

## DWM Assignment 4

1. (i) Page rank algorithm is used by google to rank web pages in their search engine results.  
 (ii) PR technique is used by google to prioritize pages returned from search by looking at web structure. Importance of page is calculated based on no. of pages which point to it - Back links. Weighting is used to provide more importance

(iii) Page Rank is defined as  $PR(P)$

$$PR(p) = c \left( \frac{PR(i)}{N_i} + \dots + \frac{PR(n)}{N_n} \right)$$

$PR(i)$  : Page Rank for page  $i$  which points to target

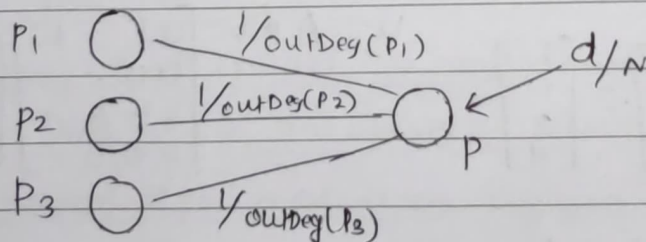
$N_i$  : no of links coming out of that page  $P$

$c$  : constant (0-1) used for normalization.

eg. Google's Page Rank

Rank of a web page depends on rank of web page pointing to it

$$PR(P) = \frac{d}{N} + (1-d) \left[ \frac{PR(P_1)}{\text{outDeg}(P_1)} + \frac{PR(P_2)}{\text{outDeg}(P_2)} + \frac{PR(P_3)}{\text{outDeg}(P_3)} \right]$$



Set  $PR \rightarrow [r_1, r_2 \dots r_n]$  these are same initial rank of  $p_1$  to  $N$ .  $d = 0.15$ ,  $dD \Rightarrow \left[ \frac{1}{N} \dots \frac{1}{N} \right]^T$

$A$  is adjacency matrix

$$do \quad PR_{i+1} \leftarrow A^T * PR_i$$

$$PR_{i+1} \leftarrow (1-d) * PR_{i-1} + d * D$$

$b \leftarrow ||PR_{i+1} - PR_i||$  , while  $b < \epsilon$  where  $\epsilon$  is small no.  
indicating convergence threshold  
return PR

2. (i) HITS (Hyperlink induced topic search) algo is a kind of rank of algo that analyzes web resources based on local link Algorithm is as follows

① Use text keyword matching to acquire root set R, which includes thousands of URL or more.

② Define S to R i.e.  $S \subseteq R$  are equal

③ To each page P in R, put hyperlinks included by R into set S put pages referring to P into set S

④ S is acquired expanded neighbourhood set. The inputs for HITS algo are a set of nodes with an adjacency matrix and a value k (no of iteration). First we need to define 2 variables

u: hubs weighted vector and v: Authority weighted vector where A: adjacency matrix of all pages

$$u = A * v \quad \text{and} \quad v = A^T * u$$

To calculate v we assume  $u=1$  initially and calculate v & then using this v we calculate u. Based on u, v we get the ranking for hub and authentication

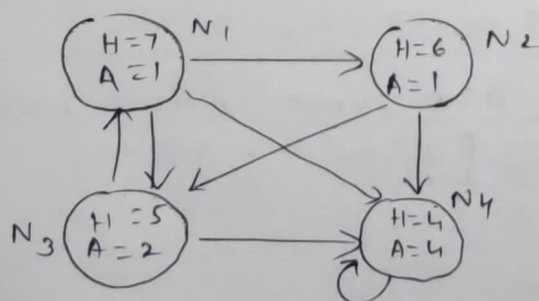
For eg A is given as

Initial weight  
 $u = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$

$$\begin{matrix} & N_1 & N_2 & N_3 & N_4 \\ \begin{matrix} N_1 \\ N_2 \\ N_3 \\ N_4 \end{matrix} & \begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \end{matrix}$$
 with nodes  $N_1, N_2, N_3, N_4, k=2$

$$v = A^T \cdot u = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 2 \\ 4 \end{bmatrix}$$

$$u = A * v = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 2 \\ 4 \end{bmatrix} = \begin{bmatrix} 7 \\ 6 \\ 5 \\ 4 \end{bmatrix}$$



For  $k=1$

| Nodes | Hub | Authority |
|-------|-----|-----------|
| N1    | 7   | 1         |
| N2    | 6   | 1         |
| N3    | 5   | 2         |
| N4    | 4   | 4         |

Hub:  $N_1, N_2, N_3, N_4$

Authority:  $N_4, N_3, \{N_2, N_1 \text{ TIE}\}$



Now for  $k=2$ , we have to use unit matrix of  $u, v$

$$|v| = \sqrt{1^2 + 1^2 + 2^2 + 4^2} = \sqrt{22}$$

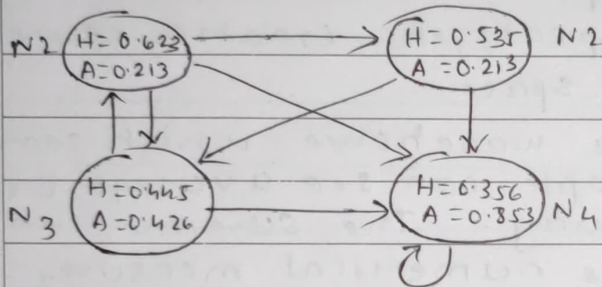
$$v' = \frac{1}{\sqrt{22}}, \frac{1}{\sqrt{22}}, \frac{2}{\sqrt{22}}, \frac{4}{\sqrt{22}}$$

$$v' = \begin{bmatrix} 0.213 \\ 0.213 \\ 0.426 \\ 0.853 \end{bmatrix}$$

similarly  $|v| = \sqrt{7^2 + 6^2 + 5^2 + 4^2} = \sqrt{126}$

$$u' = \frac{7}{\sqrt{126}}, \frac{6}{\sqrt{126}}, \frac{5}{\sqrt{126}}, \frac{4}{\sqrt{126}}$$

$$u' = \begin{bmatrix} 0.623 \\ 0.535 \\ 0.445 \\ 0.356 \end{bmatrix}$$



For  $k=2$

|    | Hub   | Authority |
|----|-------|-----------|
| N1 | 0.623 | 0.213     |
| N2 | 0.535 | 0.213     |
| N3 | 0.445 | 0.426     |
| N4 | 0.356 | 0.853     |

Ranking in same as  $k=1$

3. (i) Multimedia database consists of audio, video, images and text media. They can be stored on object-oriented databases. They are used to store complex info in a pre specified format.
- (ii) They store geographical info in spatial database. Store data in form of co-ordinates, topology, lines, polygon, etc.
- (iii) Time series databases: Time series databases contain stock exchange data and user logged activities. Handles array of nos indexed by time, date, etc.
- (iv) WWW: It is a collection of documents and resources like audio, video, text, etc which are identified by URL through web browsers, linked by HTML pages and accessible via Internet network.

### \* Issues in web mining :

Size of web is very huge and rapidly increasing. Complexity of web pages ; web pages don't follow order of traditional text documents. Web is a dynamic information source. Data such as news, stock market, sports, etc. are continuously changing. Diversity of use communities ; community is rapidly expanding and have diff. backgrounds & interest. Relevance of info : It is considered that a particular person is interested only in a small portion of the web, while the rest is not relevant to the user.

### \* Issues in spatial mining :

Accurately representing true shape and size. Representing non-continuous data eg. road lines, peaks, etc. Creating aesthetically pleasing maps. Conserving disk space.

4.] We can develop spatial data warehouse which stores highway traffic info so that people can see average or peak time traffic flow by the highway. The schema can be a star with flow and count as numerical measures. The region map measure is spatial and represents collection of spatial pointers to corresponding regions. So based on this data warehouses, as a city planner of Mumbai, different kinds of information can be mined. For eg. he can notice where max traffic is generated & broaden the roads or divert the vehicles. We can use mining to learn live traffic updates and no. of accidents at a location. Such info is useful for planning infrastructure.

5.] It is unrealistic to treat video, as a long sequence of individual stills and analyze each picture. This would result in too many pictures. In order to capture outliers from video data, it is better to treat each video clip first as a collection of actions and events in time, segment them into video shots. A shot is a group of frames where video content from one frame to adjacent ones does not change abruptly. Each frame can be analyzed using image feature extraction and other methods. A process can then sample the data quickly and analyze images in real time. The process compares features of key frame with current static images to see if there is a major difference. If so then alert an alarm that there is some unusual event occurring.