

CS47300: Web Information Search and Management

Using Graph Structure for Retrieval

Prof. Chris Clifton

24 September 2019

Material adapted from slides created by Dr. Rong Jin (formerly Michigan State, now at Alibaba)



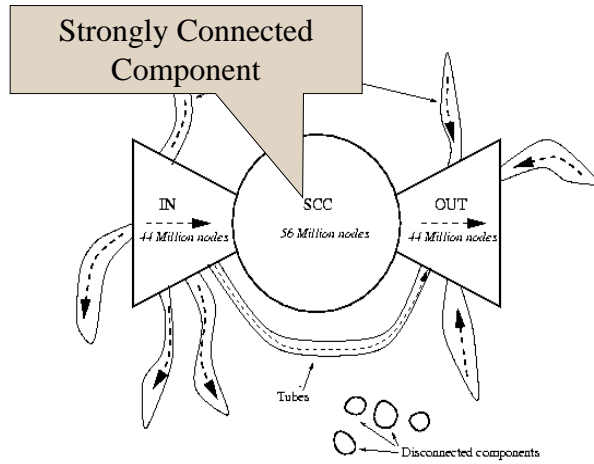
Ad-Hoc Retrieval: Beyond the Words

- Web is a graph
 - Each web site correspond to a node
 - A link from one site to another site forms a directed edge



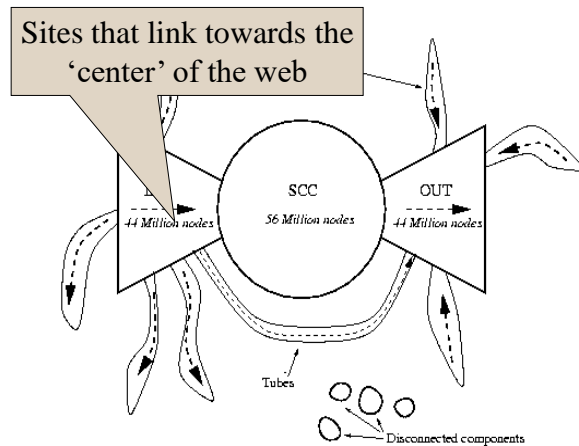
- Web is a graph
 - Each web site correspond to a node
 - A link from one site to another site forms a directed edge
- What does it look like?
 - Web is small world
 - The diameter of the web is 19
 - e.g. the average number of clicks from one web site to another is 19

Bowtie Structure



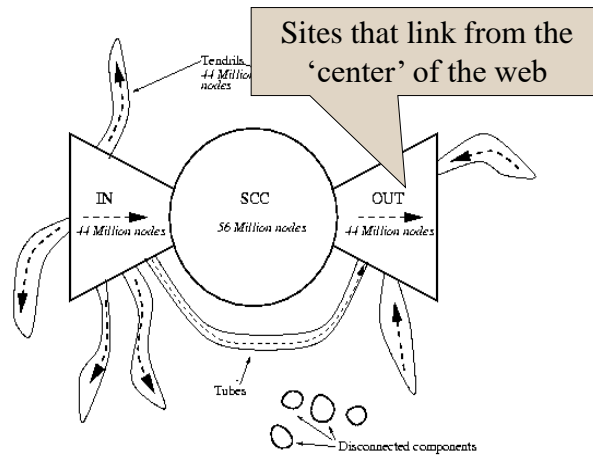
Broder et al., 2001

Bowtie Structure



Broder et al., 2001

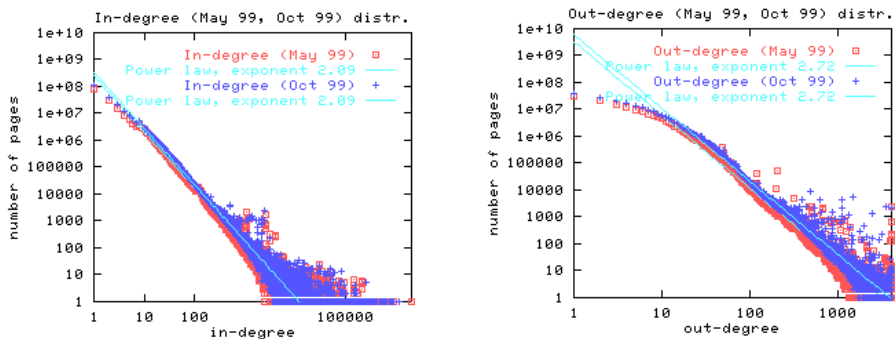
Bowtie Structure



Broder et al., 2001

Inlinks and Outlinks

- Both degrees of incoming and outgoing links follow power law



Broder et al., 2001

So what good is Link Structure?

- When you can't find something, do you:
 - Keep looking in the same place?
 - Look somewhere else?
 - Give up?
 - Ask for help?
- *Other people may already know the answer!*
 - Links: Reflect human judgement

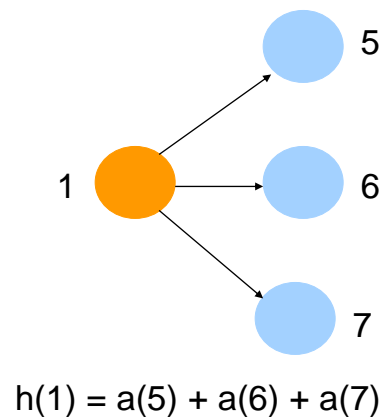
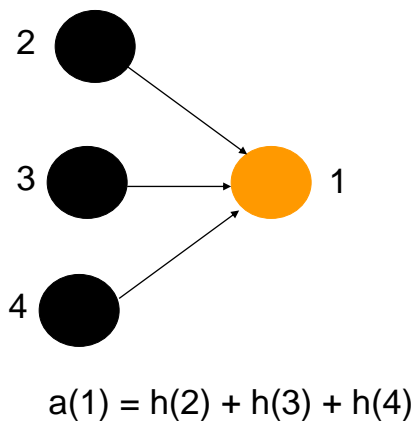
28

Early Approaches

- Basic Assumptions
 - Hyperlinks contain information about the human judgment of a site
 - The more incoming links to a site, the more it is judged important
- Bray 1996
 - The visibility of a site is measured by the number of other sites pointing to it
 - The luminosity of a site is measured by the number of other sites to which it points
 - Limitation: failure to capture the relative importance of different parents (children) sites

30

- HITS – Hypertext Induced Topic Selection
- For each vertex $v \in V$ in a subgraph of interest:
 - $a(v)$ - the authority of v
 - $h(v)$ - the hubness of v
- A site is very authoritative if it receives many citations.
 - Citation from important sites weight more than citations from less-important sites
- Hubness shows the importance of a site.
 - A good hub is a site that links to many authoritative sites



Authority and Hubness: Version 1

Recursive dependency

$$a(v) = \sum_{w \in pa[v]} h(w)$$

$$h(v) = \sum_{w \in ch[v]} a(w)$$

HubsAuthorities(G)

```

1  1 ← [1,...,1] ∈ R|V|
2  a0 ← h0 ← 1
3  t ← 1
4  repeat
5      for each v in V
6          do at(v) ← ∑w ∈ pa[v] ht-1(w)
7              ht(v) ← ∑w ∈ pa[v] at-1(w)
8  t ← t + 1
9  until || at - at-1 || + || ht - ht-1 || < ε
10 return (at, ht)

```

Problems ?

34

Authority and Hubness: Version 2

Recursive dependency

$$a(v) = \sum_{w \in pa[v]} h(w)$$

$$h(v) = \sum_{w \in ch[v]} a(w)$$

+ Normalization

$$a(v) = \frac{a(v)}{\sum_w a(w)}$$

$$h(v) = \frac{h(v)}{\sum_w h(w)}$$

HubsAuthorities(G)

```

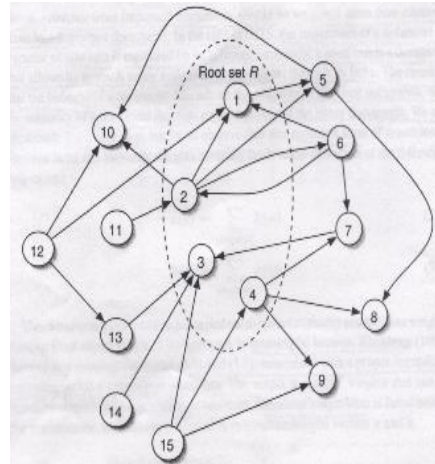
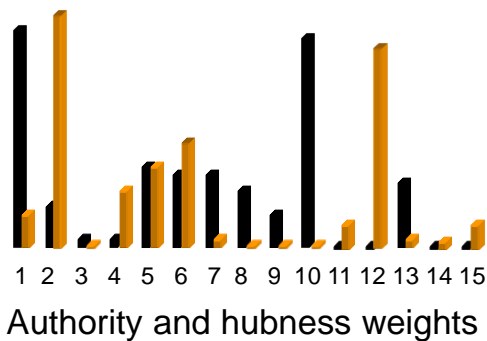
1  1 ← [1,...,1] ∈ R|V|
2  a0 ← h0 ← 1
3  t ← 1
4  repeat
5      for each v in V
6          do at(v) ← ∑w ∈ pa[v] ht-1(w)
7              ht(v) ← ∑w ∈ pa[v] at-1(w)
8          at ← at / || a ||
9          ht ← ht / || h ||
10 t ← t + 1
11 until || at - at-1 || + || ht - ht-1 || < ε
12 return (at, ht)

```

35

HITS Example Results

Authority
 Hubness



36

Authority and Hubness

- Authority score
 - Not only depends on the number of incoming links
 - But also the ‘quality’ (e.g., hubness) of the incoming links
- Hubness score
 - Not only depends on the number of outgoing links
 - But also the ‘quality’ (e.g., hubness) of the outgoing links

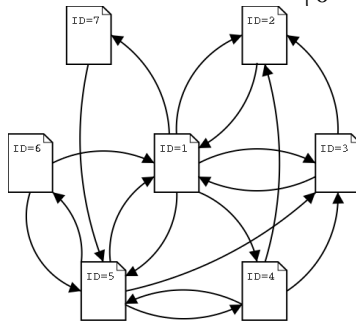
37

Authority and Hub

- Column vector \mathbf{a} : a_i is the authority score for the i -th site
- Column vector \mathbf{h} : h_i is the hub score for the i -th site

- Matrix \mathbf{M} :

$$\mathbf{M}_{i,j} = \begin{cases} 1 & \text{the } i\text{th site points to the } j\text{th site} \\ 0 & \text{otherwise} \end{cases}$$



$$\mathbf{M} = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

Authority and Hub

- Vector \mathbf{a} : a_i is the authority score for the i -th site
- Vector \mathbf{h} : h_i is the hub score for the i -th site

- Matrix \mathbf{M} :

$$\mathbf{M}_{i,j} = \begin{cases} 1 & \text{the } i\text{th site points to the } j\text{th site} \\ 0 & \text{otherwise} \end{cases}$$

- Recursive dependency:

$$a(v) \leftarrow \sum_{w \in \text{pa}[v]} h(w)$$

$$h(v) \leftarrow \sum_{w \in \text{ch}[v]} a(w)$$

Authority and Hub

- Column vector \mathbf{a} : a_i is the authority score for the i -th site
- Column vector \mathbf{h} : h_i is the hub score for the i -th site

- Matrix \mathbf{M} :

$$\mathbf{M}_{i,j} = \begin{cases} 1 & \text{the } i\text{th site points to the } j\text{th site} \\ 0 & \text{otherwise} \end{cases}$$

- Recursive dependency:

$$\begin{aligned} a(v) &\leftarrow \sum_{w \in \text{pa}[v]} h(w) \\ h(v) &\leftarrow \sum_{w \in \text{ch}[v]} a(w) \end{aligned} \quad \Rightarrow \quad \begin{aligned} \mathbf{a} &= \mathbf{M}^T \mathbf{h} \\ \mathbf{h} &= \mathbf{M} \mathbf{a} \end{aligned}$$

40

Authority and Hub

- Column vector \mathbf{a} : a_i is the authority score for the i -th site
- Column vector \mathbf{h} : h_i is the hub score for the i -th site
- Matrix \mathbf{M} :

$$\mathbf{M}_{i,j} = \begin{cases} 1 & \text{the } i\text{th site points to the } j\text{th site} \\ 0 & \text{otherwise} \end{cases}$$

- Recursive dependency:

$$\begin{aligned} a(v) &\leftarrow \sum_{w \in \text{pa}[v]} h(w) \\ h(v) &\leftarrow \sum_{w \in \text{ch}[v]} a(w) \end{aligned} \quad \Rightarrow \quad \begin{aligned} \mathbf{a}_t &= \alpha_t \mathbf{M}^T \mathbf{h}_t \\ \mathbf{h}_t &= \beta_t \mathbf{M} \mathbf{a}_t \end{aligned}$$

Normalization
Procedure

41

Authority and Hub

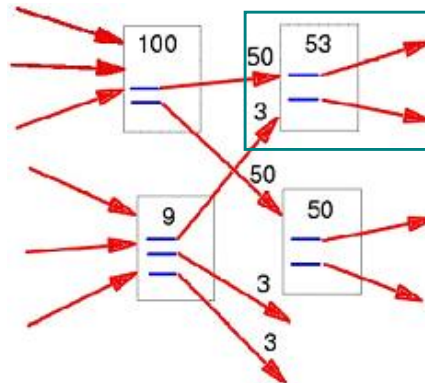
$$\left. \begin{aligned} \mathbf{a}_t &= \alpha_t \mathbf{M}^T \mathbf{h}_t \\ \mathbf{h}_t &= \beta_t \mathbf{M} \mathbf{a}_t \end{aligned} \right\} \rightarrow \begin{aligned} \mathbf{a}_t &= \alpha_t \beta_t \mathbf{M}^T \mathbf{M} \mathbf{a}_t \\ \mathbf{h}_t &= \alpha_t \beta_t \mathbf{M} \mathbf{M}^T \mathbf{h}_t \end{aligned}$$

- Apply SVD to matrix \mathbf{M}

$$\mathbf{M} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \sum_i \lambda_i \mathbf{u}_i \mathbf{v}_i^T \longrightarrow \mathbf{a} = \mathbf{u}_1, \mathbf{h} = \mathbf{v}_1$$

PageRank

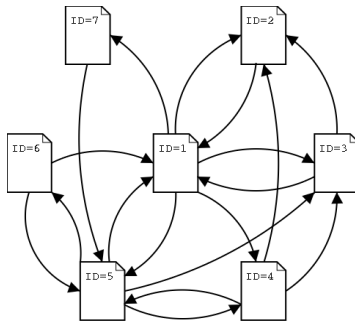
- Introduced by Page et al. (1998)
 - The weight is assigned by the rank of parents
- Difference from HITS
 - HITS separates Hubness & Authority weights
 - Page rank is proportional to its parents' rank, but inversely proportional to its parents' outdegree



Matrix Notation

$$M_{i,j} = \begin{cases} 1 & \text{the } i\text{th site points to the } j\text{th site} \\ 0 & \text{otherwise} \end{cases}$$

$$B_{i,j} = \begin{cases} \frac{1}{\sum_j M_{i,j}} & \sum_j M_{i,j} > 0 \\ 0 & \text{otherwise} \end{cases}$$



$$M = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

$$B = \begin{pmatrix} 0 & 1/5 & 1/5 & 1/5 & 1/5 & 0 & 1/5 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/3 & 1/3 & 0 & 1/3 & 0 & 0 \\ 1/4 & 0 & 1/4 & 1/4 & 0 & 1/4 & 0 \\ 1/2 & 0 & 0 & 0 & 1/2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

46

Matrix Notation

\mathbf{r} : r_i represents the rank score for the i -th web page

$$r(v) = \alpha \sum_{w \in \text{pa}[v]} \frac{r(w)}{|\text{ch}[w]|}$$

$$\mathbf{r} = \alpha \mathbf{B}^T \mathbf{r}$$

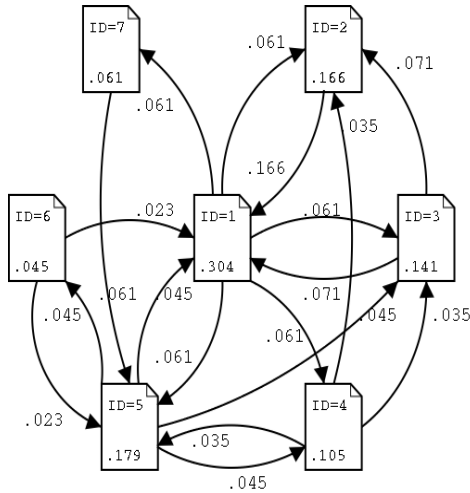
α : eigenvalue
 \mathbf{r} : eigenvector of \mathbf{B}

Finding Pagerank

→ find principle eigenvector of \mathbf{B}

47

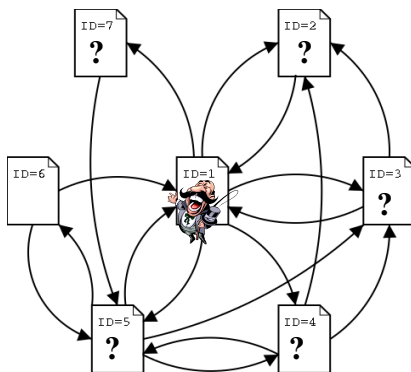
Matrix Notation



PR	ID	OutLink	InLink
0.304	1	2,3,4,5,7	2,3,5,6
0.179	5	1,3,4,6	1,4,6,7
0.166	2	1	1,3,4
0.141	3	1,2	1,4,5
0.105	4	2,3,5	1,5
0.061	7	5	1
0.045	6	1,5	5

Random Walk Model

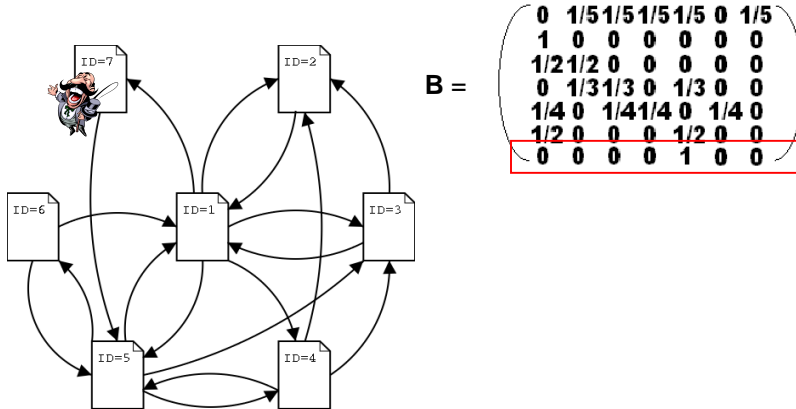
- Consider a random walk through the Web graph



$$B = \begin{pmatrix} 0 & 1/5 & 1/5 & 1/5 & 1/5 & 0 & 1/5 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/3 & 1/3 & 0 & 1/3 & 0 & 0 \\ 1/4 & 0 & 1/4 & 1/4 & 0 & 1/4 & 0 \\ 1/2 & 0 & 0 & 0 & 1/2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

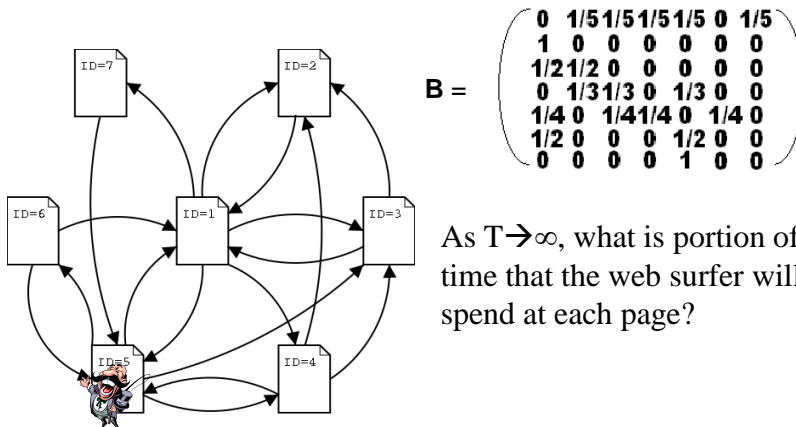
Random Walk Model

- Consider a random walk through the Web graph



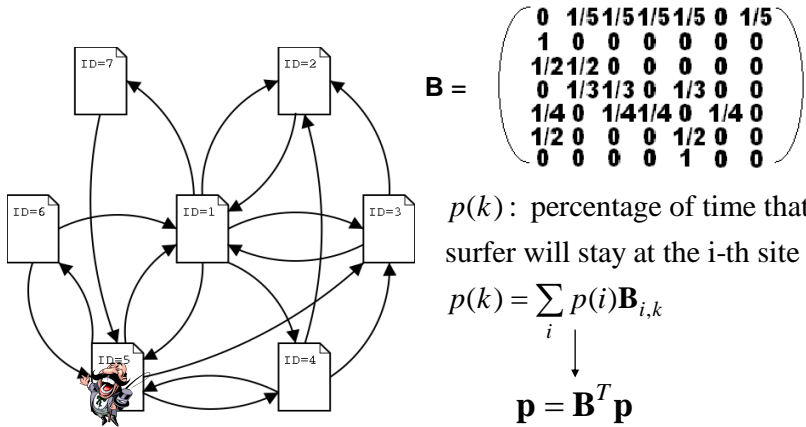
Random Walk Model

- Consider a random walk through the Web graph



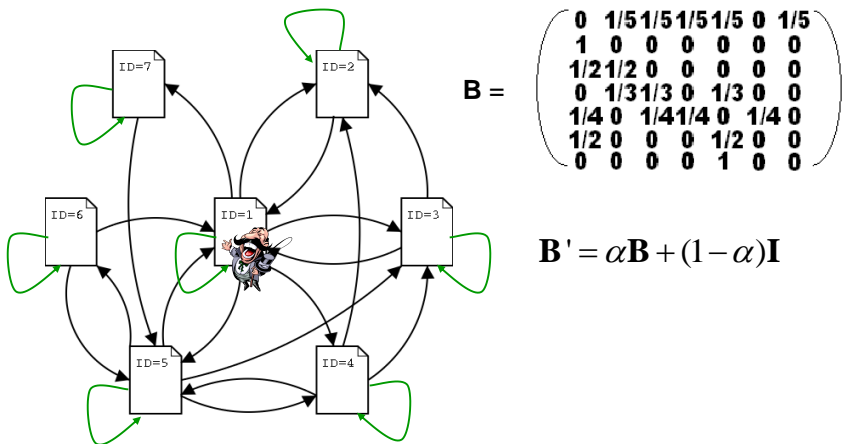
Random Walk Model

- Consider a random walk through the Web graph



Adding Self Loop

- Allow surfer to decide to stay on the same place

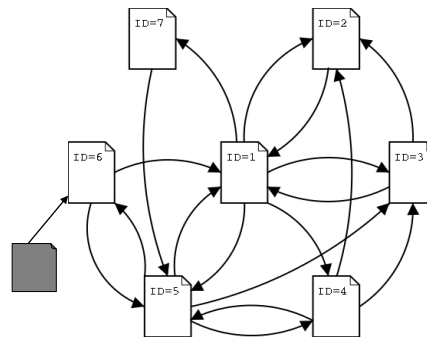


Problem

- “Rank Sink” Problem
 - Many Web pages have no inlinks
 - Results in dangling edges in the graph

$$B = \begin{pmatrix} 0 & 1/5 & 1/5 & 1/5 & 1/5 & 0 & 1/5 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/3 & 1/3 & 0 & 1/3 & 0 & 0 & 0 \\ 1/4 & 0 & 1/4 & 1/4 & 0 & 1/4 & 0 & 0 \\ 1/2 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

$$r(\text{new page}) = 0$$



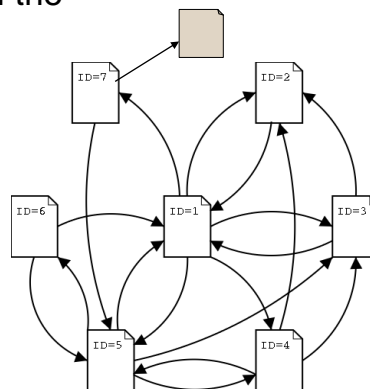
54

Problem

- “Rank Sink” Problem
 - Many Web pages have no outlinks
 - Results in dangling edges in the graph

$$B = \begin{pmatrix} 0 & 1/5 & 1/5 & 1/5 & 1/5 & 0 & 1/5 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/3 & 1/3 & 0 & 1/3 & 0 & 0 & 0 \\ 1/4 & 0 & 1/4 & 1/4 & 0 & 1/4 & 0 & 0 \\ 1/2 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$r(\text{new page}) = 1$$



55

Distribution of the Mixture Model

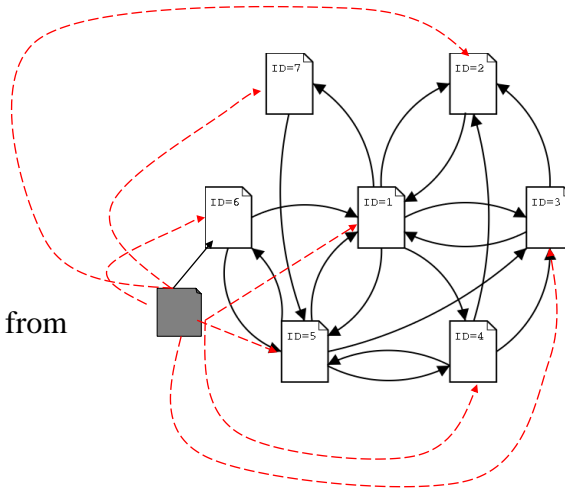
$$\mathbf{H}_{i,j} = 1/n$$

$$\mathbf{B}' = \varepsilon\mathbf{H} + (1 - \varepsilon)\mathbf{B}$$



$$\mathbf{r} = \mathbf{B}'^T \mathbf{r}$$

Prevents the page ranks from being 0 or 1



56

Stability

- Are link analysis algorithms based on eigenvectors stable?
 - Will small changes in graph result in major changes in outcomes?
- What if the connectivity of a portion of the graph is changed arbitrarily?
 - How will this affect the results of algorithms?

57

Stability of HITS

Ng et al (2001)

- A bound on the number of hyperlinks k that can be added or deleted from one page without affecting the authority or hubness weights
- It is possible to perturb a symmetric matrix by a quantity that grows as δ that produces a constant perturbation of the dominant eigenvector

$$k \leq \left(\sqrt{d + \frac{\alpha\delta}{4 + \sqrt{2}\alpha}} - \sqrt{d} \right)^2$$

$$\|\mathbf{a} - \tilde{\mathbf{a}}\|_2 \leq \alpha$$

δ : eigengap $\lambda_1 - \lambda_2$

d : maximum outdegree of G

58

Stability of PageRank

$$\|\tilde{r} - r\| \leq \frac{2 \sum_{j \in V} r(j)}{\epsilon} \quad \text{Ng et al (2001)}$$

V : the set of vertices touched by the perturbation

- The parameter ϵ of the mixture model has a stabilization role
- If the set of pages affected by the perturbation have a small rank, the overall change will also be small

59