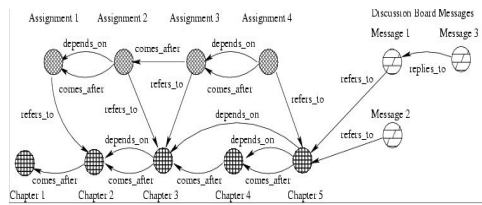


Web

Maria Luisa Sapino (BDM 2018)

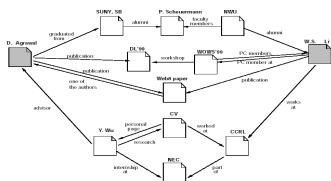
Web



Maria Luisa Sapino (BDM 2018)

Web

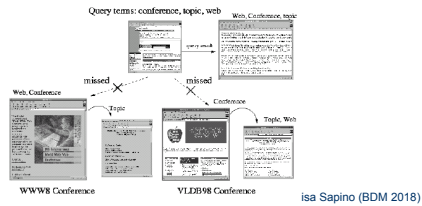
- Approach 1: use standard IR techniques to find pages that satisfy a query



Maria Luisa Sapino (BDM 2018)

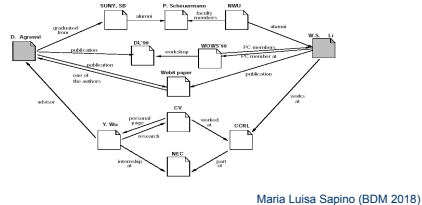
Web

- Approach 1: use standard IR techniques to find pages that satisfy a query



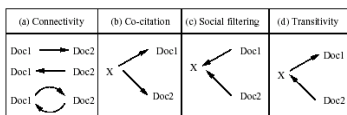
Web

- Approach 2: integrate IR techniques with structure/link analysis



Web

- Approach 2: integrate IR techniques with structure/link analysis



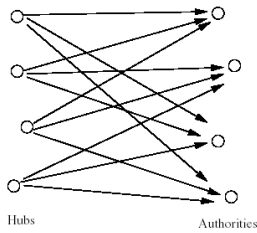
Maria Luisa Sapino (BDM 2018)

HITS algorithm

- Good pages are categorized into two types
 - Hubs: point to many pages of high quality
 - Authorities: pages of high quality

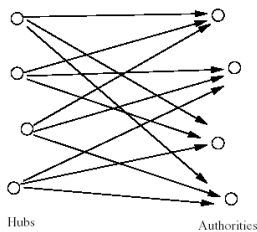
Maria Luisa Sapino (BDM 2018)

Hubs and authorities



Maria Luisa Sapino (BDM 2018)

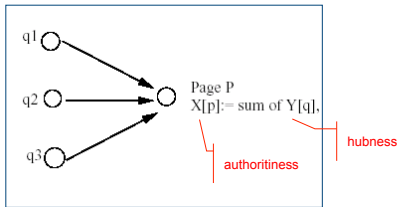
Hubs and authorities



- Good hubs should point to good authorities
- Good authorities must be pointed by good hubs.

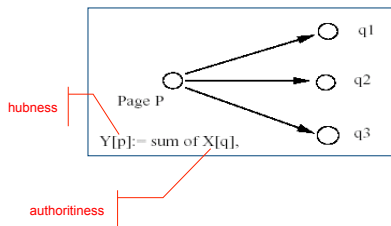
Maria Luisa Sapino (BDM 2018)

Topic distillation by iterative mutual reinforcement



Maria Luisa Sapino (BDM 2018)

Topic distillation by iterative mutual reinforcement



Maria Luisa Sapino (BDM 2018)

HITS

- Use IR to find the candidate pages

Maria Luisa Sapino (BDM 2018)

HITS

- Use IR to find the candidate pages
- Expand to include all pages which link or are linked by this core set

Maria Luisa Sapino (BDM 2018)

HITS

- Use IR to find the candidate pages
- Expand to include all pages which link or are linked by this core set
- Compute authority and hub values for all pages (iterate!!)

$$a(i) = \sum_{j \in \text{In}(i)} h(j) \quad h(i) = \sum_{j \in \text{Out}(i)} a(j)$$

Maria Luisa Sapino (BDM 2018)

HITS

- Matrix notation

$$\vec{a} = E^T \vec{h} \quad \vec{h} = E \vec{a}$$

Maria Luisa Sapino (BDM 2018)

...reminder

- Eigenvalue and eigenvector
- Given a matrix E , let c (scalar) and x (vector) be such that

$$c \vec{x} = E \vec{x}$$

Eigenvalue Eigenvector

Maria Luisa Sapino (BDM 2018)

...authorities

$$\vec{a} = E^T \vec{h}$$

Maria Luisa Sapino (BDM 2018)

...authorities

$$\vec{a} = E^T E \vec{a}$$

\vec{a} is an eigenvector of $E^T E$

Maria Luisa Sapino (BDM 2018)

...hubs

$$\vec{h} = EE^T \vec{h}$$

h is an
eigenvector of
 EE^T

Maria Luisa Sapino (BDM 2018)

HITS and LSI???

- ...reminder: SVD of E is

$$E = U M V^T$$

where

- M is a diagonal matrix
- $U^T U = I$
- $V^T V = I$

Maria Luisa Sapino (BDM 2018)

..then

$$EE^T = U M V^T V M U^T$$

Maria Luisa Sapino (BDM 2018)

..then

$$EE^T = U M V^T V M U^T$$

$$EE^T = U M M U^T$$

$$EE^T = U M^2 U^T$$

Maria Luisa Sapino (BDM 2018)

...then again

$$EE^T = U M^2 U^T$$

$$EE^T U = U M^2 U^T U = U M^2$$

Maria Luisa Sapino (BDM 2018)

...then again

$$EE^T = U M^2 U^T$$

$$EE^T U = U M^2 U^T U = U M^2$$

$$EE^T U_j = U_j m_j^2$$

Eigenvector Eigenvalue

Maria Luisa Sapino (BDM 2018)

...in fact

$$EE^T U_j = U_j m_j^2$$

Eigenvector | Eigenvector Eigenvector | Eigenvector

$$EE^T \vec{h} = \vec{h}$$

h = U_j |

Maria Luisa Sapino (BDM 2018)

...in fact

$$EE^T U_j = U_j m_j^2$$

Eigenvector | Eigenvector Eigenvector | Eigenvector

$$EE^T \vec{h} = \vec{h}$$

h = U_j |

HITS is similar to LSI, but on (source, destination) rather than (term,document) matrix

Maria Luisa Sapino (BDM 2018)

Problem

- Topic drift
 - Pages include neighbors
 - Neighbors may be good hubs, authorities; but may not have good content match

Maria Luisa Sapino (BDM 2018)

Clever system

- Use IR to find the candidate pages
- Expand to include all pages which link or are linked by this core set
- Compute authority and hub values for all pages (iterate!!)
- **Consider text next to the link!!!!**

$$a(i) = \sum_{j \in \text{in}(i)} w(j \rightarrow i) h(j) \quad h(i) = \sum_{j \in \text{out}(i)} w(i \rightarrow j) a(j)$$

Maria Luisa Sapino (BDM 2018)

Problem

- Topic drift
 - Pages include neighbors
 - Neighbors may be good hubs, authorities; but may not have good content match
- Slow
 - Iteration is not good!
 - One eigenvector computation per query

Maria Luisa Sapino (BDM 2018)

PageRank

- Random Surfer
 - Jumps from page to page with uniform probability
 - Occasionally jump to a random page with small probability (1-β)
 - If no out page, then jump to any page with equal probability

Maria Luisa Sapino (BDM 2018)

PageRank

- Random Surfer (N pages)
 - Jumps from page to page with uniform probability
 - Occasionally jump to a random page with small probability (1-β)
 - If no out page, then jump to any page with equal probability

$$\mathbf{Z} = (1 - \beta) \begin{bmatrix} \frac{1}{N} \\ \vdots \\ \frac{1}{N} \end{bmatrix}_{N \times N} + \beta \mathbf{M}$$

$M_{ij} = \begin{cases} \frac{1}{out(i)} & \text{if there is an edge from } i \text{ to } j \\ 0 & \text{otherwise} \end{cases}$

Transition matrix

Maria Luisa Sapino (BDM 2018)

PageRank

- Random Surfer (N pages)
 - Jumps from page to page with uniform probability
 - Occasionally jump to a random page with small probability (1-β)
 - If no out page, then jump to any page with equal probability

$$P(j) = \frac{1 - \beta}{N} + \beta \sum_{i \in in(j)} \frac{P(i)}{out(i)}$$

Probability that the surfer is at page j

Maria Luisa Sapino (BDM 2018)

PageRank

- Random Surfer (N pages)
 - Jumps from page to page with uniform probability
 - Occasionally jump to a random page with small probability (1-β)
 - If no out page, then jump to any page with equal probability

$$P(j) = \frac{1 - \beta}{N} + \beta \sum_{i \in in(j)} \frac{P(i)}{out(i)}$$

Probability that the surfer is at page j

Primary eigenvector of the transition matrix Z

Maria Luisa Sapino (BDM 2018)
