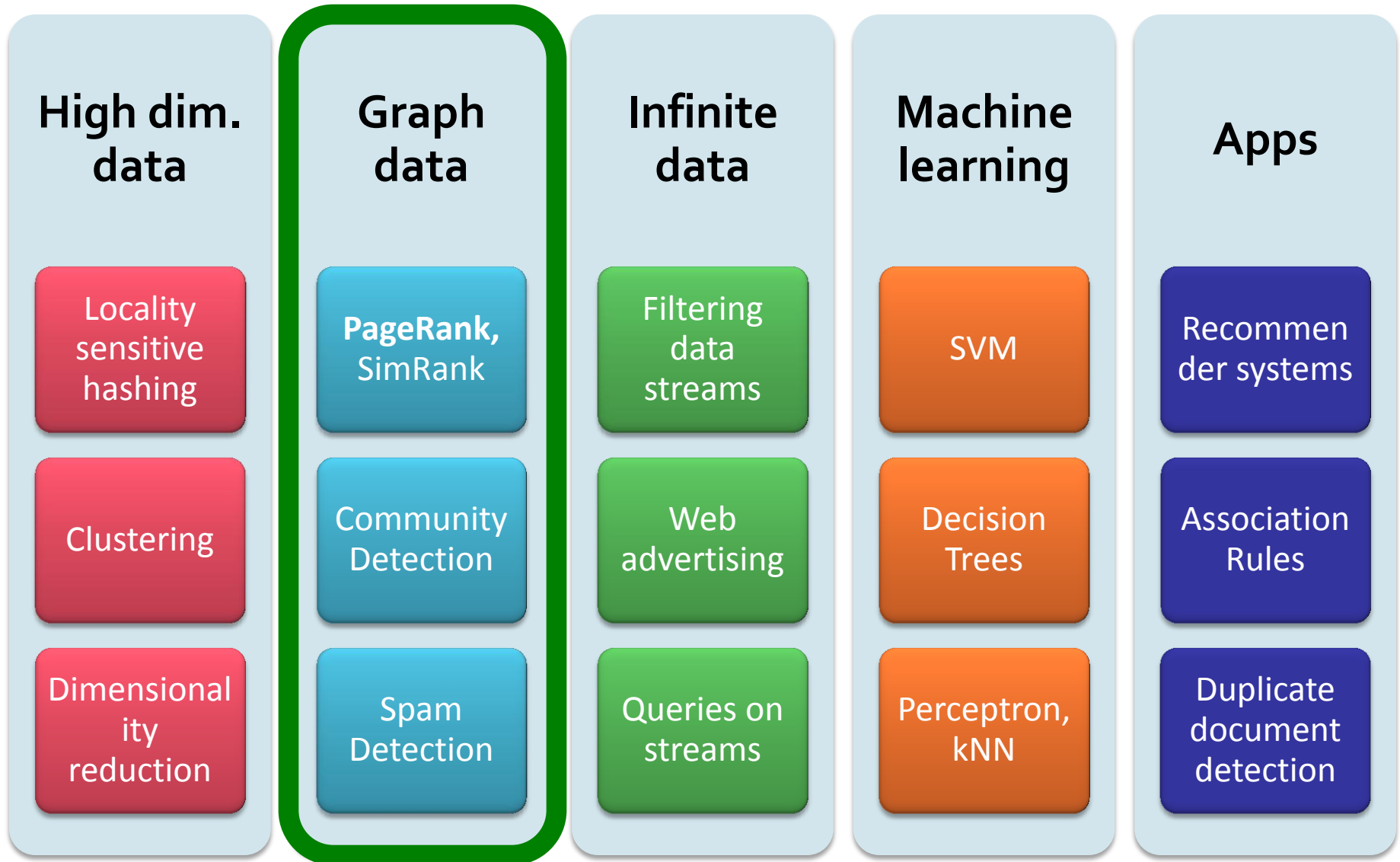


Link Analysis: PageRank and Similar Ideas

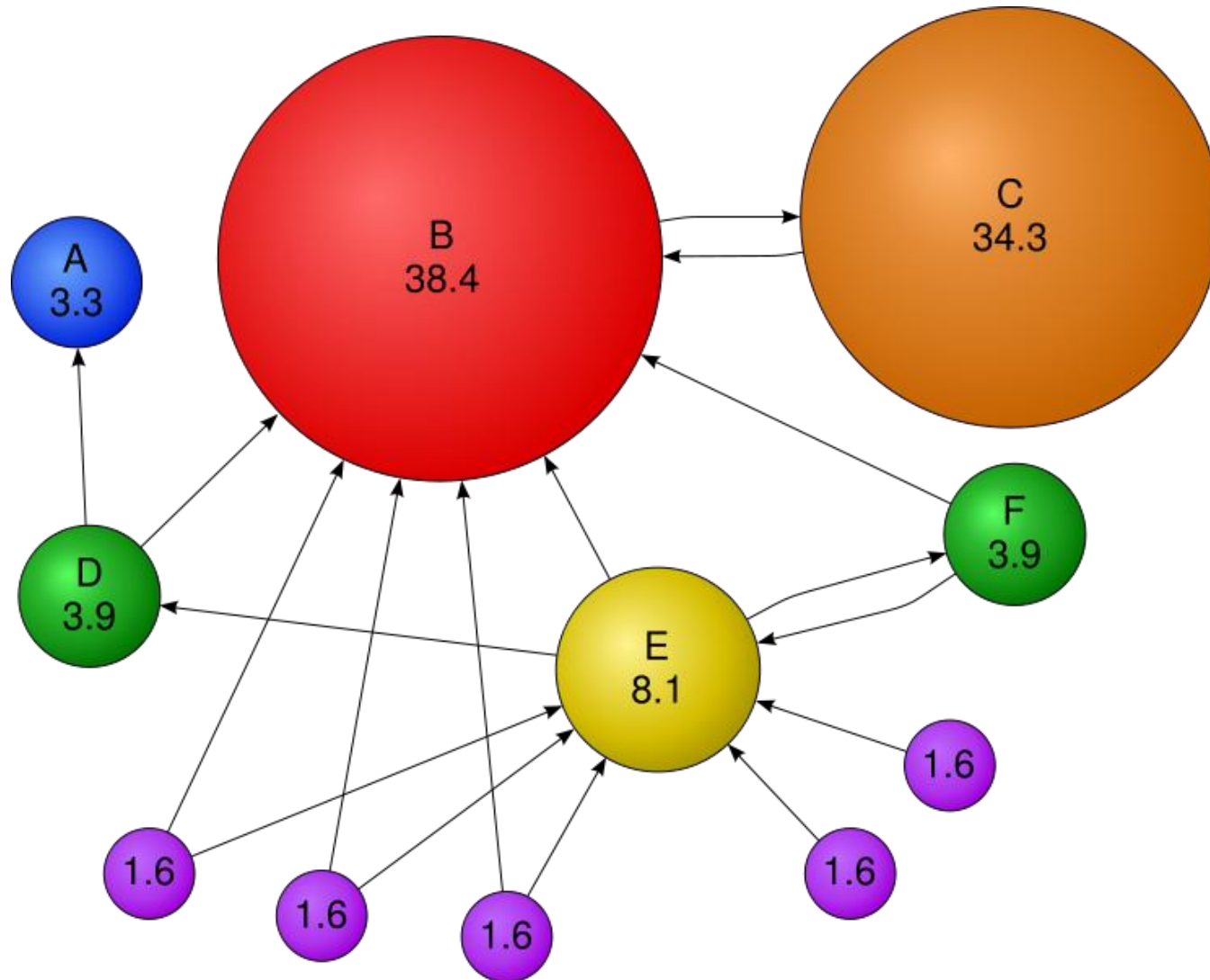
CS246: Mining Massive Datasets
Jure Leskovec, Stanford University
<http://cs246.stanford.edu>



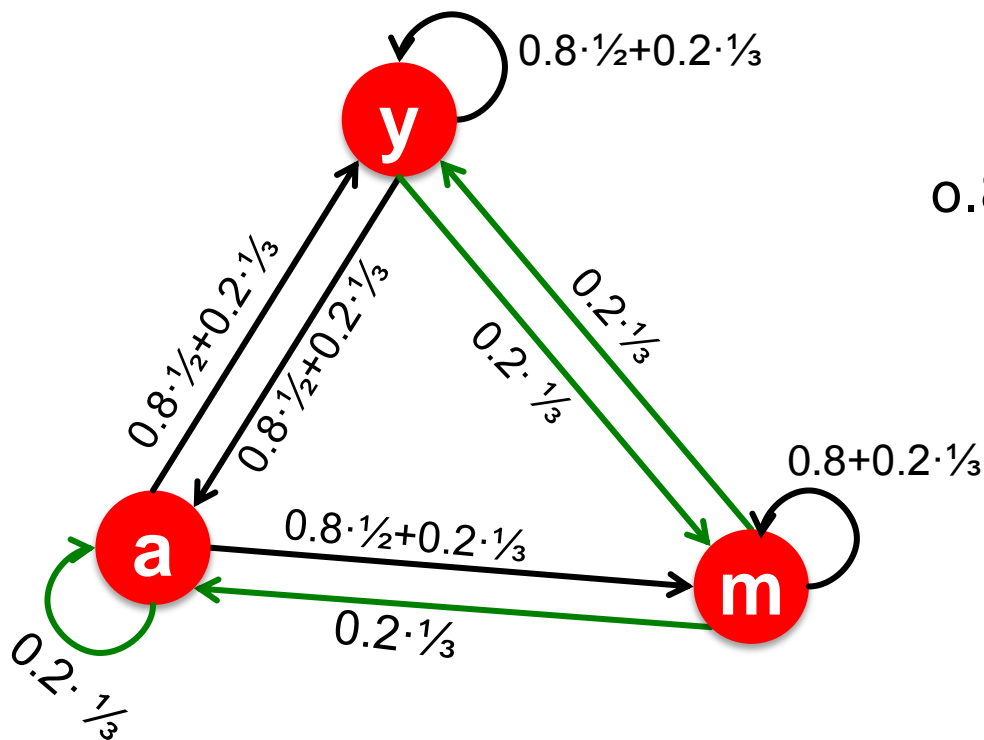
New Topic: Graph Data!



Example: PageRank Scores



Random Teleports ($\beta = 0.8$)



$$0.8 \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

y	7/15	7/15	1/15
a	7/15	1/15	1/15
m	1/15	7/15	13/15

A

y		1/3	0.33	0.24	0.26		7/33
a	=	1/3	0.20	0.20	0.18	...	5/33
m		1/3	0.46	0.52	0.56		21/33

$$\mathbf{r} = \mathbf{A} \mathbf{r}$$

Computing Page Rank

- **Key step is matrix-vector multiplication**

- $r^{\text{new}} = A \cdot r^{\text{old}}$

- Easy if we have enough main memory to hold A , r^{old} , r^{new}

- **Say $N = 1$ billion pages**

- We need 4 bytes for each entry (say)

- 2 billion entries for vectors, approx 8GB

- **Matrix A has N^2 entries**

- 10^{18} is a large number!

$$A = \beta \cdot M + (1-\beta) [1/N]_{N \times N}$$

$$A = 0.8 \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

$$= \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix}$$

Matrix Formulation

- Suppose there are N pages
- Consider page i , with d_i out-links
- We have $M_{ji} = 1/d_i$ when $i \rightarrow j$
and $M_{ji} = 0$ otherwise
- **The random teleport is equivalent to:**
 - Adding a **teleport link** from i to every other page and setting transition probability to $(1-\beta)/N$
 - Reducing the probability of following each out-link from $1/d_i$ to β/d_i
 - **Equivalent:** Tax each page a fraction $(1-\beta)$ of its score and redistribute evenly

Rearranging the Equation

- $\mathbf{r} = \mathbf{A} \cdot \mathbf{r}$, where $A_{ji} = \beta M_{ji} + \frac{1-\beta}{N}$
- $r_j = \sum_{i=1}^N A_{ji} \cdot r_i$
- $r_j = \sum_{i=1}^N \left[\beta M_{ji} + \frac{1-\beta}{N} \right] \cdot r_i$
 $= \sum_{i=1}^N \beta M_{ji} \cdot r_i + \frac{1-\beta}{N} \sum_{i=1}^N r_i$
 $= \sum_{i=1}^N \beta M_{ji} \cdot r_i + \frac{1-\beta}{N}$ since $\sum r_i = 1$
- So we get: $\mathbf{r} = \beta \mathbf{M} \cdot \mathbf{r} + \left[\frac{1-\beta}{N} \right]_N$

Note: Here we assumed \mathbf{M} has no dead-ends.

$[x]_N$... a vector of length N with all entries x

Sparse Matrix Formulation

- We just rearranged the **PageRank equation**

$$\mathbf{r} = \beta \mathbf{M} \cdot \mathbf{r} + \left[\frac{1 - \beta}{N} \right]_N$$

- where $[(1-\beta)/N]_N$ is a vector with all N entries $(1-\beta)/N$
- \mathbf{M} is a **sparse matrix!** (with no dead-ends)
 - 10 links per node, approx $10N$ entries
- So in each iteration, we need to:
 - Compute $\mathbf{r}^{\text{new}} = \beta \mathbf{M} \cdot \mathbf{r}^{\text{old}}$
 - Add a constant value $(1-\beta)/N$ to each entry in \mathbf{r}^{new}
 - **Note if \mathbf{M} contains dead-ends then $\sum_j r_j^{\text{new}} < 1$ and we also have to renormalize \mathbf{r}^{new} so that it sums to 1**

PageRank: The Complete Algorithm

- Input: Graph G and parameter β

- Directed graph G with spider traps and dead ends
- Parameter β

- Output: PageRank vector r

- **Set:** $r_j^{(0)} = \frac{1}{N}, \quad t = 1$

- **do:**

- $\forall j: r_j'^{(t)} = \sum_{i \rightarrow j} \beta \frac{r_i^{(t-1)}}{d_i}$

- $r_j'^{(t)} = 0$ if in-degree of j is 0

- **Now re-insert the leaked PageRank:**

- $\forall j: r_j^{(t)} = r_j'^{(t)} + \frac{1-S}{N}$ **where:** $S = \sum_j r_j'^{(t)}$

- $t = t + 1$

- **while** $\sum_j |r_j^{(t)} - r_j^{(t-1)}| > \varepsilon$

If the graph has no dead-ends then the amount of leaked PageRank is $1-\beta$. But since we have dead-ends the amount of leaked PageRank may be larger. We have to explicitly account for it by computing S .

Sparse Matrix Encoding

- Encode sparse matrix using only nonzero entries
 - Space proportional roughly to number of links
 - Say $10N$, or $4 \times 10 \times 1$ billion = 40GB
 - **Still won't fit in memory, but will fit on disk**

source node	degree	destination nodes
0	3	1, 5, 7
1	5	17, 64, 113, 117, 245
2	2	13, 23

Basic Algorithm: Update Step

- Assume enough RAM to fit r^{new} into memory
 - Store r^{old} and matrix M on disk
- 1 step of power-iteration is:

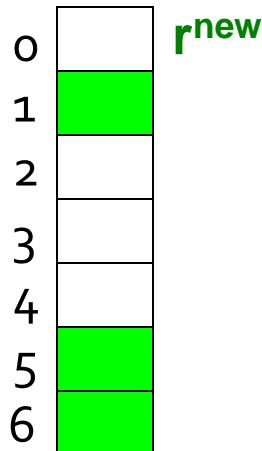
Initialize all entries of $r^{new} = (1-\beta) / N$

For each page i (of out-degree d_i):

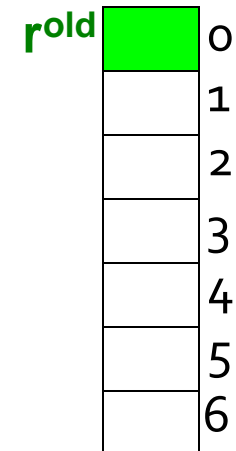
Read into memory: $i, d_i, dest_1, \dots, dest_{d_i}, r^{old}(i)$

For $j = 1 \dots d_i$

$r^{new}(dest_j) += \beta r^{old}(i) / d_i$



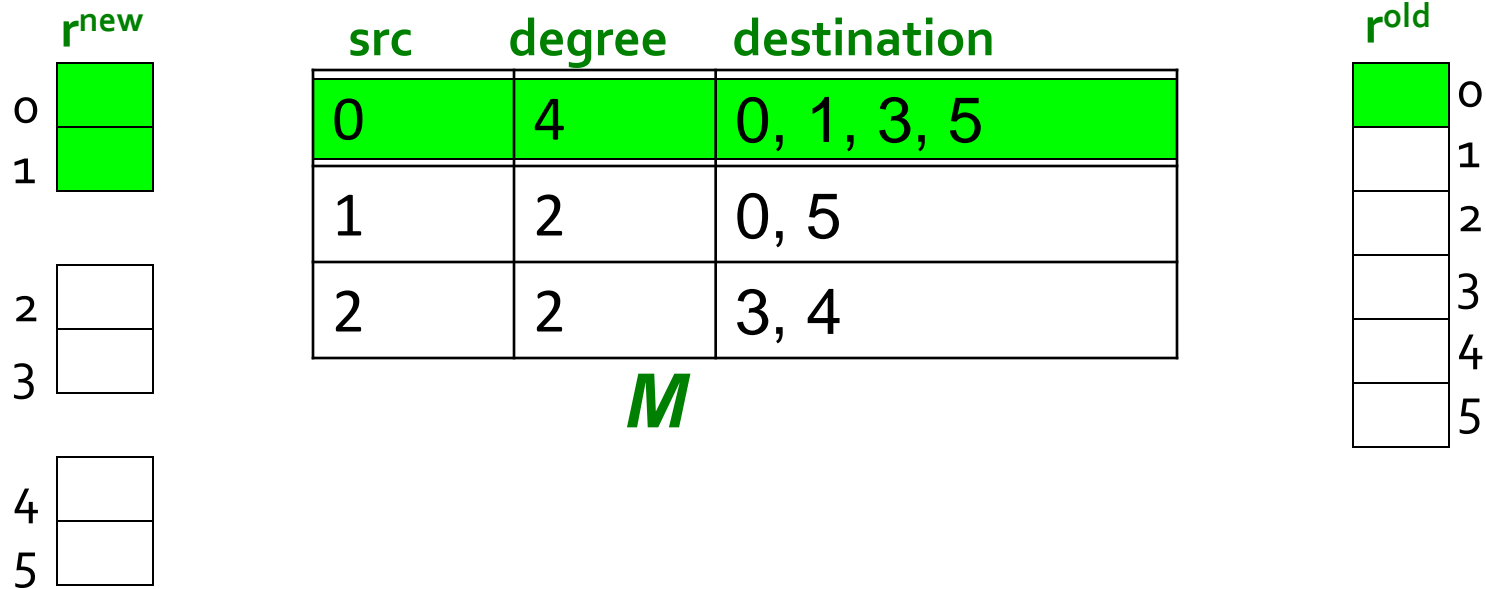
source	degree	destination
0	3	1, 5, 6
1	4	17, 64, 113, 117
2	2	13, 23



Analysis

- Assume enough RAM to fit r^{new} into memory
 - Store r^{old} and matrix M on disk
- In each iteration, we have to:
 - Read r^{old} and M
 - Write r^{new} back to disk
 - Cost per iteration of Power method:
 $= 2|r| + |M|$
- Question:
 - What if we could not even fit r^{new} in memory?

Block-based Update Algorithm



- Break r^{new} into k blocks that fit in memory
- Scan M and r^{old} once for each block

Analysis of Block Update

- **Similar to nested-loop join in databases**
 - Break r^{new} into k blocks that fit in memory
 - Scan M and r^{old} once for each block
- **Total cost:**
 - k scans of M and r^{old}
 - **Cost per iteration of Power method:**
$$k(|M| + |r|) + |r| = k|M| + (k+1)|r|$$
- **Can we do better?**
 - **Hint:** M is much bigger than r (approx 10-20x), so we must avoid reading it k times per iteration

Block-Stripe Update Algorithm

r^{new}

0	
1	

src	degree	destination
0	4	0, 1
1	3	0
2	2	1

2	
3	

0	4	3
2	2	3

4	
5	

0	4	5
1	3	5
2	2	4

r^{old}

	0
	1
	2
	3
	4
	5

Break M into stripes! Each stripe contains only destination nodes in the corresponding block of r^{new}

Block-Stripe Analysis

- Break M into stripes
 - Each stripe contains only destination nodes in the corresponding block of r^{new}
- Some additional overhead per stripe
 - But it is usually worth it
- **Cost per iteration of Power method:**
 $= |M|(1+\varepsilon) + (k+1)|r|$

Some Problems with Page Rank

- **Measures generic popularity of a page**
 - Biased against topic-specific authorities
 - **Solution:** Topic-Specific PageRank (**next**)
- **Uses a single measure of importance**
 - Other models of importance
 - **Solution:** Hubs-and-Authorities
- **Susceptible to Link spam**
 - Artificial link topographies created in order to boost page rank
 - **Solution:** TrustRank

Topic-Specific PageRank

Topic-Specific PageRank

- Instead of generic popularity, can we measure popularity within a topic?
- **Goal:** Evaluate Web pages not just according to their popularity, but by how close they are to a particular topic, e.g. “sports” or “history”
- **Allows search queries to be answered based on interests of the user**
 - **Example:** Query “Trojan” wants different pages depending on whether you are interested in sports, history and computer security

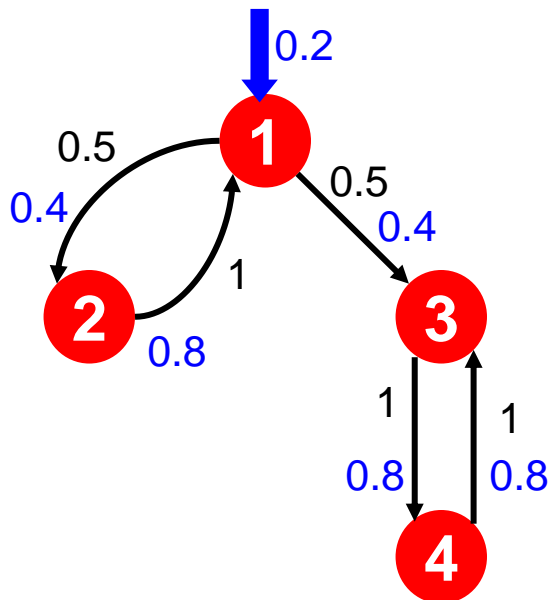
Topic-Specific PageRank

- Random walker has a small probability of teleporting at any step
- **Teleport can go to:**
 - **Standard PageRank:** Any page with equal probability
 - To avoid dead-end and spider-trap problems
 - **Topic Specific PageRank:** A topic-specific set of “relevant” pages (**teleport set**)
- **Idea: Bias the random walk**
 - When walker teleports, she pick a page from a set S
 - S contains only pages that are relevant to the topic
 - E.g., Open Directory (DMOZ) pages for a given topic/query
 - For each teleport set S , we get a different vector r_S

Matrix Formulation

- To make this work all we need is to update the teleportation part of the PageRank formulation:
 - $$A_{ij} = \begin{cases} \beta M_{ij} + (1-\beta)/|S| & \text{if } i \in S \\ \beta M_{ij} & \text{otherwise} \end{cases}$$
 - A is stochastic!
- We have weighted all pages in the teleport set S equally
 - Could also assign different weights to pages!
- Compute as for regular PageRank:
 - Multiply by M , then add a vector
 - Maintains sparseness

Example



Suppose $S = \{1\}$, $\beta = 0.8$

Node	Iteration				
	0	1	2	...	stable
1	0.25	0.4	0.28		0.294
2	0.25	0.1	0.16		0.118
3	0.25	0.3	0.32		0.327
4	0.25	0.2	0.24		0.261

$S=\{1\}$, $\beta=0.90$:

$r=[0.17, 0.07, 0.40, 0.36]$

$S=\{1\}$, $\beta=0.8$:

$r=[0.29, 0.11, 0.32, 0.26]$

$S=\{1\}$, $\beta=0.70$:

$r=[0.39, 0.14, 0.27, 0.19]$

$S=\{1,2,3,4\}$, $\beta=0.8$:

$r=[0.13, 0.10, 0.39, 0.36]$

$S=\{1,2,3\}$, $\beta=0.8$:

$r=[0.17, 0.13, 0.38, 0.30]$

$S=\{1,2\}$, $\beta=0.8$:

$r=[0.26, 0.20, 0.29, 0.23]$

$S=\{1\}$, $\beta=0.8$:

$r=[0.29, 0.11, 0.32, 0.26]$

Discovering the Topic

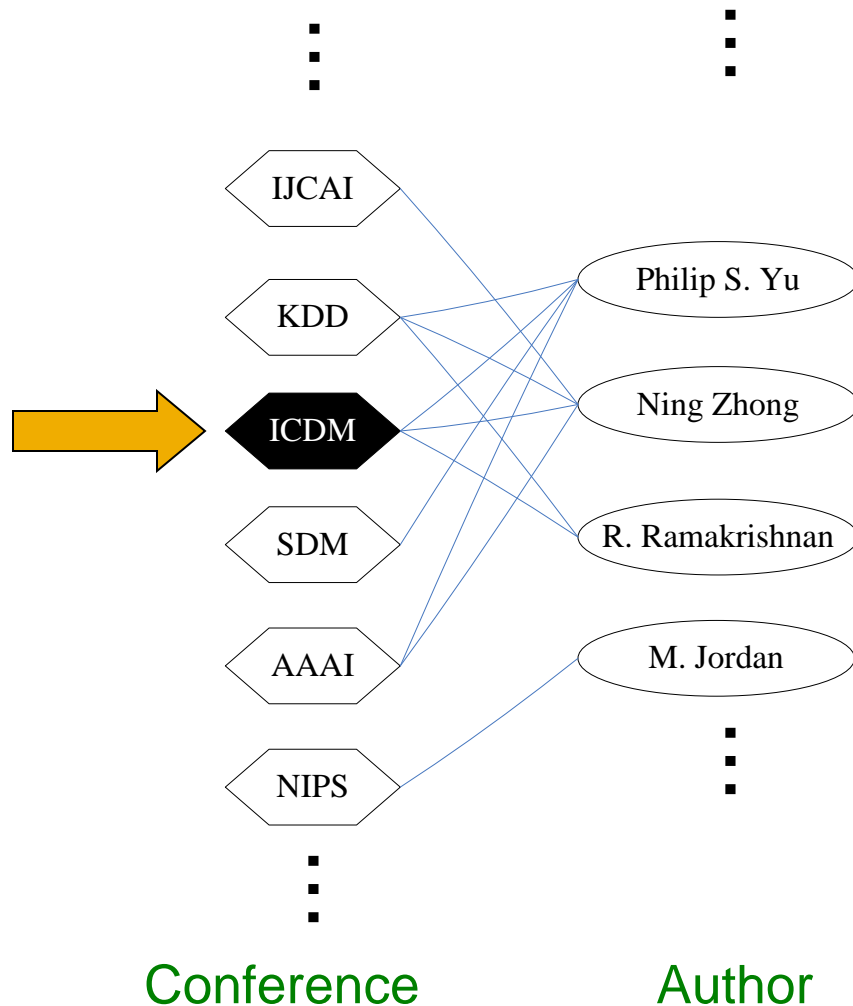
- **Create different PageRanks for different topics**
 - The 16 DMOZ top-level categories:
 - arts, business, sports,...
- **Which topic ranking to use?**
 - User can pick from a menu
 - Classify query into a topic
 - Can use the **context** of the query
 - E.g., query is launched from a web page talking about a known topic
 - History of queries e.g., “basketball” followed by “Jordan”
 - User context, e.g., user’s bookmarks, ...

SimRank: An Application of Personalized PageRank

SimRank: Idea

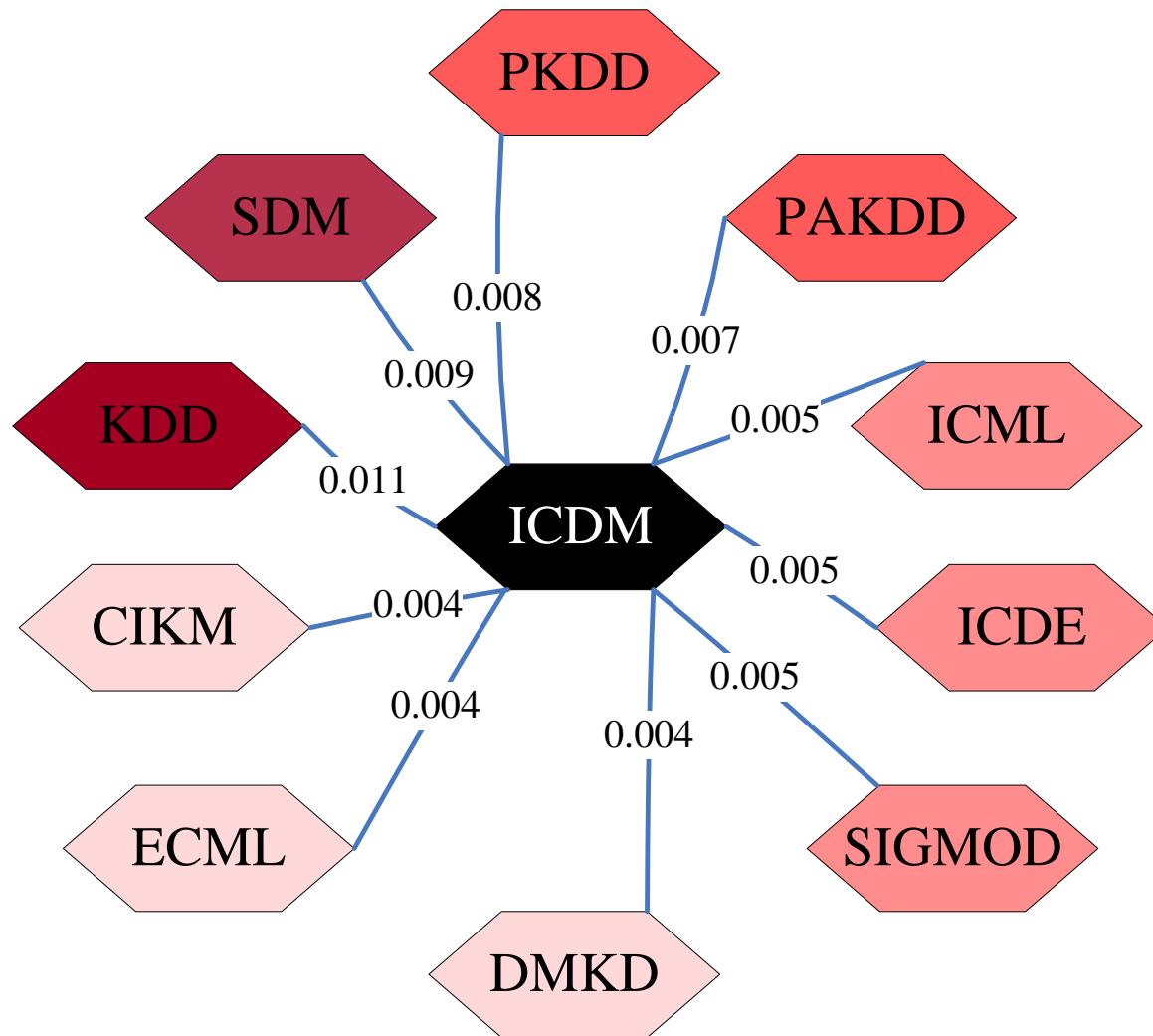
- **SimRank:** Random walks from a fixed node on k -partite graphs
- **Setting:** k -partite graph with k types of nodes
 - Example: picture nodes and tag nodes
- Do a **Random-Walk with Restarts** from node u
 - i.e., **teleport set** $S = \{u\}$
- **Resulting scores measures** similarity/proximity to node u
- **Problem:**
 - Must be done once for each node u
 - Suitable for sub-Web-scale applications

SimRank: Example



Q: What is most related conference to **ICDM**?

SimRank: Example



3 announcements:

- We will hold office hours via Google hangout:
- Sat, Sun 10am Pacific time
- HW2 is due
- HW3 is out

TrustRank: Combating the Web Spam

What is Web Spam?

- **Spamming:**

- Any deliberate action to boost a web page's position in search engine results, incommensurate with page's real value

- **Spam:**

- Web pages that are the result of spamming

- This is a very broad definition

- **SEO** industry might disagree!
- SEO = search engine optimization

- Approximately **10-15%** of web pages are spam

Web Search

- **Early search engines:**

- Crawl the Web
- Index pages by the words they contained
- Respond to search queries (lists of words) with the pages containing those words

- **Early page ranking:**

- Attempt to order pages matching a search query by “importance”
- **First search engines considered:**
 - (1) Number of times query words appeared
 - (2) Prominence of word position, e.g. title, header

First Spammers

- As people began to use search engines to find things on the Web, those with commercial interests tried to **exploit search engines** to bring people to their own site – whether they wanted to be there or not
- **Example:**
 - Shirt-seller might pretend to be about “movies”
- **Techniques for achieving high relevance/importance for a web page**

First Spammers: Term Spam

- **How do you make your page appear to be about movies?**
 - **(1)** Add the word movie 1,000 times to your page
 - Set text color to the background color, so only search engines would see it
 - **(2)** Or, run the query “movie” on your target search engine
 - See what page came first in the listings
 - Copy it into your page, make it “invisible”
- **These and similar techniques are term spam**

Google's Solution to Term Spam

- Believe what people say about you, rather than what you say about yourself
 - Use words in the anchor text (words that appear underlined to represent the link) and its surrounding text
- PageRank as a tool to measure the “importance” of Web pages

The screenshot shows a Google search interface with the query "miserable failure" entered in the search bar. The search results are displayed under the "Web" tab, showing the first 10 results. The top result is the "Biography of President George W. Bush" from the official White House website. The second result is "Welcome to MichaelMoore.com!", the official site of the gadfly of corporations. The third result is a BBC News article titled "BBC NEWS | Americas | 'Miserable failure' links to Bush". The fourth result is "Google's (and Inktomi's) Miserable Failure", which discusses how a search for "miserable failure" on Google brings up the official George W. Bush biography from the US White House web site, dismissed by Google as not a ...

Web Images Groups News Froogle Local more »
Google miserable failure Search Advanced Search Preferences

Web Results 1 - 10 of about 969,000 for miserable failure. (0.06 seconds)

[Biography of President George W. Bush](#)
Biography of the president from the official White House web site.
[www.whitehouse.gov/president/gwbbio.html](#) - 29k - [Cached](#) - [Similar pages](#)
[Past Presidents](#) - [Kids Only](#) - [Current News](#) - [President](#)
[More results from www.whitehouse.gov »](#)

[Welcome to MichaelMoore.com!](#)
Official site of the gadfly of corporations, creator of the film Roger and Me and the television show The Awful Truth. Includes mailing list, message board, ...
[www.michaelmoore.com/](#) - 35k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)

[BBC NEWS | Americas | 'Miserable failure' links to Bush](#)
Web users manipulate a popular search engine so an unflattering description leads to the president's page.
[news.bbc.co.uk/2/hi/americas/3298443.stm](#) - 31k - [Cached](#) - [Similar pages](#)

[Google's \(and Inktomi's\) Miserable Failure](#)
A search for miserable failure on Google brings up the official George W. Bush biography from the US White House web site. Dismissed by Google as not a ...
[searchenginewatch.com/sereport/article.php/3296101](#) - 45k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)

Why It Works?

- **Our hypothetical shirt-seller loses**
 - Saying he is about movies doesn't help, because others don't say he is about movies
 - His page isn't very important, so it won't be ranked high for shirts or movies
- **Example:**
 - Shirt-seller creates 1,000 pages, each links to his with "movie" in the anchor text
 - These pages have no links in, so they get little PageRank
 - So the shirt-seller can't beat truly important movie pages, like IMDB



SPAM FARMING

Google vs. Spammers: Round 2!

- Once Google became the dominant search engine, spammers began to work out ways to fool Google
- **Spam farms** were developed to concentrate PageRank on a single page
- **Link spam:**
 - Creating link structures that boost PageRank of a particular page



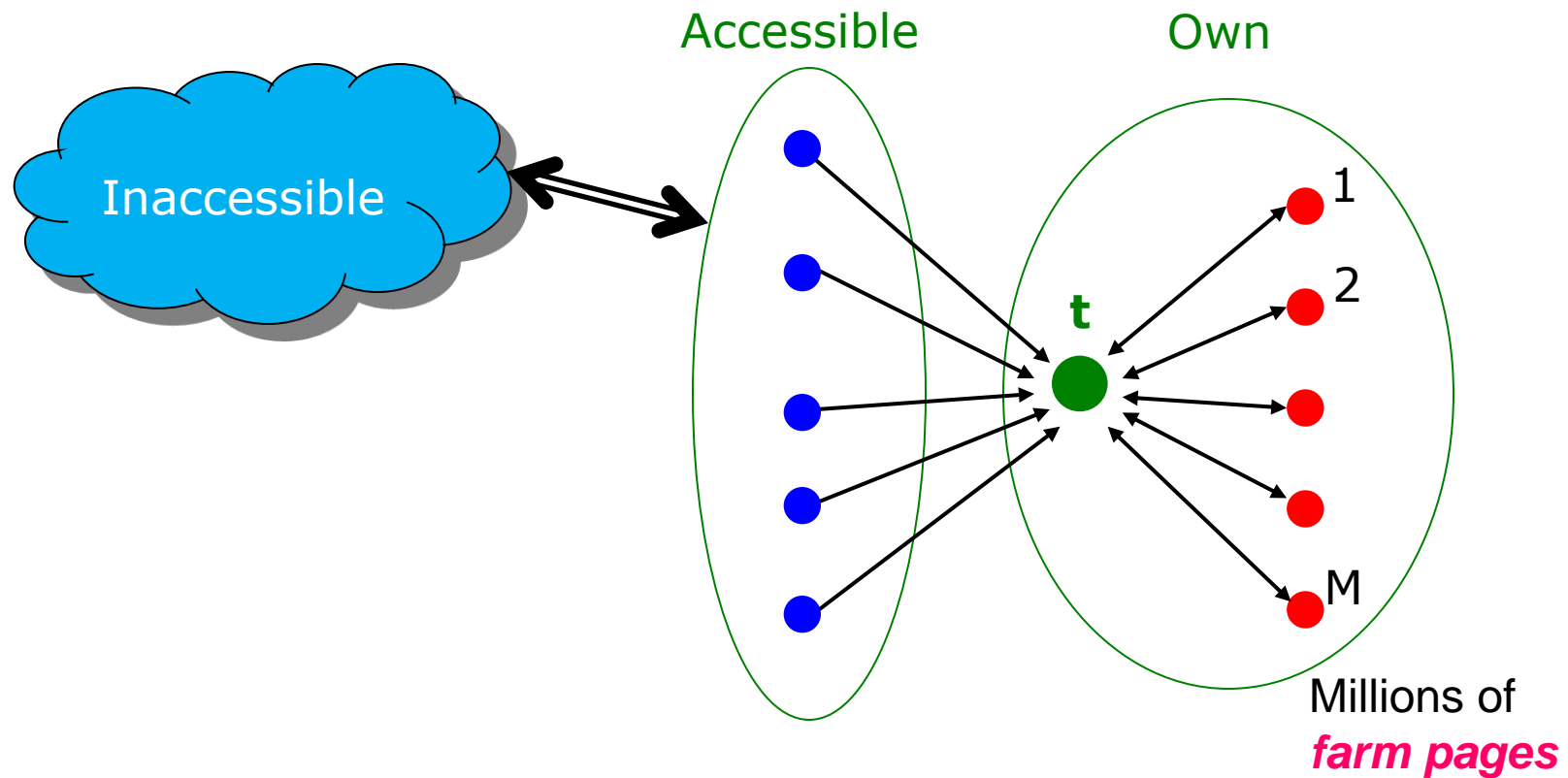
Link Spamming

- **Three kinds of web pages from a spammer's point of view**
 - **Inaccessible pages**
 - **Accessible pages**
 - e.g., blog comments pages
 - spammer can post links to his pages
 - **Own pages**
 - Completely controlled by spammer
 - May span multiple domain names

Link Farms

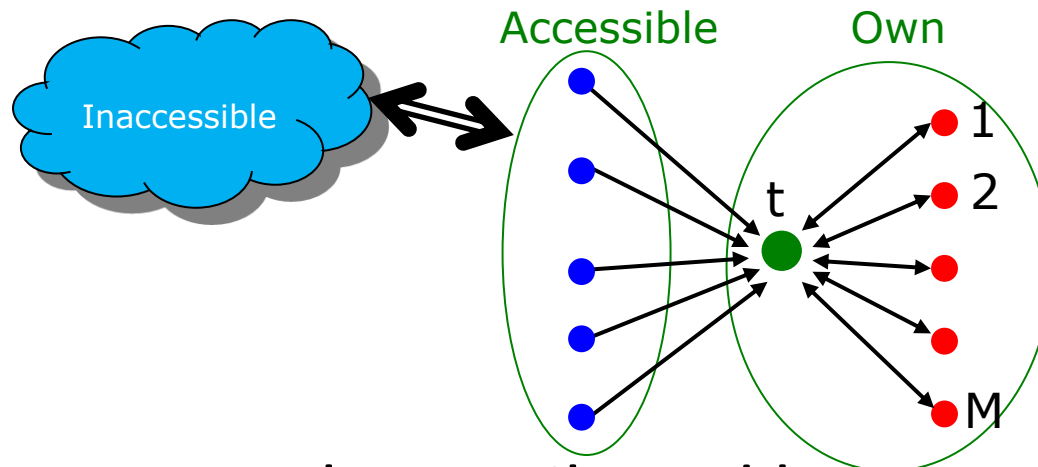
- **Spammer's goal:**
 - Maximize the PageRank of target page t
- **Technique:**
 - Get as many links from accessible pages as possible to target page t
 - Construct “link farm” to get PageRank multiplier effect

Link Farms



One of the most common and effective organizations for a link farm

Analysis



N ...# pages on the web
 M ...# of pages spammer owns

- x : PageRank contributed by accessible pages
- y : PageRank of target page t

- Rank of each “farm” page = $\frac{\beta y}{M} + \frac{1-\beta}{N}$

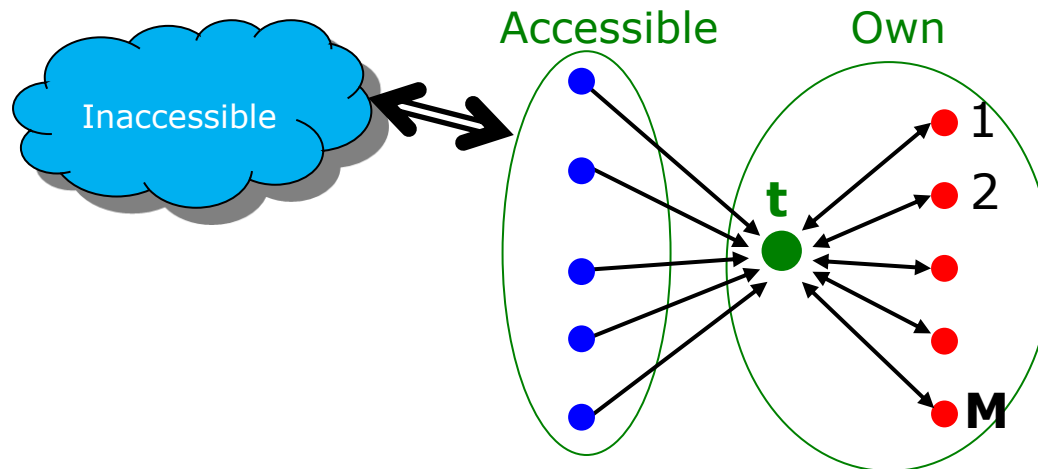
- $$y = x + \beta M \left[\frac{\beta y}{M} + \frac{1-\beta}{N} \right] + \frac{1-\beta}{N}$$

$$= x + \beta^2 y + \frac{\beta(1-\beta)M}{N} + \frac{1-\beta}{N}$$

Very small; ignore
 Now we solve for y

- $$y = \frac{x}{1-\beta^2} + c \frac{M}{N} \quad \text{where } c = \frac{\beta}{1+\beta}$$

Analysis



N ...# pages on the web
 M ...# of pages spammer owns

- $y = \frac{x}{1-\beta^2} + c \frac{M}{N}$ where $c = \frac{\beta}{1+\beta}$
- For $\beta = 0.85$, $1/(1-\beta^2) = 3.6$
- Multiplier effect for “acquired” PageRank
- By making M large, we can make y as large as we want

TrustRank: Combating the Web Spam

Combating Spam

■ Combating term spam

- Analyze text using statistical methods
- Similar to email spam filtering
- Also useful: Detecting approximate duplicate pages

■ Combating link spam

- Detection and blacklisting of structures that look like spam farms
 - Leads to another war – hiding and detecting spam farms
- **TrustRank** = topic-specific PageRank with a teleport set of “trusted” pages
 - **Example:** .edu domains, similar domains for non-US schools

TrustRank: Idea

- **Basic principle: Approximate isolation**
 - It is rare for a “good” page to point to a “bad” (spam) page
- Sample a set of **seed pages** from the web
- Have an **oracle (human)** to identify the good pages and the spam pages in the seed set
 - **Expensive task**, so we must make seed set as small as possible

Trust Propagation

- Call the subset of seed pages that are identified as **good** the **trusted pages**
- Perform a topic-sensitive PageRank with **teleport set = trusted pages**
 - **Propagate trust through links:**
 - Each page gets a trust value between **0** and **1**
- **Use a threshold value and mark all pages below the trust threshold as spam**

Why is it a good idea?

- **Trust attenuation:**

- The degree of trust conferred by a trusted page decreases with the distance in the graph

- **Trust splitting:**

- The larger the number of out-links from a page, the less scrutiny the page author gives each out-link
- Trust is **split** across out-links

Picking the Seed Set

- **Two conflicting considerations:**
 - Human has to inspect each seed page, so seed set must be as small as possible
 - Must ensure every **good page** gets adequate trust rank, so need make all good pages reachable from seed set by short paths

Approaches to Picking Seed Set

- Suppose we want to pick a seed set of k pages
- **How to do that?**
- **(1) PageRank:**
 - Pick the top k pages by PageRank
 - Theory is that you can't get a bad page's rank really high
- **(2) Use trusted domains** whose membership is controlled, like .edu, .mil, .gov

Spam Mass

- In the **TrustRank** model, we start with good pages and propagate trust
- **Complementary view:**
What fraction of a page's PageRank comes from **spam** pages?
- In practice, we don't know all the spam pages, so we need to estimate

Spam Mass Estimation

- r_p = PageRank of page p
- r_p^+ = page rank of p with teleport into **trusted** pages only
- **Then:** What fraction of a page's PageRank comes from **spam** pages?

$$r_p^- = r_p - r_p^+$$

- **Spam mass of p** = $\frac{r_p^-}{r_p}$

HITS: Hubs and Authorities

Hubs and Authorities

- **HITS** (Hypertext-Induced Topic Selection)
 - Is a measure of importance of pages or documents, similar to PageRank
 - Proposed at around same time as PageRank ('98)
- **Goal**: Imagine we want to find good newspapers
 - Don't just find newspapers. Find “experts” – people who link in a coordinated way to good newspapers
- **Idea: Links as votes**
 - Page is more important if it has more links
 - In-coming links? Out-going links?

Finding newspapers

■ Hubs and Authorities

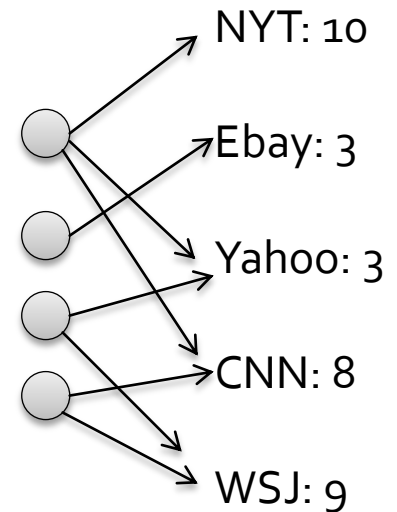
Each page has 2 scores:

■ Quality as an expert (**hub**):

- Total sum of votes of pages pointed to

■ Quality as an content (**authority**):

- Total sum of votes of experts



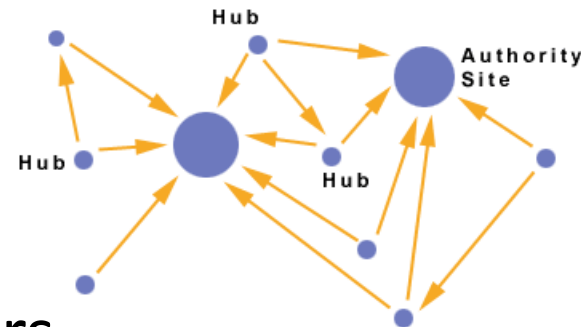
■ Principle of repeated improvement

Hubs and Authorities

Interesting pages fall into two classes:

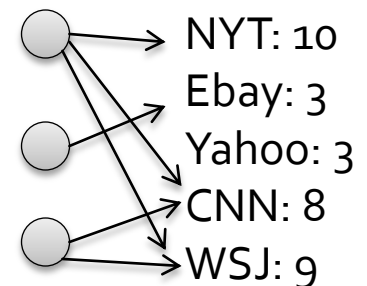
1. **Authorities** are pages containing useful information

- Newspaper home pages
- Course home pages
- Home pages of auto manufacturers

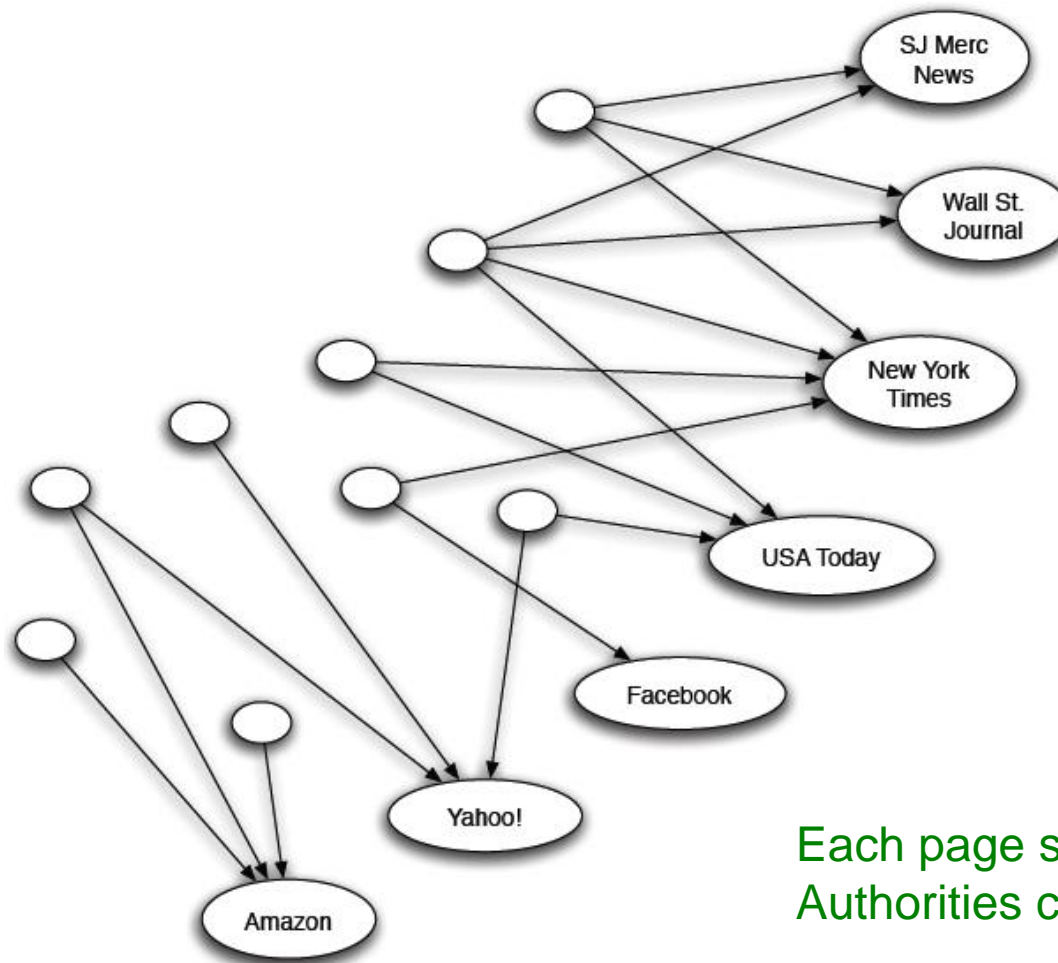


2. **Hubs** are pages that link to authorities

- List of newspapers
- Course bulletin
- List of US auto manufacturers



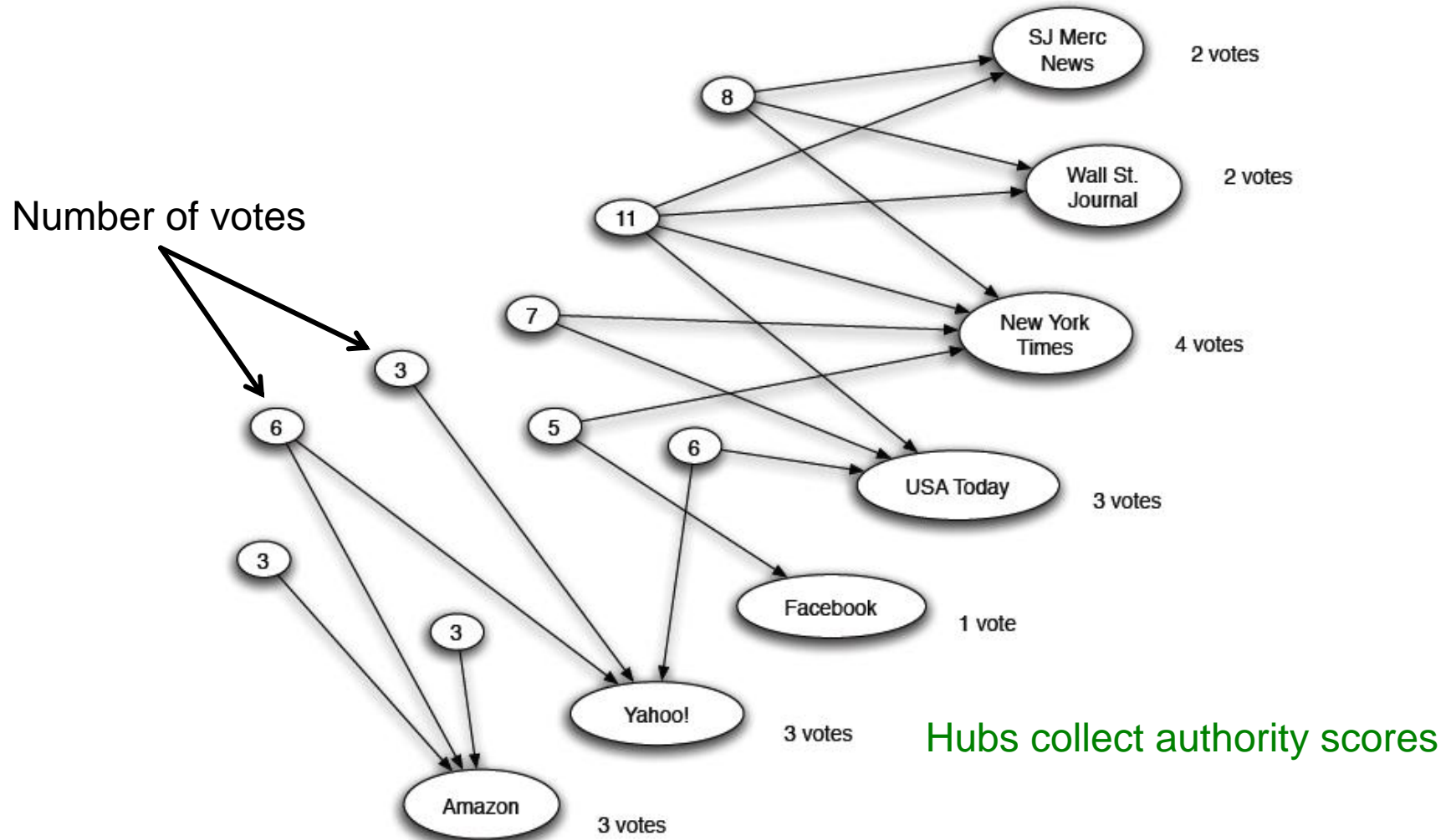
Counting in-links: Authority



Each page starts with hub score 1
Authorities collect their votes

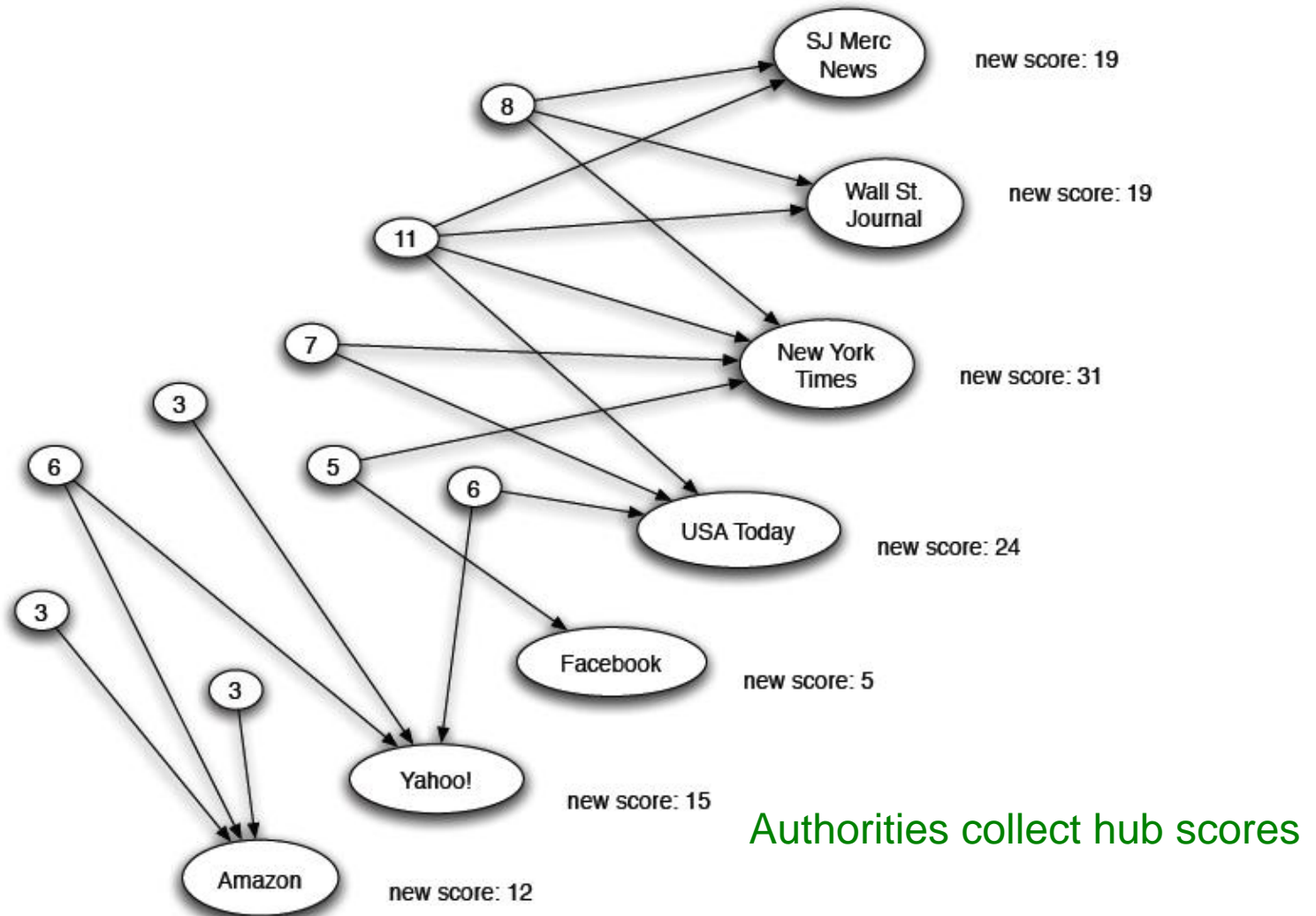
(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

Expert Quality: Hub



(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

Reweighting



(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

Mutually Recursive Definition

- A good hub links to many good authorities
- A good authority is linked from many good hubs
- Model using two scores for each node:
 - **Hub** score and **Authority** score
 - Represented as vectors h and a

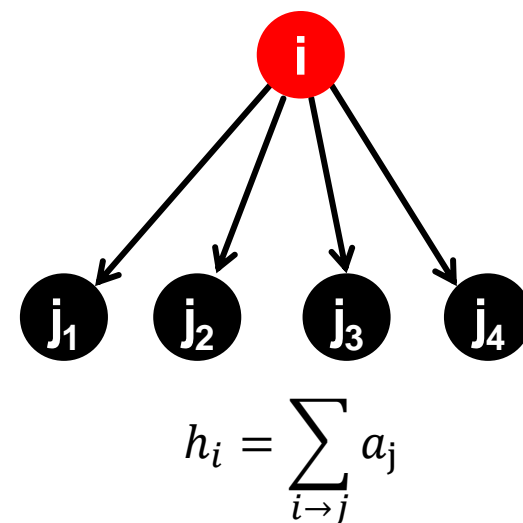
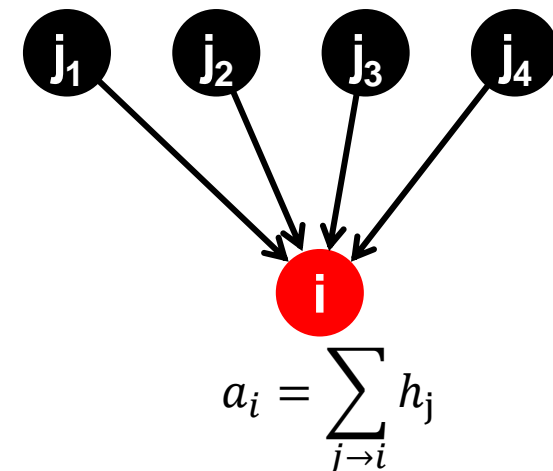
Hubs and Authorities

- Each page i has 2 scores:

- Authority score: a_i
- Hub score: h_i

HITS algorithm:

- Initialize: $a_i = 1, h_i = 1$
- Then keep iterating:
 - $\forall i$: Authority: $a_i = \sum_{j \rightarrow i} h_j$
 - $\forall i$: Hub: $h_i = \sum_{i \rightarrow j} a_j$
 - $\forall i$: normalize: $\sum_j a_j = 1, \sum_j h_j = 1$



Transition Matrix A

- **HITS converges to a single stable point**

- Slightly change the notation:

- Vector $a = (a_1, \dots, a_n)$, $h = (h_1, \dots, h_n)$

- Adjacency matrix ($n \times n$): $A_{ij} = 1$ if $i \rightarrow j$

- **Then:**

$$h_i = \sum_{i \rightarrow j} a_j \Leftrightarrow h_i = \sum_j A_{ij} a_j$$

- So: $h = A a$

- And likewise: $a = A^T h$

Hub and Authority Equations

- The **hub** score of page i is proportional to the sum of the **authority** scores of the pages it links to: $h = \lambda A a$
 - Constant λ is a scale factor, $\lambda = 1 / \sum h_i$
- The **authority** score of page i is proportional to the sum of the **hub** scores of the pages it is linked from: $a = \mu A^T h$
 - Constant μ is scale factor, $\mu = 1 / \sum a_i$

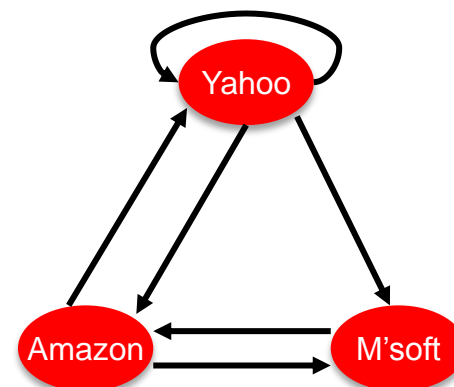
Iterative algorithm

- **The HITS algorithm:**
 - Initialize \mathbf{h} , \mathbf{a} to all 1's
 - Repeat:
 - $\mathbf{h} = \mathbf{A} \mathbf{a}$
 - Scale \mathbf{h} so that its sums to 1.0
 - $\mathbf{a} = \mathbf{A}^T \mathbf{h}$
 - Scale \mathbf{a} so that its sums to 1.0
 - Until \mathbf{h} , \mathbf{a} converge (i.e., change very little)

Example

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$



$$\begin{array}{lcl} a(\text{yahoo}) & = & 1 \quad 1 \quad 1 \quad 1 \quad \dots \quad 1 \\ a(\text{amazon}) & = & 1 \quad 1 \quad 4/5 \quad 0.75 \quad \dots \quad 0.732 \\ a(\text{m'soft}) & = & 1 \quad 1 \quad 1 \quad 1 \quad \dots \quad 1 \end{array}$$

$$\begin{array}{lcl} h(\text{yahoo}) & = & 1 \quad 1 \quad 1 \quad 1 \quad \dots \quad 1.000 \\ h(\text{amazon}) & = & 1 \quad 2/3 \quad 0.71 \quad 0.73 \quad \dots \quad 0.732 \\ h(\text{m'soft}) & = & 1 \quad 1/3 \quad 0.29 \quad 0.27 \quad \dots \quad 0.268 \end{array}$$

Hubs and Authorities

- **HITS algorithm in new notation:**

- Set: $a = h = 1^n$

- **Repeat:**

- $h = A a, \quad a = A^T h$

- Normalize

- Then: $a = A^T \underbrace{(A a)}_{\substack{\text{new } h \\ \text{new } a}}$

- Thus, in $2k$ steps:

$$a = (A^T A)^k a$$

$$h = (A A^T)^k h$$

a is being updated (in 2 steps):

$$A^T(A a) = (A^T A) a$$

h is updated (in 2 steps):

$$A(A^T h) = (A A^T) h$$

Repeated matrix powering

Existence and Uniqueness

- $h = \lambda A a$
 - $a = \mu A^T h$
 - $h = \lambda \mu A A^T h$
 - $a = \lambda \mu A^T A a$
- $\lambda = 1 / \sum h_i$
 $\mu = 1 / \sum a_i$
- Under reasonable assumptions about A , the HITS iterative algorithm **converges to vectors h^* and a^*** :
 - h^* is the **principal eigenvector** of matrix $A A^T$
 - a^* is the **principal eigenvector** of matrix $A^T A$

PageRank and HITS

- PageRank and HITS are two solutions to the same problem:
 - What is the value of an in-link from u to v ?
 - In the PageRank model, the value of the link depends on the links into u
 - In the HITS model, it depends on the value of the other links out of u
- The destinies of PageRank and HITS post-1998 were very different