

Link Analysis

Web Ranking

- Documents on the web are first ranked according to their relevance vrs the query
- Additional ranking methods are needed to cope with huge amount of information
- Additional ranking methods:
 - Classification (manual, automatic)
 - **Link Analysis** (today's lesson)

Why link analysis?

- The web is **not** just a collection of documents – its hyperlinks are important!
- A link from page *A* to page *B* may indicate:
 - *A* is related to *B*, or
 - *A* is recommending, citing, voting for or endorsing *B*
- Links are either
 - referential – *click here and get back **home***, or
 - Informational – *click here to get more **detail***
- Links affect the ranking of web pages and thus have **commercial value**.
- The idea of using links is somehow “borrowed” by citation analysis

A is related to B, B is referential



DIPARTIMENTO DI INFORMATICA



Cerca

DIPARTIMENTO

STRUTTURE

DIDATTICA

RICERCA

NOTIZIE

Home

HOME



IN EVIDENZA

- ▶ Lectio Magistralis
- ▶ Notizie
- ▶ Riconoscimenti
- ▶ Seminari
- ▶ Dicono di noi...
- ▶ Incontri con le aziende

CHI SIAMO

DOVE SIAMO

GOVERNO

UFFICI

PERSONE

RICONOSCIMENTI

DICONO DI NOI...

BANDI

TRASPARENZA

BENVENUTI NEL SITO DEL DIPARTIMENTO DI INFORMATICA



WIRED

SERVIZI

- ▶ Webmail
- ▶ Modulistica
- ▶ U-Gov
- ▶ Infostud
- ▶ Visualizzazione aule

A is citing B, B is informational

- GOVERNO
- UFFICI
- PERSONE
- RICONOSCIMENTI
- DICONO DI NOI...
- BANDI
- TRASPARENZA



- Riconoscimenti
- Seminari
- Dicono di noi...
- Incontri con le aziende

BENVENUTI NEL SITO DEL DIPARTIMENTO DI INFORMATICA



CONGRATULAZIONI!


Dal 1 marzo Chiara Petrioli è diventata professore ordinario nel nostro dipartimento. Chiara ha ottenuto la promozione direttamente dal Ministero...

WIRED

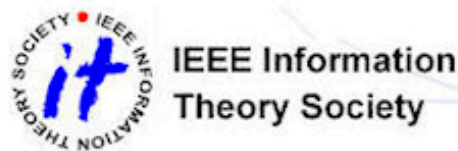
CDROID VS ANDROID

Wired.it, l'edizione italiana della nota rivista sull'innovazione digitale e tecnologica, ha dedicato un articolo...

- ### SERVIZI
- Webmail
 - Modulistica
 - U-Gov
 - Infostud
 - Visualizzazione aule

Seguiteci sul nostro gruppo 

Informatica@Sapienza



SHANNON AWARD 2014

Il prof. János Körner è stato insignito del



PRIMI IN ITALIA!

Si è da poco concluso l'esercizio di

A is recommending B, B is informational

- GOVERNO
- UFFICI
- PERSONE
- RICONOSCIMENTI
- DICONO DI NOI...
- BANDI
- TRASPARENZA



- ▶ Riconoscimenti
- ▶ Seminari
- ▶ Dicono di noi...
- ▶ Incontri con le aziende

BENVENUTI NEL SITO DEL DIPARTIMENTO DI INFORMATICA



CONGRATULAZIONI!

Dal 1 marzo Chiara Petrioli è diventata professore ordinario nel nostro dipartimento. Chiara ha ottenuto la promozione direttamente dal Ministero...

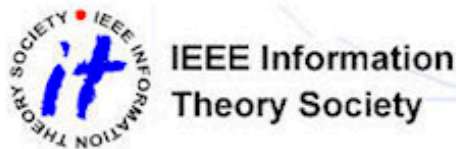
WIRED

CDROID VS ANDROID

Wired.it, l'edizione italiana della nota rivista sull'innovazione digitale e tecnologica, ha dedicato un articolo...

- ### SERVIZI
- ▶ Webmail
 - ▶ Modulistica
 - ▶ U-Gov
 - ▶ Infostud
 - ▶ Visualizzazione aule

Seguiteci sul nostro
gruppo 
informatica@Sapienza



SHANNON AWARD 2014

Il prof. János Körner è stato insignito del



PRIMI IN ITALIA!

Si è da poco concluso l'esercizio di

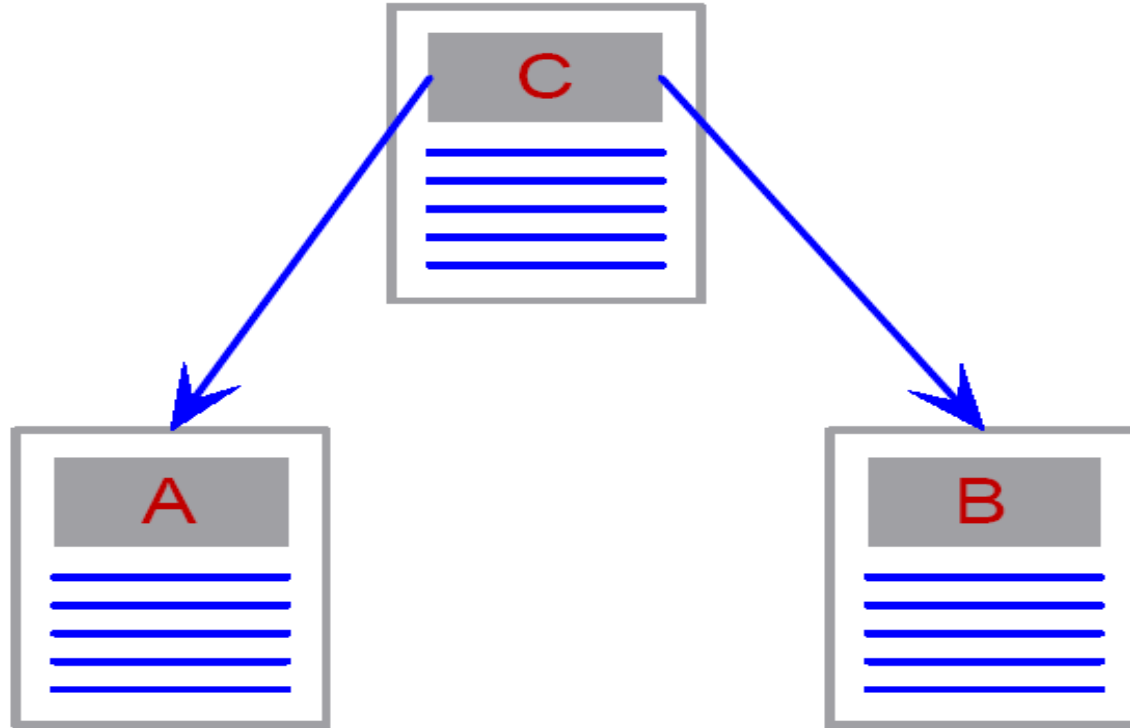


Citation Analysis

- The **impact factor** of a journal = A/B
 - A is the number of **current year citations** to articles appearing in the journal during previous two years.
 - B is the **number of articles** published in the journal during previous two years.

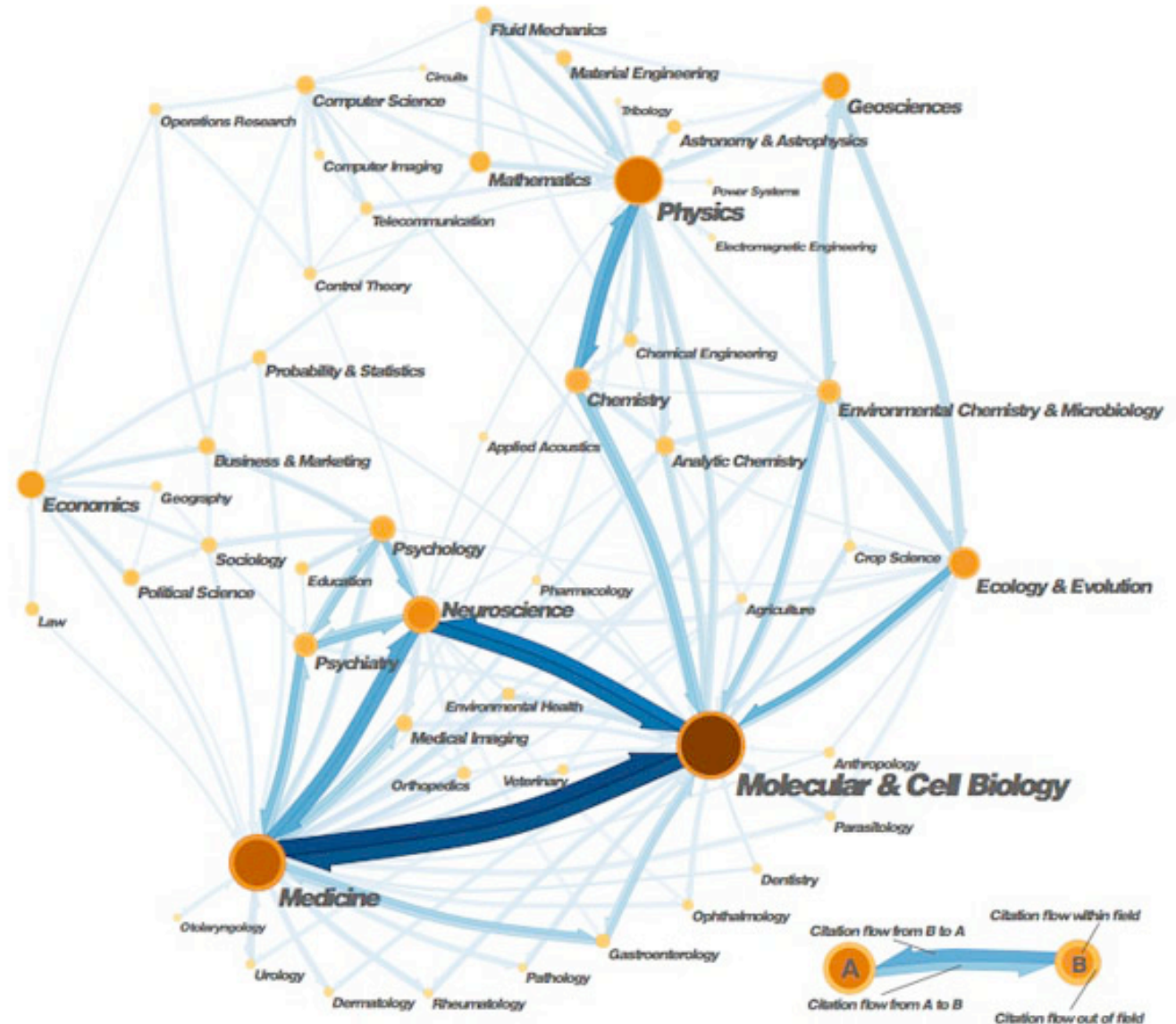
Journal Title (AI)	Impact Factor (2004)
J. Mach. Learn. Res.	5.952
IEEE T. Pattern Anal.	4.352
IEEE T. Evolut. Comp.	3.688
Artif. Intell.	3.570
Mach. Learn.	3.258

Co-Citation



- *A* and *B* are co-cited by *C*, implying that they are related or associated.
- The strength of co-citation between *A* and *B* is the number of times they are co-cited.

Clusters from Co-Citation Graph



Citations vs. Links

- *Web links are a bit different than citations:*
 - *Many links are navigational.*
 - *Many pages with high out-degree are **portals**, not content providers.*
 - *Not all links are endorsements (e.g. pointers to “fake” conferences).*
 - *Company websites don’t point to their competitors.*

However, the general idea that

“many citations = authority”

has been borrowed in link analysis

HITS and Page Rank: algebra that you need

- Eigenvector, eigenvalue and eigen-decomposition of normal matrixes
- **Iterative methods** 

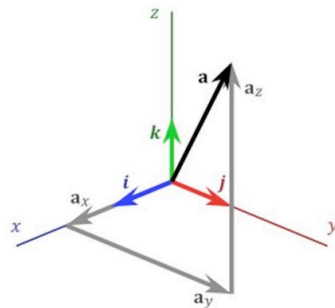
Iterative methods

- A mathematical procedure that generates a sequence of improving approximate solutions for a class of problems
- General formulation: $\mathbf{x}^{t+1} = \mathbf{A}\mathbf{x}^t$ where \mathbf{x} is a vector and t an iteration
- Iterative methods converge under specific hypotheses for matrix \mathbf{A} .
- **Condition a:** \mathbf{A} is square, real and symmetric
 - A real symmetric matrix is also **normal** and it exists a decomposition $\mathbf{U}\mathbf{\Delta}\mathbf{U}^{-1}$ such that $\mathbf{\Delta} = \text{diag}(\lambda_1, \dots, \lambda_n)$ and $\lambda_1 > \dots > \lambda_n$
 - Under these conditions, the method converges (as shown later)
- **Condition b:** \mathbf{A} is square, stochastic and irreducible
 - In a stochastic matrix, either $\sum_i (a_{ij}) = 1$ (right stochastic) or $\sum_j (a_{ij}) = 1$ (left stochastic)
 - A $n \times n$ matrix is **reducible** if indices $1, 2, \dots, n$ can be divided into two disjoint nonempty sets i_1, i_2, \dots, i_μ and j_1, j_2, \dots, j_ν such that $a_{i_\alpha j_\beta} = 0$ for $\alpha = 1, 2, \dots, \mu$ and $\beta = 1, 2, \dots, \nu$ (equivalent to say that subsumed graph has disconnected components)
- Conditions a and b are not the ONLY conditions for convergence of iterative methods, **but those we need here**

Geometric or graph interpretation of matrixes

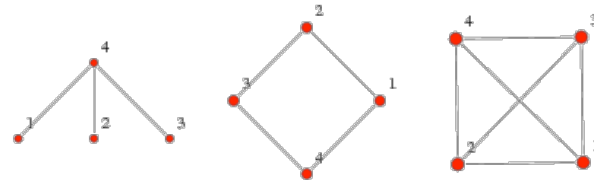
Geometric interpretation

- a_{ij} are coordinates of column vectors of the matrix on the carthesian axes $i=1..n$
- Ax is a linear transformation: if A is **normal**, $\lambda_1 > .. > \lambda_n$ and there exist an othonormal space defined by A 's eigenvectors on which x is projected.



Graph representation

- A matrix is a **weighted graph**, a_{ij} represent the weight of an edge between nodes i and j
- Irreducible matrix= the subsumed graph is connected



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

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

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

Link Analysis

- **HITS (Hyperlink Induced Topic Search)**
Jon Kleinberg
- Page Rank *Larry Page, Sergei Brin*

Hyperlink Induced Topic Search (HITS)

- Or Hypertext-Induced Topic Search(HITS) developed by Jon Kleinberg, while visiting IBM Almaden
- IBM expanded HITS into Clever, a web search engine for advertising (no longer an active project)
- However, HITS still used in many graph-based applications (e.g. social networks)

Main concept of the algorithm

- HITS stands for **Hypertext Induced Topic Search**.
- HITS is search query dependent.
- When the user issues a search query,
 - HITS first expands the list of “relevant” (according to, e.g. vector space model) pages returned by a search engine
 - Next, it produces two rankings of the expanded set of pages, **authority ranking** and **hub ranking**.

Main concept of the algorithm-cont.

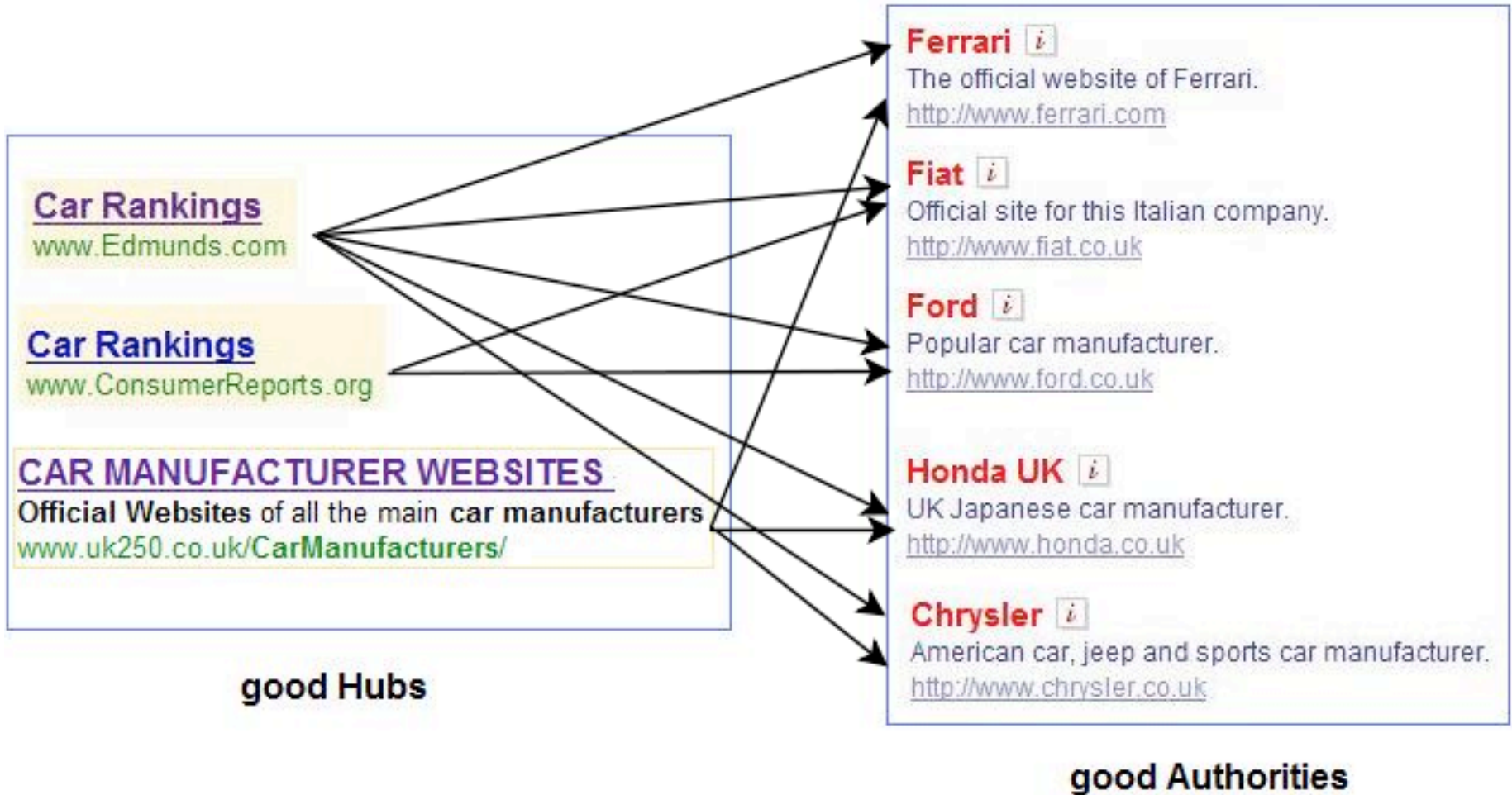
Authority: A authority is a page with many incoming links (**in-links, back-links**).

- The idea is that the page may have good or authoritative content on some topic and thus many people trust it and link to it.

Hub: A hub is a page with many out-links.

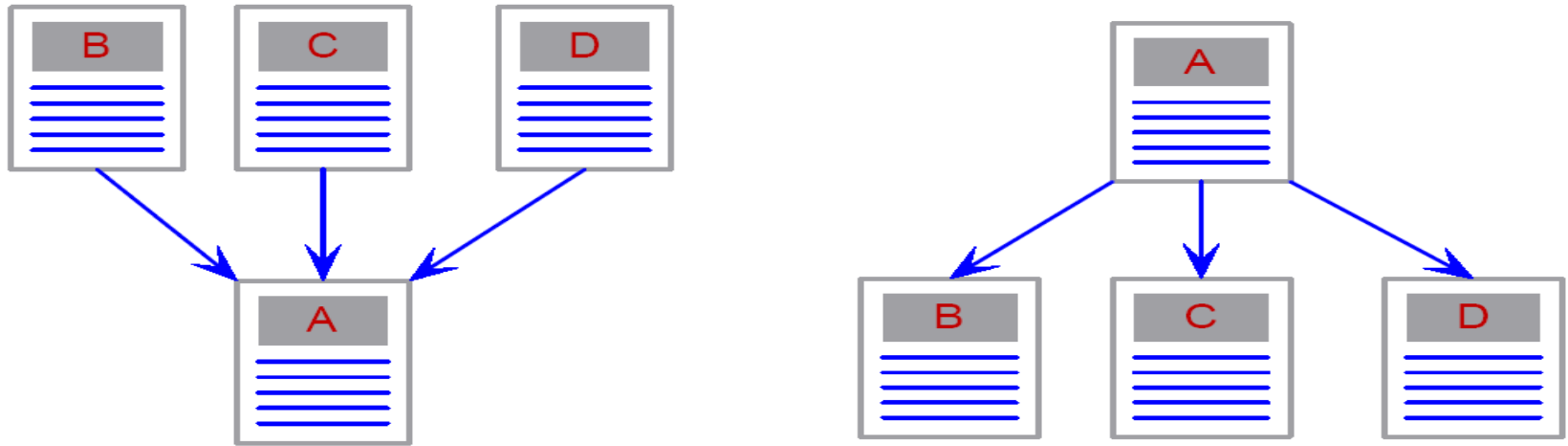
- The page serves as an **organizer of the information** on a particular topic and points to many good authority pages on the topic (e.g. a portal).

Example



Query: **Top automobile makers**

HITS – Hubs and Authorities –



- **A** on the left is an **authority**
- **A** on the right is a **hub**

Description of HITS

- A good hub points to many good authorities, and
- A good authority is pointed to by many good hubs.
- Authorities and hubs have a **mutual reinforcement relationship**. The figure shows some densely linked authorities and hubs (a **bipartite sub-graph**).

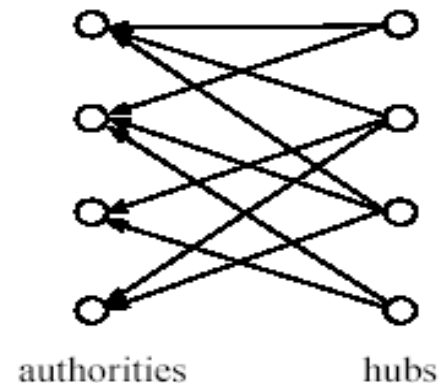


Fig. 8. A densely linked set of authorities and hubs

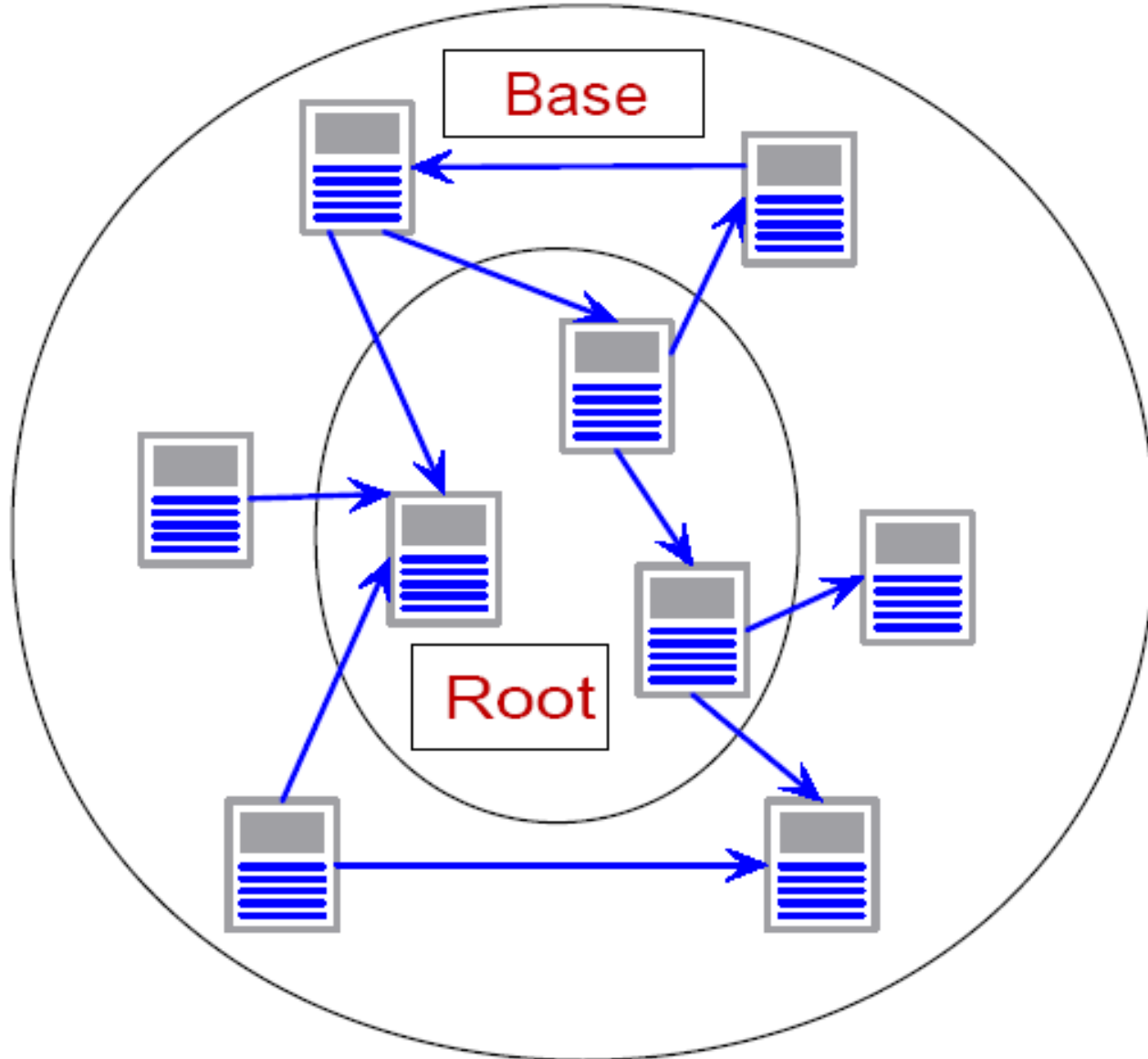
Hubs and Authorities: two steps

- First Step:
 - Constructing a **focused subgraph** of the WWW, based on a user's query
- Second step:
 - **Iteratively** calculate authority weight and hub weight for each page in the subgraph

The HITS algorithm: focused graph

- Given a broad search query, q , HITS collects a set of pages as follows:
 - It sends the query q to a search engine.
 - It then collects t ($t = 200$ is used in the HITS paper) highest ranked pages. This set is called the **root** set W .
 - It then grows W by including any page **pointed to** by a page in W and any page **that points** to a page in W . This gives a larger set S , **base set**.

Expanding the Root Set

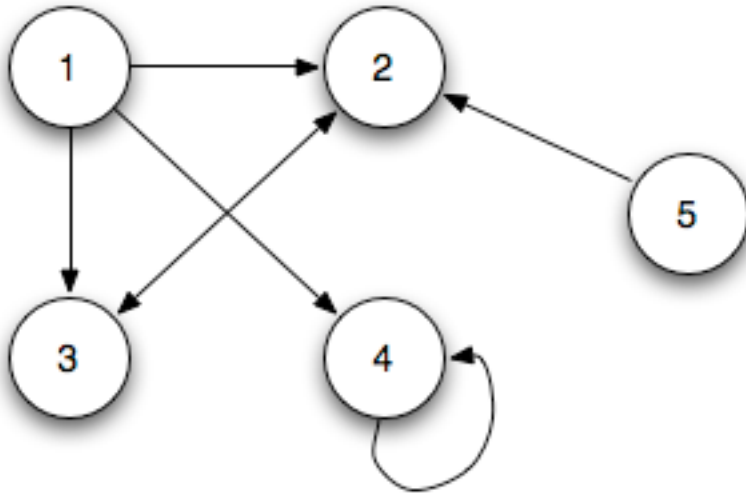


The link graph G

- HITS works on the pages in $S = B \cup R$, and assigns every page in S an **authority score** and a **hub score**.
- Let the number of pages in S be n .
- We use $G = (V, E)$ to denote the **hyperlink graph** of S . (V nodes, E edges). (NOTE: by construction, this is a CONNECTED graph)
- We use L to denote the **adjacency matrix** of the graph.

$$L_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E \\ 0 & \text{otherwise} \end{cases} \quad i \rightarrow j$$

Adjacency matrix (directed graph)



	1	2	3	4	5
1	0	1	1	1	0
2	0	0	1	0	0
3	0	1	0	0	0
4	0	0	0	1	0
5	0	1	0	0	0

The HITS algorithm (cont'd)

- Let the authority score of the page i be $a(i)$, and the hub score of page i be $h(i)$.
- The **mutual reinforcing relationship** of the two scores is represented as follows:

$$a(i) = \sum_{(j,i) \in E} h(j)$$

$$h(i) = \sum_{(i,j) \in E} a(j)$$

Remember: E is the set of edges of the derived hyperlink graph

HITS in matrix form

- We use \mathbf{a} to denote the column **vector** with all the authority scores,

$$\mathbf{a} = (a(1), a(2), \dots, a(n))^T, \text{ and}$$

- use \mathbf{h} to denote the column **vector** with all the "hub scores",

$$\mathbf{h} = (h(1), h(2), \dots, h(n))^T,$$

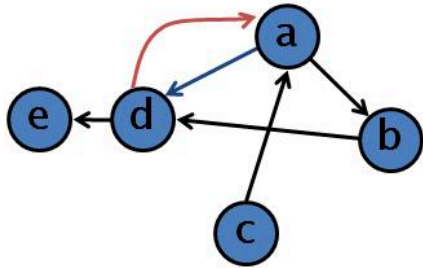
- Then, we can express previous formulas in matrix form as:

$$\mathbf{a} = \mathbf{L}^T \mathbf{h} \quad (\text{I step})$$

$$\mathbf{h} = \mathbf{L} \mathbf{a} \quad (\text{O step})$$

$$\text{normalize: } \mathbf{a} = \mathbf{a} / \|\mathbf{a}\| \quad \mathbf{h} = \mathbf{h} / \|\mathbf{h}\|$$

- It is an equivalent formulation wrt the sum in previous formula is **for all j linked to i** in E and \mathbf{L} has 1 where there is a link between i and j since $a(i) = \sum_{(j,i) \in E} h(j)$ 27



	a	b	c	d	e
a	0	1	0	1	0
b	0	0	0	1	0
c	1	0	0	0	0
d	1	0	0	0	1
e	0	0	0	0	0

$$\begin{pmatrix} h_a \\ h_b \\ h_c \\ h_d \\ h_e \end{pmatrix} = \begin{pmatrix} 01010 \\ 00010 \\ 10000 \\ 10001 \\ 00000 \end{pmatrix} \times \begin{pmatrix} a_a \\ a_b \\ a_c \\ a_d \\ a_e \end{pmatrix}$$



$$\begin{aligned} h_a &= a_b + a_d \\ h_b &= a_d \\ h_c &= a_a \\ h_d &= a_b + a_e \\ h_e &= 0 \end{aligned}$$

Computation of HITS

- The computation of authority scores and hub scores uses **power iteration** iterative method.
- If we use \mathbf{a}^k and \mathbf{h}^k to denote authority and hub vectors at the k_{th} iteration, the iterations for generating the final (stationary) solutions are:

$$\begin{array}{l} \mathbf{a}^k = \mathbf{L}^T \mathbf{h}^{k-1} \\ \mathbf{h}^{k-1} = \mathbf{L} \mathbf{a}^{k-1} \end{array} \quad \longrightarrow \quad \begin{array}{l} \mathbf{a}^k = \mathbf{L}^T \mathbf{L} \mathbf{a}^{k-1} \\ \mathbf{h}^k = \mathbf{L} \mathbf{L}^T \mathbf{h}^{k-1} \end{array}$$

Example (simple algorithm)

- 2nd Iteration
- I Step
- O Step



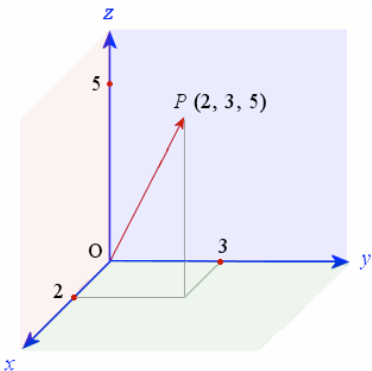
The HITS algorithm (with normalization)

- $\underline{h}^{(0)} := (1, 1, \dots)$
- $k := 1$
- *Until convergence, do:*
 - $\underline{a}^{(k)} := L^T \underline{h}^{(k-1)} = L^T L \underline{a}^{(k-1)}$ (*update a*)
 - $\underline{h}^{(k)} := L \underline{a}^{(k)} = L L^T \underline{h}^{(k-1)}$ (*update h*)
 - $\underline{a}^{(k)} := \underline{a}^{(k)} / \|\underline{a}^{(k)}\|$ and $\underline{h}^{(k)} := \underline{h}^{(k)} / \|\underline{h}^{(k)}\|$ (*normalize*)

Does it converge to a stationary solution?

Digression: many vector notations, don't be confused!

- \underline{X} , \vec{X} and \mathbf{X} (underline, arrow and bold) are all valid *intensional* notations for vectors!!
- $X: (x_1, x_2, \dots, x_n)$ is an *extensional* notation (shows all coordinates of the vector) and is either column or row:



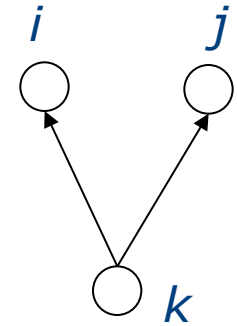
- *Finally, we have the **graphic** notation*

$$\begin{pmatrix} x_1 \\ x_2 \\ \square \\ \square \\ x_{n-1} \\ x_n \end{pmatrix}$$

Back to HITS: meaning of the $L L^T$ and $L^T L$ matrixes

- L is the adjacency matrix of the graph
- $L^T L$ is the authority matrix:

$$A = L^T L = \left(\begin{array}{c} \vdots \\ \text{--- } A_{ij} \text{ ---} \\ \vdots \end{array} \right) = \underbrace{\left(\begin{array}{c} \vdots \\ \text{--- } L_{ik}^T \text{ ---} \\ \vdots \end{array} \right)}_{L^T} \underbrace{\left(\begin{array}{c} \vdots \\ \text{--- } L_{kj} \\ \vdots \end{array} \right)}_L$$



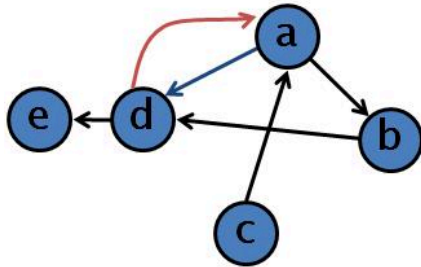
L_{kj} column “means” that j is pointed by all non-zero k ;

L_{ik}^T “means” that i is pointed by all non-zero k

$$A_{ij} = \sum_{k=1}^n L_{ik}^T L_{kj} = \sum_{k=1}^n L_{ki} L_{kj}$$

A_{ij} is the number of **co-citations**, the Number of nodes pointing to both i and j

..is this something you have already seen???????



	a	b	c	d	e
a	0	1	0	1	0
b	0	0	0	1	0
c	1	0	0	0	0
d	1	0	0	0	1
e	0	0	0	0	0

$$L^T \times L = \begin{pmatrix} 00110 \\ 10000 \\ 00000 \\ 11000 \\ 00010 \end{pmatrix} \times \begin{pmatrix} 01010 \\ 00010 \\ 10000 \\ 10001 \\ 00000 \end{pmatrix} = \begin{pmatrix} 20001 \\ 01010 \\ 00000 \\ 01020 \\ 10001 \end{pmatrix}$$

A_{ij} is the number of nodes pointing both i and j .

*For example, $A_{ae} = 1$ since only node **d** points to both **a** and **e***

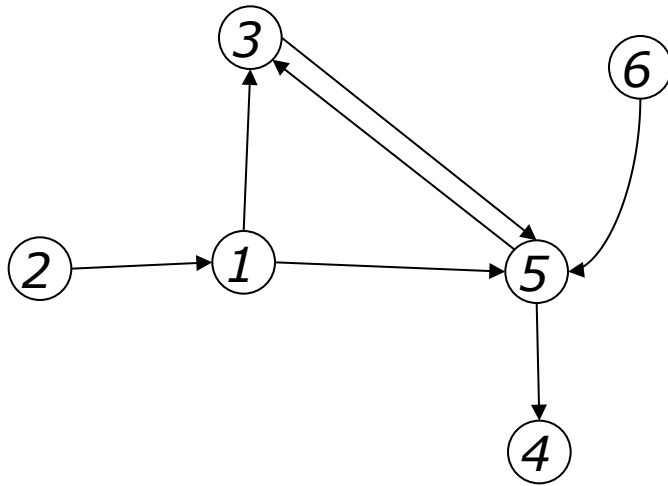
Proof of convergence of HITS (power method)

- Since matrix $A (=LL^T)$ is **square and symmetric**, it has an eigen-decomposition $U\Lambda U^{-1}$ where $\lambda_1, \lambda_2, \dots, \lambda_k$ are then **eigenvalues** and $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_k|$ (note: they are eigenvalues of A and **singular values** of L !!)
- $\underline{x}_1, \dots, \underline{x}_k$ are the **eigenvectors** of A and they form an orthonormal basis (e.g. : $\alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2 + \dots + \alpha_k \underline{x}_k = 0$ iff $\alpha_1 = \dots = \alpha_k = 0$, since $\underline{a}^{(k)} := L^T L \underline{a}^{(k-1)} = A \underline{a}^{(k-1)}$)
- A generic vector \underline{v}_0 (in our case, either $\underline{h}^{(0)}$ or $\underline{a}^{(0)}$) can be re-written as:
 - $\underline{v}^0 = \alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2 + \dots + \alpha_k \underline{x}_k$ (its projection on the orthonormal space of A)
- Hence (let $\forall i: \underline{x}_i, \lambda_i$ eigenvalue, eigenvector of $A, A = \lambda_i \underline{x}_i$)
 - $\underline{h}^1 = A \underline{h}^0 = \alpha_1 A \underline{x}_1 + \alpha_2 A \underline{x}_2 + \dots + \alpha_k A \underline{x}_k = \alpha_1 \lambda_1 \underline{x}_1 + \alpha_2 \lambda_2 \underline{x}_2 + \dots + \alpha_k \lambda_k \underline{x}_k =$
 $\lambda_1 [\alpha_1 \underline{x}_1 + \alpha_2 (\lambda_2 / \lambda_1) \underline{x}_2 + \dots + \alpha_k (\lambda_k / \lambda_1) \underline{x}_k]$
- And in general:
 - $\underline{h}^m = A \underline{h}^{m-1} = A^m \underline{h}^0 = \alpha_1 A^m \underline{x}_1 + \alpha_2 A^m \underline{x}_2 + \dots + \alpha_k A^m \underline{x}_k = \alpha_1 \lambda_1^m \underline{x}_1 + \alpha_2 \lambda_2^m \underline{x}_2 + \dots +$
 $\alpha_k \lambda_k^m \underline{x}_k = \lambda_1^m [\alpha_1 \underline{x}_1 + \alpha_2 (\lambda_2 / \lambda_1)^m \underline{x}_2 + \dots + \alpha_k (\lambda_k / \lambda_1)^m \underline{x}_k]$
- Since $|\lambda_i / \lambda_1| < 1, i = 2, 3, \dots, n$, we get:

$$\lim_{m \rightarrow \infty} \frac{1}{\lambda_1^m} \underline{h}^m = \lim_{m \rightarrow \infty} \frac{1}{\lambda_1^m} A^m \underline{h}^0 = \alpha_1 \underline{x}_1$$

Speed of convergence depends on λ_2 / λ_1 and on initial choice of \underline{h}^0

HITS: Example (1)



$$L = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \begin{pmatrix} 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \end{matrix}$$

$$L^T L = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

$$LL^T = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \begin{pmatrix} 2 & 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 2 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{pmatrix} \end{matrix}$$

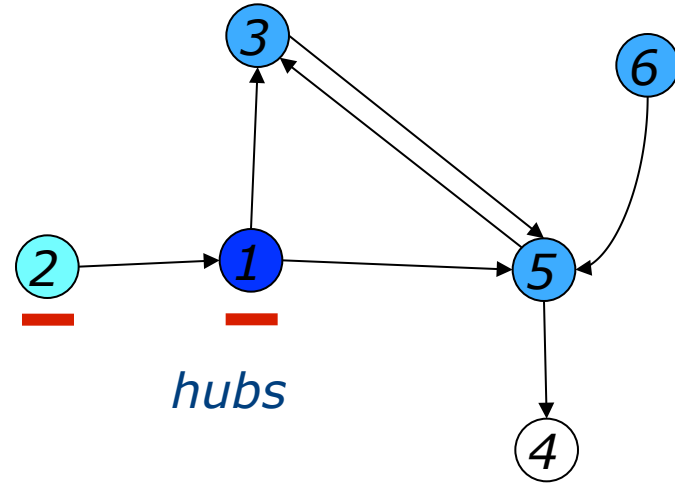
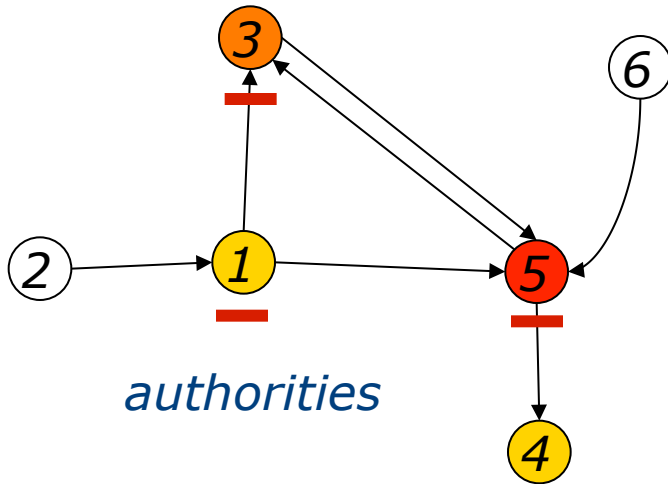
Ex: 3 and 4 are “co-cited” by 5

3 and 1 co-cite 5”

$$-\underline{a}^{(1)} := L^T \underline{h}^{(0)}$$

$$-\underline{h}^{(1)} := L \underline{a}^{(1)}$$

HITS: Example (2)



$$\begin{pmatrix} \underline{a}_1 \\ 0 \\ 0.516 \\ 0.258 \\ 0.775 \\ 0 \end{pmatrix} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \begin{pmatrix} 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \end{matrix} T \begin{pmatrix} \underline{h}_0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

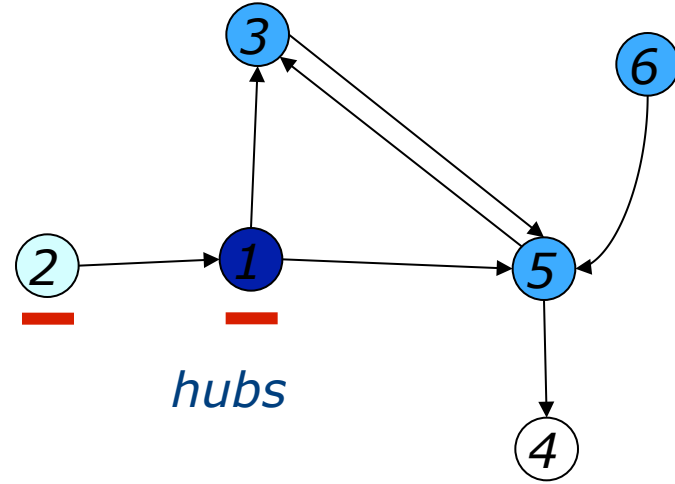
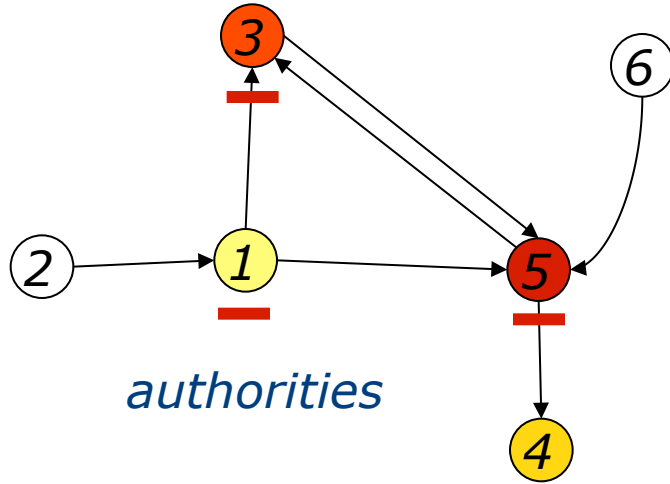
$$\begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \begin{pmatrix} 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \end{matrix} \begin{pmatrix} \underline{a}_1 \\ 0.258 \\ 0 \\ 0.516 \\ 0.258 \\ 0.775 \\ 0 \end{pmatrix} = \begin{pmatrix} \underline{h}_1 \\ 0.687 \\ 0.137 \\ 0.412 \\ 0 \\ 0.412 \\ 0.412 \end{pmatrix}$$

(NOTE: normalization step is not shown, however results are normalized)

$$-\underline{a}^{(2)} := L^T \underline{h}^{(1)}$$

$$-\underline{h}^{(2)} := L \underline{a}^{(2)}$$

HITS: Example (3)

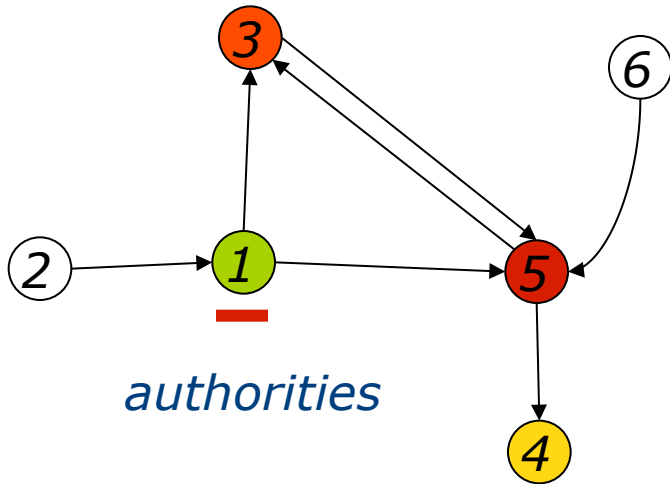


	1	2	3	4	5	6	\underline{h}_1	\underline{a}_2		1	2	3	4	5	6	\underline{a}_2	\underline{h}_2								
1	⎛	0	0	1	0	1	0	⎝	⎞	0.687	=	⎝	0.072	⎞	⎞	0.072	⎝	0.706							
2		1	0	0	0	0	0			0.137			0			2		1	0	0	0	0	0	0	0.037
3		0	0	0	0	1	0			0.412			0.573			3		0	0	0	0	1	0	0.573	0.409
4		0	0	0	0	0	0			0			0.215			4		0	0	0	0	0	0	0.215	0
5		0	0	1	1	0	0			0.412			0.788			5		0	0	1	1	0	0	0.788	0.409
6		0	0	0	0	1	0			0.412			0			6		0	0	0	0	1	0	0	0.409

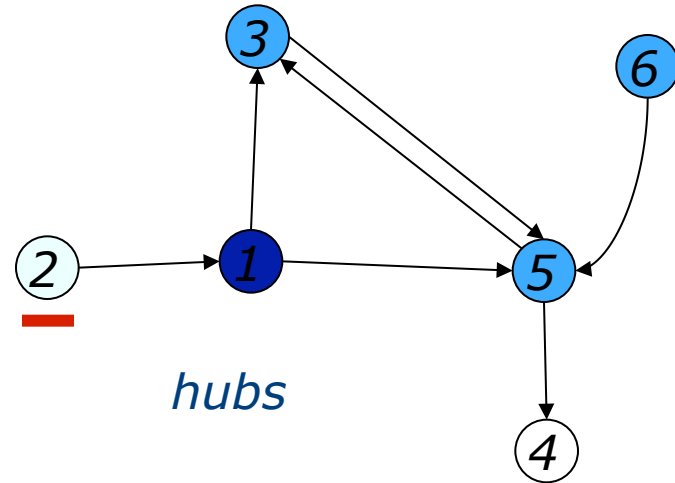
HITS: Example (4)

$$-\underline{a}^{(3)} := L^T \underline{h}^{(2)}$$

$$-\underline{h}^{(3)} := L \underline{a}^{(3)}$$



authorities



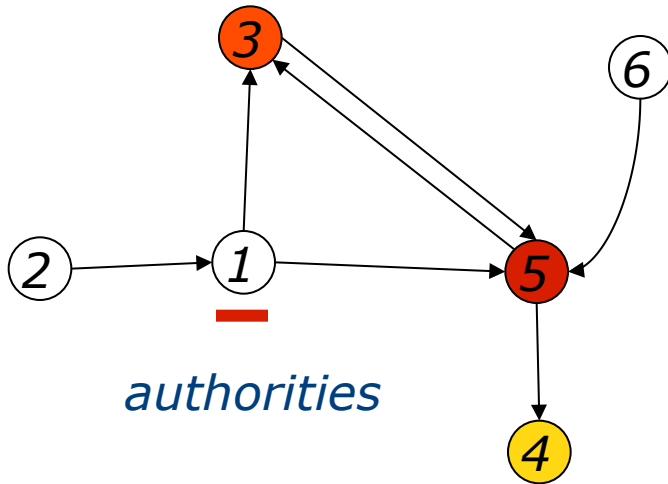
hubs

	1	2	3	4	5	6	\underline{h}_2	\underline{a}_3		1	2	3	4	5	6	\underline{a}_3	\underline{h}_3
1	0	0	1	0	1	0	0.706	0.019	1	0	0	1	0	1	0	0.019	0.707
2	1	0	0	0	0	0	0.037	0	2	1	0	0	0	0	0	0	0.001
3	0	0	0	0	1	0	0.409	0.577	3	0	0	0	0	1	0	0.577	0.408
4	0	0	0	0	0	0	0	0.212	4	0	0	0	0	0	0	0.212	0
5	0	0	1	1	0	0	0.409	0.789	5	0	0	1	1	0	0	0.789	0.408
6	0	0	0	0	1	0	0.409	0	6	0	0	0	0	1	0	0	0.408

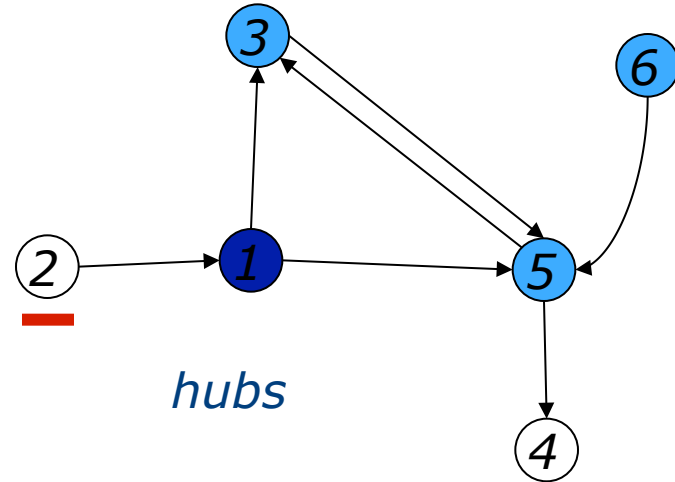
$$-\underline{a}^{(4)} := L^T \underline{h}^{(3)}$$

$$-\underline{h}^{(4)} := L \underline{a}^{(4)}$$

HITS: Esempio (5)



authorities



hubs

	1	2	3	4	5	6	\underline{h}_3	\underline{a}_4		1	2	3	4	5	6	\underline{a}_4	\underline{h}_4
1	0	0	1	0	1	0	0.707	0	1	0	0	1	0	1	0	0	0.707
2	1	0	0	0	0	0	0.001	0	2	1	0	0	0	0	0	0	0
3	0	0	0	0	1	0	0.408	0.577	3	0	0	0	0	1	0	0.577	0.408
4	0	0	0	0	0	0	0	0.211	4	0	0	0	0	0	0	0.211	0
5	0	0	1	1	0	0	0.408	0.789	5	0	0	1	1	0	0	0.789	0.408
6	0	0	0	0	1	0	0.408	0	6	0	0	0	0	1	0	0	0.408

Strengths and weaknesses of HITS

- **Strength**: its ability to rank pages according to the query topic, which may be able to provide more relevant authority and hub pages.
- **Weaknesses**:
 - **It is easily spammed**. It is in fact quite easy to influence HITS since adding out-links in one's own page is so easy.
 - **Topic drift**. Many pages in the expanded set may not be on topic.
 - **Inefficiency at query time**: The query time evaluation is slow. Collecting the root set, expanding it and performing eigenvector computation are all expensive operations

Applications of HITS

- Search engine querying (speed is an issue)
- Finding web communities.
- Finding related pages.
- Populating categories in web directories.
- Citation analysis
- Social network analysis

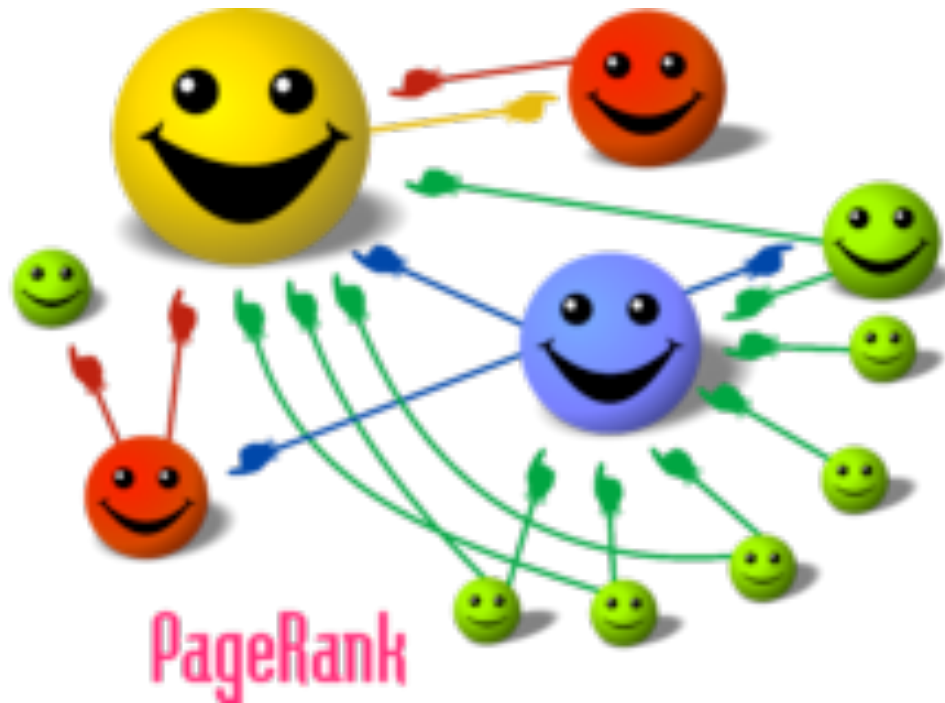
Link Analysis

- HITS (Hyperlink Induced Topic Search)
Jon Kleinberg
- **Page Rank** *Larry Page, Sergei Brin*

Page Rank

- Ranks pages by authority.
- Applied to the **entire web** rather than a local neighborhood of pages surrounding the results of a query.
- **Not query-dependent**
- It is the Google algorithm for ranking pages

PageRank----Idea



Every page has some number of out-links and in-links

PageRank----Idea

Two cases PageRank is interesting:

- 1. Web pages vary greatly in terms of the number of backlinks (in-links) they have. For example, the Netscape home page has 62,804 backlinks compared to most pages which have just a few backlinks. **Generally, highly linked pages are more “important” than pages with few links.***

EUGENE GARFIELD, FRANCIS NARIN, PAGERANK: THE THEORETICAL BASES OF THE GOOGLE SEARCH ENGINE

PageRank----Idea

2. In-links coming from important pages convey more importance to a page. For example, if a web page has a link off the Yahoo home page, it may be just one link but it is a very important one.

A page has high rank if the sum of the ranks of its incoming links is high. This covers both the case when a page has many in-links and when a page has a few highly ranked in-links.

PageRank----Definition

u,v: a web page

F_u: set of pages that u points to

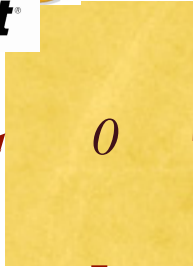
B_u: set of pages that point to u (backlinks or inlinks)

N_u=|F_u|: the number of links outgoing from u (outlinks)

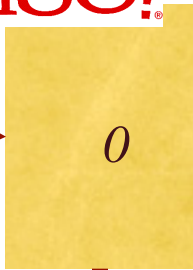
c: a factor used for normalization

$$r(u) = c \sum_{v \in B_u} \frac{r(v)}{N_v}$$

The equation is recursive, but it may be computed by starting with any set of ranks and iterating the computation until it converges.

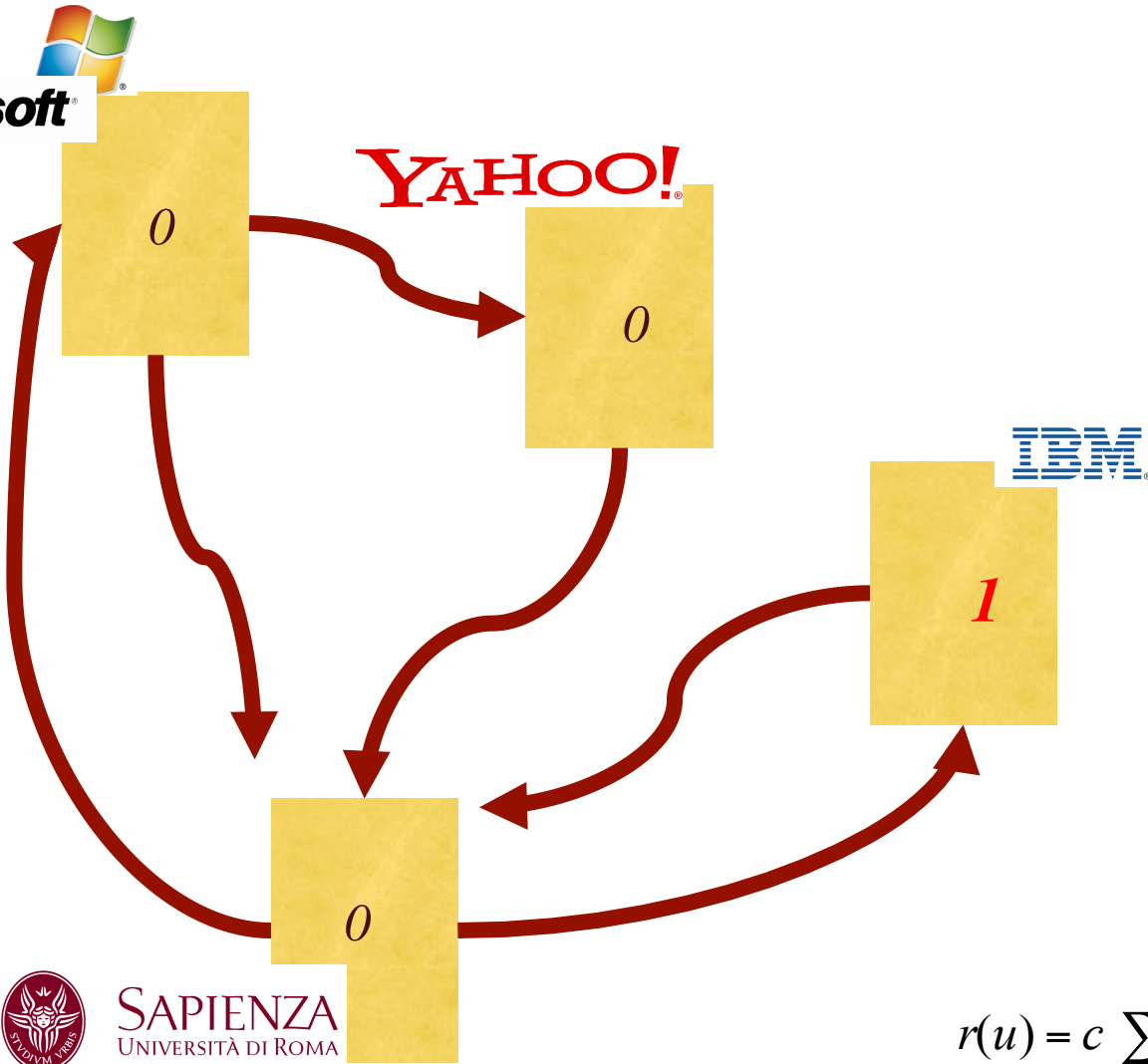


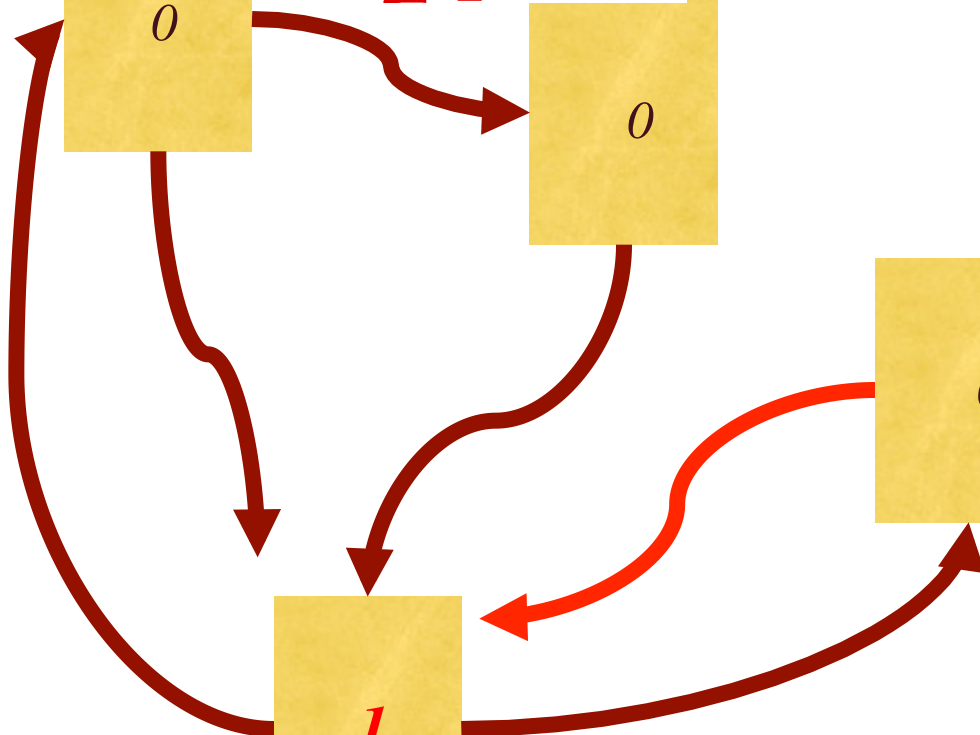
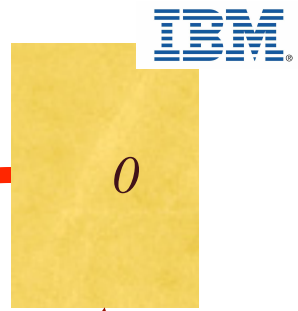
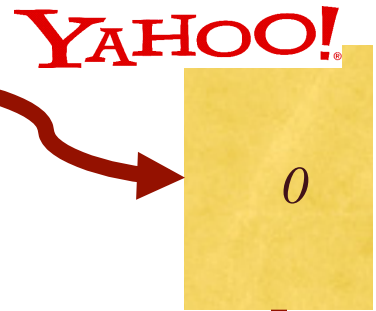
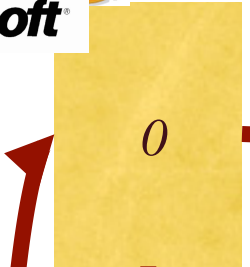
YAHOO!



SAPIENZA
UNIVERSITÀ DI ROMA

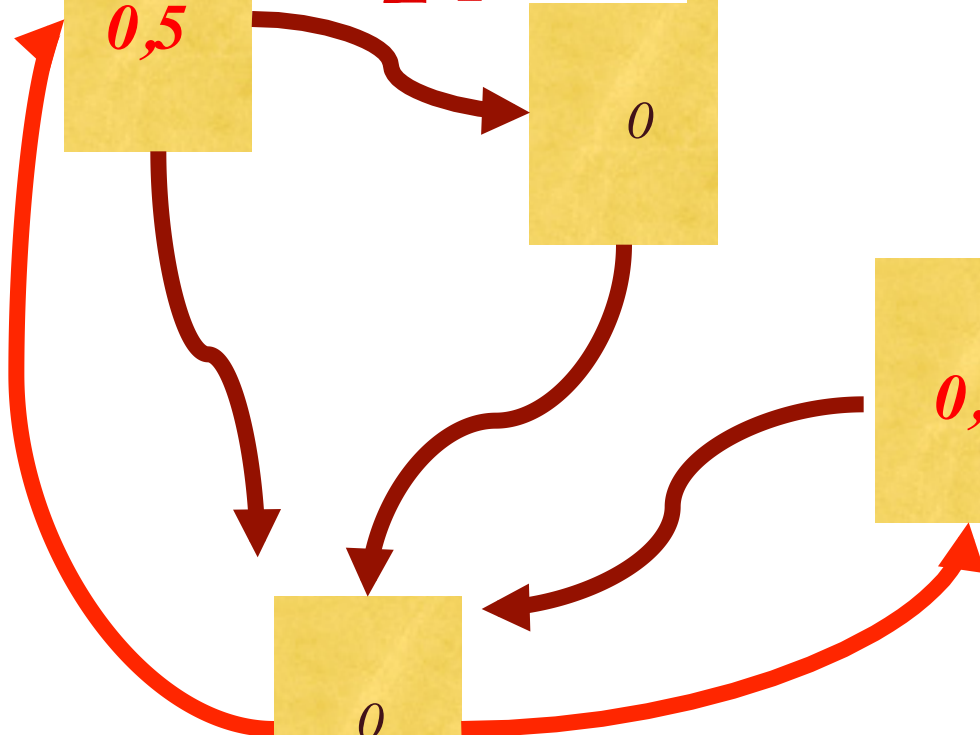
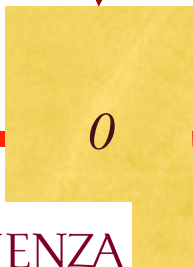
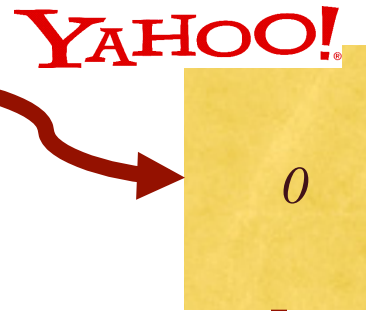
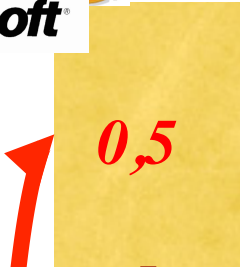
$$r(u) = c \sum_{v \in B_u} \frac{r(v)}{N_v}$$





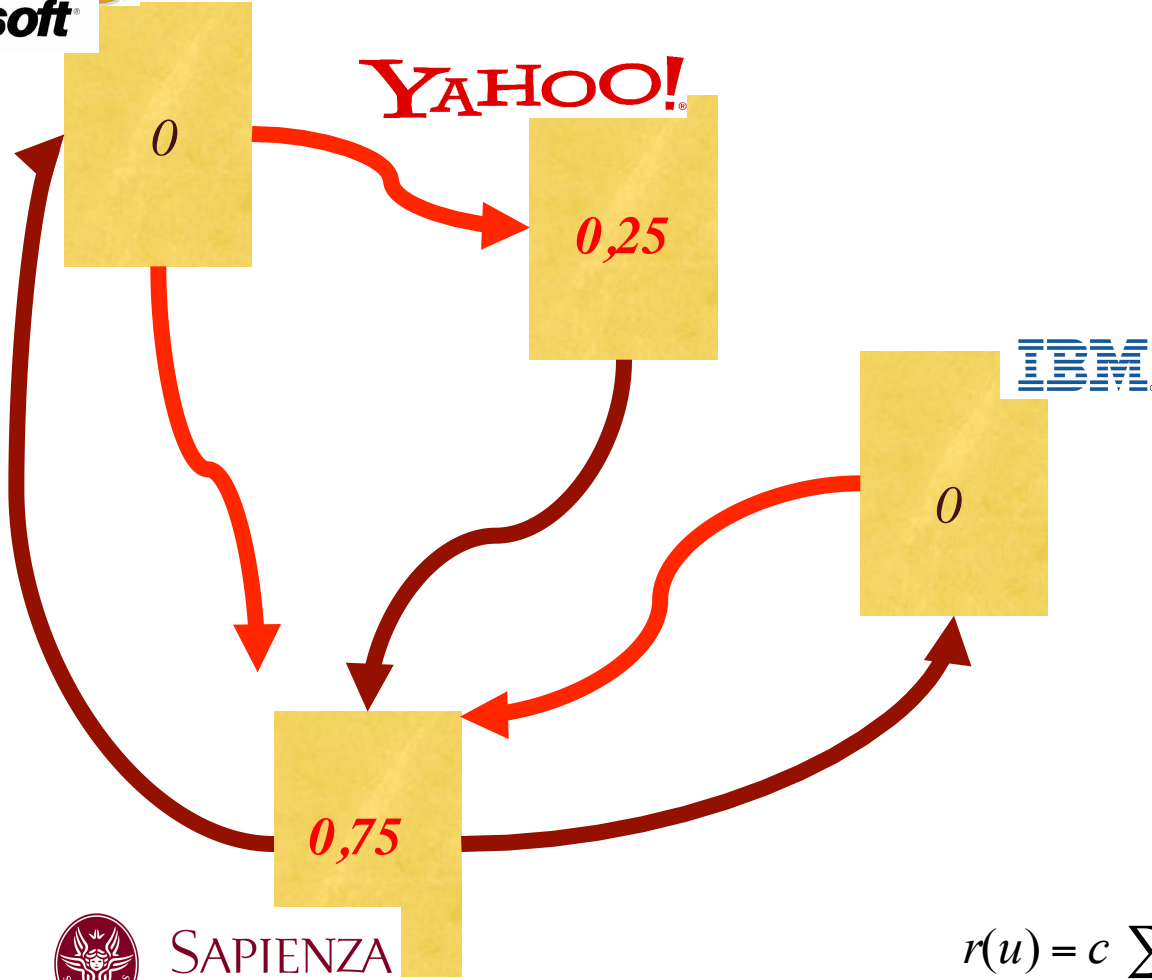
SAPIENZA
UNIVERSITÀ DI ROMA

$$r(u) = c \sum_{v \in B_u} \frac{r(v)}{N_v}$$



SAPIENZA
UNIVERSITÀ DI ROMA

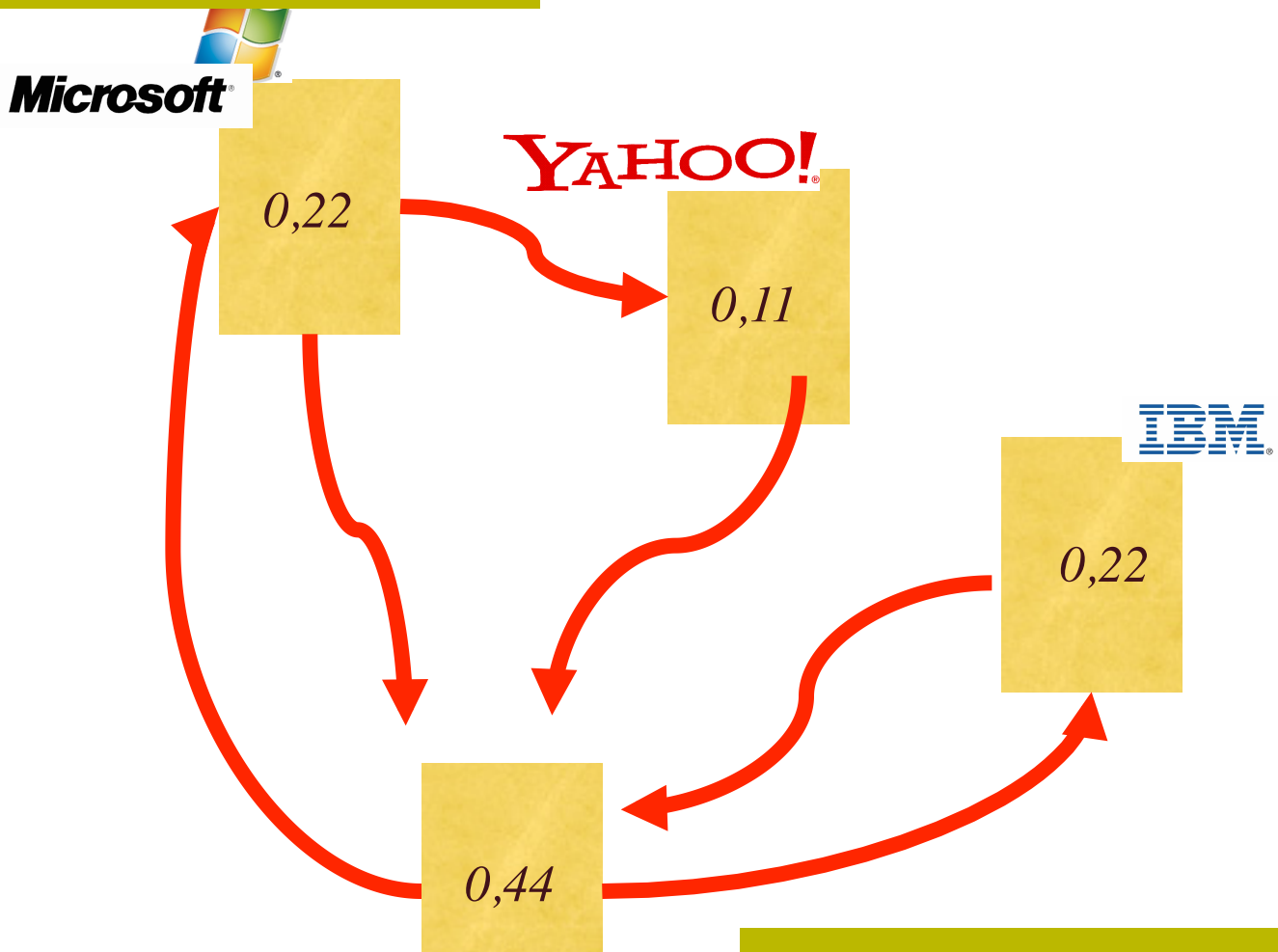
$$r(u) = c \sum_{v \in B_u} \frac{r(v)}{N_v}$$



SAPIENZA
UNIVERSITÀ DI ROMA

$$r(u) = c \sum_{v \in B_u} \frac{r(v)}{N_v}$$

After several iterations..



Why stops here?

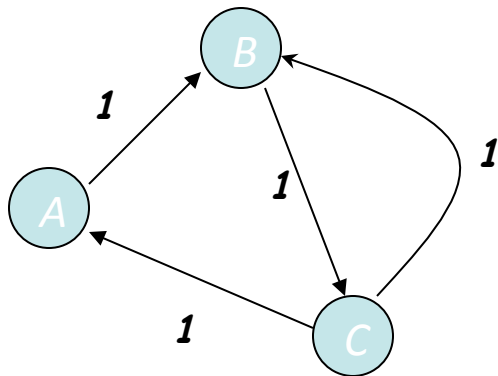
A probabilistic interpretation of PageRank

- *The definition corresponds to the probability distribution of a **random walk** on the web graphs.*
- *First note: we can write $r(u) = c \sum_{v \in B_u} \frac{r(v)}{N_v}$ in matrix iterative form as $\mathbf{r}^t = P \mathbf{r}^{t-1}$*
- *P (transition matrix) is **left stochastic** since $\sum_j p_{ij} = 1$ (column j in P corresponds to the outlinks of node j , and for k outlinks, each weights $1/k$)*

An example

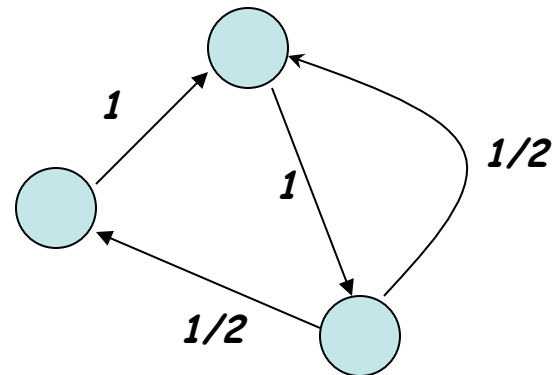
0	1	0
0	0	1
1	1	0

Adjacency matrix A



0	1	0
0	0	1
1/2	1/2	0

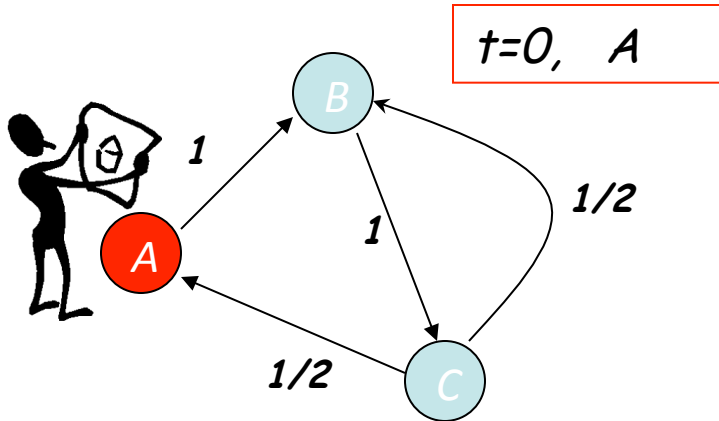
Transition matrix P



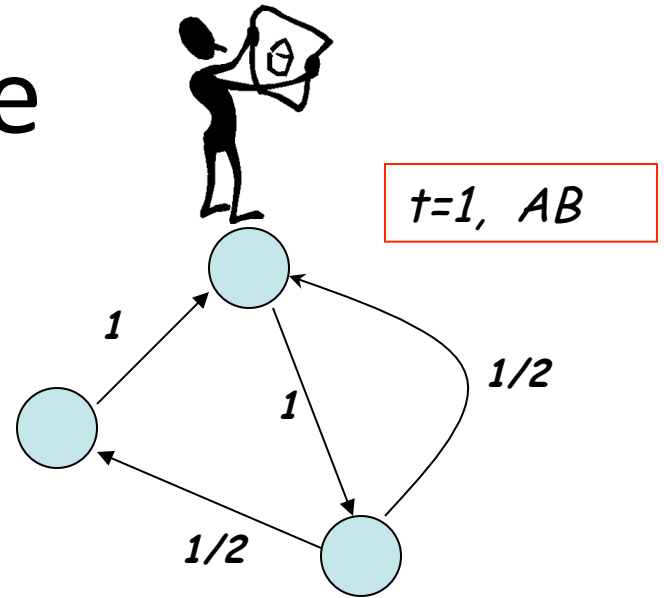
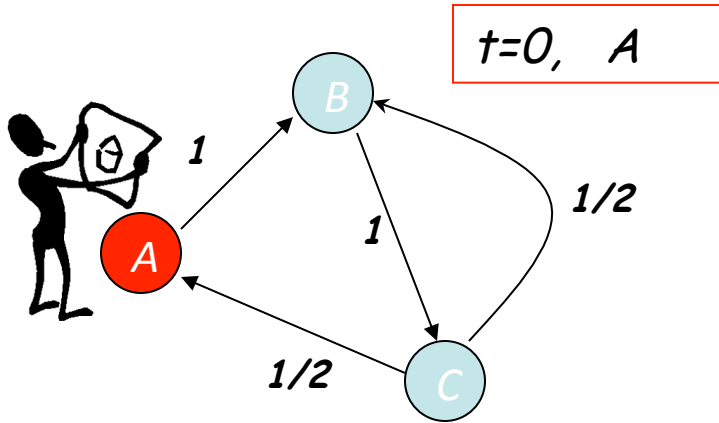
What is a Random Walk?

- Given a graph and a starting point (node), we select a neighbor (= a pointed node) of it **at random**, and move to this neighbor;
- Then we select (at random) a neighbor of this node and move to it, and so on;
- The (random) sequence of nodes selected this way is a **random walk** on the graph
- In our case, if the walker is on node **j** at time **t**, it has **1/k** probability of jumping on any of its hyperlinked nodes at time **t+1**

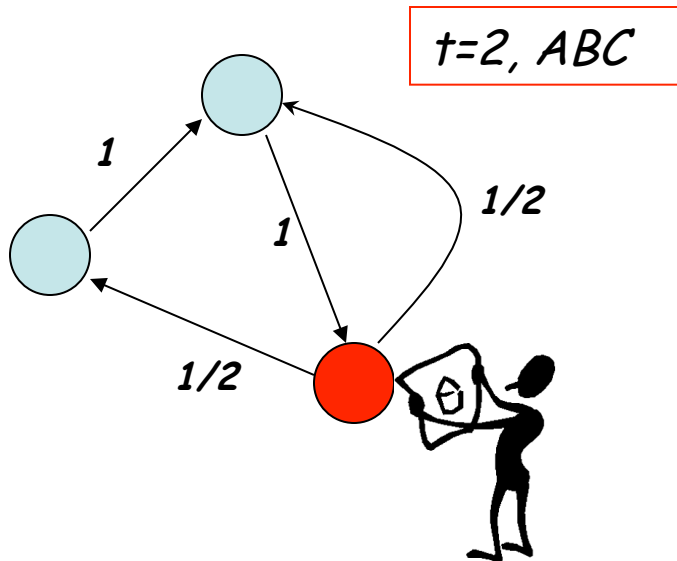
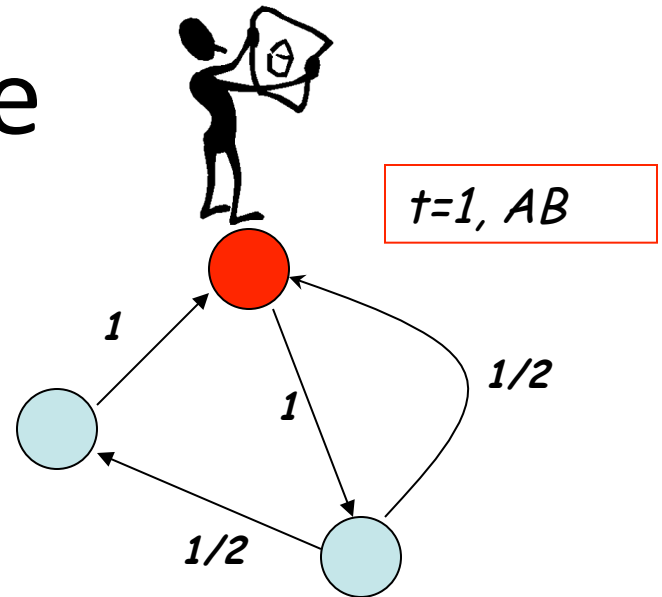
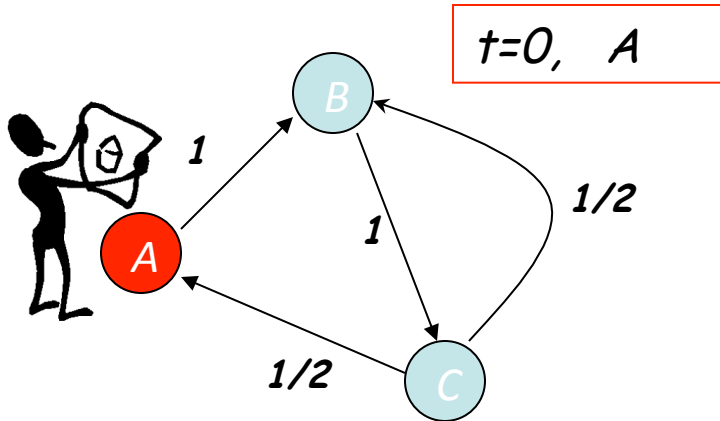
An example



An example

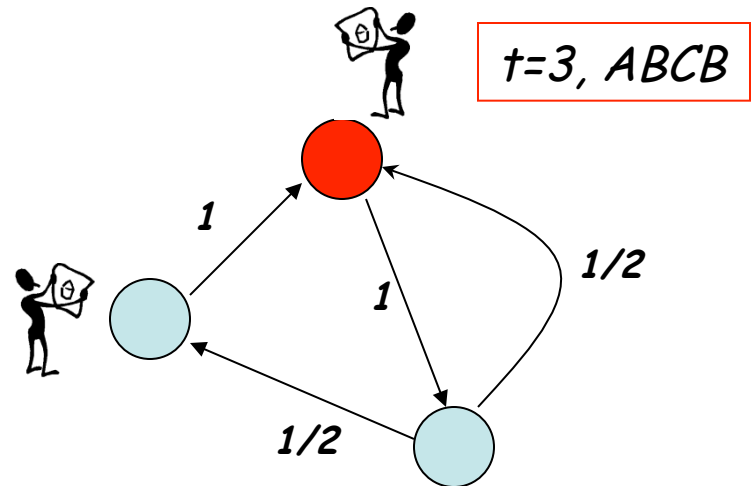
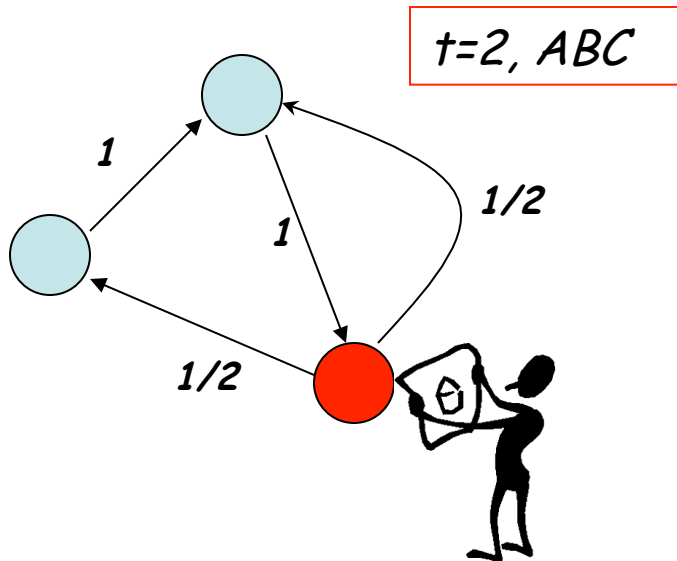
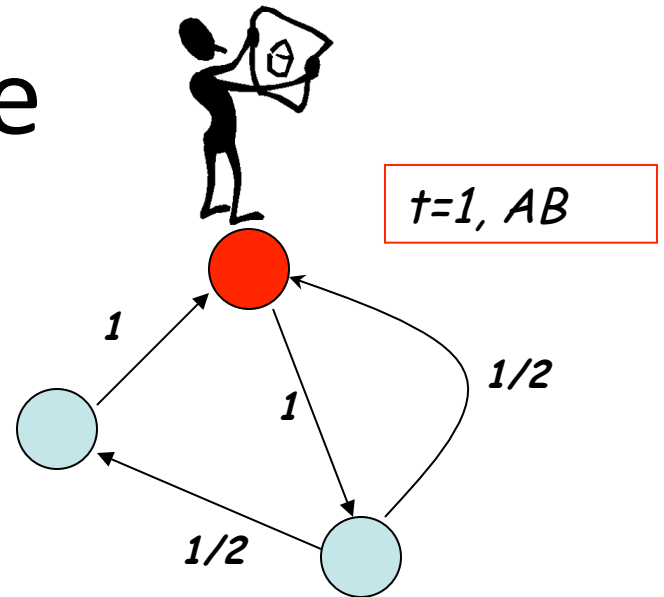
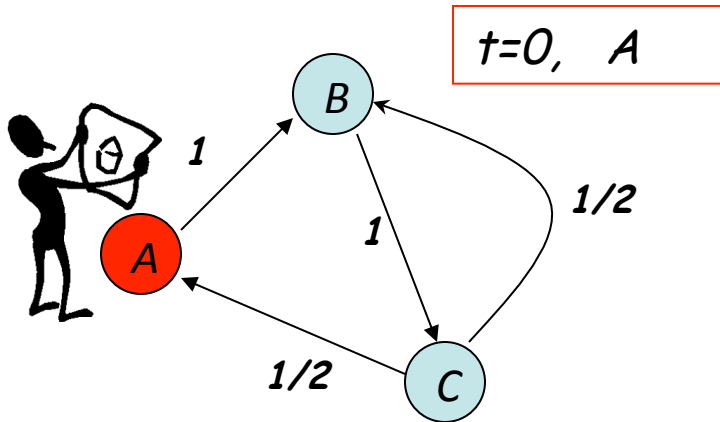


An example



In node C the random walker has a 0.5 probability of jumping to node B and 0.5 of jumping to node A

An example



Probabilistic interpretation

- N total number of web pages
- k outgoing links from page j
- P **Transition matrix** with elements:

$$P_{ij} = \begin{cases} 1/k & \text{if } i \rightarrow j \\ 0 & \\ P_{ii} > 0 \quad \forall i & \end{cases} \quad \sum_j P_{ij} = 1$$

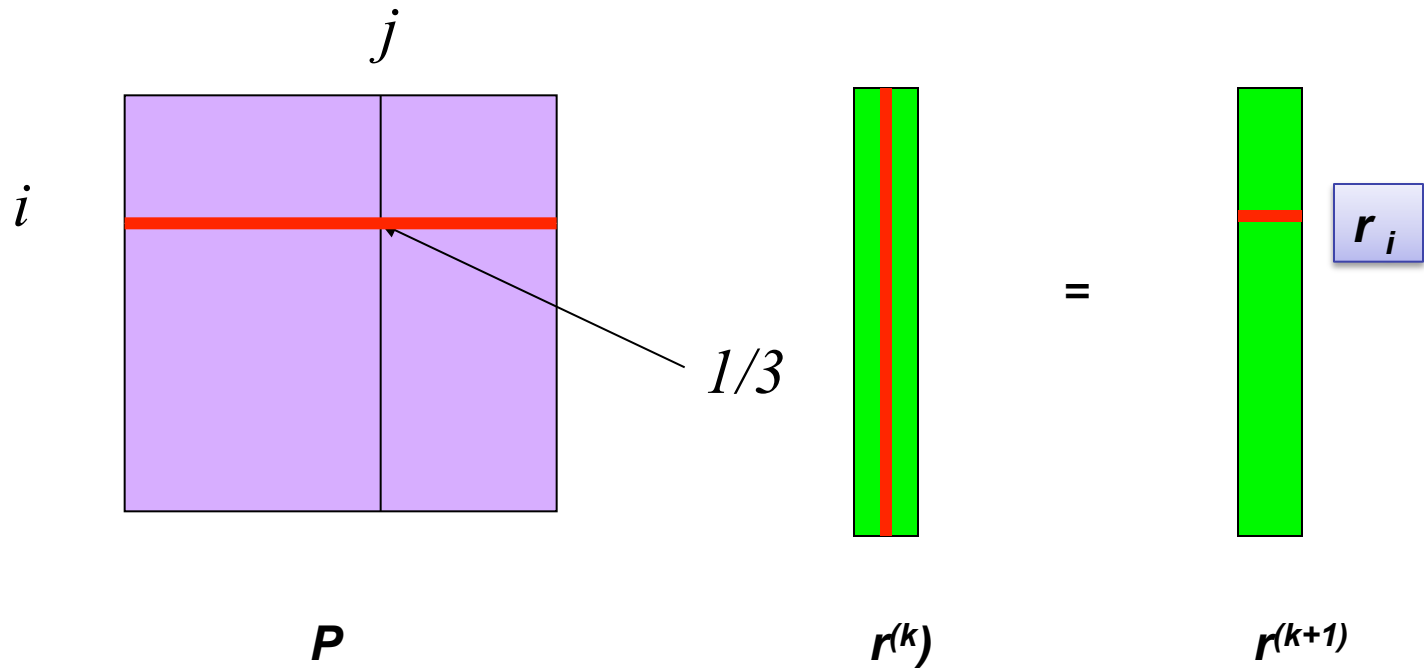
- The PageRank formulation can be written

as:

$$r^{(k)} = P \cdot r^{(k-1)} \quad P \text{ is a left stochastic matrix, as we anticipated}$$

Example

Suppose page j links to 3 pages, including i



The new value of r_i (the page rank of node i) at iteration k is obtained by summing the r_j of all pages pointing to i , multiplied by their p_{ij} value: $r_i^{(k+1)} = \sum_j p_{ij} r_j^{(k)}$ where p_{ij} is the (uniform) probability of jumping from j to i .

How to compute the vector \mathbf{r} of page ranks?

- The random surfer (or random walks) model can be represented using **Markov Chains**

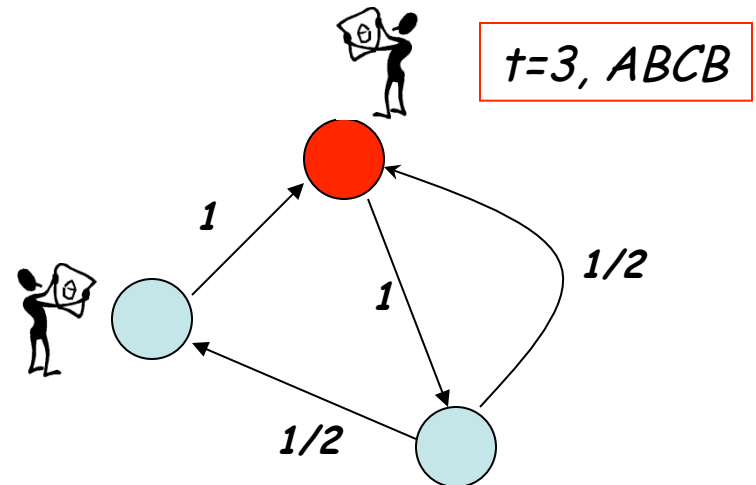
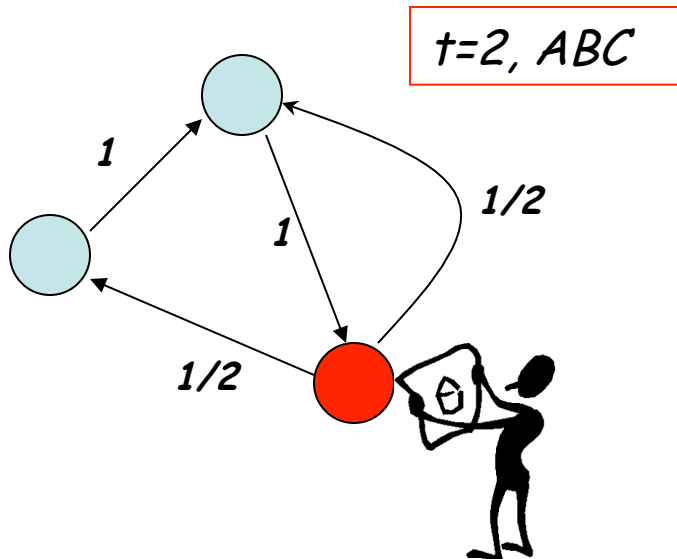
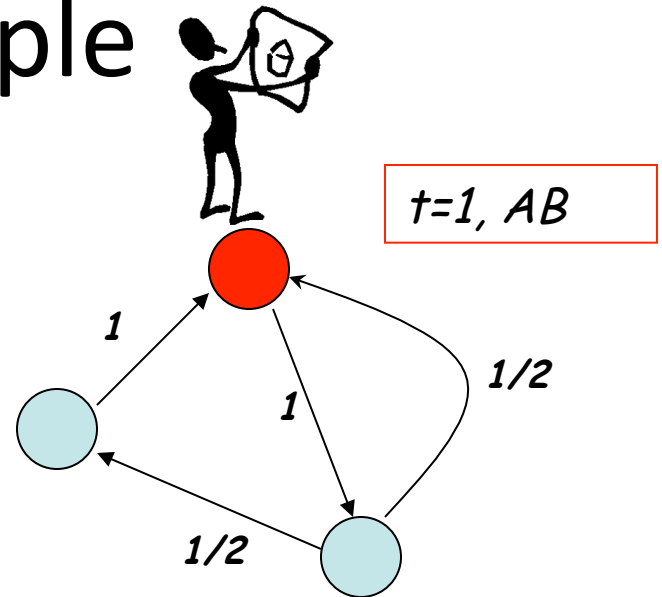
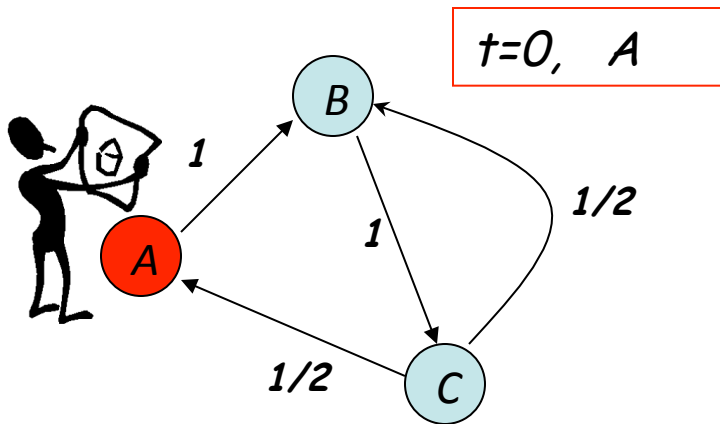
Markov Chains (1)

- A Markov Chain consists in N states (let S the set of possible states), and a matrix of transition probabilities $N \times N$, \mathbf{P} .
- At each step, the system is precisely in one state (states are web pages, in our case).
- For $1 \leq i, j \leq N$, $P(s_i \rightarrow s_j) = P_{ij}$ is the probability of jumping to s_j , given we are in s_i .
- Furthermore, if X_k is the random variable indicating the state \mathbf{s} reached at time t_k (X gets values in S), then:

$$P(X_k / X_1, X_2, \dots, X_{k-1}) = P(X_k / X_{k-1})$$

- The value of X_k at time k **depends only on the value of the random variable at time $k-1$!** (This is the basic property of Markov Chains: 1 state memory!)
- Which means: the probability of being on page s_j at time k only depends on the page s_i on which the surfer was at time $k-1$

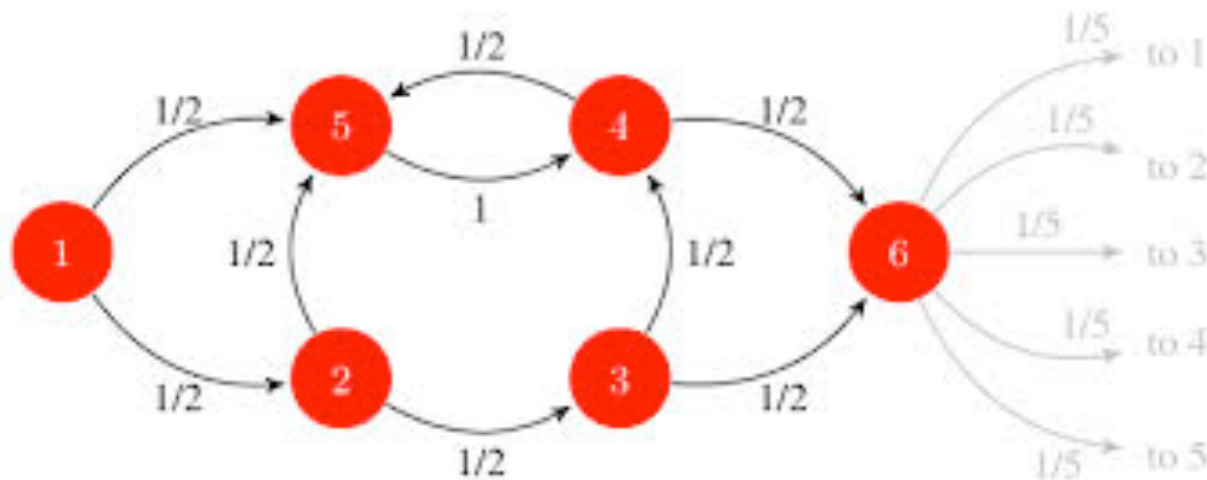
Back to previous example



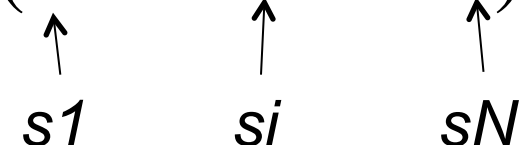
$P(X_4=B/X_1=A, X_2=B, X_3=C) = P(X_4=B/X_3=C)$ in a Markov Model!

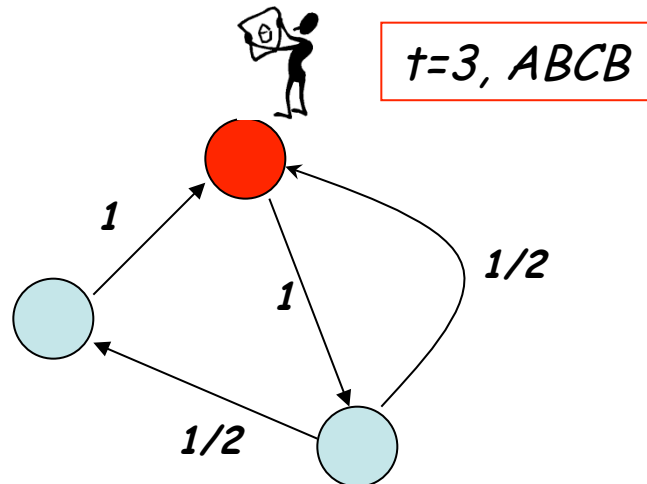
Markov chains (2)

- We also have that: $\sum_{j=1}^n P_{ij} = 1$.
- **Markov Chains are a model of random walks!**



Probability Vectors

- Let $\mathbf{x}^{(t)} = (x_1, \dots, x_N)$ be an S -dimensional vector indicating the state reached at time t
- Ex: $(000\dots 1\dots 000)$ means we are in state s_i .




Example:

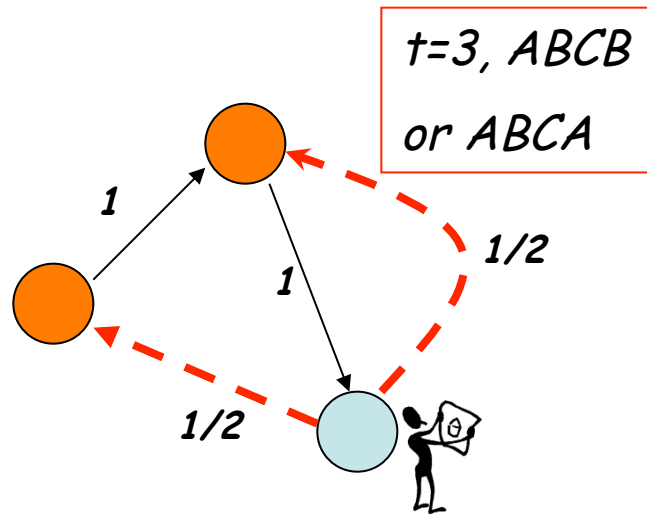
$X^{(3)}=(0,1,0)$ since the walker in $t=3$ jumps to state B

Probability Vectors

Since **we are modeling a stochastic process**, we can define a vector of probabilities $\mathbf{r}^{(t)} = (P(t, s_1), \dots, P(t, s_N)) = (r_1^t, \dots, r_N^t)$, indicating that at step t the random walk will bring to state s_i with probability r_i^t ,

Example:

$$\sum_{i=1}^N P(t, s_i) = \sum_{i=1}^N r_i^t = 1.$$



$X^{(3)} = (0.5, 0.5, 0)$ since the walker in $t=2$ has a 0.5 probability of jumping To B and 0.5 of jumping to A in $t=3$

Difference between P_{ij} and $P(t,s)$

- P_{ij} is the probability of jumping in state (page) j when we are in state (page) i : those probabilities *are known and uniform* (all equal to $1/k$ for k outlinks) for any starting node i
- $P(t,s_i)=r_i^t$ is the probability of being in state i when starting in some initial state and after t jumps δ (=at time t).
- We wish to compute **stationary values r_i** for r_i^t !! These **are the Page Ranks**: *the asymptotic probabilities of being in a given page for a walker who start in a page at random and travels the web graph at random.*

Ergodic Markov Chains

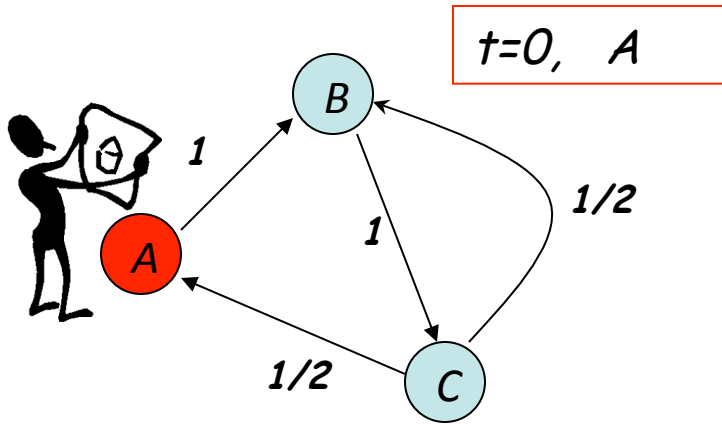
- The Random Walk is modeled with Markov Chains. Stationary probabilities can be computed if the process is ERGODIC
- A Markov Chain is **ergodic** if:
 - There is a path between any pair of states (= adjacency matrix is **irreducible**, \rightarrow the graph is connected)
 - Starting from any state, after a finite transition time T_0 , the probability to reach any other state in a finite time $T > T_0$ is always different from zero.
 - Note: **not true for the web graph!** Since not connected. Will see how to cope with this

Ergodic Chains

- If a Markov Chain is ergodic, every state has a **stationary probability** of being visited, regardless of the initial state of the random walker .
 - The vector $\mathbf{r}^{(t)}$ of state probabilities converges to a **stationary vector** \mathbf{r} as $t \rightarrow \infty$

Computing State Probability Vector

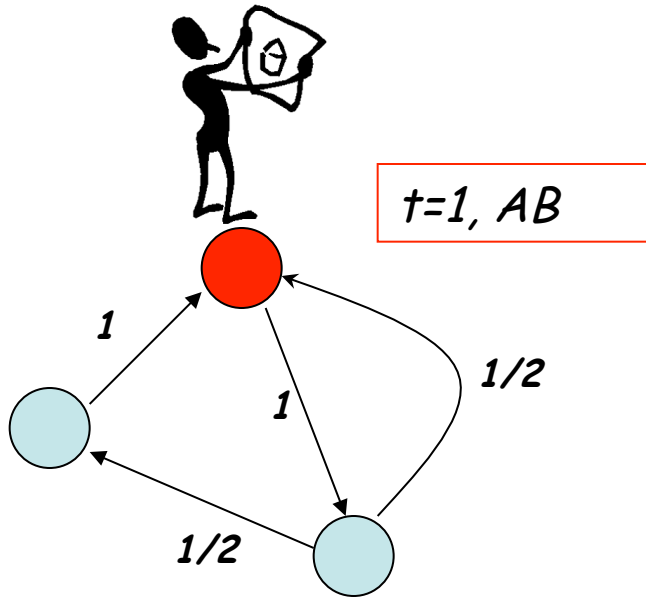
- If $\mathbf{r}^{(t)} = (r_1, \dots, r_N)$ is the probability vector in step t , how would it change after the next jump?
- The adjacency matrix \mathbf{P} tells us where we are likely to jump from any state (*since it has all transition probabilities from s_i to the other linked states*):
- Therefore, from $\mathbf{r}^{(t)}$, the probability of next state $\mathbf{r}^{(t+1)}$ is computed according to: $(\mathbf{r}^{(t+1)} = \mathbf{P}\mathbf{r}^{(t)})$
- Even under the random walk model, *we obtain again our iterative formulation!*



$$(100) \times \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

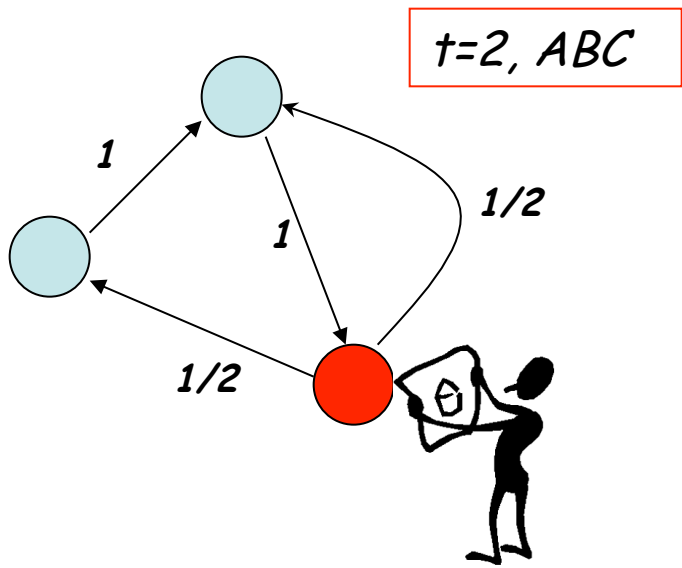
↑
 $r^{(0)}$

↑
 $r^{(1)}$

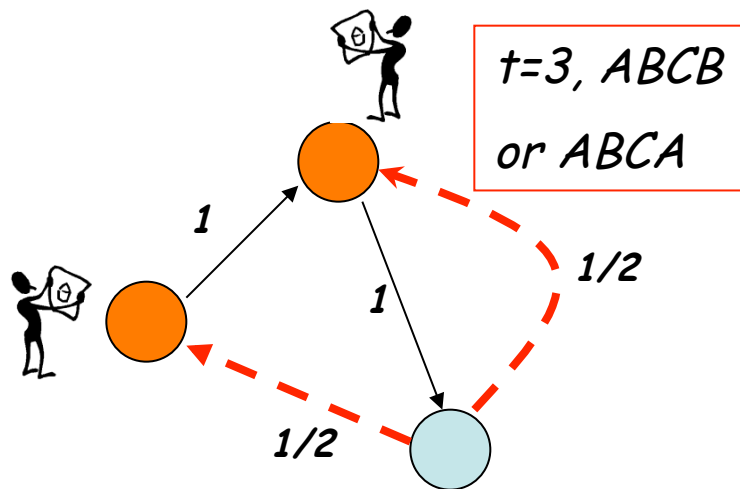


$$(010) \times \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

\uparrow $r(1)$
 \uparrow $r(2)$



$$\begin{matrix}
 \uparrow \\
 r^{(2)}
 \end{matrix}
 (001) \times
 \begin{pmatrix}
 0 & 1 & 0 \\
 0 & 0 & 1 \\
 \frac{1}{2} & \frac{1}{2} & 0
 \end{pmatrix}
 =
 \begin{pmatrix}
 1 \\
 \frac{1}{2} \\
 1 \\
 \frac{1}{2} \\
 0
 \end{pmatrix}
 \begin{matrix}
 \uparrow \\
 r^{(3)}
 \end{matrix}$$



$$\begin{matrix}
 \begin{matrix} \uparrow \\ r^{(3)} \end{matrix} & \begin{pmatrix} 1 & 1 \\ 2 & 2 \end{pmatrix} 0 & \times & \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix} & = & \begin{matrix} \begin{matrix} \uparrow \\ r^{(4)} \end{matrix} & \begin{pmatrix} 0 \\ 1 \\ 2 \\ 1 \\ 1 \\ 2 \end{pmatrix}
 \end{matrix}$$

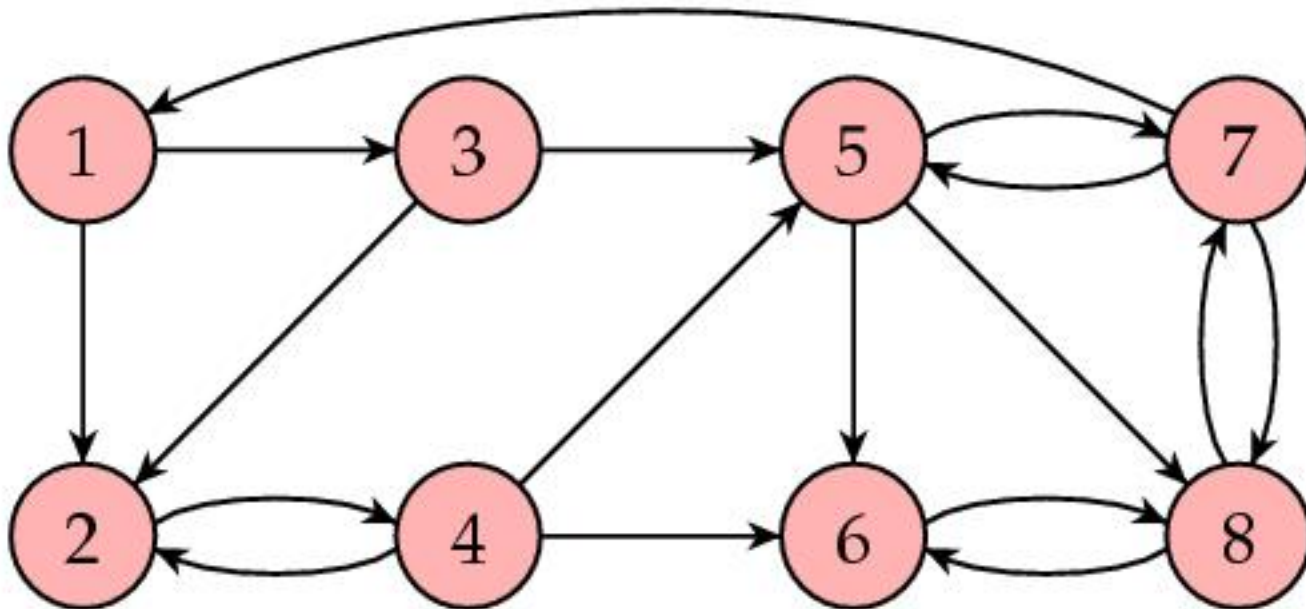
Computing Stationary Probability Vector

- If the process is ergodic, \mathbf{r}^t will eventually converge to a vector r such that $r=Pr$
- **Since P is a matrix and r is a vector, what kind of vector is r ?? In other terms, if $r=Pr$ holds, what vector is r ??**

Again: the Power method!

- $\underline{\mathbf{r}}^{(k+1)} = \mathbf{P}\underline{\mathbf{r}}^{(k)}$
- The sequence of vectors \mathbf{r}^k converge to the stationary vector \mathbf{r} if \mathbf{P} is stochastic and irreducible
- To compute \mathbf{r} we use the same method as for HITS
- $\underline{\mathbf{x}}^{(k+1)} = \mathbf{P}\underline{\mathbf{x}}^{(k)} = \mathbf{P}^k \underline{\mathbf{x}}^{(0)} = \lambda_1^k \left[\alpha_1 r_1 + \alpha_2 \left(\frac{\lambda_2}{\lambda_1}\right)^k r_2 + \dots + \alpha_n \left(\frac{\lambda_n}{\lambda_1}\right)^k r_n \right]$
- The method converges provided there is a dominant (principal) eigenvector
- Since the stationary condition is: $\mathbf{r} = \mathbf{P}\mathbf{r}$, \mathbf{r} is the principal eigenvector r_1 of \mathbf{P} and $\lambda_1 = 1$
- The principal eigenvalue is 1 (because \mathbf{P} is stochastic)

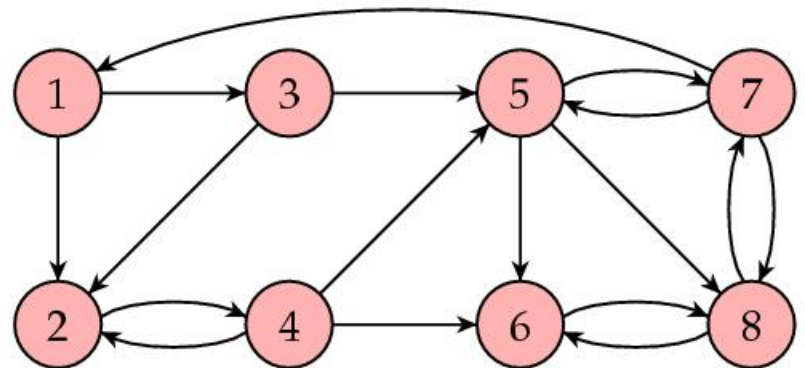
Example



The normalized adjacency matrix P

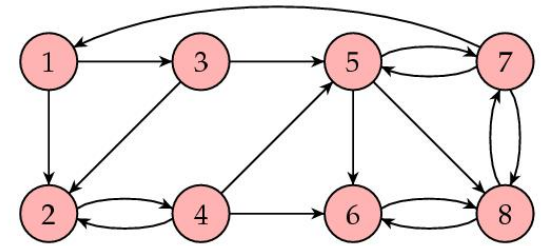
$$P = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 1/3 & 0 \\ 1/2 & 0 & 1/2 & 1/3 & 0 & 0 & 0 & 0 \\ 1/2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1/2 & 1/3 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 1/3 & 1/3 & 0 & 0 & 1/2 \\ 0 & 0 & 0 & 0 & 1/3 & 0 & 0 & 1/2 \\ 0 & 0 & 0 & 0 & 1/3 & 1 & 1/3 & 0 \end{bmatrix}$$

$1/N(u_i)$



Iterations

$$\underline{x}^{(k+1)} = P \underline{x}^{(k)}$$



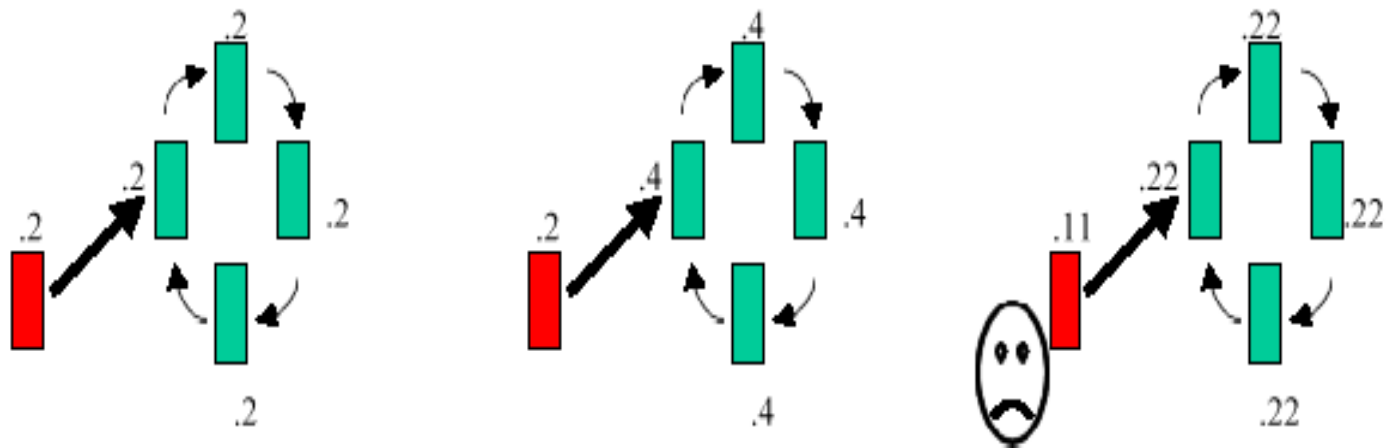
x^0	x^1	x^2	x^3	x^4	x^{60}	x^{611}	
1	0	0	0	0.0278	...	0.06	0.06
0	0.5	0.25	0.1667	0.0833	...	0.0675	0.0675
0	0.5	0	0	0	...	0.03	0.03
0	0	0.5	0.25	0.1667	...	0.0675	0.0675
0	0	0.25	0.1667	0.1111	...	0.0975	0.0975
0	0	0	0.25	0.1806	...	0.2025	0.2025
0	0	0	0.0833	0.0972	...	0.18	0.18
0	0	0	0.0833	0.3333	...	0.295	0.295

Problem with our PageRank formulation

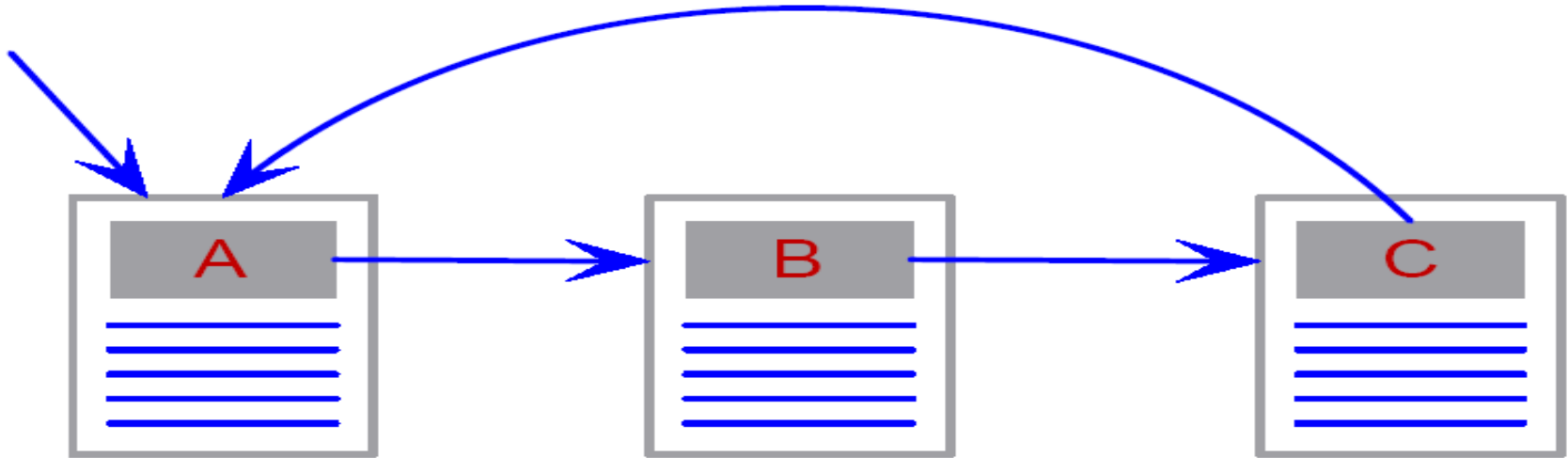
- Markov process converges under condition of ergodicity. Iterative computation (power method) converges if matrix P is irreducible. Two ways of saying the same thing!
 - As we said, these conditions are satisfied if the graph is fully connected and not deeply cyclic, **which is not the case** for the web graph
 - What causes the problem in practice?

Rank sink

If a group of web pages point to each other but to no other page, during the iteration, this loop will accumulate rank but never distribute any rank.



Teleporting

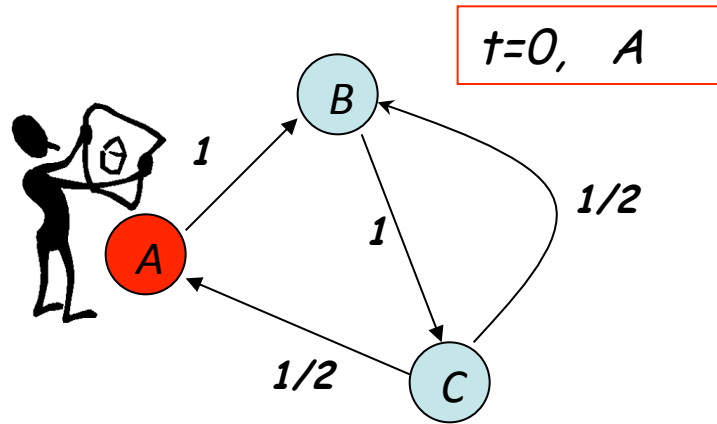


- Problem: Pages in a loop (or in a disconnected component) accumulate rank but do not distribute it to the rest of the graph.
- Solution: **Teleportation**, i.e. with a certain small probability the surfer can jump to any other web page (to which it is not connected) to get out of the loop.

Definition modified (with teleporting)

$$r(u) = c \sum_{v \in B_u} \frac{r(v)}{N_v} + (1 - c)E(u)$$

- *E(u) is some vector of probabilities over the set of web pages (for example uniform prob., favorite page etc.) that corresponds to a source of rank.*
- *c is called the **dumping factor** (also denoted with d)*
- *E(u) can be thought as if the random surfer “gets bored” periodically to travel from one page to another adjacent page, and “flies” to a different page (even though not connected), so he is not kept in a loop forever.*



NOTE: needs Rank-normalization step to get $\sum r_i = 1$

$$0.9 \times (100) \times \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix} + 0.1 \times \begin{pmatrix} 1 \\ 3 \\ \frac{1}{3} \\ 3 \\ \frac{1}{3} \\ 3 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 0.033 \\ 0.033 \\ 0.033 \end{pmatrix} = \begin{pmatrix} 0.033 \\ 1.033 \\ 0.033 \end{pmatrix}$$

$$c \sum_{v \in B_u} \frac{r(v)}{N_v} + (1-c)E(u) = r(u)$$

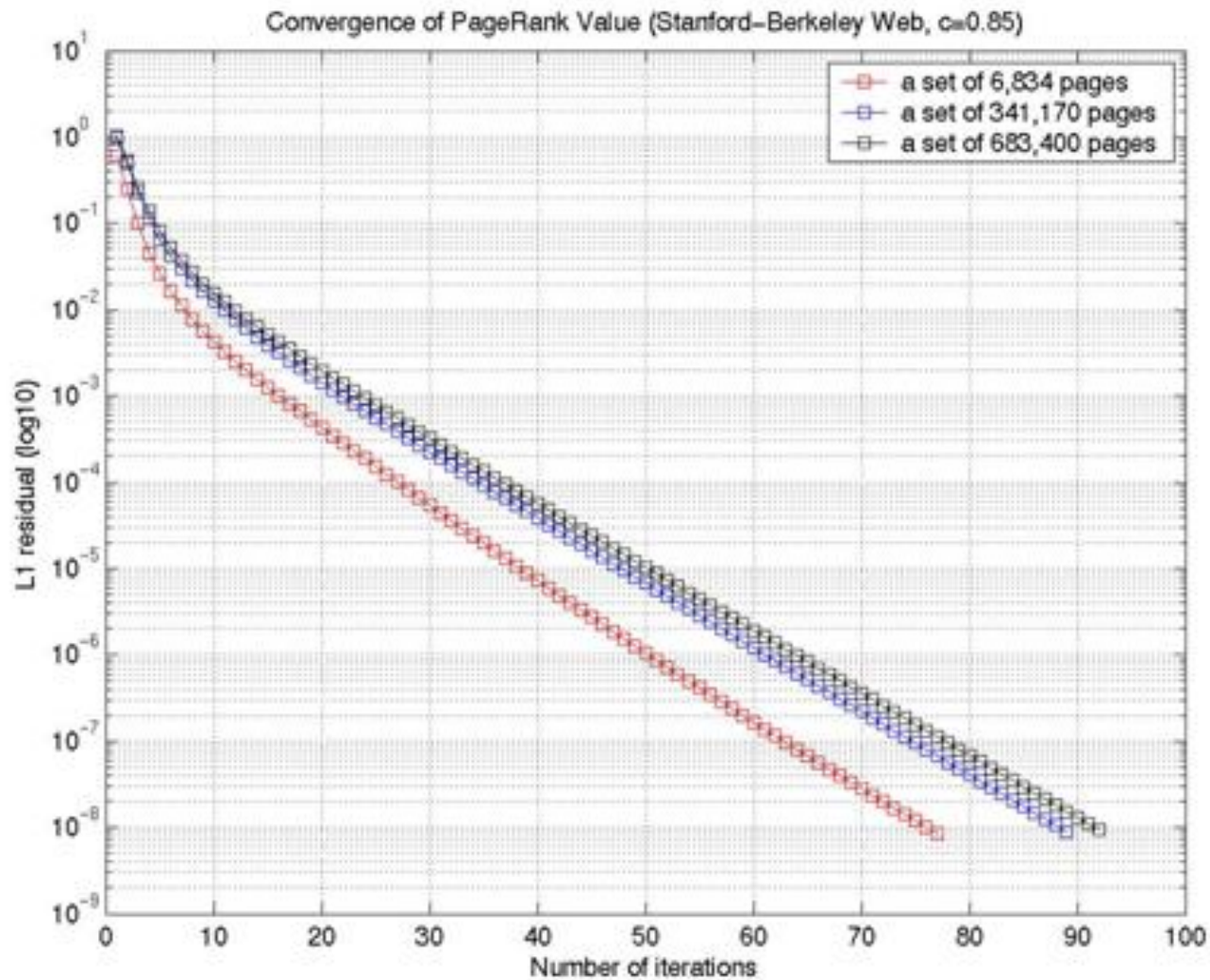


Figure 5: PageRank convergence as a function of the size of the web graph

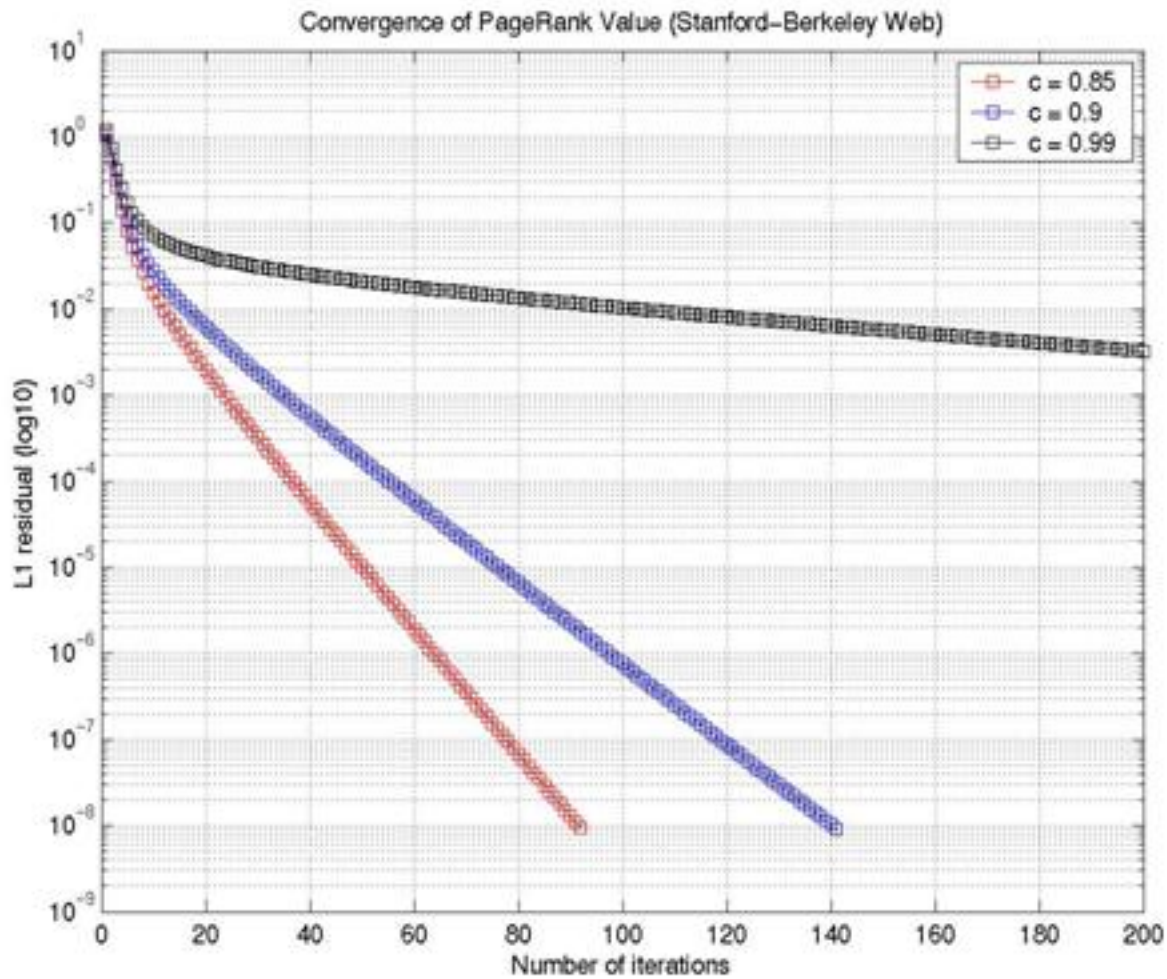


Figure 6: PageRank convergence as a function of C

Note: "c" is the dumping factor

Teleporting is a great “trick”

- This solves:
 - Sink problem
 - Disconnectedness of the web graph
 - Converges fast if *set of initial values* $r^{(0)}$ is chosen appropriately. In algebraic terms, the initial vector must have a non-zero component in the direction of the principal eigenvector (else it will never move in that direction)
- But we still have problems:
 1. Computing Page Rank for all web pages is computationally very intensive, plus needs frequent updates (web is dynamic)
 2. Does not take into account the specific query
 3. Easy to fool (less than HITS: less easy to be cited than cite!)

The Largest Matrix Computation in the World

- Computing PageRank can be done via matrix multiplication, where the matrix has billions rows and columns.
- The matrix is sparse as average number of outlinks is between 7 and 8.
- Setting $c = 0.85$ or above requires at most 100 iterations to convergence.
- Researchers still trying to speed-up the computation (“Big Data” problem, but you have a “Big Data” course).

Monte Carlo Methods in Computing PageRank

- Rather than following a single long random walk, the random surfer can follow many “sampled” random walks (threads).
- Each walk starts at a random page and either teleports with probability c or continues choosing a connected link with uniform probability.
- The PageRank of a page is the proportion of times a “sample” random walk ended at that page.

Another variant: Personalised PageRank

$$R(u) = c \sum_{v \in B_u} \frac{R(v)}{N_v} + (1 - c)v$$

- Change $cE(v)$ with $c\mathbf{v}$
- Instead of teleporting uniformly to any page we *bias* the jump to prefer some pages over others.
 - E.g. v is 1 for “your home page” and 0 otherwise.
 - E.g. v prefers the topics you are interested in.

Weblogs influence on PageRank

- A weblog (or blog) is a frequently updated web site on a particular topic, made up of entries in reverse chronological order.
- Blogs are a rich source of links, and therefore **their links influence PageRank.**
- Although there might be attempts to influence google rankings by unfair behaviour, this is severely punished by Google downgrading the page rank of the “bomber” – so that the page will never be shown to users

Summary

- Link analysis one of the main mechanisms to rank web pages (others are content-based methods and personalization)
- PageRank is the most well know and used, HITS is used but more in social networks (due to query-time computation delay)