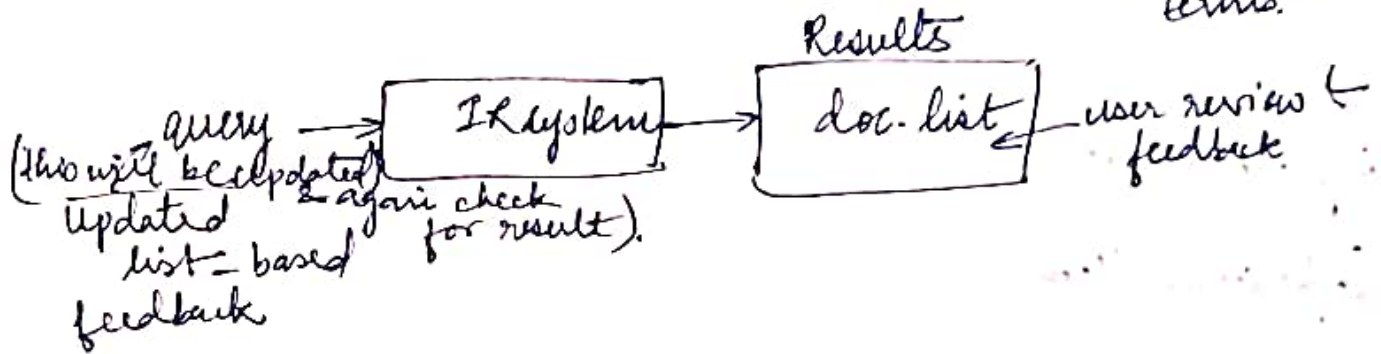


- 1) a rose is a
- 2) rose is a rose
- 3) is a rose is
- 4) a rose is a
- 5) rose is a rose

08/10/24

Query Expansion - Technique to increase the recall.

(To improve the query string so that more relevance doc. are retrieved)
 Relevance feedback - which retrieves the relevant data
 Add more terms /
 Delete unnecessary terms.



JACCARD CO-EFFICIENT

$$J(A, B) = \frac{A \cap B}{A \cup B}$$

$$J[S(d_1), S(d_2)] = \frac{S(d_1) \cap S(d_2)}{S(d_1) \cup S(d_2)} \geq \tau \quad \begin{matrix} \nearrow d_1 \& d_2 \\ \text{are} \\ \text{similar} \end{matrix}$$

can find out the similarity.

$K=4$

slings

$S(d_1)$

- ① This is a test
- ② is a test for
- ③ a test for IR

$d_1 \rightarrow$ This is a test for IR
 $d_2 \rightarrow$ This is a test for IR Assessment
 threshold = 0.7

$S(d_2)$

- ① This is a test
- ② is a test for
- ③ a test for IR
- ④ test for IR assessment

08/10/24

$$J(s(d_1), s(d_2)) = \frac{3}{4} = 0.75 > \text{thr. (0.7)}$$

$\therefore d_1$ & d_2 are similar.

Q.

$s(d_1)$	$s(d_2)$	$s(d_3)$
1	0	1
1	1	0
1	1	1
1	1	0

identify which docs are similar if threshold for $J.C = 0.6$.

$$J(s(d_1), s(d_2)) = \frac{4}{5} = 0.8 > \text{thr.} - \text{similar}$$

$$J(s(d_2), s(d_3)) = \frac{2}{5} = 0.4 < \text{thr.} - \text{Not similar.}$$

$$J(s(d_3), s(d_1)) = \frac{3}{5} = 0.6 \geq \text{thr.} - \text{similar}$$

Q.

$s(d_1)$	$s(d_2)$	$s(d_3)$	$s(d_4)$
1	0	1	1
1	1	0	0
1	1	1	1
1	1	1	1
1	1	0	1

$$\underline{\text{thr} = 0.8}$$

08.10.24

402

$$= \frac{41}{21.21} \\ = \frac{172.31}{21.21} \\ = 6.1$$

 $J(d_1, d_2)$

d_i	d_j	$J(s(d_i), s(d_j))$
d_1	d_2	$4/5 = 0.8 \geq \frac{1}{2} \checkmark$
d_2	d_3	$2/5 = 0.4 < \frac{1}{2} \times$
d_3	d_4	$3/4 = 0.75 \leq \frac{1}{2} \times$
d_4	d_1	$4/5 = 0.8 \geq \frac{1}{2} \checkmark$
d_1	d_3	$3/5 = 0.6 < \frac{1}{2} \times$
d_2	d_4	$3/5 = 0.6 < \frac{1}{2} \times$
d_2	d_4	

d_1 & d_2 } similar.
 d_4 & d_1 }