

# Random Walks on Graphs

Based on materials

by J. Leskovec, A. Rajaraman, J. Ullman:  
Mining of Massive Datasets, <http://www.mmds.org>

And

by LalaAdamic and Purnamrita Sarkar

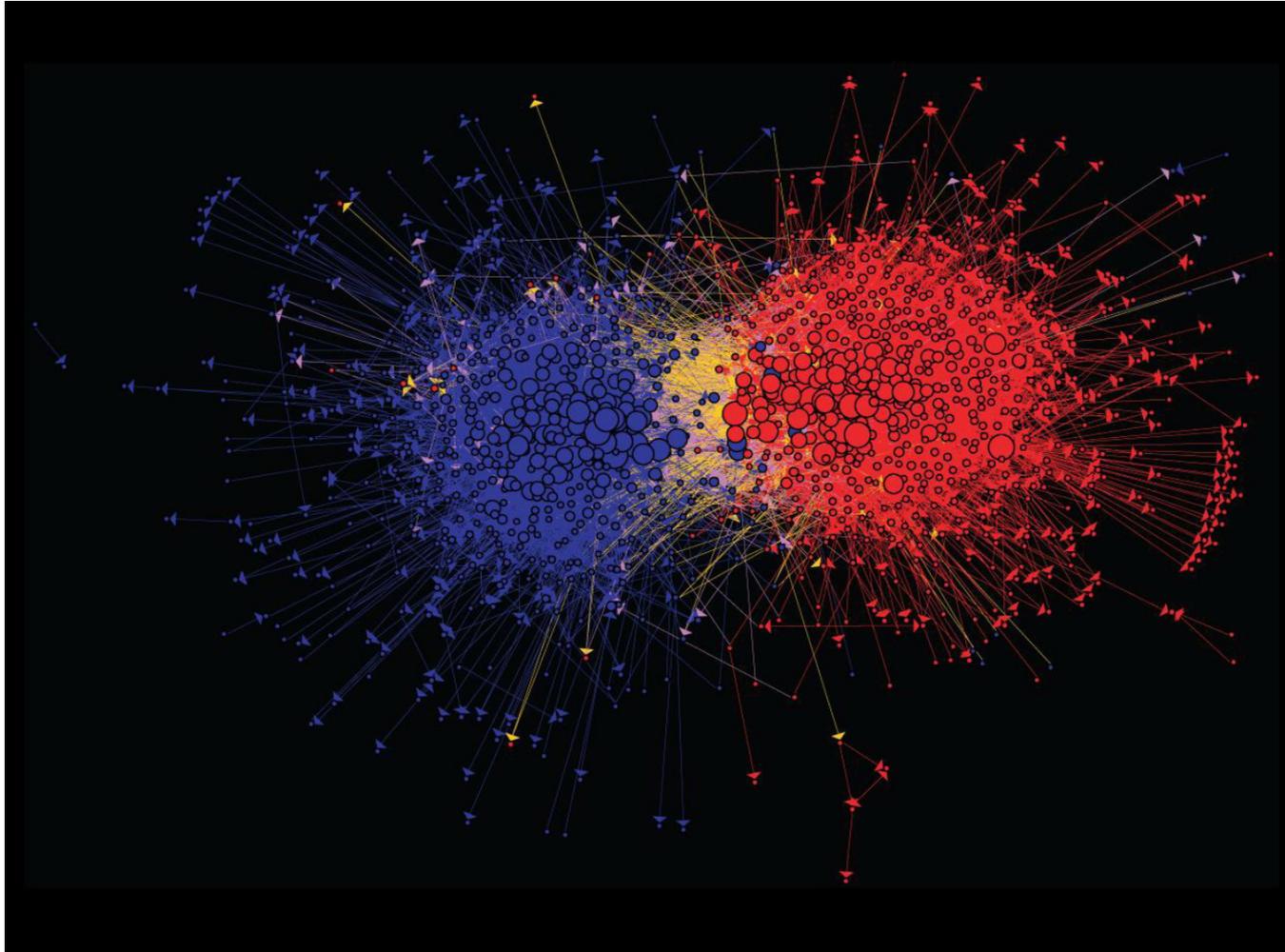
# Graph Data: Social Networks



## Facebook social graph

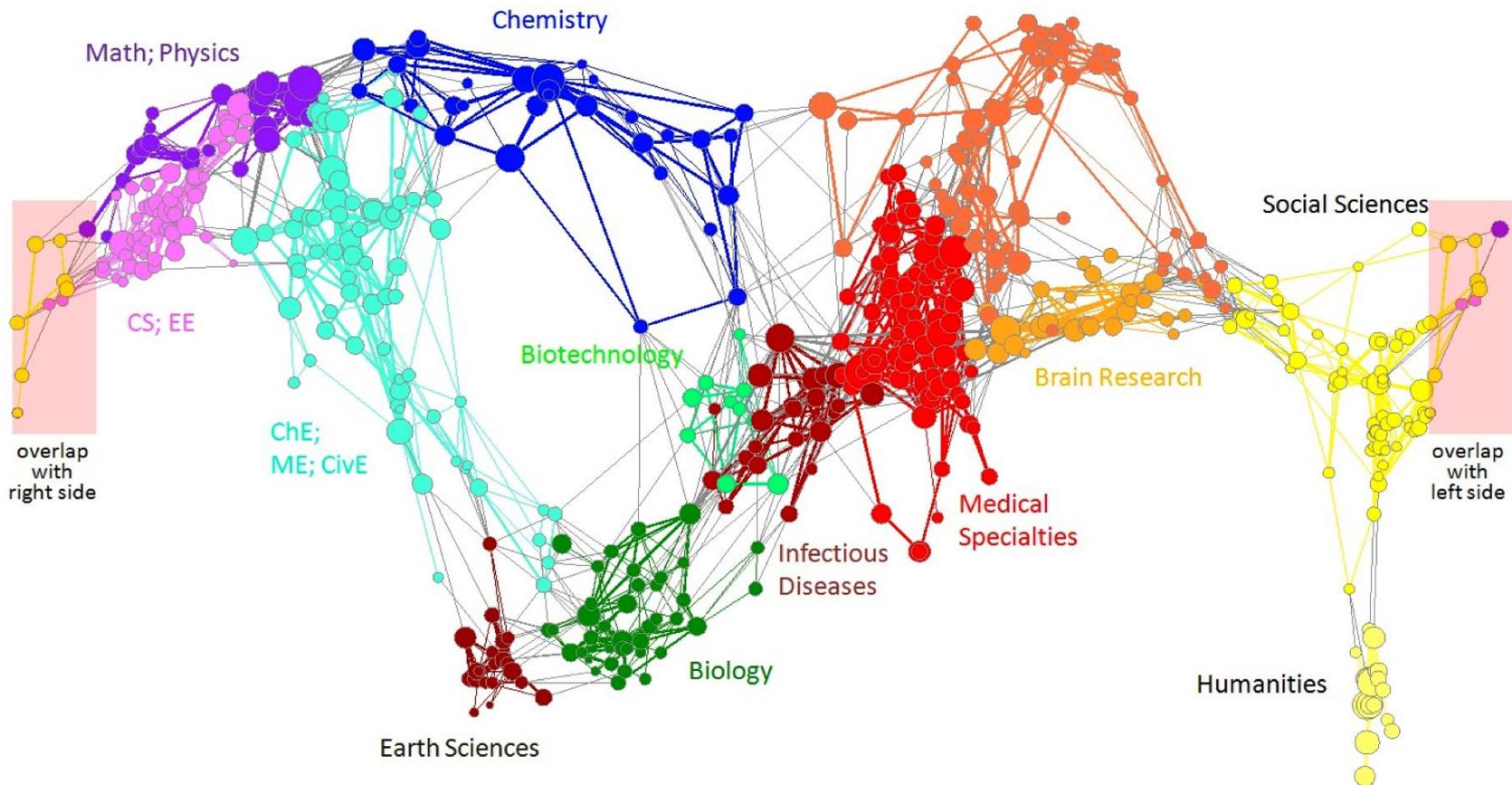
4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]

# Graph Data: Media Networks



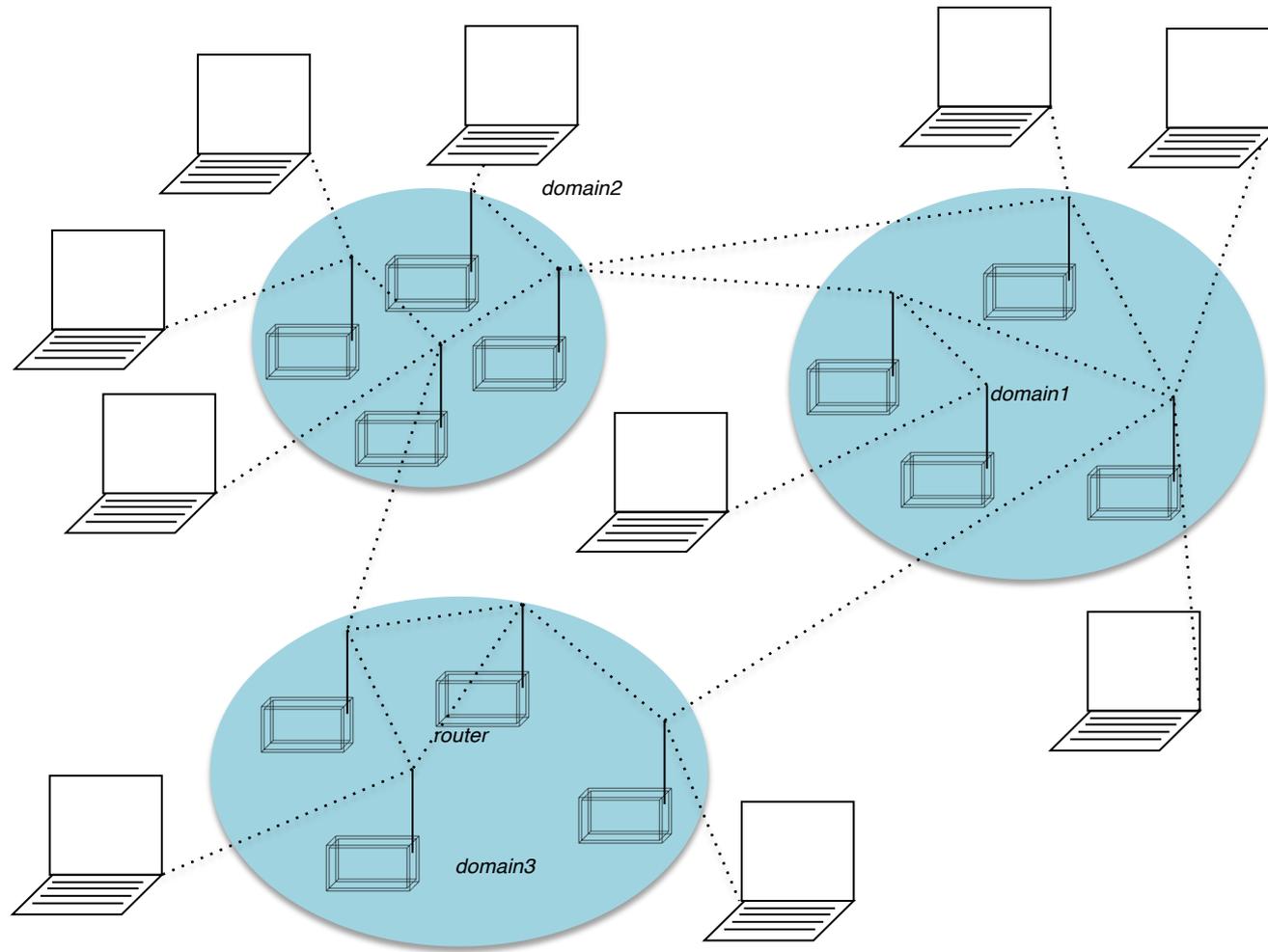
**Connections between political blogs**  
Polarization of the network [Adamic-Glance, 2005]

# Graph Data: Information Nets



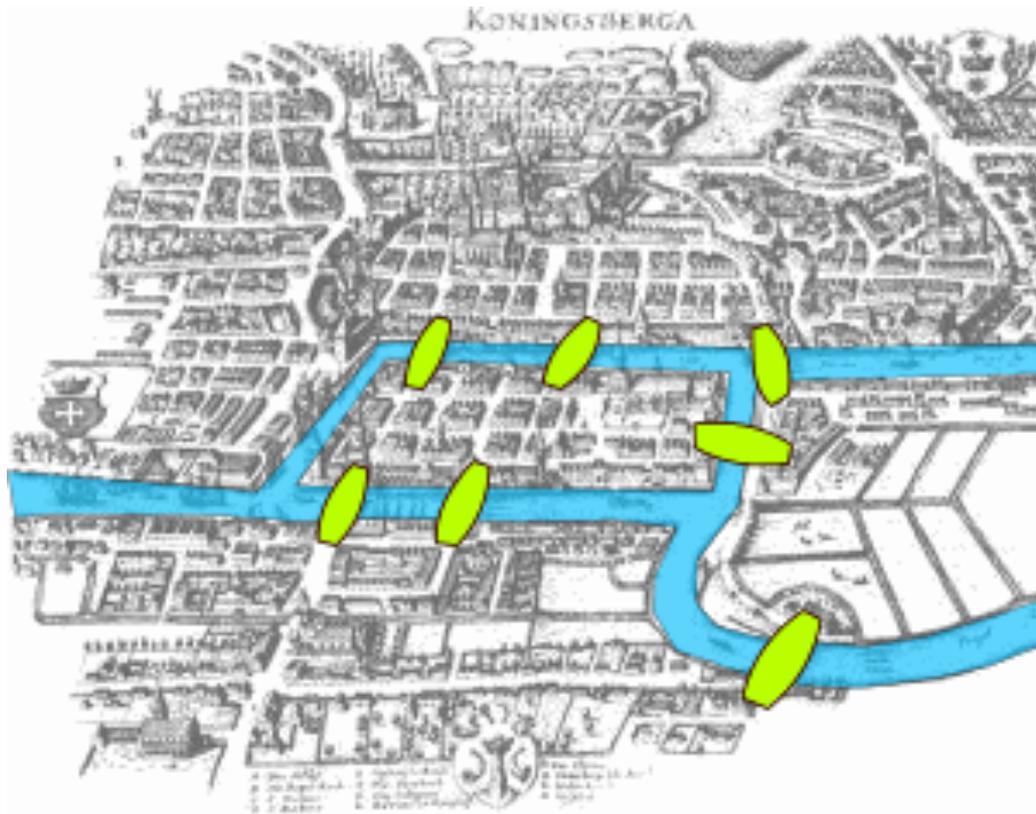
**Citation networks and Maps of science**  
[Börner et al., 2012]

# Graph Data: Communication Nets



# Internet

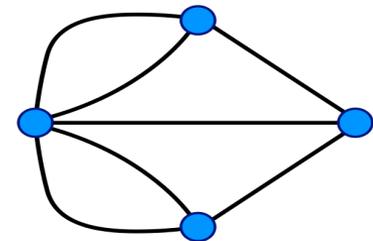
# Graph Data: Technological Networks



## Seven Bridges of Königsberg

[Euler, 1735]

Return to the starting point by traveling each link of the graph once and only once.



# Web as a Graph

- **Web as a directed graph:**
  - **Nodes: Webpages**
  - **Edges: Hyperlinks**

I teach a  
class on  
Networks.

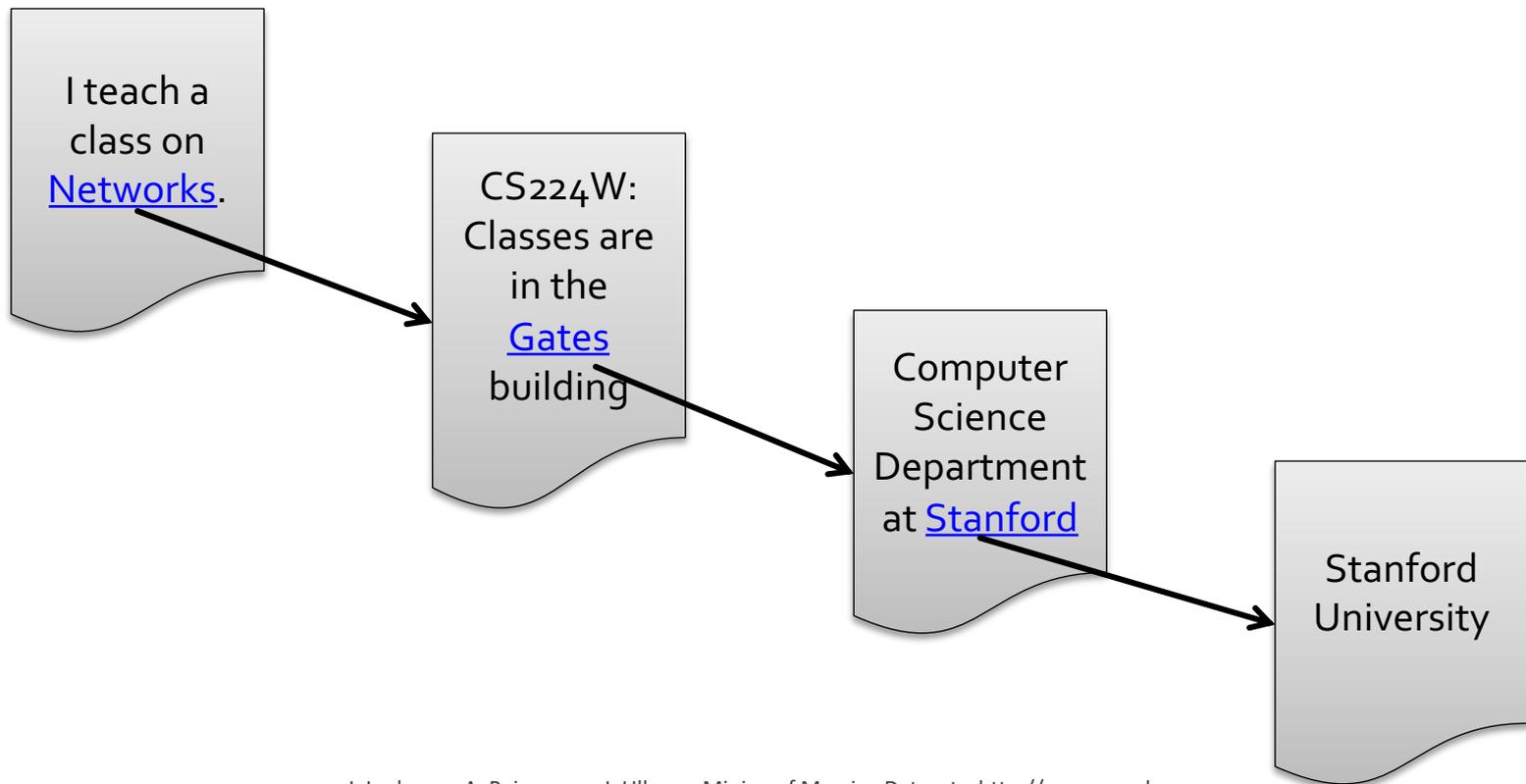
CS224W:  
Classes are  
in the  
Gates  
building

Computer  
Science  
Department  
at Stanford

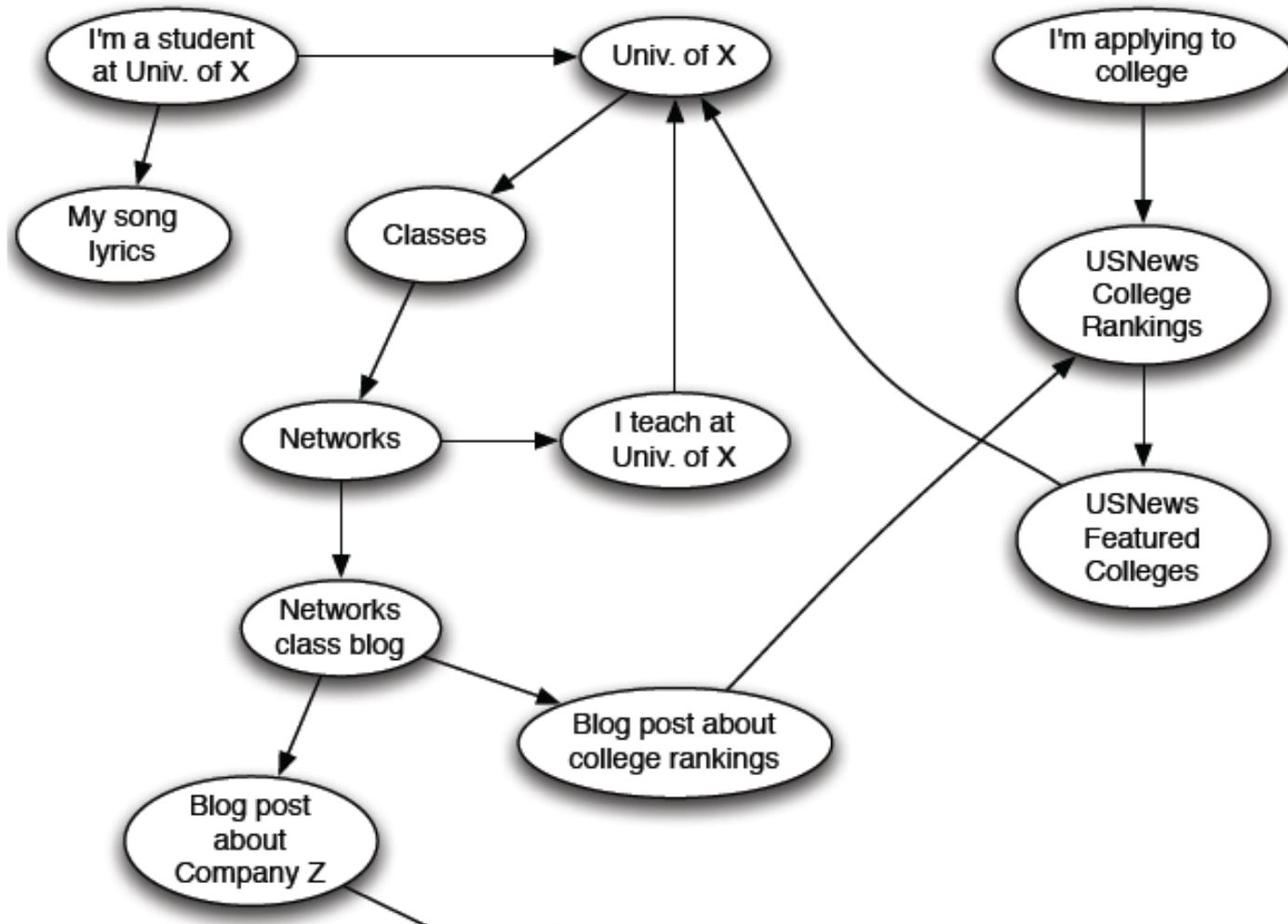
Stanford  
University

# Web as a Graph

- **Web as a directed graph:**
  - **Nodes: Webpages**
  - **Edges: Hyperlinks**



# Web as a Directed Graph



# Broad Question

- **How to organize the Web?**
- **First try: Human curated Web directories**
  - Yahoo, DMOZ, LookSmart
- **Second try: Web Search**
  - **Information Retrieval** investigates:  
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.



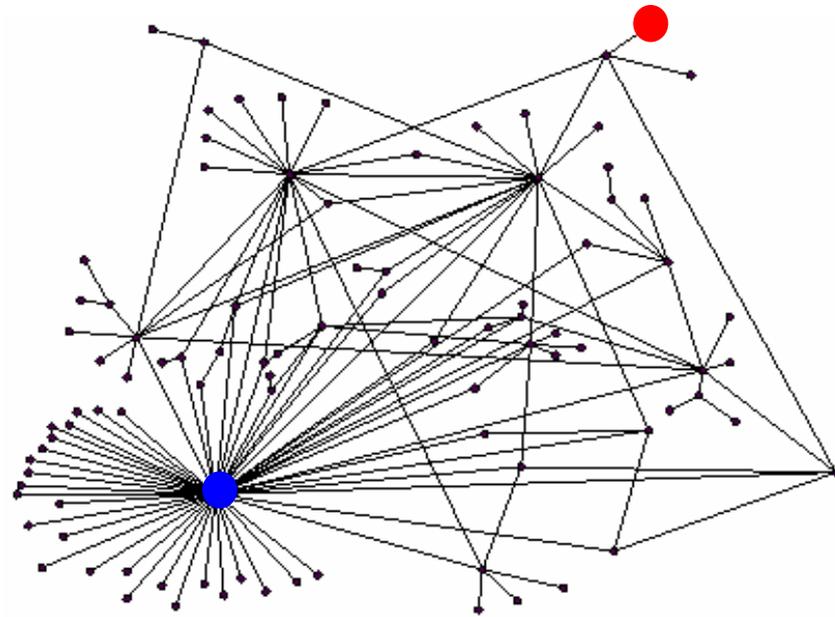
# Web Search: 2 Challenges

## 2 challenges of web search:

- (1) Web contains many sources of information  
Who to “trust”?
  - **Trick:** Trustworthy pages may point to each other!
- (2) What is the “best” answer to query “newspaper”?
  - No single right answer
  - **Trick:** 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”
- There is large diversity in the web-graph node connectivity.  
**Let's rank the pages by the link structure!**



# Link Analysis Algorithms

- We will cover the following **Link Analysis approaches** for computing **importances** of nodes in a graph:
  - Page Rank
  - Topic-Specific (Personalized) Page Rank
  - Web Spam Detection Algorithms

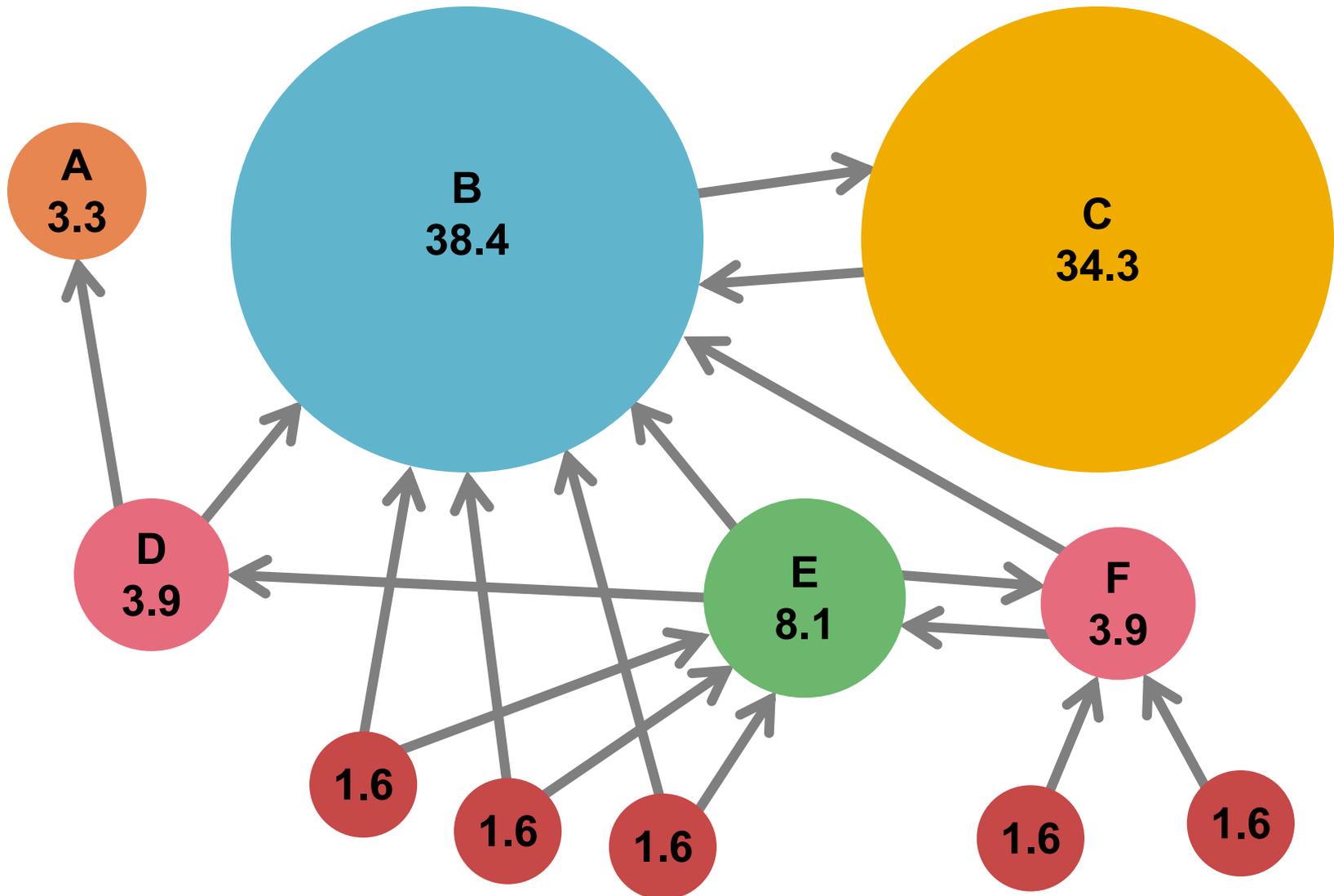
# PageRank: The “Flow” Formulation

---

# Links as Votes

- **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](http://www.stanford.edu) has 23,400 in-links
- **Are all in-links are equal?**
  - **Links from important pages count more**
  - Recursive question!

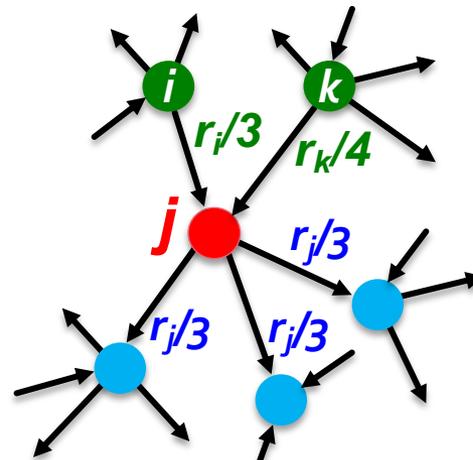
# Example: PageRank Scores



# Simple Recursive Formulation

- Each link's vote is proportional to the **importance** of its source page
- If page  $j$  with importance  $r_j$  has  $n$  out-links, each link gets  $r_j/n$  votes
- Page  $j$ 's own importance is the sum of the votes on its in-links

$$r_j = r_i/3 + r_k/4$$



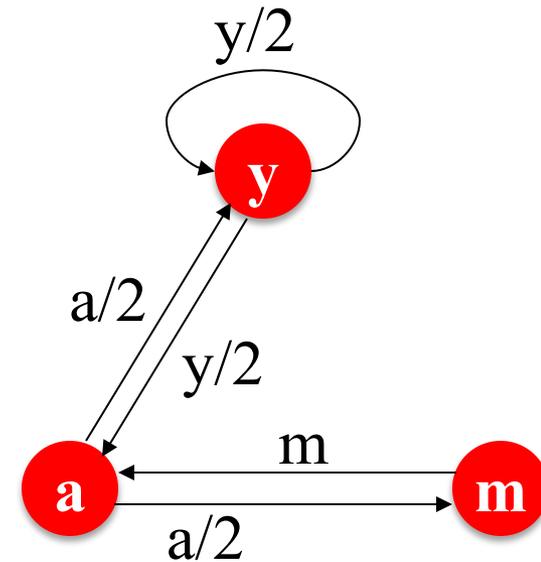
# PageRank: The “Flow” Model

- A “vote” from an important page is worth more
- A page is important if it is pointed to by other important pages
- Define a “rank”  $r_j$  for page  $j$

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

$d_i$  ... out-degree of node  $i$

The web in 1839



“Flow” equations:

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

# Solving the Flow Equations

- **3 equations, 3 unknowns, no constants**

Flow equations:

$$\mathbf{r}_y = \mathbf{r}_y / 2 + \mathbf{r}_a / 2$$
$$\mathbf{r}_a = \mathbf{r}_y / 2 + \mathbf{r}_m$$
$$\mathbf{r}_m = \mathbf{r}_a / 2$$

- No unique solution
- All solutions equivalent modulo the scale factor
- **Additional constraint forces uniqueness:**
  - $\mathbf{r}_y + \mathbf{r}_a + \mathbf{r}_m = 1$
  - **Solution:**  $\mathbf{r}_y = \frac{2}{5}$ ,  $\mathbf{r}_a = \frac{2}{5}$ ,  $\mathbf{r}_m = \frac{1}{5}$
- **Gaussian elimination method works for small examples, but we need a better method for large web-size graphs**
- **We need a new formulation!**

# PageRank: Matrix Formulation

- **Stochastic adjacency matrix  $M$**

- Let page  $i$  has  $d_i$  out-links

- If  $i \rightarrow j$ , then  $M_{ji} = \frac{1}{d_i}$  else  $M_{ji} = 0$

- $M$  is a **column stochastic matrix**

- Columns sum to 1

- **Rank vector  $r$** : vector with an entry per page

- $r_i$  is the importance score of page  $i$

- $\sum_i r_i = 1$

- **The flow equations can be written**

$$r = M \cdot r$$

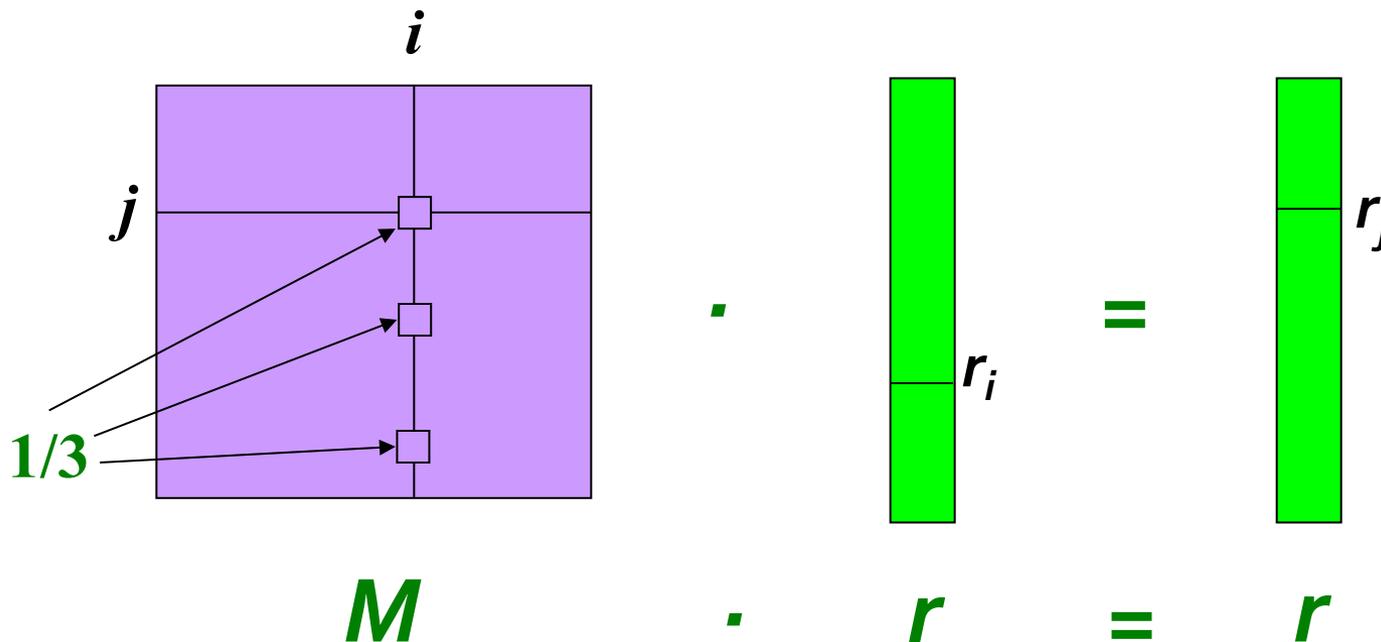
$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

# Example

- Remember the flow equation:  $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- Flow equation in the matrix form

$$M \cdot r = r$$

- Suppose page  $i$  links to 3 pages, including  $j$



# Eigenvector Formulation

- The flow equations can be written

$$r = M \cdot r$$

- So the rank vector  $r$  is an eigenvector of the stochastic web matrix  $M$

- In fact, its first or principal eigenvector, with corresponding eigenvalue  $1$

- Largest eigenvalue of  $M$  is  $1$  since  $M$  is column stochastic (with non-negative entries)

- *We know  $r$  is unit length and each column of  $M$  sums to one, so  $Mr \leq 1$*

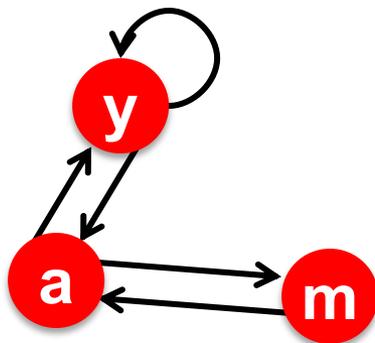
- We can now efficiently solve for  $r$ !

The method is called Power iteration

**NOTE:**  $x$  is an eigenvector with the corresponding eigenvalue  $\lambda$  if:

$$Ax = \lambda x$$

# Example: Flow Equations & M



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

$$r = M \cdot r$$

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

$$r_a = r_y / 2 + r_m$$

$$r_m = r_a / 2$$

$$\begin{array}{|c|} \hline y \\ \hline a \\ \hline m \\ \hline \end{array} = \begin{array}{|ccc|} \hline 1/2 & 1/2 & 0 \\ \hline 1/2 & 0 & 1 \\ \hline 0 & 1/2 & 0 \\ \hline \end{array} \begin{array}{|c|} \hline y \\ \hline a \\ \hline m \\ \hline \end{array}$$

# Power Iteration Method

- Given a web graph with  $n$  nodes, where the nodes are pages and edges are hyperlinks
- **Power iteration:** a simple iterative scheme

- Suppose there are  $N$  web pages

- Initialize:  $\mathbf{r}^{(0)} = [1/N, \dots, 1/N]^T$

- Iterate:  $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$

- Stop when  $\|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}\|_1 < \varepsilon$

$\|\mathbf{x}\|_1 = \sum_{1 \leq i \leq N} |x_i|$  is the **L1** norm

Can use any other vector norm, e.g., Euclidean

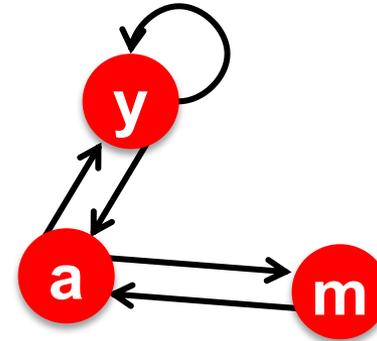
$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

$d_i$  .... out-degree of node  $i$

# PageRank: How to solve?

## ■ Power Iteration:

- Set  $r_j = 1/N$
- **1:**  $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- **2:**  $r = r'$
- 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:

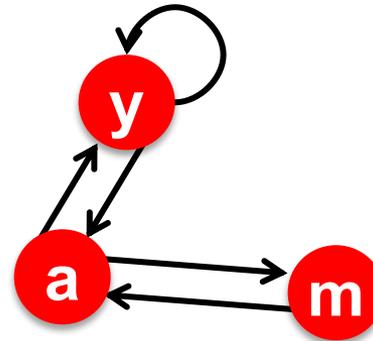
$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{matrix} 1/3 \\ 1/3 \\ 1/3 \end{matrix}$$

Iteration 0, 1, 2, ...

# PageRank: How to solve?

## ■ Power Iteration:

- Set  $r_j = 1/N$
- **1:**  $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- **2:**  $r = r'$
- 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, ...

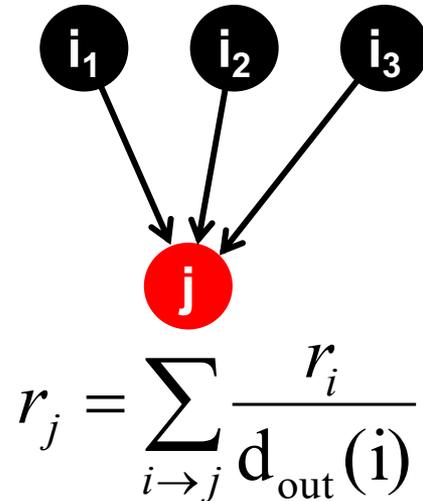
# Random Walk Interpretation

- **Imagine a random web surfer:**

- At any time  $t$ , surfer is on some page  $i$
- At time  $t + 1$ , the surfer follows an out-link from  $i$  uniformly at random
- Ends up on some page  $j$  linked from  $i$
- Process repeats indefinitely

- **Let:**

- $\mathbf{p}(t)$  ... vector whose  $i^{\text{th}}$  coordinate is the prob. that the surfer is at page  $i$  at time  $t$
- So,  $\mathbf{p}(t)$  is a probability distribution over pages

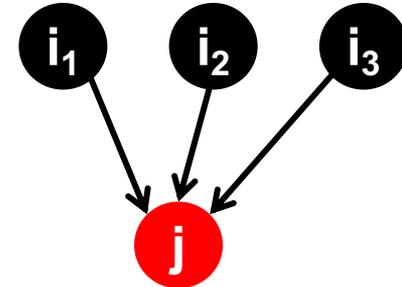


# The Stationary Distribution

- **Where is the surfer at time  $t+1$ ?**

- Follows a link uniformly at random

$$\mathbf{p}(t+1) = \mathbf{M} \cdot \mathbf{p}(t)$$



$$p(t+1) = M \cdot p(t)$$

- Suppose the random walk reaches a state

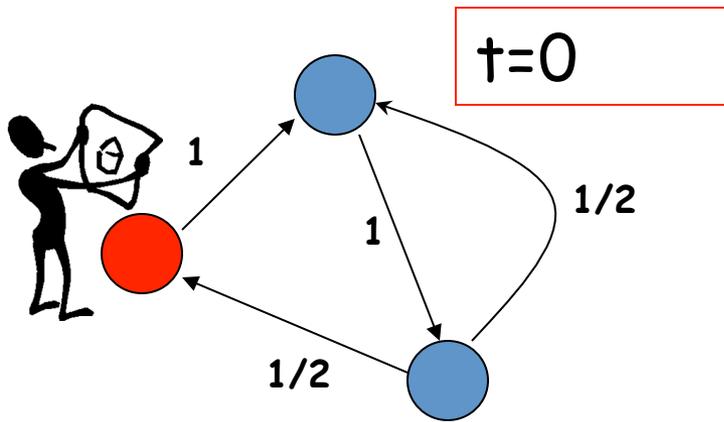
$$\mathbf{p}(t+1) = \mathbf{M} \cdot \mathbf{p}(t) = \mathbf{p}(t)$$

then  $\mathbf{p}(t)$  is **stationary distribution** of a random walk

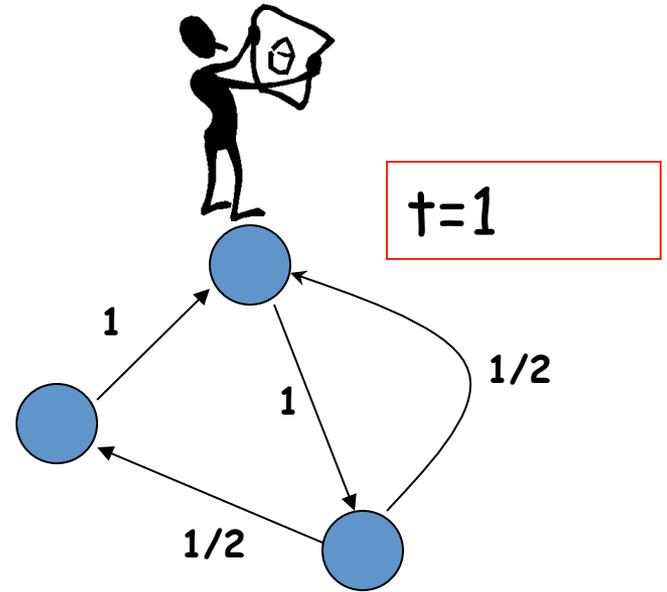
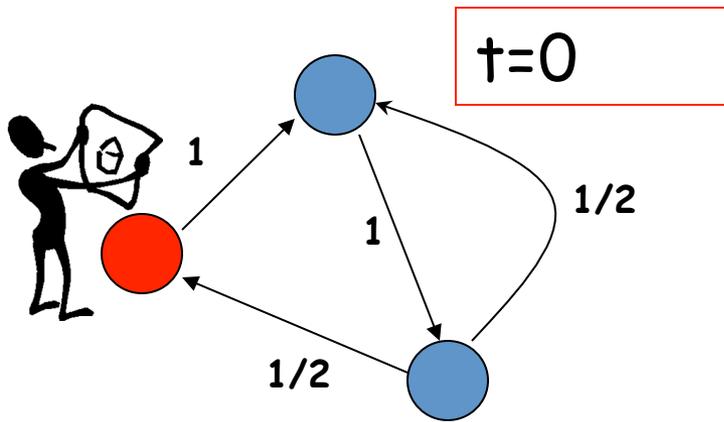
- **Our original rank vector  $\mathbf{r}$  satisfies  $\mathbf{r} = \mathbf{M} \cdot \mathbf{r}$**

- **So,  $\mathbf{r}$  is a stationary distribution for the random walk**

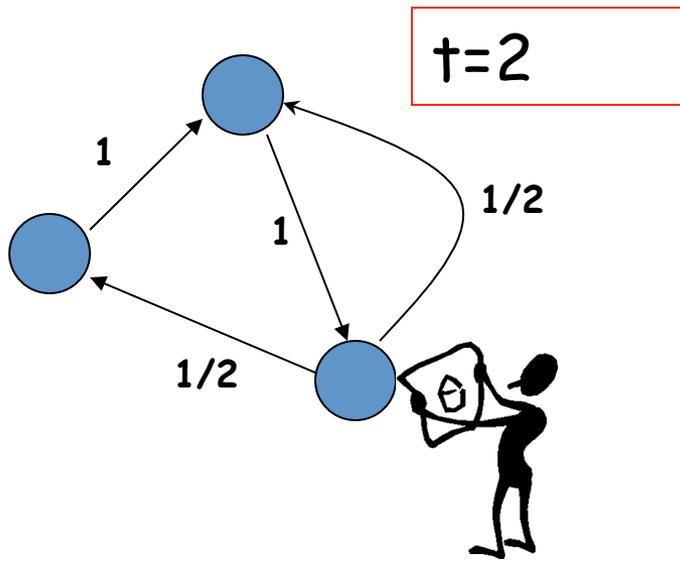
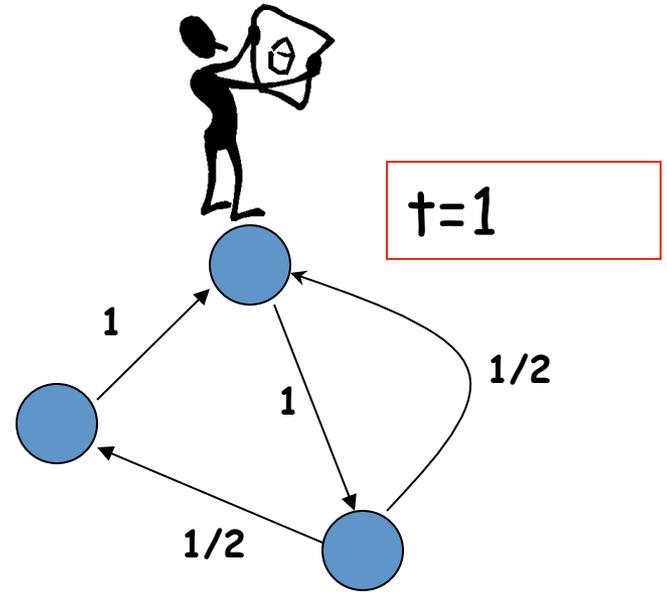
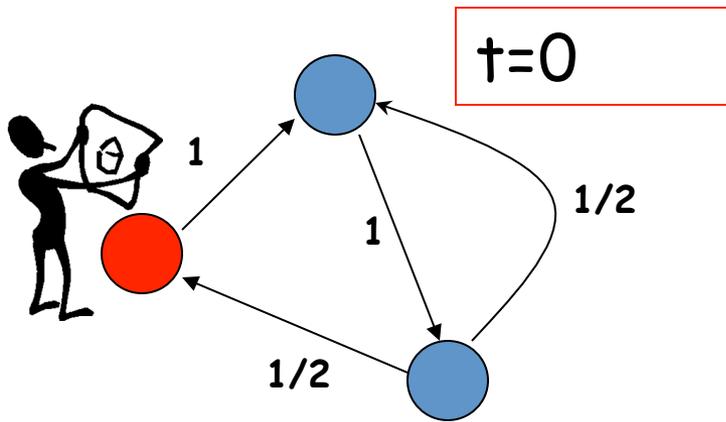
# What is a random walk



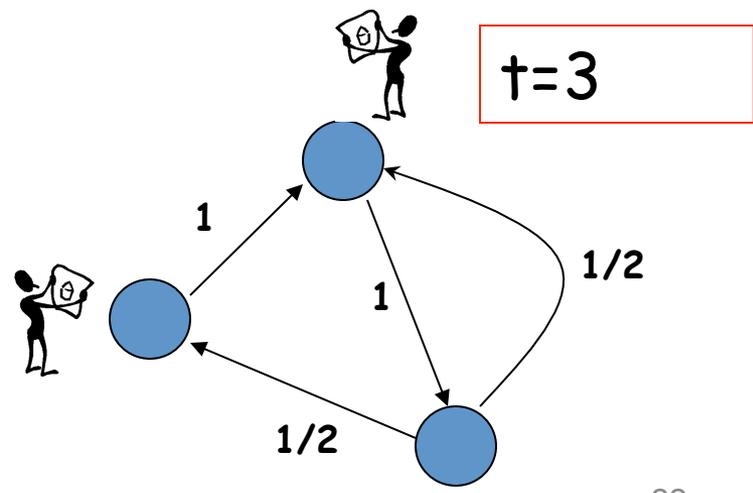
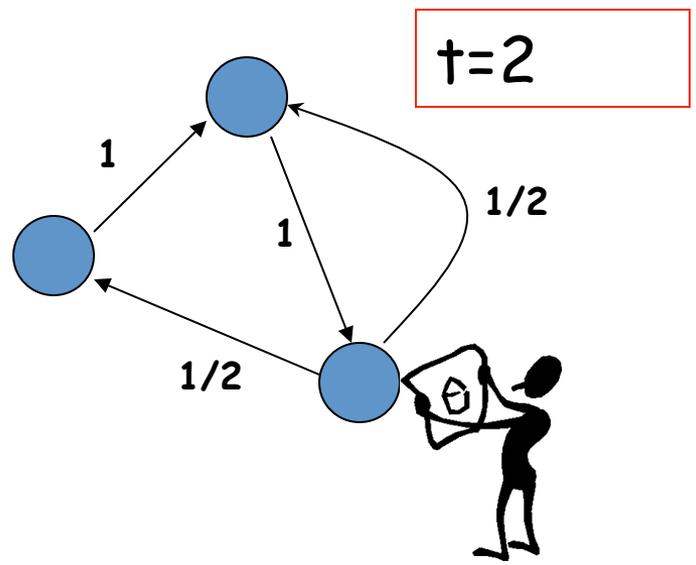
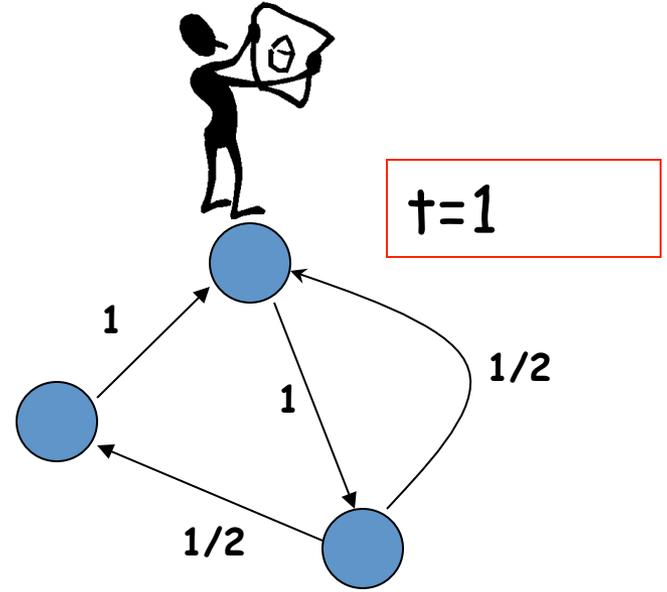
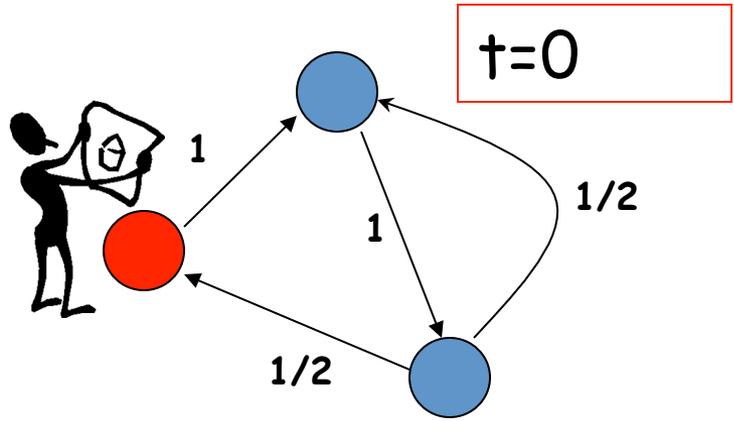
# What is a random walk



# What is a random walk



# What is a random walk



# Existence and Uniqueness

- A central result from the theory of random walks (a.k.a. Markov processes):

For graphs that satisfy **certain conditions**, the **stationary distribution is unique** and eventually will be reached no matter what the initial probability distribution at time  $t = 0$

# PageRank: The Google Formulation

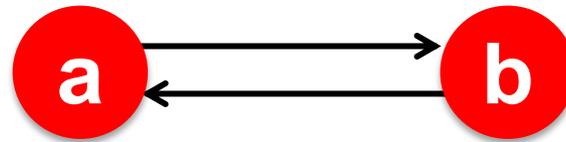
---

# PageRank: Three Questions

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

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

# Does this converge?



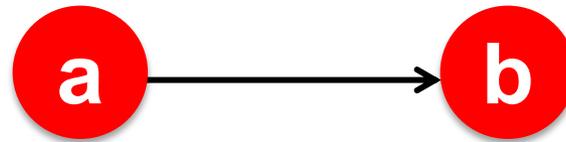
$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

## ■ Example:

$$\begin{array}{l} r_a \\ r_b \end{array} = \begin{array}{cccc} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{array}$$

Iteration 0, 1, 2, ...

# Does it converge to what we want?



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

## ■ Example:

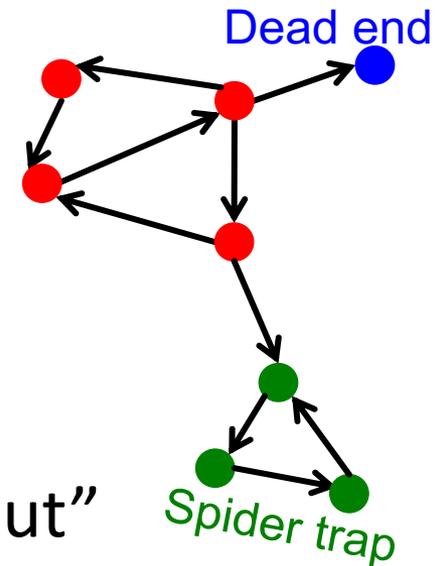
$$\begin{array}{l} r_a \\ r_b \end{array} = \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{array}$$

Iteration 0, 1, 2, ...

# PageRank: Problems

## 2 problems:

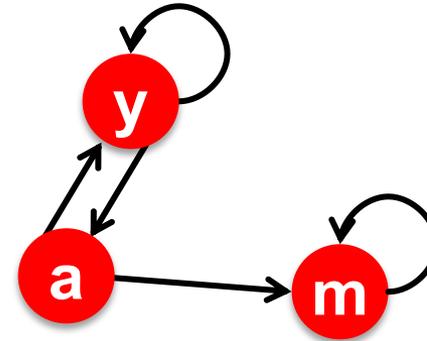
- **(1)** Some pages are **dead ends** (have no out-links)
  - Random walk has “nowhere” to go to
  - Such pages cause importance to “leak out”
- **(2) Spider traps:**  
(all out-links are within the group)
  - Random walked gets “stuck” in a trap
  - And eventually spider traps absorb all importance



# Problem: Spider Traps

## Power Iteration:

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



m is a spider trap

	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:

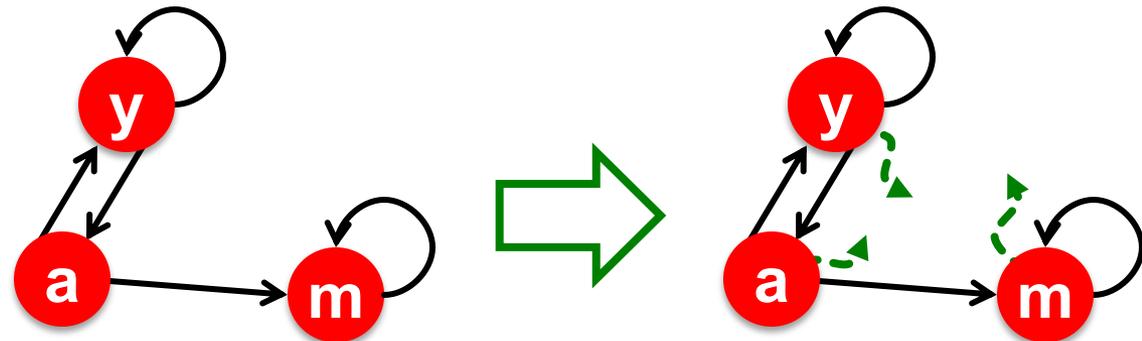
$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 3/6 & 7/12 & 16/24 & & 1 \end{matrix}$$

Iteration 0, 1, 2, ...

All the PageRank score gets "trapped" in node m.

# Solution: Teleports!

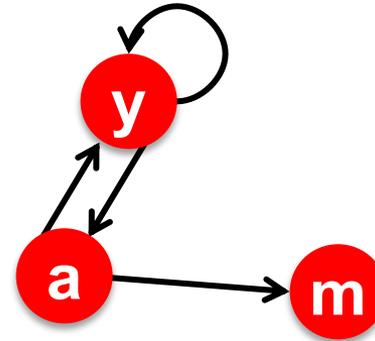
- **The Google solution for spider traps: At each time step, the random surfer has two options**
  - With prob.  $\beta$ , follow a link at random
  - With prob.  $1-\beta$ , jump to some random page
  - Common values for  $\beta$  are in the range 0.8 to 0.9
- **Surfer will teleport out of spider trap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



# Problem: Dead Ends

## ■ Power Iteration:

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

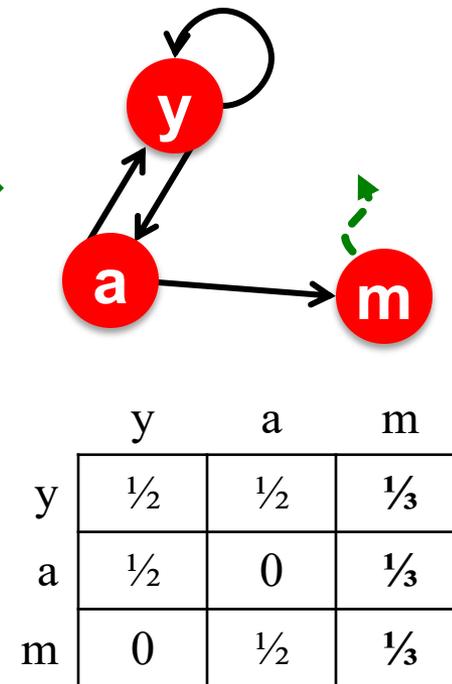
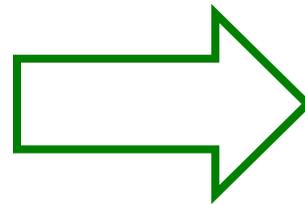
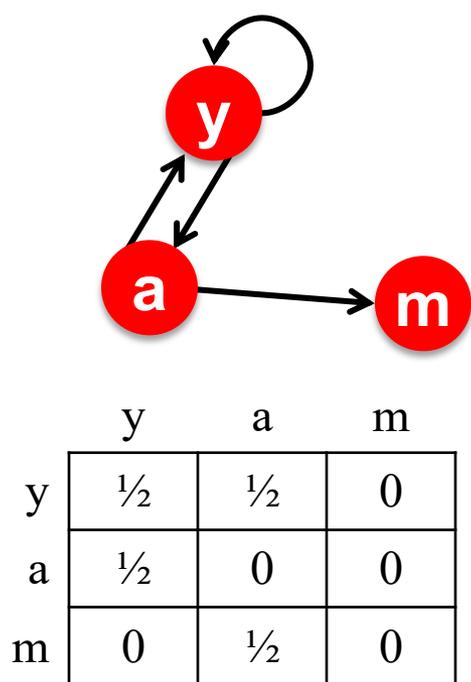
$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & & 0 \end{matrix}$$

Iteration 0, 1, 2, ...

Here the PageRank “leaks” out since the matrix is not stochastic.

# Solution: Always Teleport!

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



# Why Teleports Solve the Problem?

Why are dead-ends and spider traps a problem and **why do teleports solve the problem?**

- **Spider-traps** are not a problem, but with traps PageRank scores are **not** what we want
  - **Solution:** Never get stuck in a spider trap by teleporting out of it in a finite number of steps
- **Dead-ends** are a problem
  - The matrix is not column stochastic so our initial assumptions are not met
  - **Solution:** Make matrix column stochastic by always teleporting when there is nowhere else to go

# Solution: Random Teleports

- Google's solution that does it all:

At each step, random surfer has two options:

- With probability  $\beta$ , 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  $i$

This formulation assumes that  $M$  has no dead ends. We can either preprocess matrix  $M$  to remove all dead ends or explicitly follow random teleport links with probability 1.0 from dead-ends.

# The Google Matrix

- **PageRank equation** [Brin-Page, '98]

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

- **The Google Matrix  $A$ :**

[ $1/N$ ] $_{N \times N}$ ... $N$  by  $N$  matrix  
where all entries are  $1/N$

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

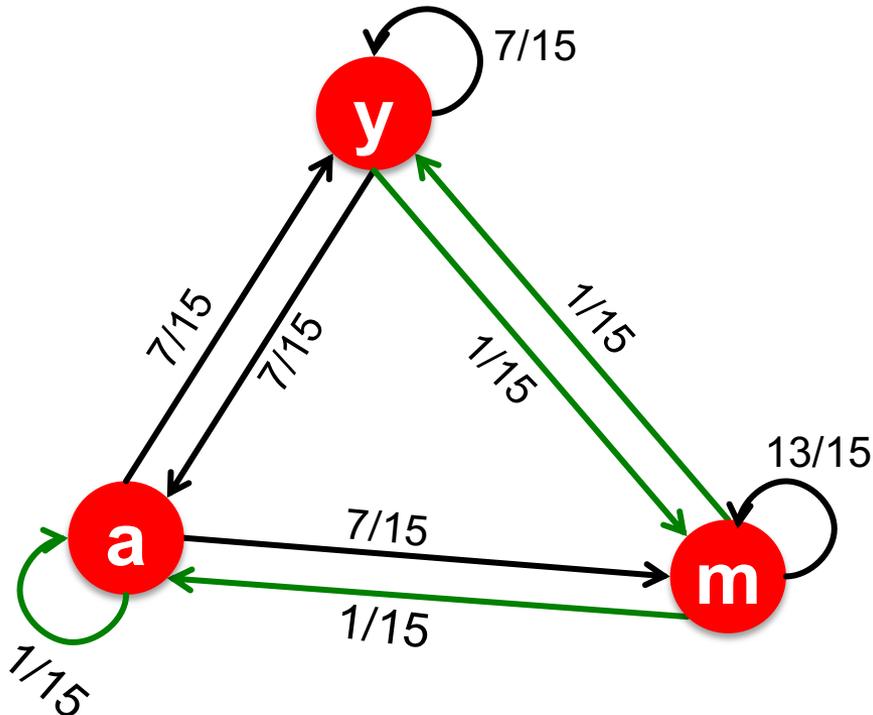
- **We have a recursive problem:  $r = A \cdot r$**

**And the Power method still works!**

- **What is  $\beta$ ?**

- In practice  $\beta = 0.8, 0.9$  (make 5 steps on avg., jump)

# Random Teleports ( $\beta = 0.8$ )



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

$$\begin{matrix} y \\ a \\ m \end{matrix} \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix}$$

**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