# Data Mining and Web Algorithms

- Why web search algorithms( information Retrieval systems)
- Different Retrieval Systems
  - Boolean Retrieval
  - Ranked Retrieval
  - Link Analysis
  - Ranking algorithms
- Web Crawlers
- Web Caching Algorithms
- Recommendation Systems

Original PageRank algorithm :-

→ It is an iterative algorithm that starts with gram.

$$PR(A) = (1-d) + d\left[\frac{PR(T_1)}{C(T_1)} + \frac{PR(T_2)}{C(T_2)} \cdots \frac{P(Tn)}{C(Tn)}\right]$$

where,

→ PR(A) = PageRank of page A    | the probability for random surfer jumping to a page via always (1-d). |

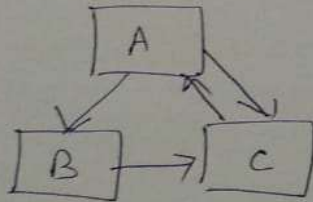→ PR(T_i) = PageRank of Pages $T_i$ which link
          to page A (inbounds link)

→ C(T_i) = No of outbound links on page $T_i$

→ d = damping factor lies B/w 0 & 1.

→ According to Page and Brin (who have proposed)

d = 0.85.

Eg



No of Nodes = 3

Let us assume initial are 1.

$$PR(A) = (1-d) + d\left[\frac{PR(C)}{C(C)}\right]$$

$$= (0.15) + (0.15^-)\left[\frac{PR(C)}{1}\right] \quad\text{———①}$$

$$PR(B) = (1-d) + d\left[\frac{PR(A)}{C(A)}\right]$$

$$= (1-0.15) + (0.15)\left[\frac{PR(A)}{2}\right] \quad\text{———②}$$

(X)

$$PR(c) = (0.5) + (0.5) \left[ \frac{PR(A)}{C(A)} + \frac{PR(B)}{C(B)} \right]$$

$$= 0.5 + 0.5 \left[ \frac{PR(A)}{2} + \frac{PR(B)}{1} \right] \quad —③$$

Let

$$\boxed{PR(A) = P(B) = P(c) = 1} \quad \text{initial} \quad \boxed{0^{th} \text{ iteration}}$$

$$PR(A) = 0.5 + 0.5 = 1 \qquad \qquad \boxed{I^{st} \text{ iteration}}$$

$$PR(B) = 0.75$$

$$PR(c) = 0.5 + 0.5 \left[ 0.5 + 0.75 \right]$$

$$= 0.5 + 0.5 * 1.25$$

$$= 0.5 + \cancel{0.615} = 1.125$$

Sum of PageRank of all the pages = No of web pages

$$= 1 + .75 + 1.125 = 2.875 \, ??$$

It has to repeated (iteratively) until we get the same (or almost same) value of page Rank.

Another notation of pagerank is :-

$$PR(A) = \frac{(1-d)}{N} + d \left[ \frac{PR(T1)}{C(T1)} + \frac{PR(T2)}{C(T2)} \cdots \frac{PR(Tn)}{C(Tn)} \right]$$

Then sum of all pages will be one.

(7)

$$PR(C) = (0.5) + (0.5)\left[\frac{PR(A)}{C(A)} + \frac{PR(B)}{C(B)}\right]$$

$$= 0.5 + 0.5\left[\frac{PR(A)}{2} + \frac{PR(B)}{1}\right] \quad\text{---}\quad ③$$

Let

| $PR(A) = P(B) = P(C) = 1$ | initial | 0th iteration |

$PR(A) = 0.5 + 0.5 = 1$ 

$\left.\begin{array}{l}\text{I}^{st}\text{ iteration}\end{array}\right.$

$PR(B) = 0.75$

$PR(C) = 0.5 + 0.5\left[0.5 + .75\right]$

$= 0.5 + 0.5 * 1.25$

$= 0.5 + 3.625 = 1.125$

Sum of Pagerank of all the pages = No of web pages

$= 1 + .75 + 1.125 = 2.875$

It has to repeated (iteratively) until we get the same (or almost same) values of page Rank.

Another notation of pagerank n° :-

$$PR(A) = \frac{(1-d)}{N} + d\left[\frac{PR(T1)}{C(T1)} + \frac{PR(T2)}{C(T2)} \cdots \frac{PR(Tn)}{C(Tn)}\right]$$
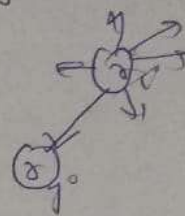
Then sum of all pages will be one.

Determining the pagerank with above equations becomes difficult if we have numerous equations in large web. [Secondly, it assumes that every node has non-zero outdegree].

## Power Iteration Method :-

Consider web as a graph (nodes as pages and edges are hyperlinks :-

Now the earlier equation # :-
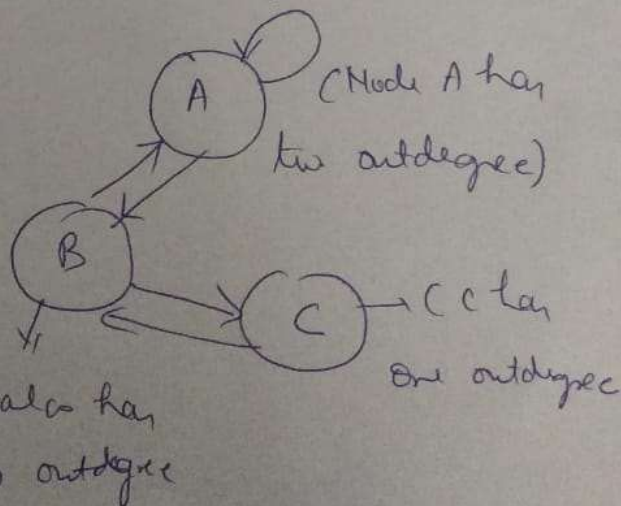
(Page Rank of Page $j$ at $(t+1)$th iteration)

$$r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i} \quad \text{(can be)}$$

written in the matrix form as

$$r^{(t+1)} = M \cdot r^{(t)}$$

Here, M is the matrix of the order $N \times N$ if N is No of nodes in the webgraph.

$$M = \begin{array}{c} \\ A \\ B \\ C \end{array} \begin{array}{ccc} A & B & C \\ \left[\begin{array}{ccc} 0\,½ & ½ & 0 \\ ½ & 0 & 1 \\ 0 & ½ & 0 \end{array}\right] \end{array}$$

(Node A has two outdegree)

B also has two outdegree

(C has one outdegree)

④ __Algo:__—  $\rightarrow$ or $r^0 = \left[\frac{1}{N}, \frac{1}{N}, \frac{1}{N}, \cdots \frac{1}{N}\right]^T$ ⟲(Teitspan)

(1) Initialize: $r^0 = \left(\frac{1}{N}\right)$ for all the pages

(2) Iterate: $r^{(t+1)} = M \cdot r^{(t)}$

(3) Stop: when $|r^{(t+1)} - r^{(t)}| \leq \epsilon$ or the

Repeated      number are same for last two iterations.

$$r^0 = \left[\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right]^T = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}$$

$$r^{t+1} = M \cdot r^{t}$$

$\rightarrow$ M is column-matrix
$\rightarrow$ it is stochastic i.e. sum of
all column matrix is 1
in each column

$$r^1 = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}$$
$$\quad\quad 3 \times 3 \quad\quad 3\times 1$$

$$= \begin{bmatrix} \frac{1}{2}\times\frac{1}{3} + \frac{1}{2}\times\frac{1}{3} \\ \frac{1}{2}\times\frac{1}{3} + \frac{1}{?} \\ 0 + \frac{1}{2}\times\frac{1}{3} + 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} \\ \frac{?}{6} \\ \frac{1}{6} \end{bmatrix}$$

$$r^2 = M \cdot r^1$$
$$= \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{3} \\ \frac{2}{6} \\ \frac{1}{6} \end{bmatrix} = \begin{bmatrix} \frac{1}{6} + \frac{1}{4} + 0 \\ \frac{1}{6} + 0 + \frac{1}{6} \\ \frac{1}{2}\times\frac{2}{6} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{5}{12} \\ \frac{1}{3} \\ \frac{1}{6} \end{bmatrix}$$

After some iteration it will converge.

(10)

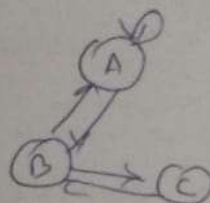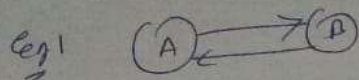In case of random surfer model, damping factor / teleporting factor $d$ or $\beta$ or $\alpha$ is given there.

matrix equation become

$$\boxed{r^{(t+1)} = (\alpha \cdot M)(r^{t})}$$

typically, $\alpha = 0.8 \mid 0.9$ [or in Between]

Example :-

$$M = \begin{array}{c} \\ A \\ B \\ C \end{array} \begin{array}{ccc} A & B & C \\ \left[\begin{array}{ccc} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{array}\right] \end{array}$$



$$r^{0} = \left[\begin{array}{ccc} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{array}\right]^{T}$$

$$r^{1} = 0.8 \left[\begin{array}{ccc} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{array}\right] \left[\begin{array}{c} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{array}\right]$$

$$= \left[\begin{array}{ccc} .4 & .4 & 0 \\ .4 & 0 & .8 \\ 0 & .4 & 0 \end{array}\right] \left[\begin{array}{c} 1/3 \\ 1/3 \\ 1/3 \end{array}\right] = \left[\begin{array}{c} .8/3 \\ 1.2/3 \\ .4/3 \end{array}\right]$$

Do it iteratively to get the result (or convergence condition is met.

⑪

## Problems :-

① "Spider Trap" Problem (~~State Problem~~)

eg1  (A) ⇄ (B)

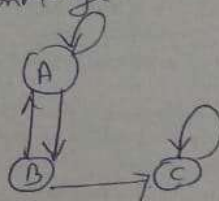→ Outgoing links are in the group. i.e. in above example A can not go to any other page and vice versa.

if we intialize the page rank with

$$r_a = 1 \quad , \; 0 \quad , \; 1 \quad , \; 0$$
$$r_b = 0 \quad , \; 1 \quad , \; 0 \quad , \; 1 \quad ---$$

It will never converge.

eg 2: C is a spider trap as it can not go anywhere else except c.

→ This problem can be solved with teleporting.

② "Dead end" Problem (Dangling Nodes / No outgoing links)

→ It happens when a node has no outgoing link. Then its importance can not be transferred.
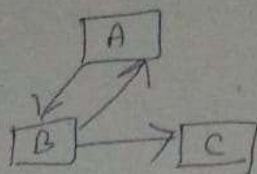
(A) ⟶ (B)

$$\left. \begin{array}{l} r_A = 1 \quad , \; 0 \; , \; 0 \; , \; 0 \\ r_B = 0 \quad , \; 1 \; , \; 0 \; , \; 0 \end{array} \right\} \text{score is zero after some iteration}$$

→ This problem can be solved by replacing zero entries with (N)

(2)

Eg

| | A | B | C | |
|---|---|---|---|---|
| A | 0 | 1/2 | 0 | = M |
| B | 1 | 0 | 0 | |
| C | 0 | 1/2 | 0 | |



So remove the dead end problem by replacing
3rd column with $\frac{1}{N}$ i.e. $\frac{1}{3}$.

$$
\begin{array}{c}
 & A & B & C \\
A & \begin{bmatrix} 0 & 1/2 & 0\ 1/3 \\ 1 & 0 & 1/3 \\ 0 & 1/2 & 1/3 \end{bmatrix} & & = M \\
B & & & \\
C & & &
\end{array}
$$

Google Matrix equation is :-

$$A = \beta \cdot M + (1-\beta)\left(\frac{1}{N}\right)_{N \times N}$$

then calculate rank iteratively by following

equation :-

$$r^{t+1} = A \cdot r^{t}$$

$\beta$ = damping / teleporting factor

M = Matrix (column stochastic matrix)

A = Google Matrix

Reference :- www.mmds.org. (Stanford University)

Mining of Massive datasets : Leskovec, Rajasaraman, & Ulman.

# Hubs and Authorities
### HITS (Hyperlink - Induced Topic Search)

→ It is developed by John Kleinberg in 1997.

→ Every web page is assigned two scores, one is called __Hub score__ and other is its authority __Score__

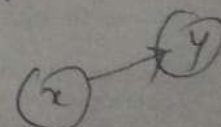→ For any query, we compute two ranked lists i.e hub score list & authority score list

→ A good hub page is one which points to many good authorities & a good authority page is one which is pointedtoby many good hub pages.
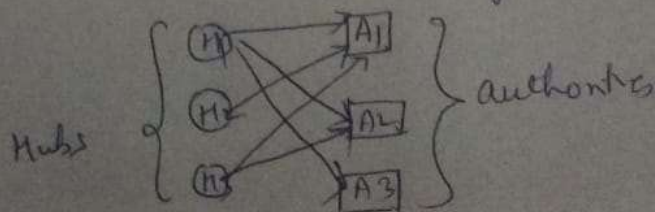
→ Let us take a subset of webgraph and then compute its hub & authority score.

for a webpage $x$, $x$ has a hyperlink that points to $y$.

Hub score of node $x = h(x) = \sum_{x \to y} a(y)$

authority score of node $y = a(x) = \sum_{y \to x} h(y)$

## Algorithm :-

1. Assign each node an authority and hub score of webgraph.

   i.e intially $h(x) = a(x) = 1$ or
   $$h(x) = a(x) = \frac{1}{\sqrt{M}}$$

2. Apply the authority update rule : each node's authority score is the sum of hub score of each node that it points to.
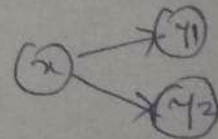   $$a(x) = \sum_{y \to x} h(y)$$

   

3. Apply the hub update rule : each node is hub score is sum of authority score of each node that it points to
   $$h(x) = \sum_{x \to y} a(y)$$

   

4. Normalize authority and hub score :
   $$auth(j^\circ) = \frac{auth(j^\circ)}{\sum_{i \in M} auth(i^\circ)}$$

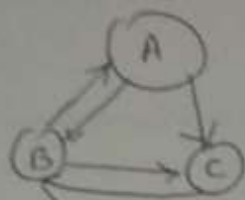   $$hub(j) = \frac{hub(j^\circ)}{\sum_{i \in M} hub(i^\circ)}$$

5. Repeat 2,3,4 until convergence condition is met. i.e. score become constant.
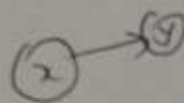
Example :-

Let $a_{i^0} = h_{i^0} = 1$



Step 2

| | old A | old H | | New A | New H |
|---|---|---|---|---|---|
| A | 1 | 1 | | 1 (B) | 2 (B & C) |
| B | 1 | 1 | | 2 (A & C) | 2 (A & C) |
| C | 1 | 1 | | 2 (A & B) | 1 (A) |

$$\text{authority } (x) = \sum_{y \to x} h(y)$$
(indegree)



$$hub (x) = \sum_{x \to y} a(y)$$



Normalize :- authority $(A) = \dfrac{A(A)}{A(A) + A(B) + A(C)} = \dfrac{1}{5}$

Similarly for other node

| new A | new H | | new A | new H |
|---|---|---|---|---|
| $\frac{1}{5}$ | $\frac{2}{5}$ | $\frac{H(B)}{\frac{2}{5}}$ | $\frac{4}{5}$ | $\frac{A(B)+A(C)}{\frac{4}{5}}$ |
| $\frac{2}{5}$ | $\frac{2}{5}$ | $\frac{H(A)+H(C)}{\frac{3}{5}}$ | $\frac{3}{5}$ | $\frac{A(A)+A(C)}{\frac{3}{5}}$ |
| $\frac{2}{5}$ | $\frac{1}{5}$ | $\frac{H(A)+H(B)}{\frac{4}{5}}$ | $\frac{2}{5}$ | $\frac{A(B)}{\frac{2}{5}}$ |

and now normalize and do it iteratively till stopping condition is met

Matrix Notation :-

Let vector $a = (a_1, a_2 - a_n)$ be a authority vector
$h = (h_1, h_2 - - h_n)$ be a hub vector

$n =$ no of nodes in the web graph

Step1 : initialize $a_i^0 = h_i^0 = \frac{1}{\sqrt{n}}$ or

$a_i^0 = h_i^0 = 1$

Step2 : Repeat until stopping condition is met (or it converges)

(i) $h = A \cdot a \implies h = (A A^T) \cdot h$

(ii) $a = A^T \cdot h \implies a = (A^T A) \cdot a$

(iii) Normalize $a$ & $h$ . (Use $L_1$ or $L_2$ norm)

Step3 :- If convergence criterion is met, stop else go to step2.

$L_2$ Norm is :-
$$\sum_{i \in N} (a_i^0)^2 = 1 \quad \& \quad \sum_{i \in N} (h_i^0)^2 = 1$$

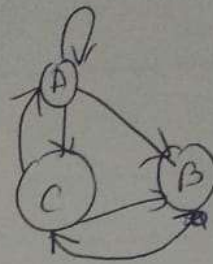Hence, at every iteration, normalize each value of $a$ & $h$ by

$$a_i^{(t)} = \frac{a_i^{0(t)}}{\sqrt{\sum_{j=1}^{n} a_j^0{}^2}} \quad \text{and} \quad h_i^{(t)} = \frac{h_i^{0(t)}}{\sqrt{\sum_{j=1}^{n} h_j^0{}^2}}$$

$$a_i^{0(t)} = \frac{a_i^0{}^{(t)}}{\sqrt{a_1^2 + a_2^2 + - a_n^2}} \quad \text{and} \quad h_i^{0(t)} = \frac{h_i^{0(t)}}{\sqrt{h_1^2 + h_2^2 - - h_n^2}}$$

Example :-

$$[A] = \begin{array}{c}A\\B\\C\end{array}\begin{bmatrix}\overset{A}{1} & \overset{B}{1} & \overset{C}{1}\\ 0 & 0 & 1\\ 1 & 1 & 0\end{bmatrix}$$

$$a = (A^T \times A \cdot a)$$
$$h = (A \times A^T \cdot h)$$

$$A^T A = \begin{bmatrix}2 & 2 & 1\\ 2 & 2 & 1\\ 1 & 1 & 2\end{bmatrix}$$

$$A A^T = \begin{bmatrix}3 & 1 & 2\\ 1 & 1 & 0\\ 2 & 0 & 2\end{bmatrix}$$

$$a = [A^T A] \cdot a$$
$$= \begin{bmatrix}2 & 2 & 1\\ 2 & 2 & 1\\ 1 & 1 & 2\end{bmatrix}\begin{bmatrix}1\\1\\1\end{bmatrix} = \begin{bmatrix}5\\5\\4\end{bmatrix}$$

$$h = [A A^T] h = \begin{bmatrix}3 & 1 & 2\\ 1 & 1 & 0\\ 1 & 1 & 2\end{bmatrix}\begin{bmatrix}1\\1\\1\end{bmatrix} = \begin{bmatrix}6\\2\\4\end{bmatrix}$$

next iteration

$$a = \begin{bmatrix}24\\24\\16\end{bmatrix}$$
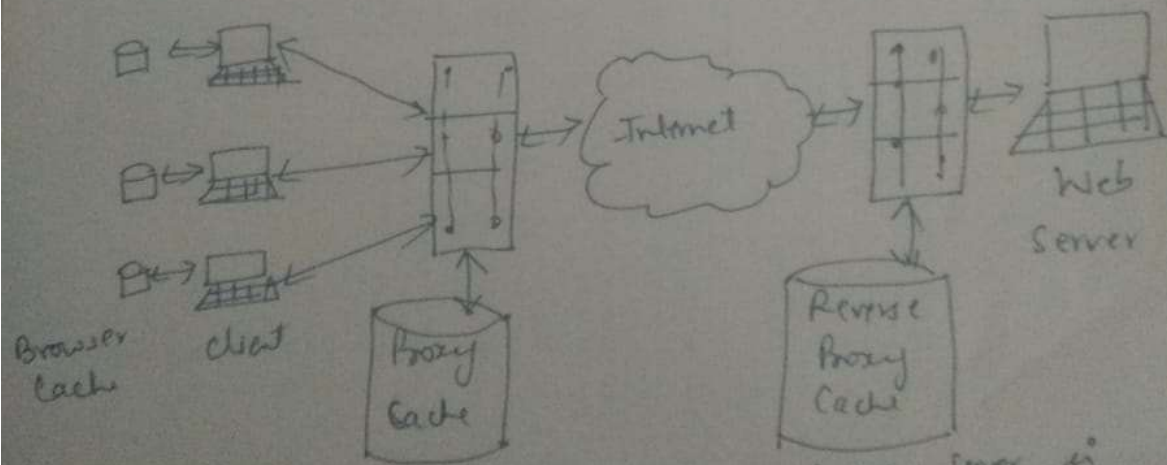
$$h = \begin{bmatrix}28\\8\\20\end{bmatrix}$$

and so on till it converges.

## World Wide Web Cacheing (www)

→ Web Caching (www) provides an efficient remedy to the latency problem.

→ It brings the documents closer to clients in less latency (time).

→ Caching can be deployed at various points in the internet Eg.
- o client browser
- o Proxy server /Reverse proxy server



Browser        client
cache

Proxy
Cache

Reverse
Proxy
Cache

Web
Server

<u>Proxy Cache</u> :- It is kept in proxy servers. A proxy server is a computer that is often placed near a gateway to the internet (as shown in figure) and provides shared cache to the clients.

→ Cache mechanism help to reduce network and server load by minimizing transmission of documents which are on demand.

→ However, every cache is naturally limited by its size. Whenever the cache is full and the proxy needs to ~~set~~ cache a new document, then cache replacement policies have to be incorporated.

Cache replacement policies tries to increase the hit ratio of cache. Here is some parameters defined :-

(i) **Cache hit** : when data is found in cache.

(ii) **Cache Miss** :- When data is not found in cache.

(iii) **Hit Ratio** :- % of times, data is found in
(Hit rate)       the cache.

(iv) **Miss ratio** :- % of times, data is not found
(Miss Rate)     in the cache

| Miss rate = 1- hit rate (hit ratio) |

(v) **Miss Penality** :- time to move data from server to cache.

→ Normally, caching's popular types are application caches and memory caches for their ability to speed up certain responses.

→ Web caching is the core design feature of HTTP protocol to minimize traffic while improving the responsiveness of the system.

What can be cached :-

→ Logos and brand images

→ Non-Rotating images (navigation icons)

→ Style sheet

→ Java Script files

→ Media files

→ Downloadable content

What not to be cached :-

→ Assets related to sensitive data

→ Content which is user-specific and frequently changed.

## Cache Replacement :-

→ It refers to the process that takes place when cache becomes full and old objects must be removed to make space for new ones.

→ Now, question is, which has to be removed?

→ To answer this question, cache researchers and developers have proposed different replacement algorithms discussed below :-

### (i) LRU (Least Recently Used)

→ Most popular replacement algorithm used by web caches.

→ It removes the objects that have not been accessed for the longest time.

→ LRU would consider time-since-reference as the only parameter.

→ In practice, web caches almost always use a variant called LRU-threshold. The "threshold" refers to object size. Objects, size larger than the threshold size are not cached.

→ Caching numerous small objects results in a higher hit ratio

(ii) FIFO (first in, first out)

→ It is simple to implement
→ Objects are purged in the same order they added
→ It is not popular because it gives lower hit ratio than LRU.

iv) LRV (Lowest Relative Value)

It expels the object that has "lowest utility value". In LRV, the utility of a document is calculated adaptively on the basis of data readily available to proxy server. It performs better than LRU. To compute the utility, it uses <u>cost, size and last reference time</u>. ~~to calculate~~. The cost of a document includes the time and processing overhead associated

with retrieving the document from original /
Proxy server.

(IV) Optimal Algorithm (OPT)

⇒ The document/page that will not be
used for the longest period of time in
the future is replaced. It involves the
knowledge of future requents to predict
which item in the cache will be needed
again.

⇒ It has lowest miss rate, but it
is difficult to implement.

(V) Randomized Algorithms :-

⇒ These algorithms expels a document drawn
randomly from the cache.

⇒ It means it uses the concept of
Probability (concept of randomness)

⇒ <u>Example</u>:-

<u>RANDOM</u> :-

→ It randomly selects M documents from the cache.

→ Select the least useful document among them for replacement based on the given policy (existing algo- like LRU, FIFO etc.)

→ The next M least useful documents out of M are kept for the next iteration.

→ Iterations are continued until there is enough space in the cache to accomdate the remaining documents.

# Web Crawler

Web search engine ( web- IR) can be visualized in following figure [1]. In Web, getting the content of the documents takes longer time because of latency, Web search engines must crawl their documents. A web crawler does following operations:

- Initialize queue with URLs of known seed pages
- Repeat
    - Take URL from queue
    - Fetch and parse page
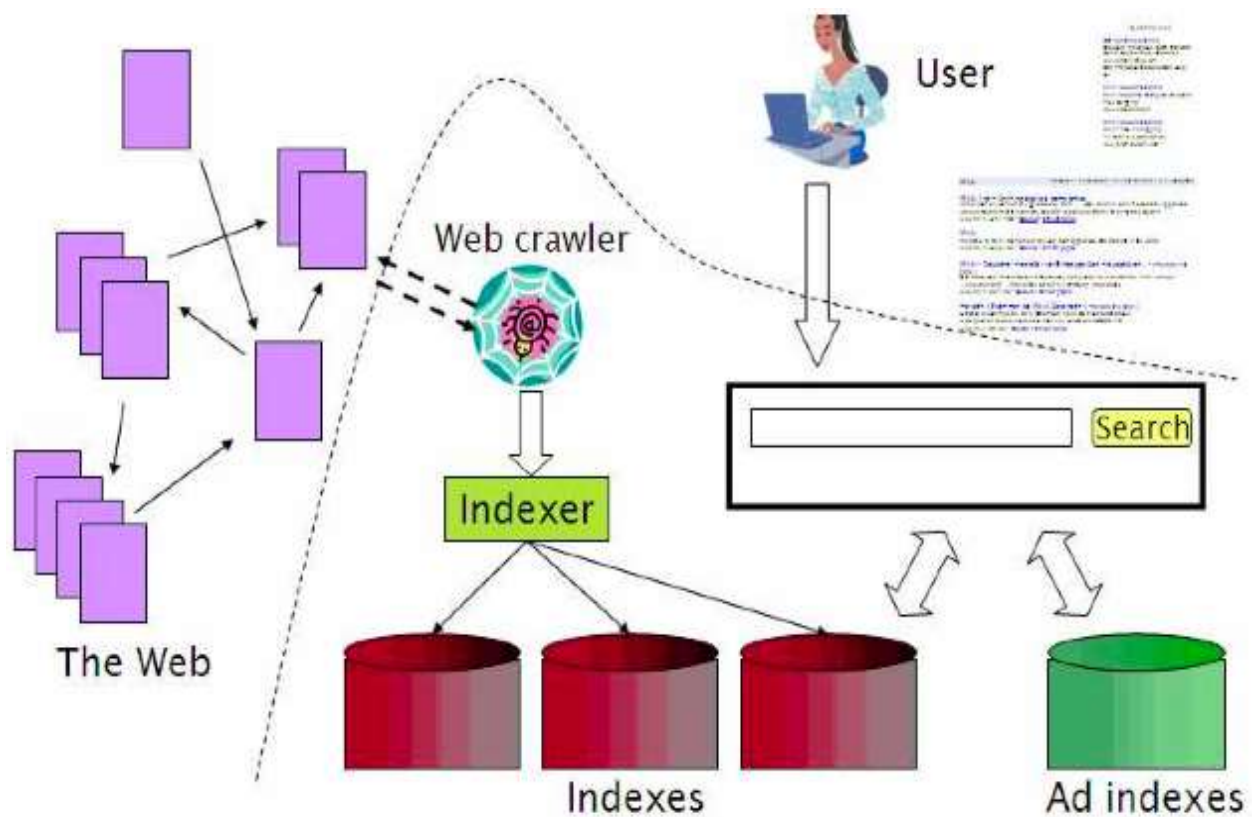    - Extract URLs from page
    - Add URLs to queue

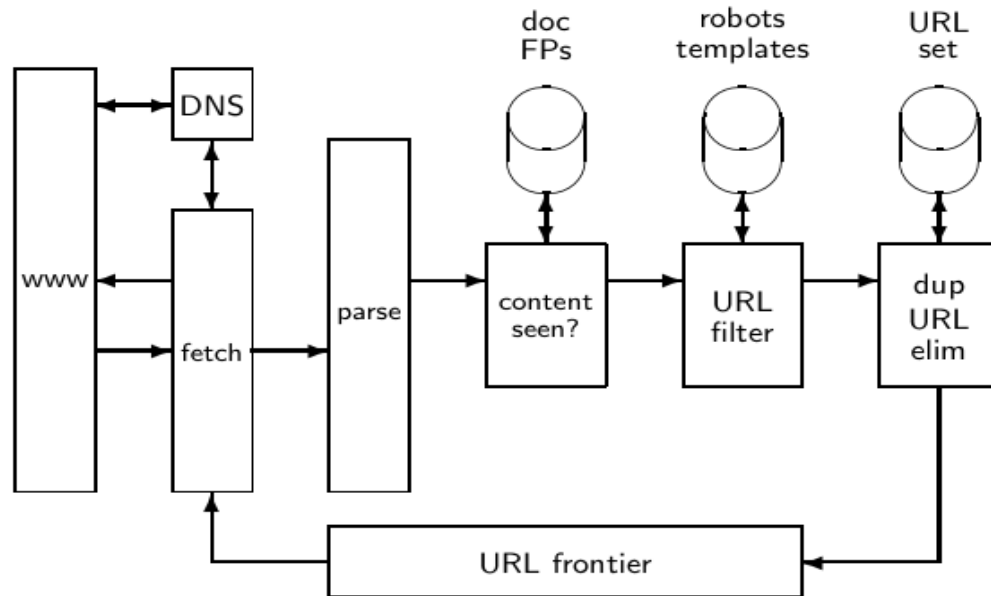Figure 1: The various components of a web search engine

**Web Crawler Architecture**



Figure 2: Web Crawler Architecture

Processing steps of crawler:

1. Pick a URL from the frontier
2. Fetch the document at the URL
3. Parse the URL
   a. Extract links from it to other docs (URLs)
4. Check if URL has content already seen
   a. If not, add to indexes
5. For each extracted URL
   a. Ensure it passes certain URL filter tests
   b. Check if it is already in the frontier (duplicate URL elimination)

## References:

[1] Manning, Christopher D., Prabhakar Raghavan, and Hinrich Schütze. Introduction to information retrieval. Cambridge university press, 2008.

[2] Leskovec, Jure, Anand Rajaraman, and Jeffrey David Ullman. Mining of massive data sets. Cambridge university press, 2020.

[3] https://www.mmds.org

[4] Brin, S.; Page, L. *(1998).* "The anatomy of a large-scale hypertextual Web search engine" (PDF). *Computer Networks and ISDN Systems. **30** (1–7): 107–117.*
[5] Wessels, Duane. Web caching. " O'Reilly Media, Inc.", 2001.

[6] Cao, Pei, and Sandy Irani. "Cost-aware www proxy caching algorithms." Usenix symposium on internet technologies and systems. Vol. 12. No. 97. 1997.

[7] Balamash, Abdullah, and Marwan Krunz. "An overview of web caching replacement algorithms." IEEE Communications Surveys & Tutorials 6.2 (2004): 44-56.