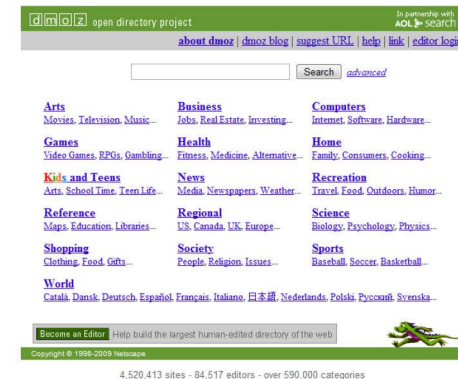


How to Organize the Web?

- **How to organize the Web?**
- **First try: Human curated Web directories**
 - Yahoo, DMOZ, LookSmart
- **Second try: Web Search**
 - **Information Retrieval** attempts to find relevant docs in a small and trusted set
 - Newspaper articles, Patents, etc.
 - **But:** Web is **huge**, full of untrusted documents, random things, web spam, etc.
 - **So we need a good way to rank webpages!**



Web Search: 2 Challenges

Two challenges of web search:

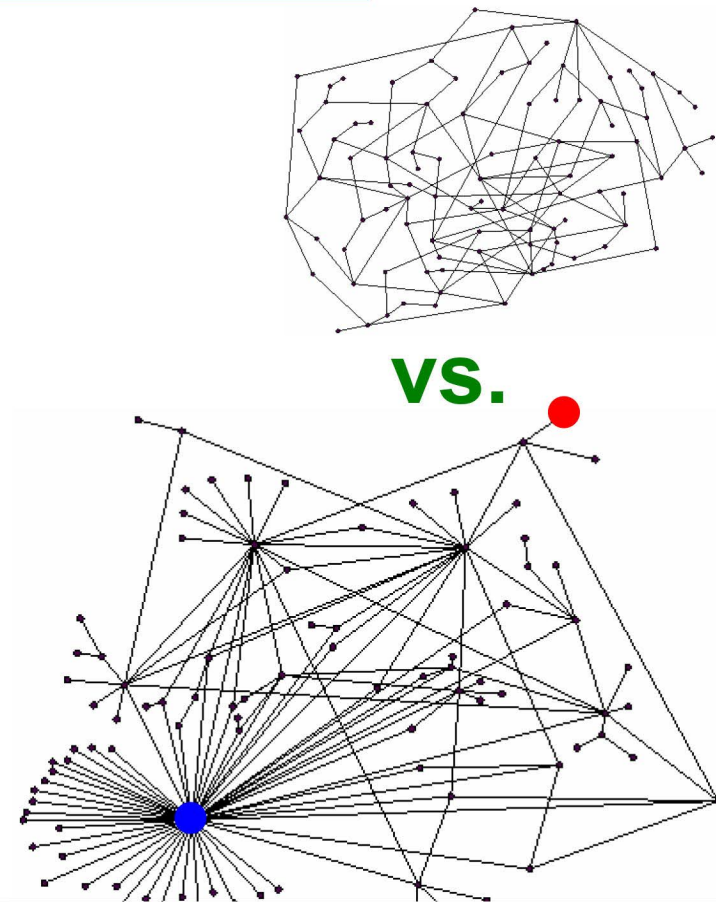
- (1) Web contains many sources of information
Who to “trust”?
 - **Insight:** Trustworthy pages may point to each other!
- (2) What is the “best” answer to query
“newspaper”?
 - No single right answer
 - **Insight:** Pages that actually know about newspapers might all be pointing to many newspapers

Ranking Nodes on the Graph

- All web pages are not equally “important”

www.joe-schmoe.com vs. www.stanford.edu

- **We already know:**
There is large diversity in the web-graph node connectivity.
- **So, let’s rank the pages using the web graph link structure!**



Link Analysis Algorithms

- We will cover the following Link Analysis approaches to computing importance of nodes in a graph:
 - Hubs and Authorities (HITS)
 - Page Rank
 - Topic-Specific (Personalized) Page Rank

Sidenote: Various notions of node centrality: Node u

- Degree centrality = degree of u
- Betweenness centrality = #shortest paths passing through u
- Closeness centrality = avg. length of shortest paths from u to all other nodes of the network
- Eigenvector centrality = like PageRank

Hubs and Authorities

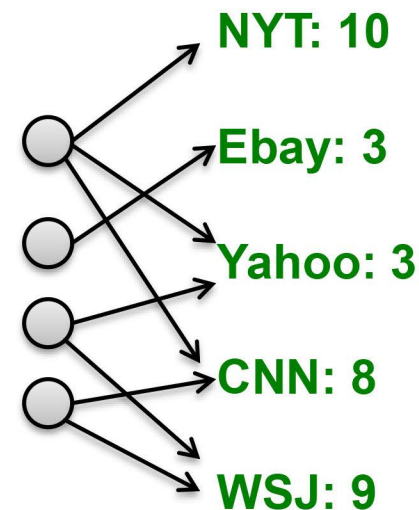
Link Analysis

- **Goal** (back to the newspaper example):
 - Don't just find newspapers. Find "experts" – pages that 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?

- **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**

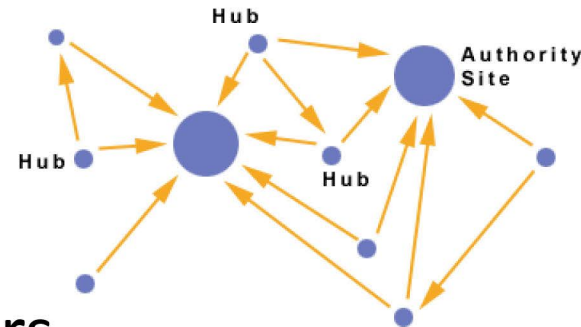


Hubs and Authorities

Interesting pages fall into two classes:

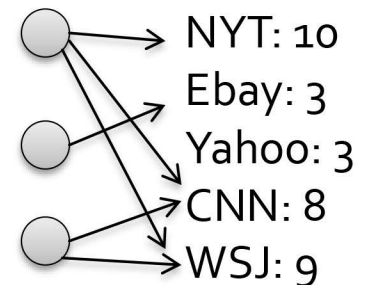
1. **Authorities** are pages containing useful information

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

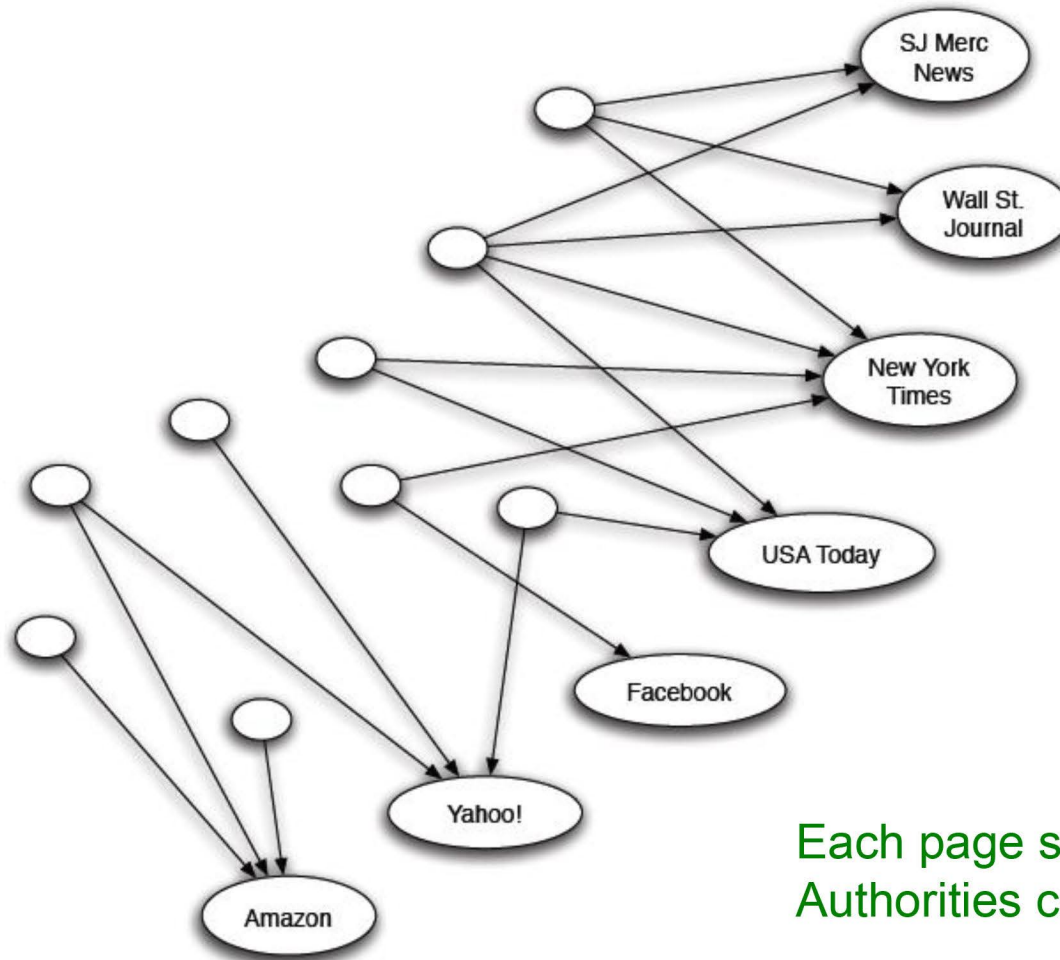


2. **Hubs** are pages that link to authorities

- List of newspapers
- Course bulletin
- List of U.S. auto manufacturers



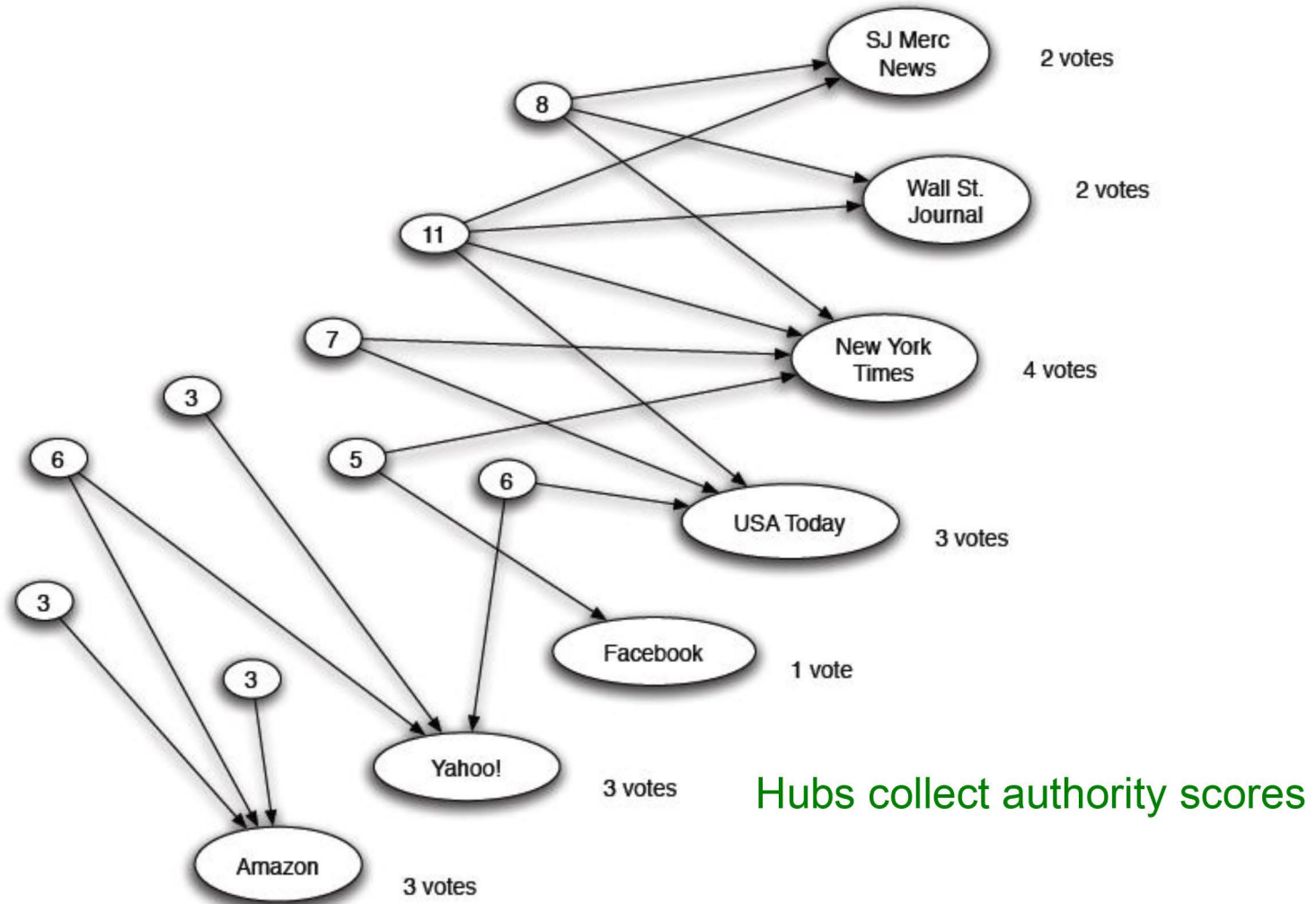
Counting in-links: Authority



Each page starts with **hub score 1**
Authorities collect their votes

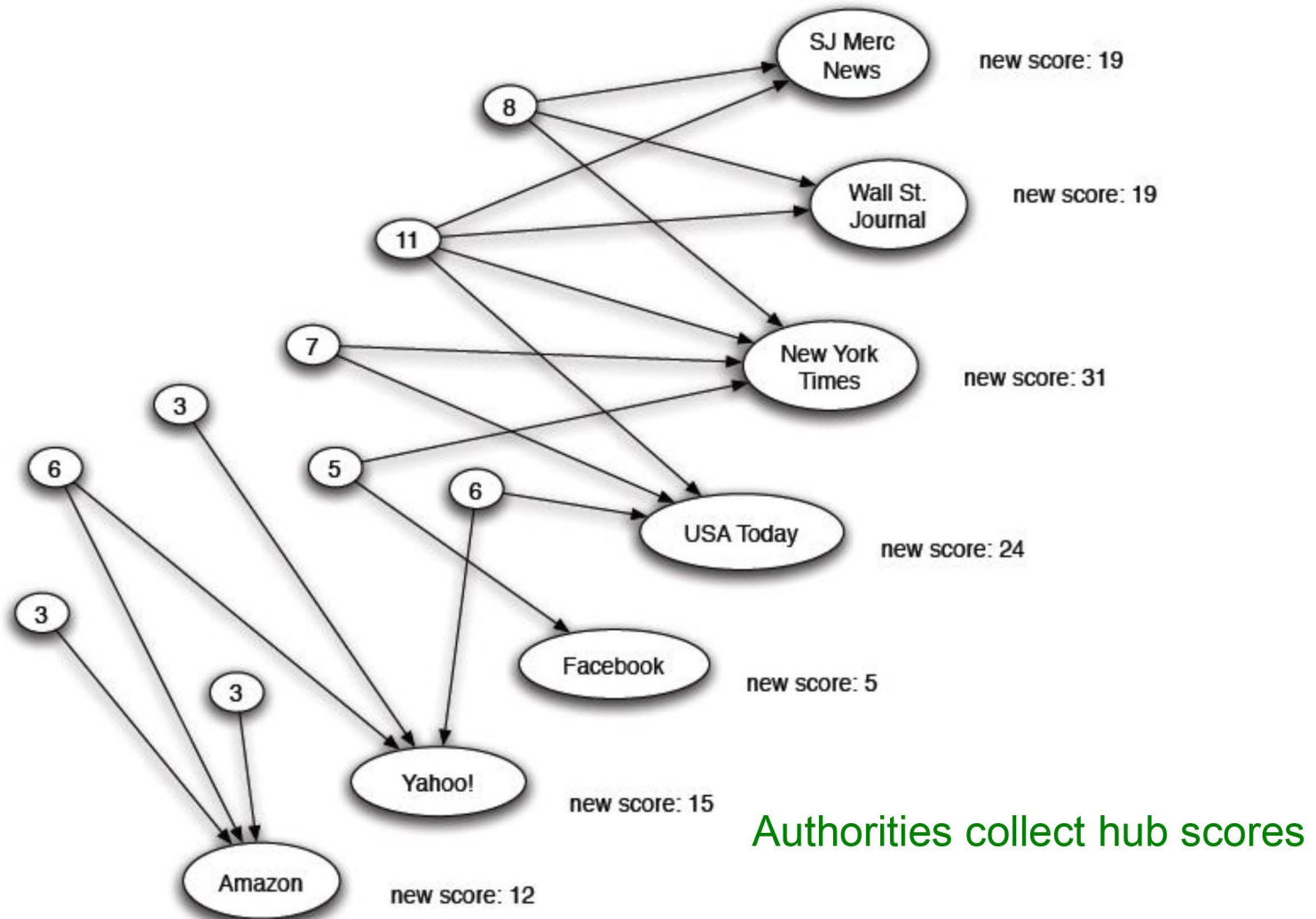
(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and the authority score)

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
 - Note a self-reinforcing recursive definition
- Model using two scores for each node:
 - Hub score and Authority score
 - Represented as vectors \mathbf{h} and \mathbf{a} , where the i -th element is the hub/authority score of the i -th node

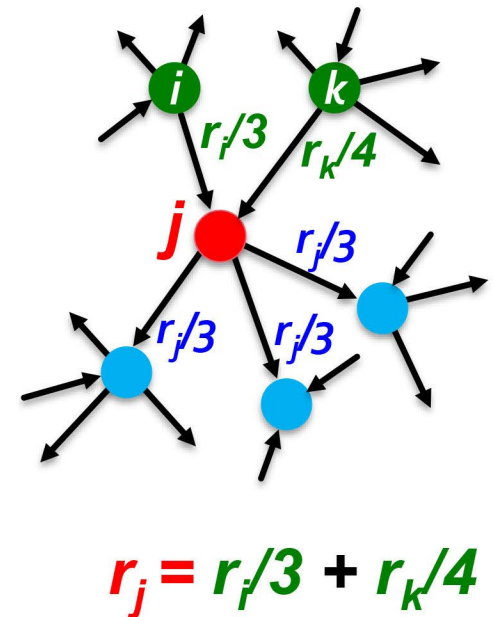
PageRank

Links as Votes

- **Still the same idea: Links as votes**
 - Page is more important if it has more links
 - In-coming links? Out-going links?
- **Think of in-links as votes:**
 - www.stanford.edu has 23,400 in-links
 - www.joe-schmoe.com has 1 in-link
- **Are all in-links equal?**
 - Links from important pages count more
 - Recursive question!

PageRank: The “Flow” Model

- A “vote” from an important page is worth more:
 - Each link’s vote is proportional to the **importance** of its source page
 - If page i with importance r_i has d_i out-links, each link gets r_i / d_i votes
 - Page j ’s own importance r_j is the sum of the votes on its in-links



PageRank: The “Flow” Model

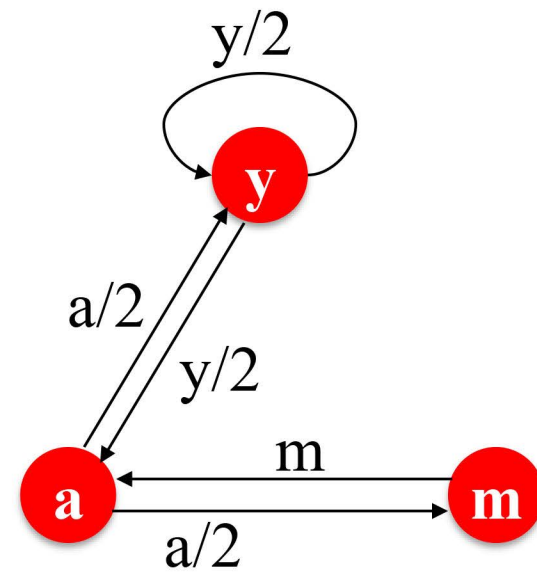
- A page is important if it is pointed to by other important pages
- Define a “rank” r_j for node j

$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

d_i ... out-degree of node i

You might wonder: Let’s just use Gaussian elimination to solve this system of linear equations. Bad idea!

The web in 1839



“Flow” equations:

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2 + r_m$$

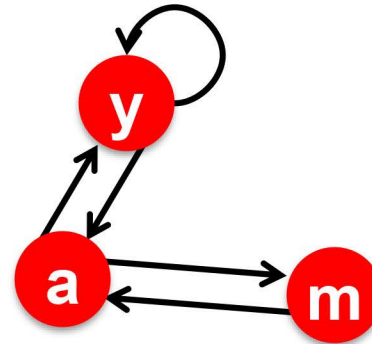
$$r_m = r_a/2$$

PageRank: How to solve?

PageRank: How to solve?

■ Power Iteration:

- Set $r_j \leftarrow 1/N$
- **1:** $r'_j \leftarrow \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- **2:** $r \leftarrow r'$
- If $|r - r'| > \varepsilon$: goto **1**



	y	a	m
y	1/2	1/2	0
a	1/2	0	1
m	0	1/2	0

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

■ Example:

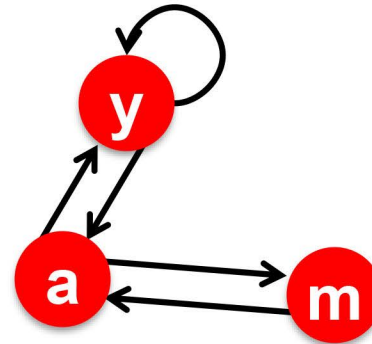
$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{matrix} 1/3 \\ 1/3 \\ 1/3 \end{matrix}$$

Iteration 0, 1, 2, ...

PageRank: How to solve?

■ Power Iteration:

- Set $r_j \leftarrow 1/N$
- **1:** $r'_j \leftarrow \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- **2:** $r \leftarrow r'$
- If $|r - r'| > \varepsilon$: goto **1**



	y	a	m
y	1/2	1/2	0
a	1/2	0	1
m	0	1/2	0

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

■ Example:

$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{matrix} 1/3 & 1/3 & 5/12 & 9/24 & & 6/15 \\ 1/3 & 3/6 & 1/3 & 11/24 & \dots & 6/15 \\ 1/3 & 1/6 & 3/12 & 1/6 & & 3/15 \end{matrix}$$

Iteration 0, 1, 2, ...

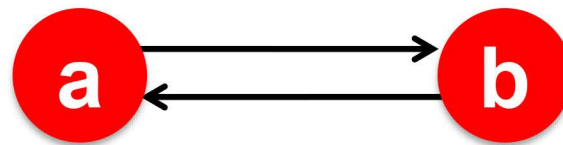
PageRank: Three Questions

$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i} \quad \text{or equivalently} \quad r = Mr$$

- Does this converge?
- Does it converge to what we want?
- Are results reasonable?

Does this converge?

- The “Spider trap” problem:



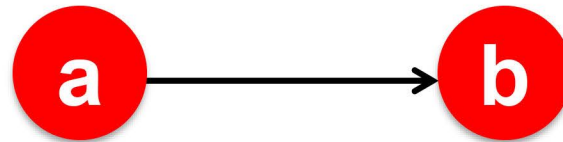
$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

- Example:

	Iteration: 0,	1,	2,	3...
r_a	1	0	1	0
r_b	0	1	0	1

Does it converge to what we want?

- The “Dead end” problem:



$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

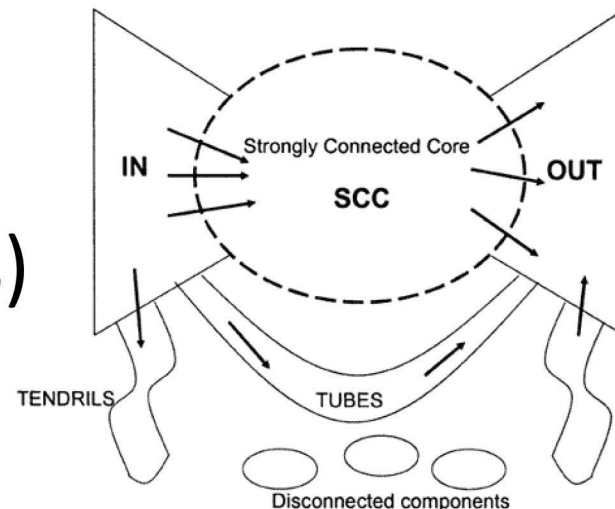
- Example:

	Iteration: 0,	1,	2,	3...
r_a	1	0	0	0
r_b	0	1	0	0

RageRank: Problems

2 problems:

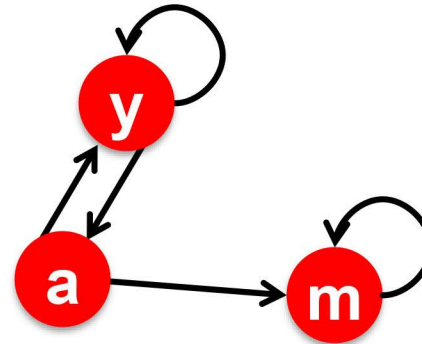
- **(1)** Some pages are **dead ends** (have no out-links)
 - Such pages cause importance to “leak out”
- **(2) Spider traps** (all out-links are within the group)
 - Eventually spider traps absorb all importance



Problem: Spider Traps

■ Power Iteration:

- Set $r_j = \frac{1}{N}$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
 - And iterate



	y	a	m
y	1/2	1/2	0
a	1/2	0	0
m	0	1/2	1

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2$$

$$r_m = r_a/2 + r_m$$

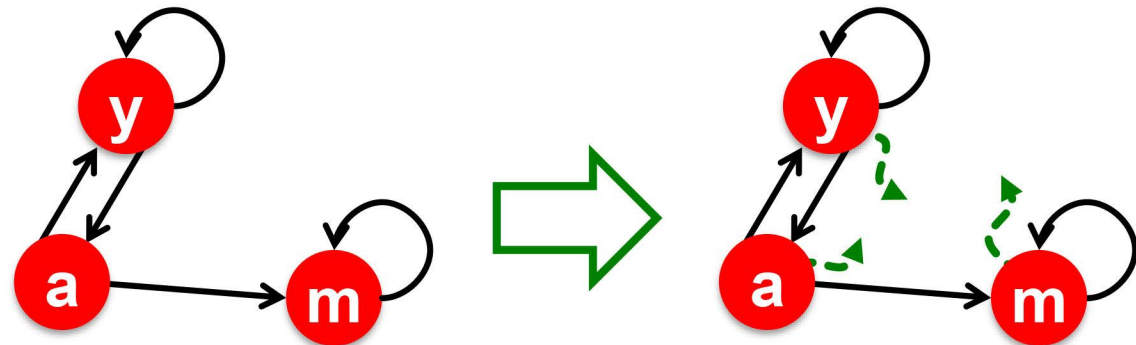
■ Example:

$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{array}{c|c|c|c|c|c} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 3/6 & 7/12 & 16/24 & & 1 \end{array}$$

Iteration 0, 1, 2, ...

Solution: Random Teleports

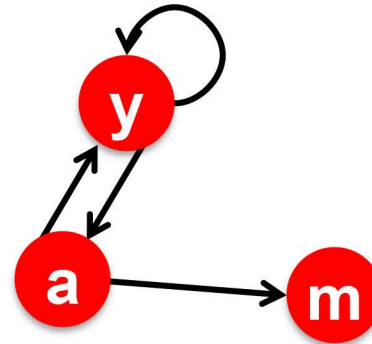
- The Google solution for spider traps: **At each time step, the random surfer has two options**
 - With prob. β , follow a link at random
 - With prob. $1-\beta$, jump to a random page
 - Common values for β are in the range 0.8 to 0.9
- **Surfer will teleport out of spider trap within a few time steps**



Problem: Dead Ends

■ Power Iteration:

- Set $r_j = \frac{1}{N}$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
 - And iterate



	y	a	m
y	1/2	1/2	0
a	1/2	0	0
m	0	1/2	0

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2$$

$$r_m = r_a/2$$

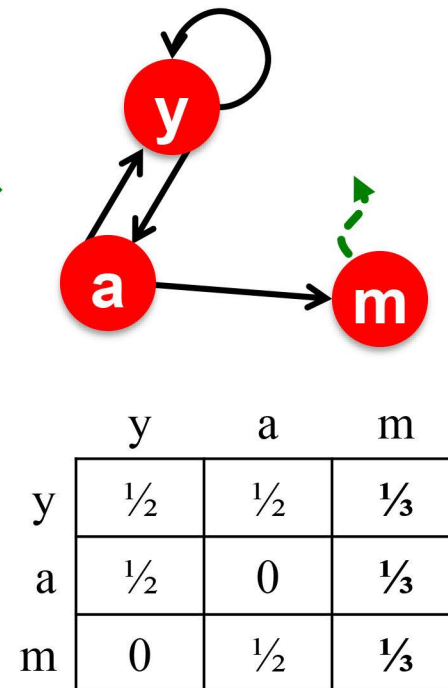
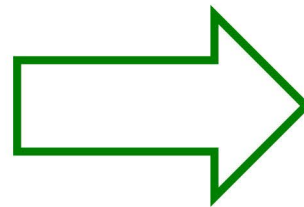
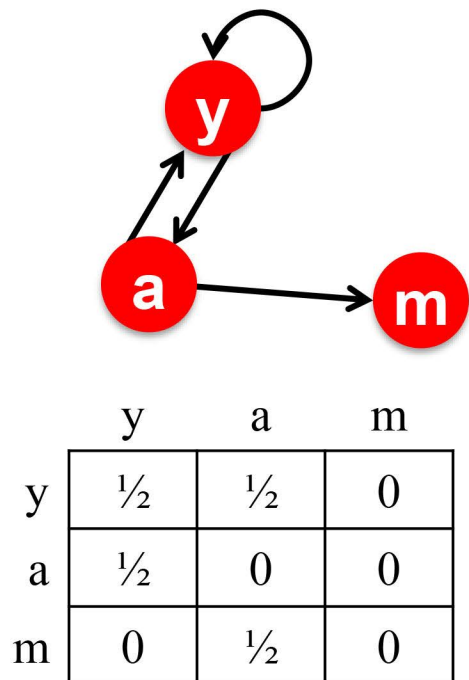
■ Example:

$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{array}{c|c|c|c|c|c}
 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\
 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\
 1/3 & 1/6 & 1/12 & 2/24 & & 0
 \end{array}$$

Iteration 0, 1, 2, ...

Solution: Always Teleport

- **Teleports:** Follow random teleport links with probability **1.0** from dead-ends
 - Adjust matrix accordingly



Final PageRank Equation

- **Google's solution:** At each step, random surfer has two options:
 - With probability β , follow a link at random
 - With probability $1-\beta$, jump to some random page
- **PageRank equation** [Brin-Page, '98]

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{n}$$

d_i ... out-degree of node i

The above formulation assumes that M has no dead ends. We can either preprocess matrix M (bad!) or explicitly follow random teleport links with probability 1.0 from dead-ends. See P. Berkhin, *A Survey on PageRank Computing*, Internet Mathematics, 2005.

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