

LC029 정보검색

2022 12 1

Chapter 21 : Link Analysis

Link Analysis

- Citation Analysis

- The number of citations is an indicative of authority.^가
- Quantify the influence of papers by analyzing the pattern of citations amongst them.

- Link Analysis for web search

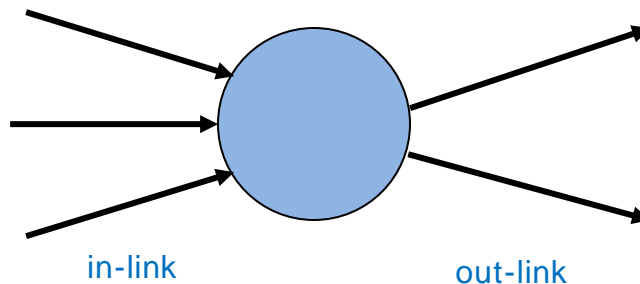
static

- Useful for ranking web search results.
- Useful to determine what page(s) to crawl next.

Link Analysis

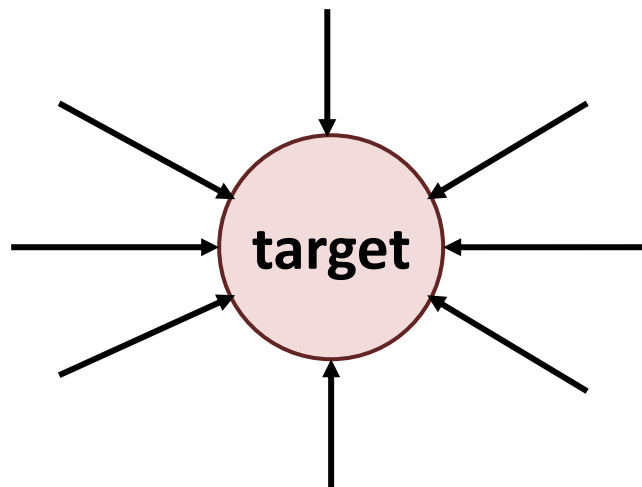
- Simple Approach

- Use link counts as simple measures of popularity.
- Undirected popularity: : in + out
 - score = the number of in-links and out-links ($3+2=5$).
- Directed popularity: : in-link
 - Score = number of in-links (3).



Link Analysis

- Query Processing
 - Retrieve all pages satisfying the query.
 - Order retrieved documents by their link popularity.
- Problem : Link Spam : in-link
 - Set up multiple web pages pointing to a target web page.



Web as a Graph

Web Graph

- Web documents have hyperlinks between them as a directed graph.

a.html

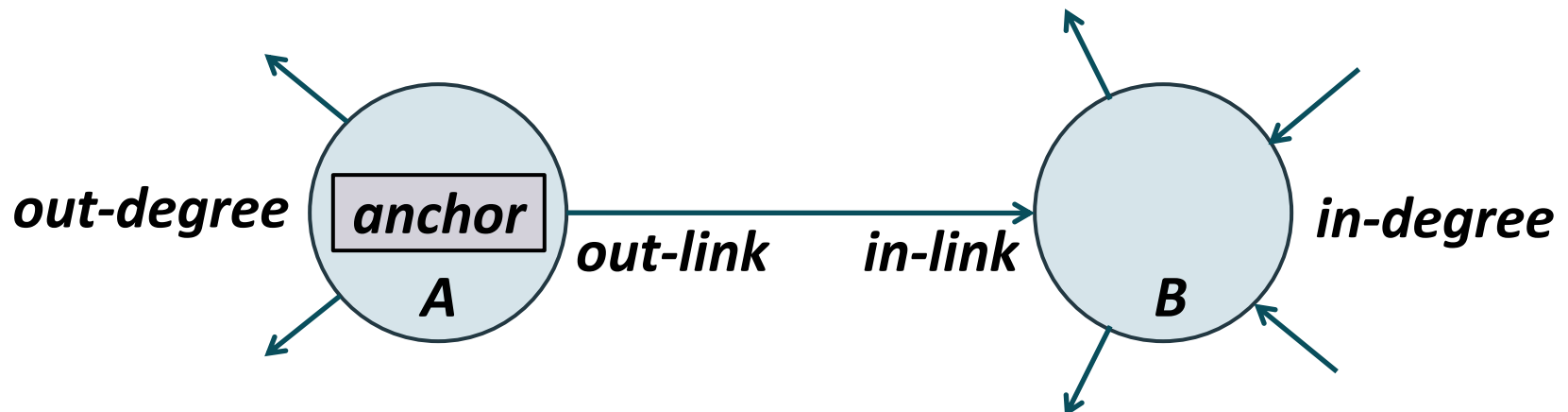
```
<A href="b.html">Samsung</A>
```

anchor text

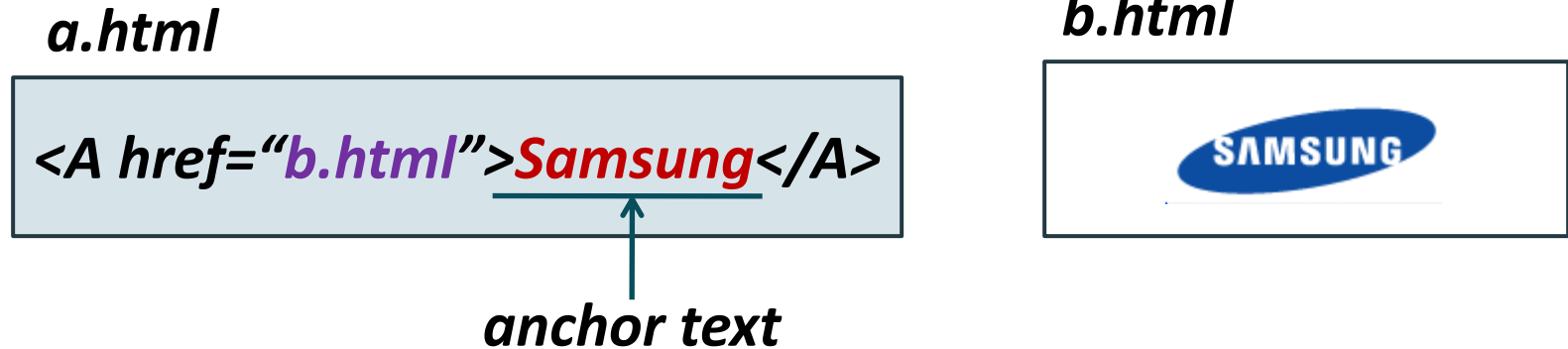
b.html



.....
..... Samsung
.....



The Web as a Directed Graph



Assumption 1: The **anchor text** of the hyperlink describes the target page (textual context).

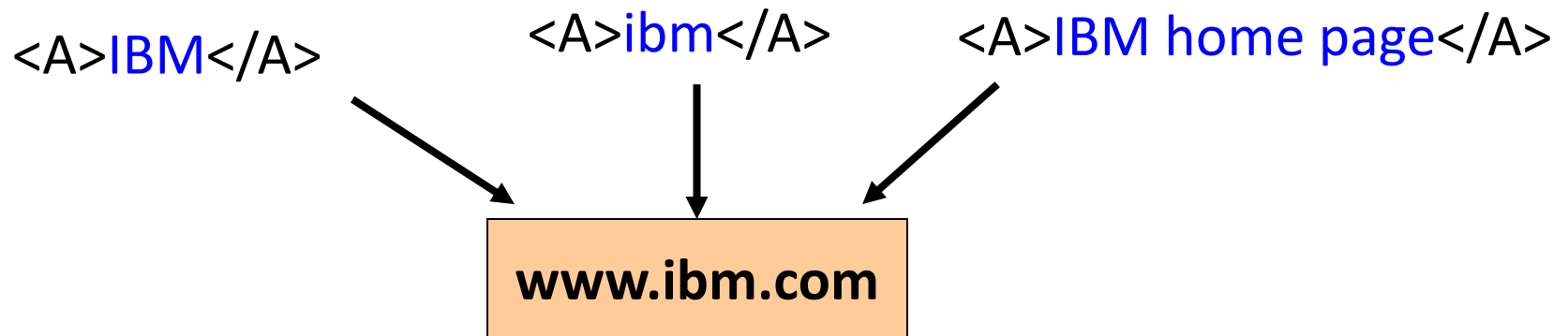
Assumption 2: A hyperlink between pages denotes author perceived relevance (quality signal).

가 1 :
가 2 :

가 () .

Anchor Text

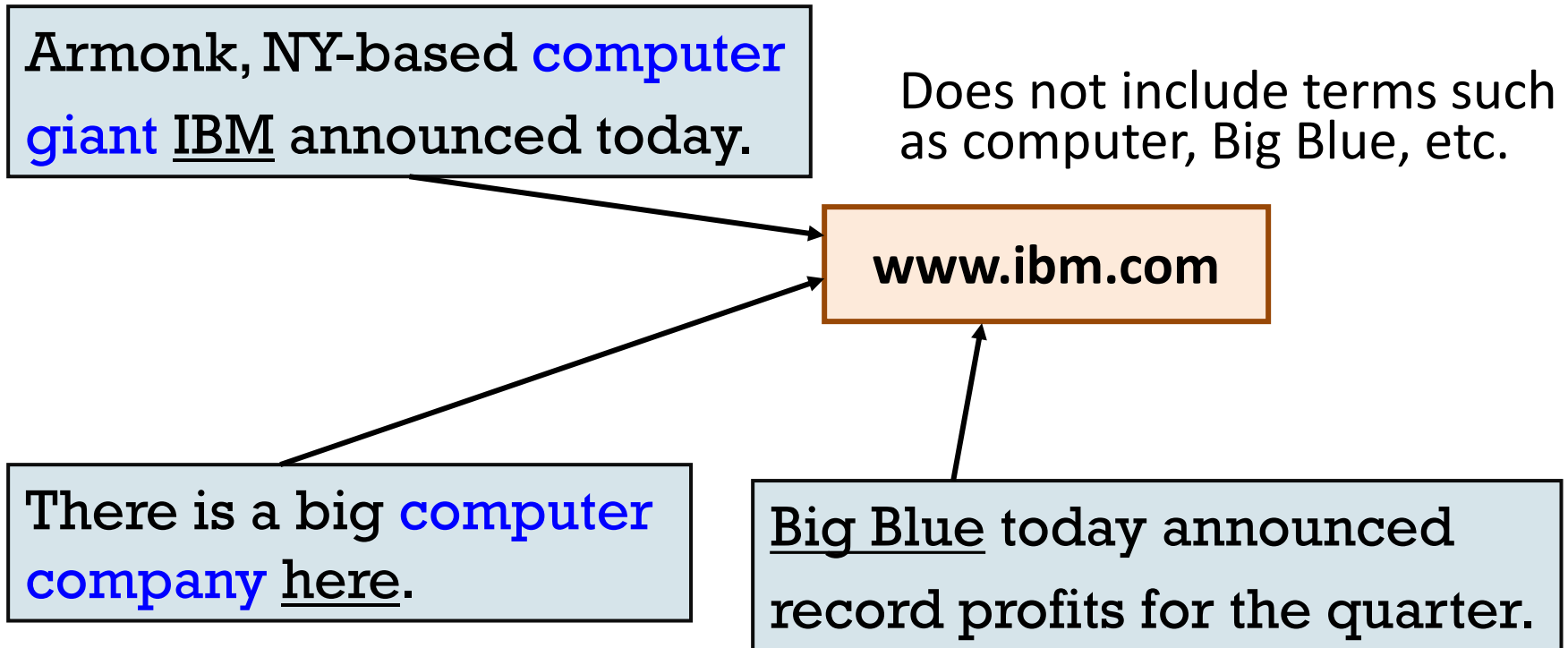
- A million pieces of anchor text with “IBM” send a strong signal.



- For **IBM** how to distinguish between:
 - IBM’s home page
 - Rival’s spam page

Anchor Text

- When indexing a document D , include anchor text or text surrounding anchor text.



Anchor Text

www.ibm.com



There is a big **computer company** here.

Armonk, NY-based **computer giant** IBM announced today.

Big Blue today announced record profits for the quarter.

term 1



...

www.ibm.com

...

term 2



...

www.ibm.com

...

term n



...

www.ibm.com

...

computer company



...

www.ibm.com

...

computer giant



...

www.ibm.com

...

big blue



...

www.ibm.com

...

가

Indexing Anchor Text

- Score anchor text with weight 가
 - Depending on the authority of the anchor page's website
e.g. if we were to assume that content from **cnn.com** or **yahoo.com** is authoritative, then trust the anchor text from them
 - Based on frequency with a penalty for terms that occur very often
e.g. penalty for terms such as **click, here, ...**

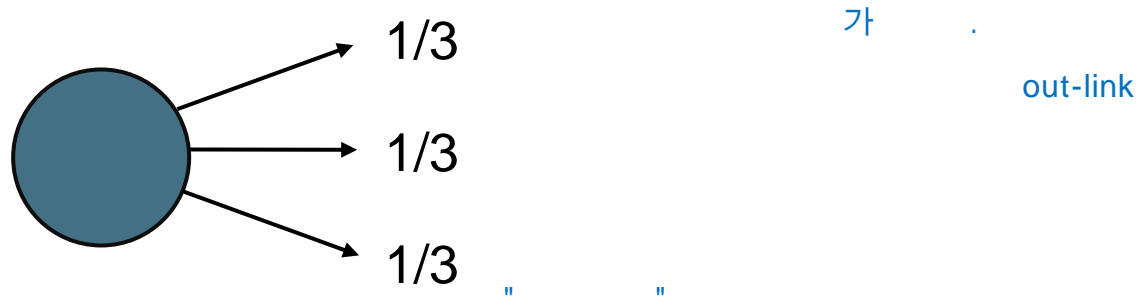
PageRank Scoring

PageRank Scoring

- PageRank
 - Score between 0 and 1, determined by link analysis.
 - The score of a page will depend on the link structure of the web graph.
 - The score is a query-independent measure of the static quality of each web page. static
 - The score is combined with other scores such as cosine similarity, term proximity, etc.
 - The composite score is used to provide a ranked list of results for the query.
 - PageRank is used in Google, but so are many other clever heuristics.

PageRank Scoring

- Imagine a surfer doing a random walk on web pages:
 - Start at a random page.
 - At each step, go out of the current page along one of the out-links on that page, equiprobably.



- “In the steady state”, each page has a long-term visit rate.
→ This rate is used as the **PageRank**.

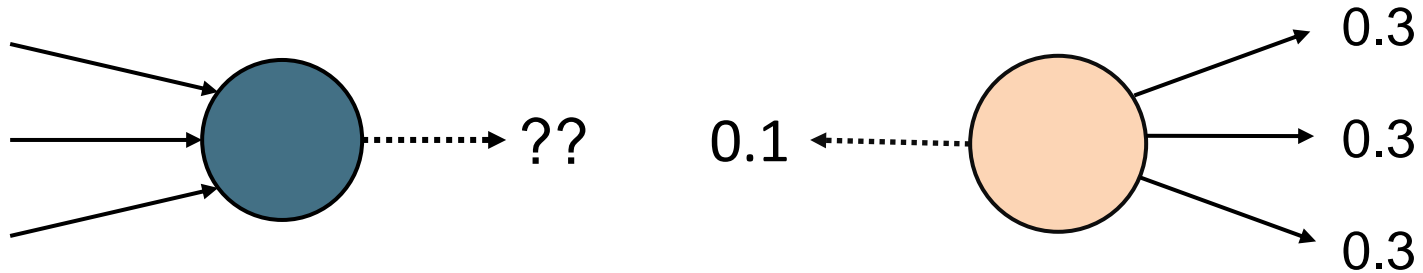
가

가

(가 가)

Teleport Operation

- The web is full of dead-ends. in out 가
 - Random walk can get stuck in dead-ends.

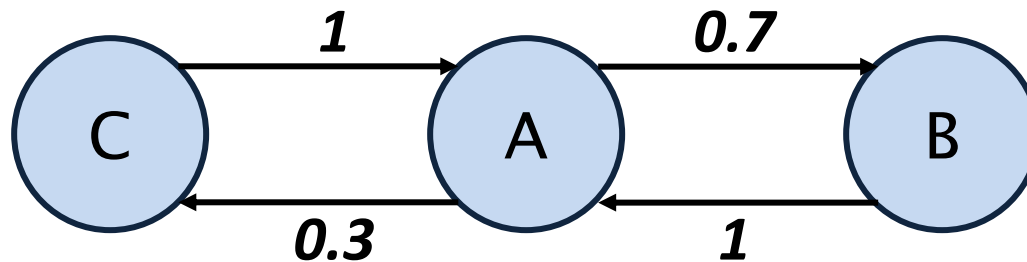


- At a **dead end**, jump to a random web page.
 - Now cannot get stuck locally.
- At any **non-dead end**, with probability α , jump to a random web page. out
 - With remaining probability $(1 - \alpha)$, go out on a random link.
 - Here α is a parameter.

Markov Chains

Markov Chain

- Consist of N **states**, with **transition probabilities** between the states.



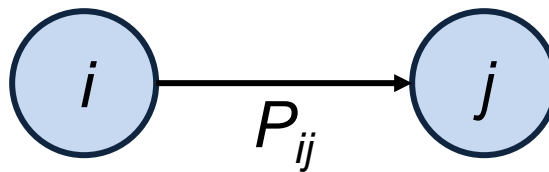
- The Markov chain can be in one of the N states at any given **time step**.
- The probability distribution of next state depends only on the current state, not on how the Markov chain arrived at the current state.

Markov Chain

- A Markov chain is characterized by an $N \times N$ transition probability matrix P , each of whose entries belongs to $[0..1]$.

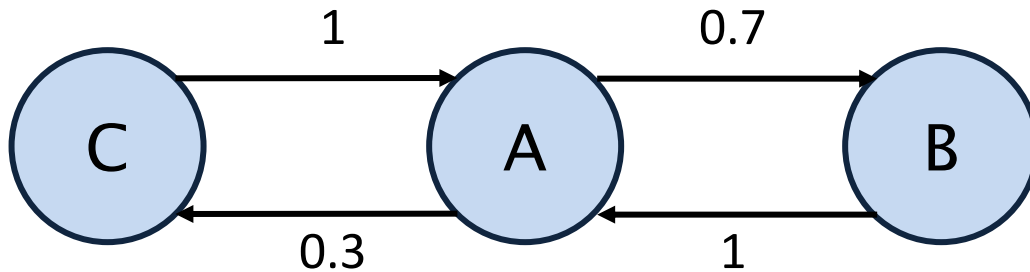
- For $1 \leq i, j \leq n$, the matrix entry P_{ij} tells us the probability of j being the next state, given we are currently in state i .

1 i, j n , P_{ij} 가 i j 가 .



Markov Chain

- Markov Chain



- Transition **Probability Matrix P**

	A	B	C
A	0	0.7	0.3
B	1	0	0
C	1	0	0

$= 1.0$

$$\forall i, \sum_{j=1}^N P_{ij} = 1$$

Probability Vector of state A
 $= (0, 0.7, 0.3)$

which tells us where we
go next from state A

Probability Vector

- A probability vector $\mathbf{x} = (x_1, \dots, x_N)$ tells us where we are at any point.

e.g. $(0, 0, 0, \dots, 1, \dots, 0, 0, 0)$ means we're in state i .

$\underset{1}{}, \dots, \underset{i}{1}, \dots, \underset{N}{0}, \dots$

- More generally, the probability vector $\mathbf{x} = (x_1, \dots, x_N)$ means that we are in state i with probability x_i .

e.g. $(0.1, 0.3, 0.2, 0, 0.1, 0.3, 0)$

Note that $\sum_{i=1}^n x_i = 1$.

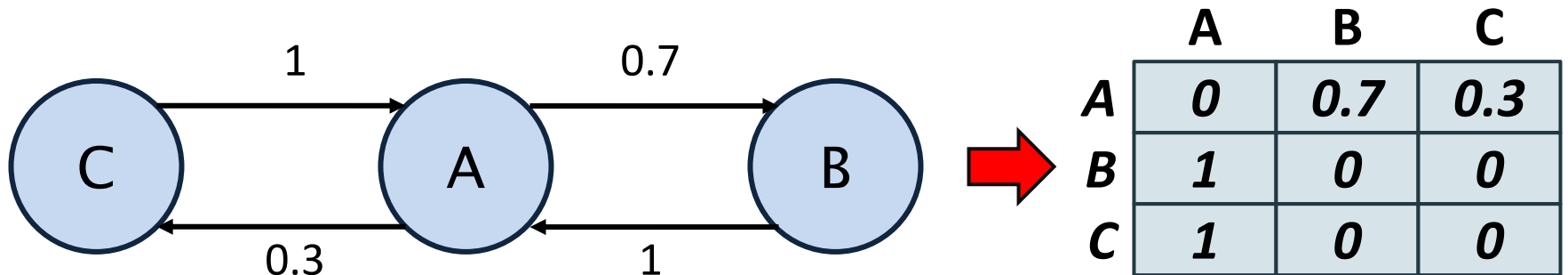
1

Change in Probability Vector

- If the probability vector is $\mathbf{x} = (x_1, \dots, x_N)$ at this step, what is it at the next step?
- Recall that row i of the transition probability matrix \mathbf{P} tells us where we go next from state i .
- So, from \mathbf{x} , our next state is distributed as \mathbf{xP} .

가 \mathbf{x}

Change in Probability Vector



transition
probability
matrix P

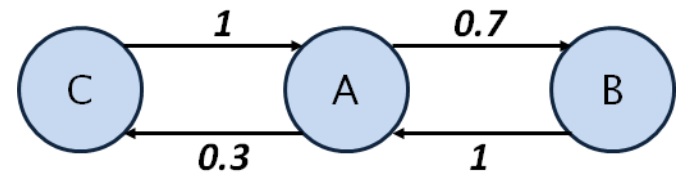
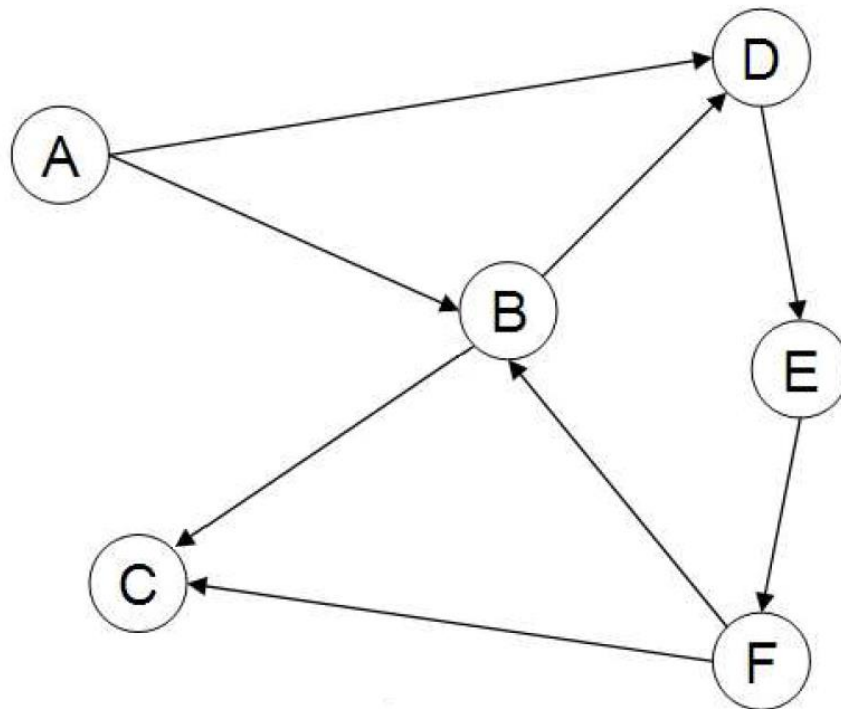
- Let $\mathbf{x} = (1, 0, 0)$. (We are in the state A.)
- Our next state is distributed as $\mathbf{x}P$.

$$(1, 0, 0) \begin{pmatrix} 0 & 0.7 & 0.3 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} = (0, 0.7, 0.3)$$

가 \mathbf{x}

Markov Chain

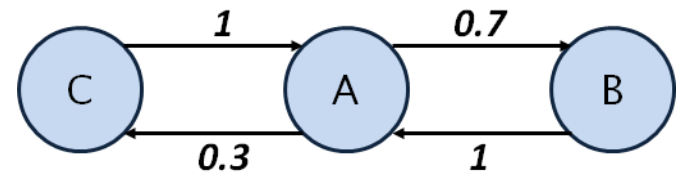
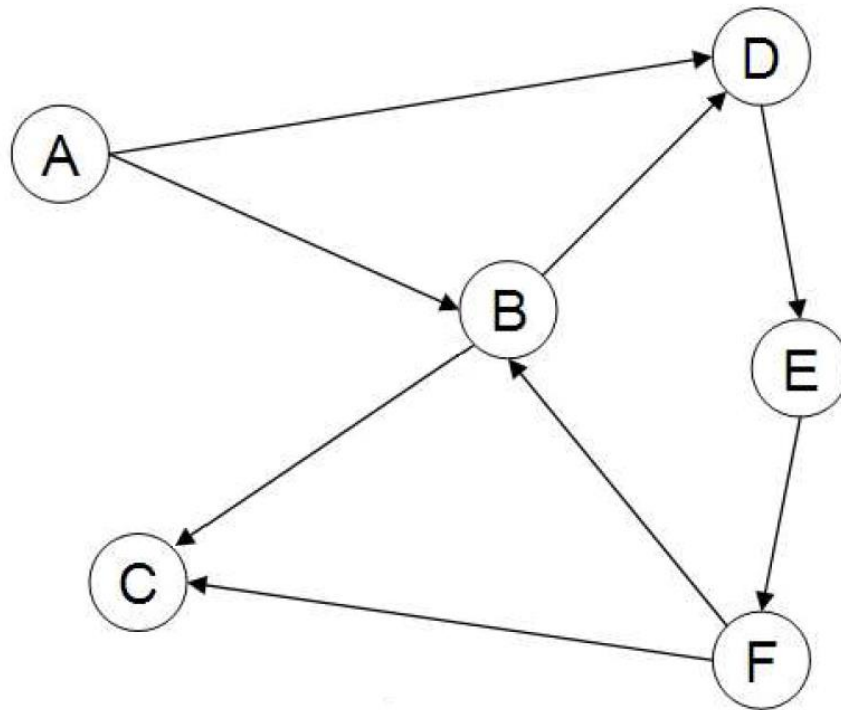
- We will use the Markov chain to represent the web.
 - Each web page corresponds to a state in the Markov chain.



← **Web Graph**

Markov Chain

- We will use the Markov chain to represent the web.
 - If we start at page i , where are we at next step?



	A	B	C
A	0	0.7	0.3
B	1	0	0
C	1	0	0

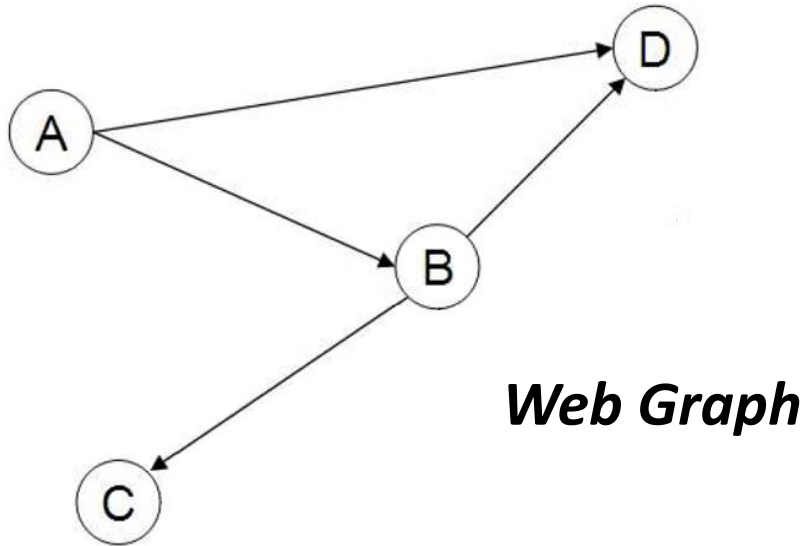
transition probability
matrix P

← **Web Graph**

Derivation of P

- Let A be an adjacency matrix of the web graph:
 $A_{ij} = 1$ if there is a hyperlink from page i to page j .
 $A_{ij} = 0$ otherwise.
- Derivation of the transition probability matrix P from A :
 - If a row of A has no 1's, replace each element by $1/N$.
 1 $1/N$ **Teleport operation**
 - Otherwise, 1
 1. Divide each 1 in A by the number of 1's in its row. 1 1
 2. Multiply the resulting matrix by $1 - \alpha$.**Random Surf**
 3. Add α/N to every entry of the resulting matrix.
 $1/N$ **Teleport operation**
- The resulting matrix is P .

Derivation of P



adjacency matrix

	A	B	C	D
A	0	1	0	1
B	0	0	1	1
C	0	0	0	0
D	0	0	0	0



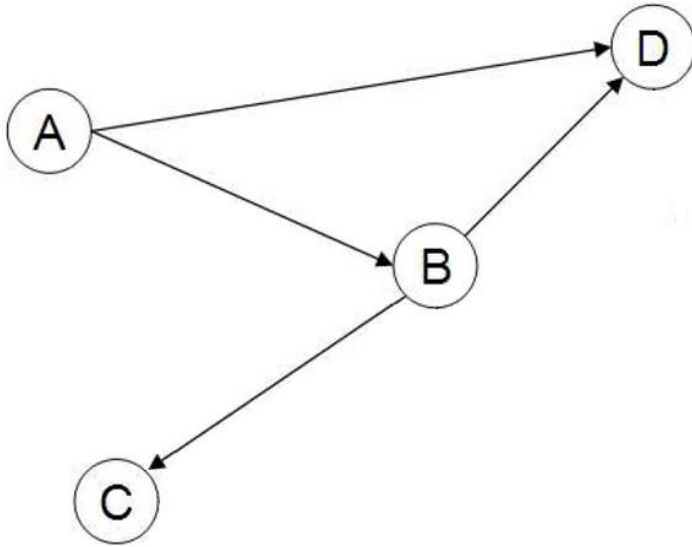
	A	B	C	D
A	.025	.475	.025	.475
B	.025	.025	.475	.475
C	.250	.250	.250	.250
D	.250	.250	.250	.250



	A	B	C	D
A	0	0.45	0	0.45
B	0	0	0.45	0.45
C	0.25	0.25	0.25	0.25
D	0.25	0.25	0.25	0.25

transition probability matrix

Derivation of P



	A	B	C	D
A	.025	.475	.025	.475
B	.025	.025	.475	.475
C	.250	.250	.250	.250
D	.250	.250	.250	.250

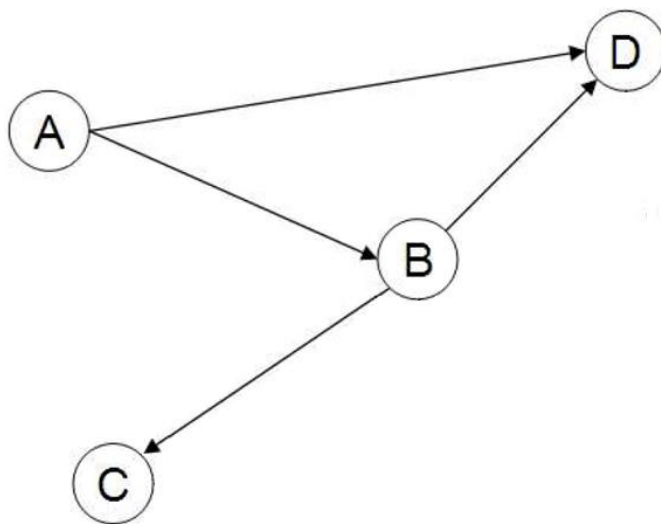
transition probability matrix

- With the given web graph, there is no path from C or D to any other node.
C, D
- The transition probability matrix, however, has transition probability > 0 from C or D to any other node.
(We applied the Teleport Operation.)

가 0

Ergodic Markov Chain

- A Markov chain is **ergodic** if
 - You have a path from any state to any other.
 - For any start state, after a finite transient time T_0 , the probability of being in any state at a fixed time $t (> T_0)$ is nonzero.



T_0
 $t (> T_0)$

	A	B	C	D
A	.025	.475	.025	.475
B	.025	.025	.475	.475
C	.250	.250	.250	.250
D	.250	.250	.250	.250

Ergodic Markov Chain

- For any ergodic Markov chain, there is a unique **long-term visit rate** for each state.
 - *Steady-state Probability Distribution.*
- Over a long time-period, we visit each state in proportion to this rate.
- It doesn't matter where we start.
- This *steady-state probability* for a state is the **PageRank** of the corresponding web page.

0

0

Steady State Probability Vector

- The **steady state probability vector** is a vector of probabilities $\mathbf{a} = (a_1, \dots, a_N)$ where a_i is the probability that we are in state i .
- If our current position is described by \mathbf{a} , then the next step is distributed as \mathbf{aP} where P is a Transition Probability Matrix.
- But \mathbf{a} is the steady state, so $\mathbf{a} = \mathbf{aP}$.
- Solving this matrix equation gives us \mathbf{a} .

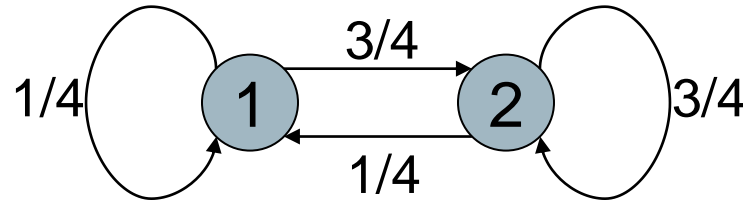
$$\mathbf{a} = \mathbf{aP}$$

(=a) a가 가

a

How do we compute this vector?

- Consider the Markov chain below.



	1	2
1	1/4	3/4
2	1/4	3/4

← *Transition
Probability
Matrix P*

- Let $\mathbf{a} = (x, y)$ where $x + y = 1.0$
 - Since $\mathbf{a} = \mathbf{aP}$, $(x, y) = (x/4 + y/4, 3x/4 + 3y/4)$.
- From (1) and (2), $x = 1/4$, $y = 3/4$.
 - Therefore, $\mathbf{a} = (1/4, 3/4)$.

How do we compute this vector?

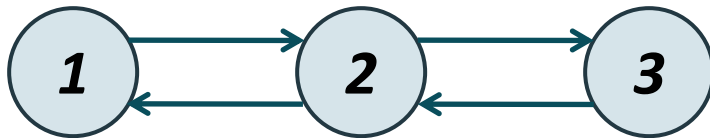
■ Power Iteration Method

- Recall, regardless of where we start, we eventually reach the steady state \mathbf{a} .
- Start with any distribution, say $\mathbf{x} = (1, 0, \dots, 0)$.
- After one step, we're at \mathbf{xP} .
- After two steps at \mathbf{xP}^2 , then \mathbf{xP}^3 and so on.
- For large k , we are at \mathbf{xP}^k where $\mathbf{xP}^k \rightarrow \mathbf{a}$.
- So, we can compute the steady-state probability vector with the multiplication of \mathbf{x} by increasing powers of \mathbf{P} until the product looks stable.



Computation of the Vector : Example

- Consider the web graph.



- Assume $\alpha = 0.5$, then $\alpha/N = 1/6$.

adjacency matrix

	1	2	3
1	0	1	0
2	1	0	1
3	0	1	0



	1	2	3
1	0	1	0
2	1/2	0	1/2
3	0	1	0



P →

	1	2	3
1	1/6	2/3	1/6
2	5/12	1/6	5/12
3	1/6	2/3	1/6

transition probability matrix

Computation of the Vector : Example

- Assume that the surfer starts in state 1, then the initial probability distribution vector $\mathbf{x}_0 = (1, 0, 0)$.
- After one step, the distribution is

$$\mathbf{x}_0 \mathbf{P} = (1, 0, 0) \begin{pmatrix} 1/6 & 2/3 & 1/6 \\ 5/12 & 1/6 & 5/12 \\ 1/6 & 2/3 & 1/6 \end{pmatrix} = (1/6, 2/3, 1/6) = \mathbf{x}_1$$

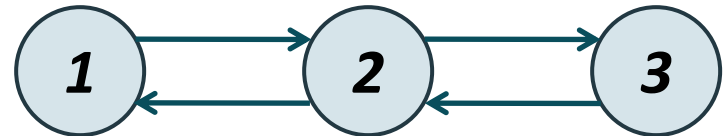
$$\mathbf{x}_1 \mathbf{P} = (1/3, 1/3, 1/3) = \mathbf{x}_2$$

$$\mathbf{x}_2 \mathbf{P} = (1/4, 1/2, 1/4) = \mathbf{x}_3$$

$$\mathbf{x}_3 \mathbf{P} = (7/24, 5/12, 7/24) = \mathbf{x}_4$$

...

$$\mathbf{x} = (5/18, 4/9, 5/18)$$



PageRank Summary

- Preprocessing:
 - Given graph of links, build a Transition Probability Matrix \mathbf{P} .
 - From \mathbf{P} , compute the Steady State Probability Vector \mathbf{a} .
 - The entry a_i is a number between 0 and 1, which is the PageRank of page i .
 - PageRank is a *query-independent* measure of the static quality of each page.
- Query processing:
 - Retrieve pages meeting query.
 - Rank them by their PageRank and other query-dependent scores such as cosine similarity, term proximity, etc.

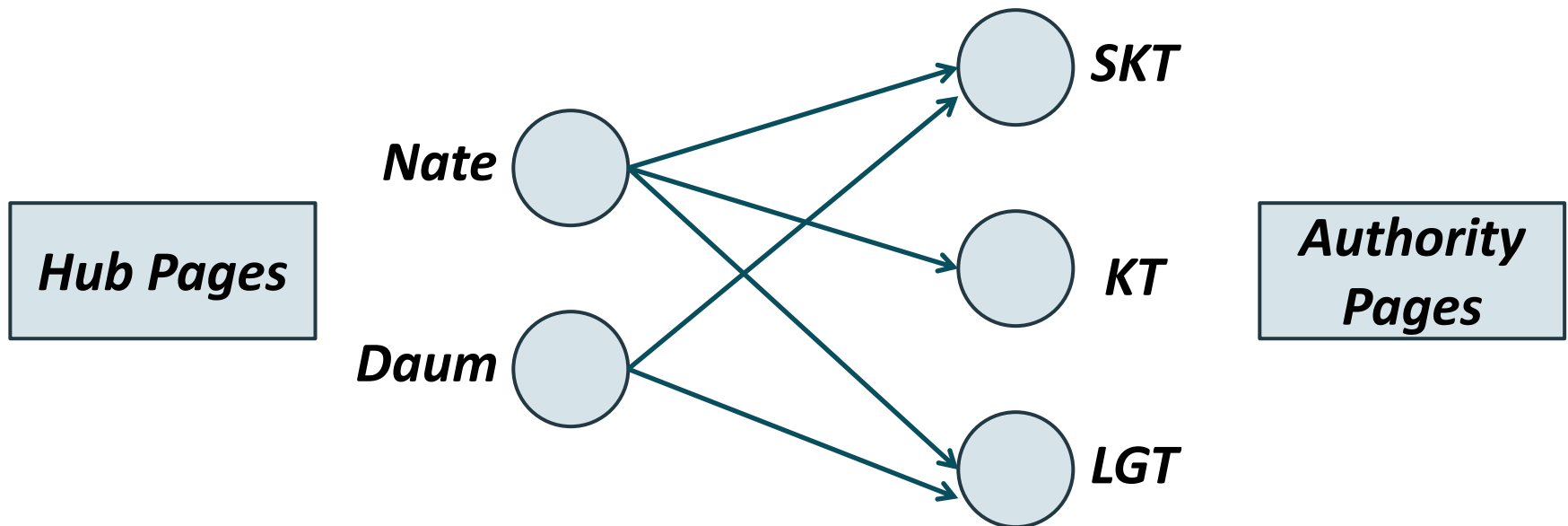
Hyperlink-Induced Topic Search

Hyperlink-Induced Topic Search (HITS)

- *In response to a query, we'd like to find two sets of inter-related pages (instead of an ordered list of pages that match the query)*
 - ***Hub pages***
that contain good lists of links on a subject.
 - ***Authority pages***
that occur recurrently on good hubs for the subject.

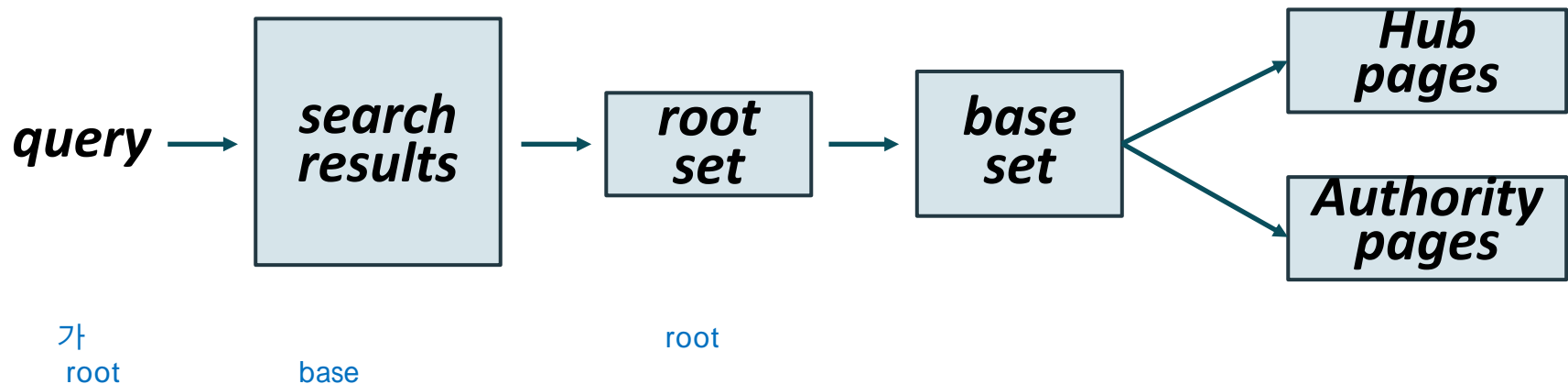
Hub Pages and Authority Pages

- A good hub page for a topic *points* to many authority pages for that topic.
- A good authority page for a topic is *pointed* to by many hub pages for that topic.



How to find Hub and Authority Pages?

- Extract from the web a **base set** of pages that *could* be good hub or authority pages.
- From these, identify a small set of hub pages and authority pages.
 - use iterative algorithm



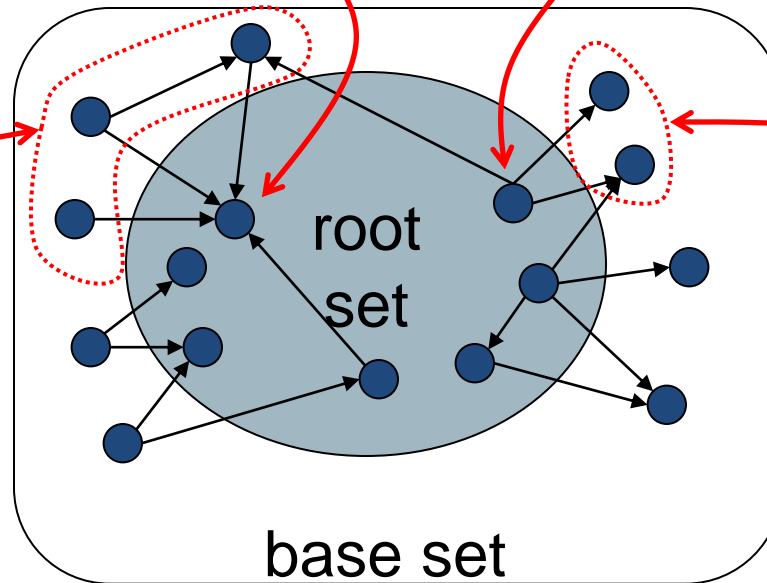
Base Set

- Given text query (say *swine flu*), retrieve pages containing *swine flu*. 가
 - Call them the **root set** of pages. root
- Add in any page that either 가
 - points to a page in the root set, or root 가
 - is pointed to by a page in the root set. root 가
- Call them the **base set**. 가 base
- Use the base set for **computing hub and authority scores**.

Base Set : Visualization

If this is good authority page

If this is good hub page



then these are good authority pages.

then these are good hub pages.

Assembling the Base Set

- **Root set** typically consists of 200-1000 pages rather than all pages matching the text query.
- **Base set** may have up to 5000 pages.
- How do you decide the base set?
 - Follow out-links by parsing root set pages (or from *connectivity server*)
 - Get in-links from a *connectivity server*

Distilling Hub and Authority Pages

- Compute, for each page x in the base set, a **hub score** $h(x)$ and an **authority score** $a(x)$.
 - Initialize: for all x , $h(x) \leftarrow 1$; $a(x) \leftarrow 1$;
 - **Iteratively update** all $h(x)$, $a(x)$;
- After iterations choose pages with
 - highest $h()$ scores as hub pages
 - highest $a()$ scores as authority pages

가 1

Iterative Update

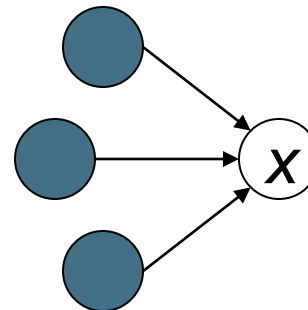
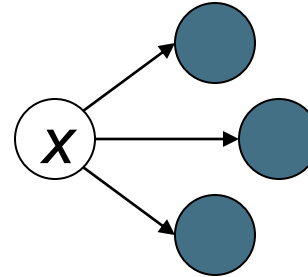
- Repeat the following updates, for all x :

$$h(x) \leftarrow \sum_{x \mapsto y} a(y)$$

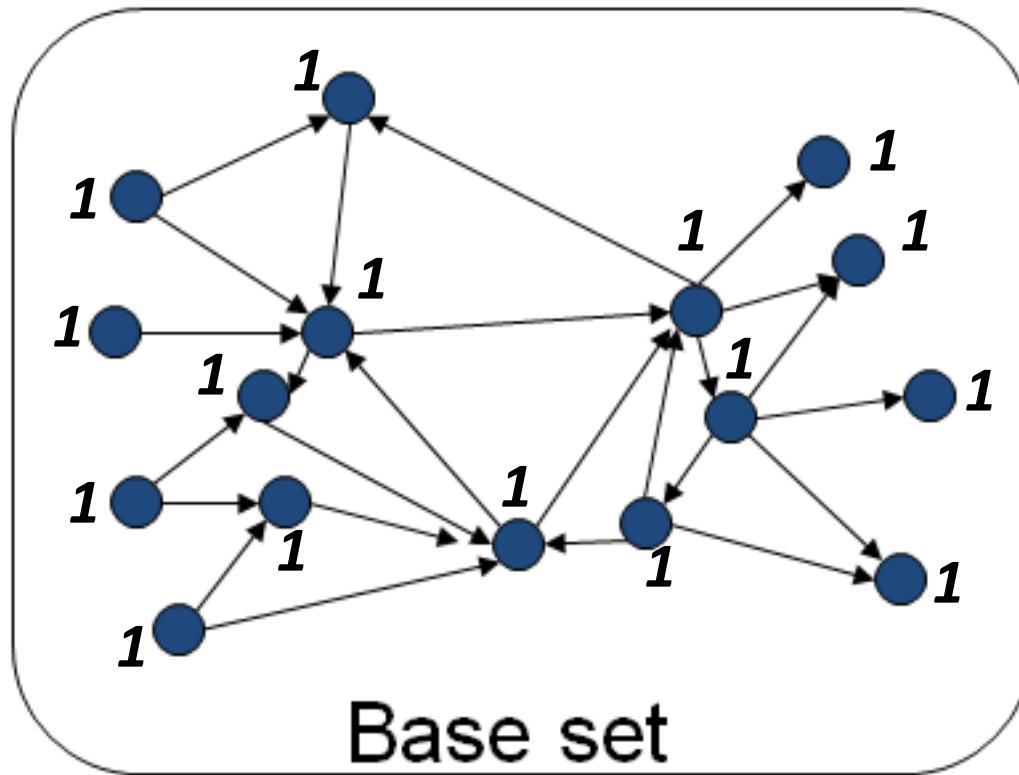
:

$$a(x) \leftarrow \sum_{y \mapsto x} h(y)$$

:



Iterative Update



Scaling Down of the Scores

- To prevent the $h()$ and $a()$ values from getting too big, scale down after each iteration.
- Scaling factor doesn't matter:
 - we only care about the *relative* values of the scores.

가

How many iterations?

- Claim
 - Relative values of scores will converge after a few iterations.
 - We only require the relative orders of the $h()$ and $a()$ scores, not their absolute values.
- Then, how many iterations?
 - In practice, at most 5 iterations will get you fairly good results.

Some interesting insights about HITS

- Frequently, the documents that emerge as top hubs and authorities include language other than the language of the query.
- This cross-language retrieval effect resulted from link analysis, with no linguistic translation.

가 .

Query: “Japan Elementary Schools”

- schools
- LINK Page-13 ← **Hub Pages**
- “ú—{,ìŠwZ
- ā%oo,,□ŠwZfz[f fy[fW
- 100 Schools Home Pages (English)
- K-12 from Japan 10/...rnet and Education)
- http://www...iglobe.ne.jp/~IKESAN
- ,l,f,j□ŠwZ,U”N,P’g•”œê
- ÒŠ—’¬—§ÒŠ—“œ□ŠwZ
- Koulutus ja oppilaitokset
- TOYODA HOMEPAGE
- Education
- Cay's Homepage(Japanese)
- —y“ì□ŠwZ,ìfz[f fy[fW
- UNIVERSITY
- %oJ—³□ŠwZ DRAGON97-TOP
- Â%oa□ŠwZ,T”N,P’gfz[f fy[fW
- ¶µ°é¼ÂÁ© ¥á¥Ë¥â¼ ¥á¥Ë¥â¼

- The American School in Japan
- The Link Page
- %oaēž—§^ä“c□ŠwZfz[f fy[fW
- Kids' Space
- ^Àéž—§^Àé¼•”□ŠwZ ← **Authority Pages**
- <{ēx³^ç’âŠw•®□ŠwZ
- KEIMEI GAKUEN Home Page (Japanese)
- Shiranuma Home Page
- fuzoku-es.fukui-u.ac.jp
- welcome to Miasa E&J school
- □“Pœ\$E%oi•ls—§’†¼□ŠwZ,ìfy
- http://www...p/~m_maru/index.html
- fukui haruyama-es HomePage
- Torisu primary school
- goo
- Yakumo Elementary,Hokkaido,Japan
- FUZOKU Home Page
- Kamishibun Elementary School...