



# Link Analysis

Shantanu Jain



# PageRank

Using inbound links as “votes”

*Credit:* Yijun Zhao, Yi Wang,  
Tan et al., Leskovec et al.

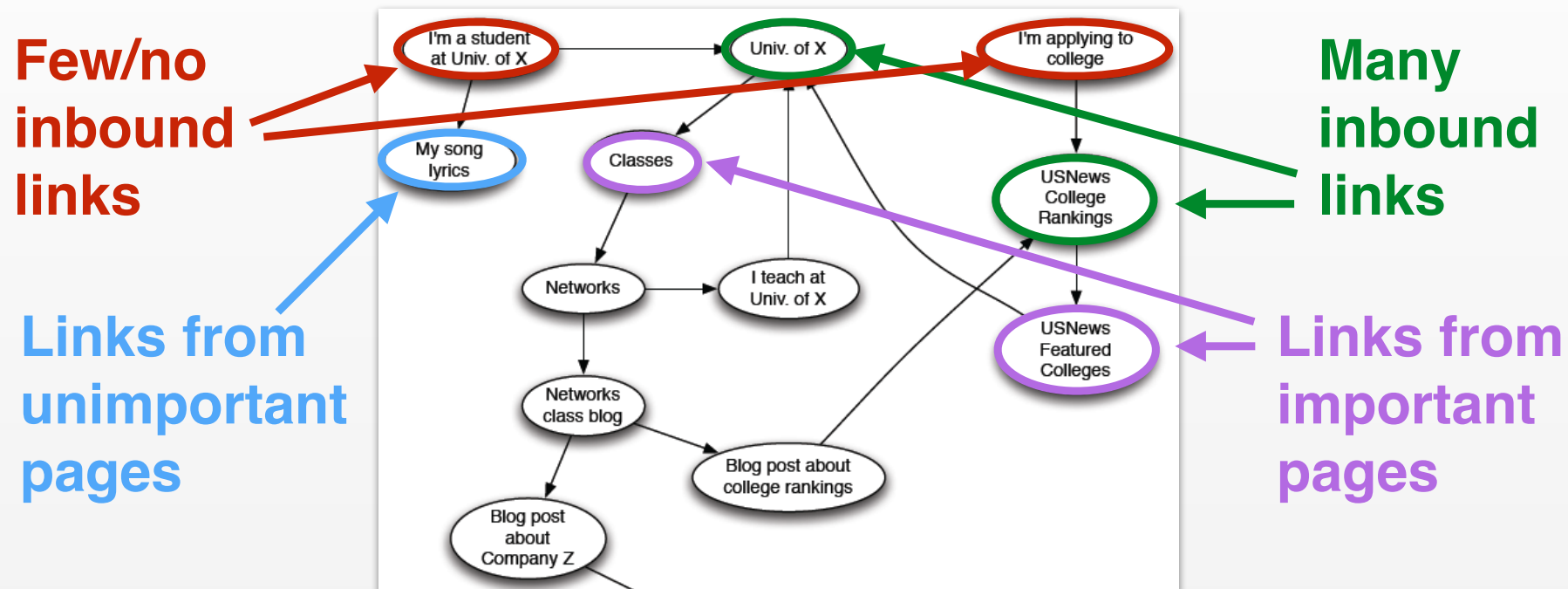
# Web search before PageRank



- Human-curated  
(e.g. Yahoo, Looksmart)
- Hand-written descriptions
- Wait time for inclusion
- Text-search  
(e.g. WebCrawler, Lycos)
- Prone to term-spam

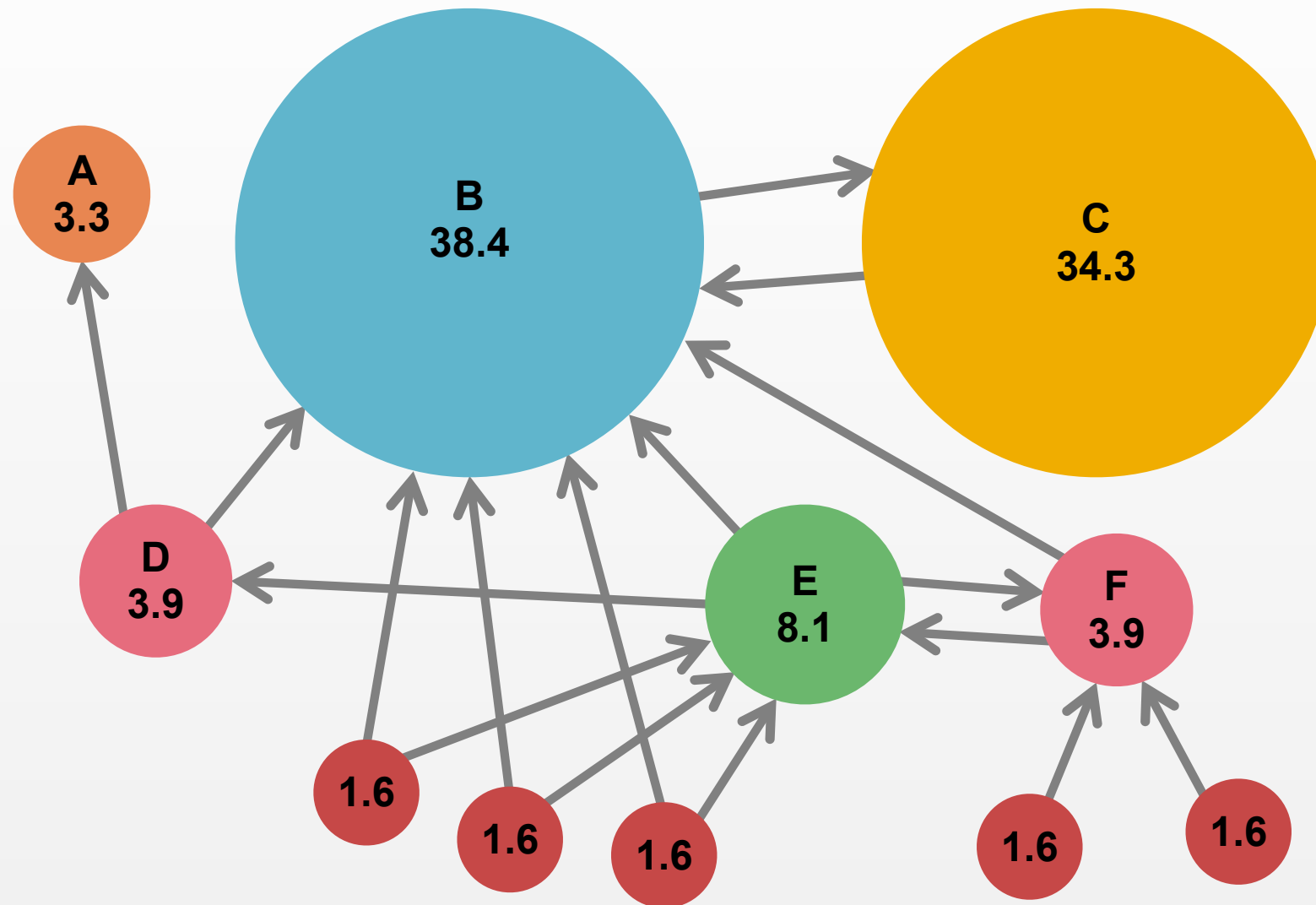
# PageRank: Links as Votes

*Not all pages are equally important*



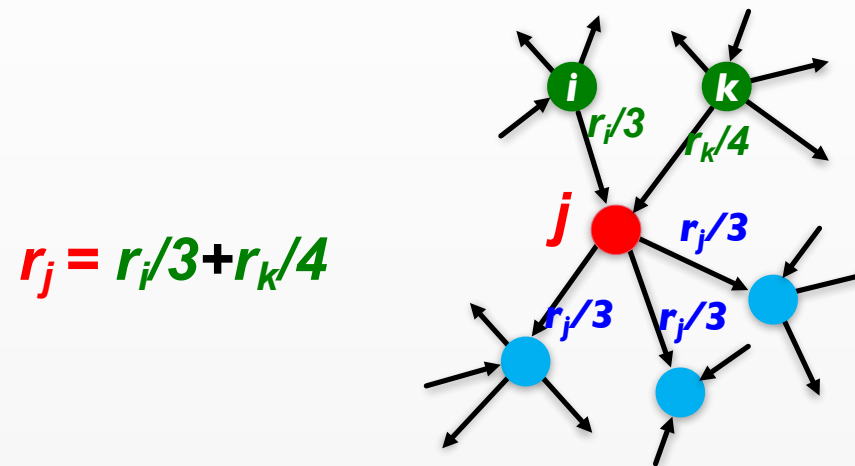
- Pages with **more inbound links** are more **important**
- Inbound **links from important pages** carry **more weight**

# Example: PageRank Scores



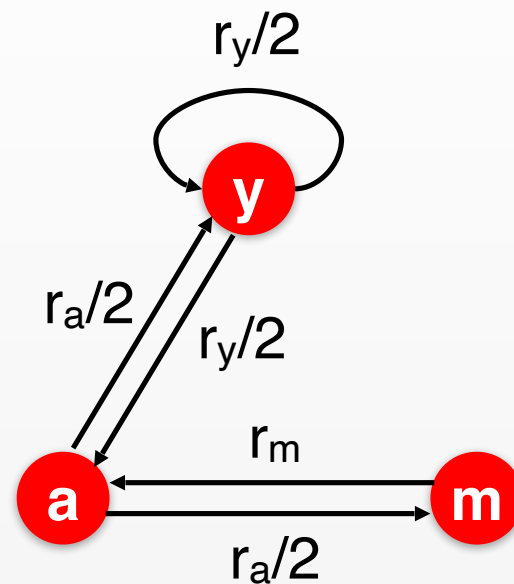
(adapted from: Mining of Massive Datasets, <http://www.mmds.org>)

# PageRank: Recursive Formulation



- A 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

# PageRank: The “Flow” Model



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

**“Flow” equations:**

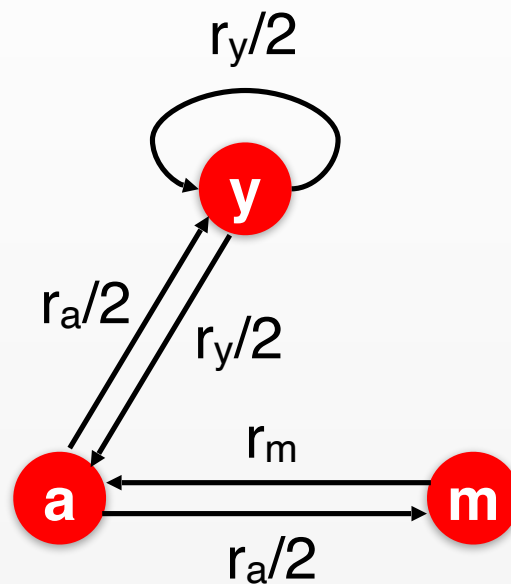
$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

- 3 equations, 3 unknowns
- However, the equations are under constrained. Has infinite solutions.
- Impose additional constraint:  $r_y + r_a + r_m = 1$  (the total sum of importances is 1)
- Solution:  $r_y = 2/5$ ,  $r_a = 2/5$ ,  $r_m = 1/5$

# PageRank: The “Flow” Model



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

**“Flow” equations:**

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

$$\mathbf{r} = \mathbf{M} \cdot \mathbf{r} \quad \begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} & \begin{matrix} y & a & m \end{matrix} \\ \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{bmatrix} & \begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} \end{matrix}$$

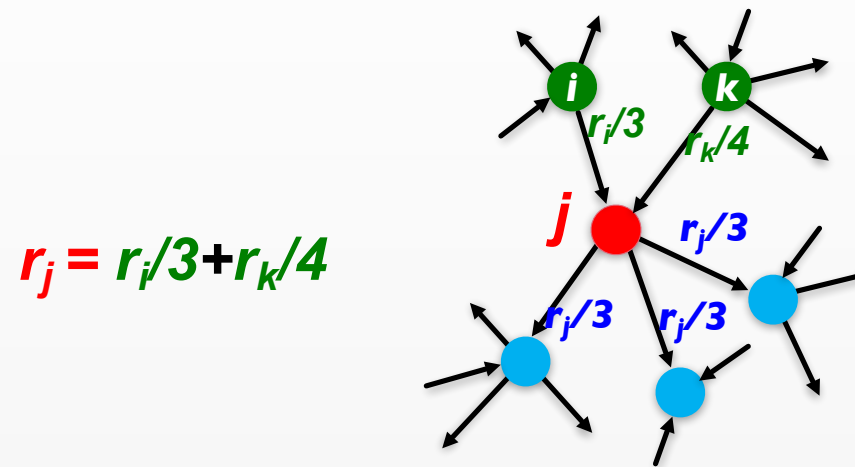
Matrix  $\mathbf{M}$  is stochastic (i.e. columns sum to one)



# PageRank: Eigenvector Problem

- PageRank: Solve for eigenvector  $r = M r$  with eigenvalue  $\lambda = 1$ .
- Normalize the eigenvector to sum to 1.
- Eigenvector with  $\lambda = 1$  is guaranteed to exist since  $M$  is a stochastic matrix (i.e. if  $a = M b$  then  $\sum a_i = \sum b_i$ )
- *Problem:* There are billions of pages on the internet. How do we solve for eigenvector with order  $10^{10}$  elements?

# Equivalent Formulation: Random Surfer



- At time  $t$  a surfer is on some page  $i$
- At time  $t+1$  the surfer follows a link to a new page at random
- Define rank  $r_i$  as fraction of time spent on page  $i$

# PageRank: Power Iteration

*Model for random Surfer:*

- At time  $t = 0$  pick a page at random
- At each subsequent time  $t$  follow an outgoing link at random

*Probabilistic interpretation:*

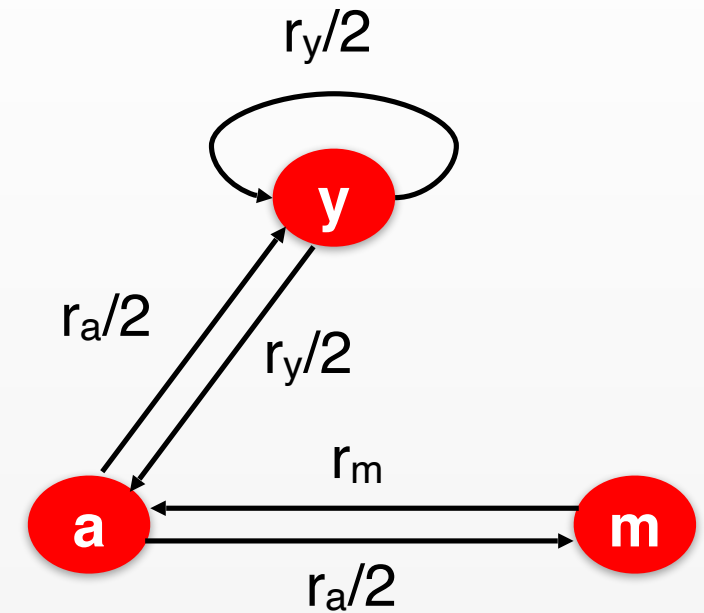
$$p(z_0 = i) = 1/N$$

$$p(z_t = i | z_{t-1} = j) = M_{ij}$$

$$p(z_t = i) = \sum_j p(z_t = i, z_{t-1} = j)$$

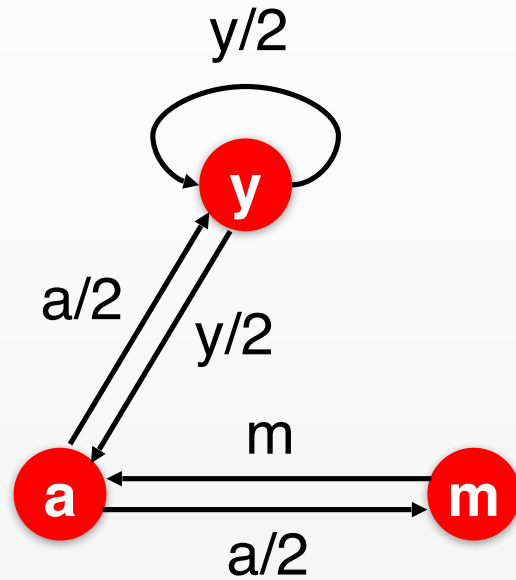
$$= \sum_j M_{ij} p(z_{t-1} = j)$$

$z_i$  is the random variable giving the index of the webpage the random surfer is at time step  $i$ .



$$M = \begin{bmatrix} y & a & m \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{bmatrix}$$

# PageRank: Power Iteration



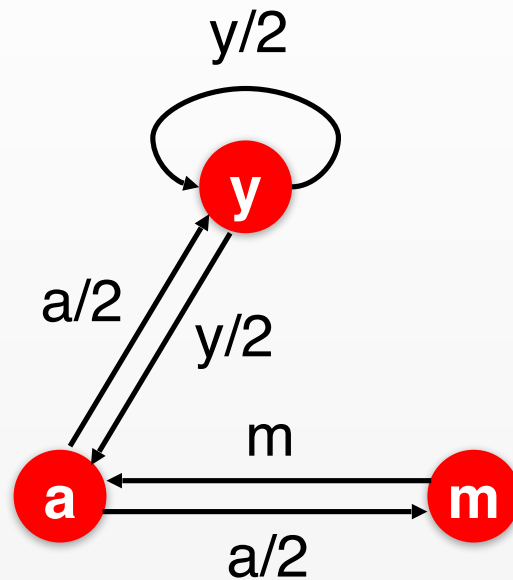
$$\mathbf{p}^t = M\mathbf{p}^{t-1} = M^t\mathbf{p}^0$$

$$\mathbf{p}^0 = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad M = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}_{p^0} \quad \begin{bmatrix} 2/6 \\ 3/6 \\ 1/6 \end{bmatrix}_{p^1 = Mp^0} \quad \begin{bmatrix} 5/12 \\ 4/12 \\ 3/12 \end{bmatrix}_{p^2 = MP^1} \quad \begin{bmatrix} 9/24 \\ 11/24 \\ 4/24 \end{bmatrix}_{p^3 = MP^2} \quad \begin{bmatrix} 20/48 \\ 17/48 \\ 11/48 \end{bmatrix}_{p^4 = MP^3} \simeq \begin{bmatrix} 2/5 \\ 2/5 \\ 1/5 \end{bmatrix}$$

$\mathbf{p}^t$  converges to  $\mathbf{r}$ . Iterate until  $\|\mathbf{p}^t - \mathbf{p}^{t-1}\| < \epsilon$

# PageRank: Convergence

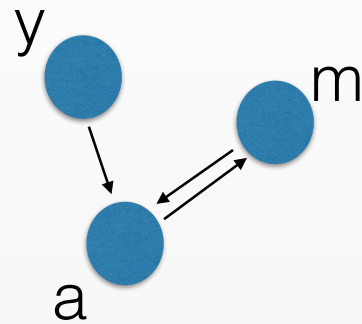


$$\mathbf{p}^t = M\mathbf{p}^{t-1} = M^t\mathbf{p}^0$$

$\mathbf{p}^t$  converges to  $\mathbf{r}$ . Iterate until  $\|\mathbf{p}^t - \mathbf{p}^{t-1}\| < \varepsilon$

Is convergence always guaranteed?  
Is the value at convergence meaningful?

# Case of non-convergence



periodicity in the network

$$\mathbf{p}^t = M\mathbf{p}^{t-1} = M^t\mathbf{p}^0$$

$$\mathbf{p}^0 = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad M = \begin{matrix} & \begin{matrix} y & a & m \end{matrix} \\ \begin{matrix} y \\ a \\ m \end{matrix} & \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \end{matrix}$$

$a$

$$\begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} \\ p^0$$

$$\begin{bmatrix} 0 \\ \frac{2}{3} \\ \frac{1}{3} \end{bmatrix} \\ p^1$$

$$\begin{bmatrix} 0 \\ \frac{1}{3} \\ \frac{2}{3} \end{bmatrix} \\ p^2$$

$$\begin{bmatrix} 0 \\ \frac{2}{3} \\ \frac{1}{3} \end{bmatrix} \\ p^3$$

$$\begin{bmatrix} 0 \\ \frac{1}{3} \\ \frac{2}{3} \end{bmatrix} \\ p^4$$

What if you started from?

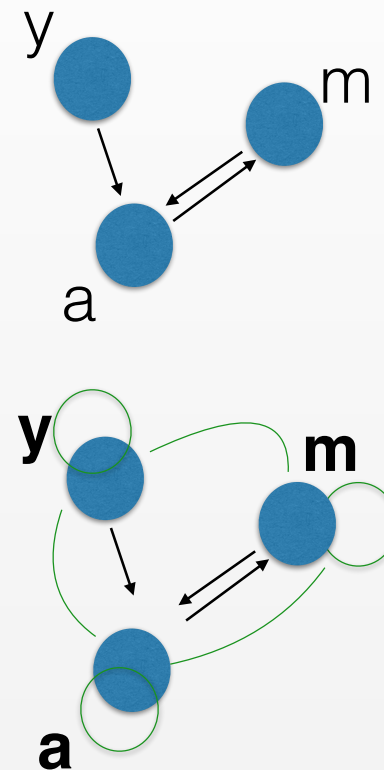
$$p^0 = \begin{bmatrix} 0 \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

# Solution: Random Teleports

Model for *teleporting* random surfer:

- At time  $t = 0$  pick a page at random
- At each subsequent time  $t$ 
  - With probability  $\beta$  follow an outgoing link at random
  - With probability  $1-\beta$  teleport to a new initial location at random

$$p(z_{t+1} = i | z_t = j) = \beta \frac{1}{d_j} + (1 - \beta) \frac{1}{N}$$



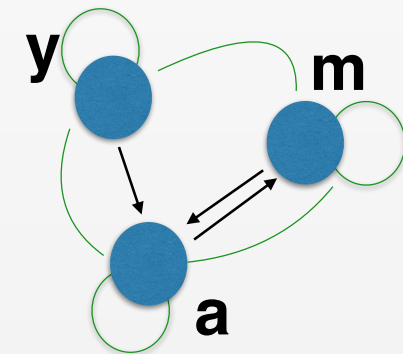
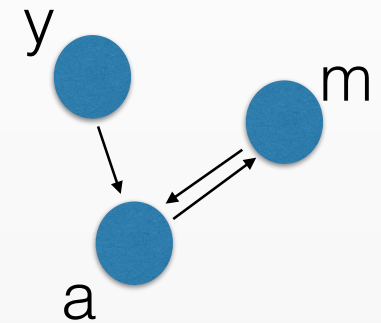
# Solution: Random Teleports

$$p(z_{t+1} = i | z_t = j) = \beta \frac{1}{d_j} + (1 - \beta) \frac{1}{N}$$

$$\tilde{M} = 0.9 \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} + 0.1 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix} = \begin{matrix} & \begin{matrix} y & a & m \end{matrix} \\ \begin{matrix} y & a & m \end{matrix} & \begin{bmatrix} .033 & .033 & .033 \\ .933 & .033 & .933 \\ .033 & .933 & .033 \end{bmatrix} \end{matrix}$$

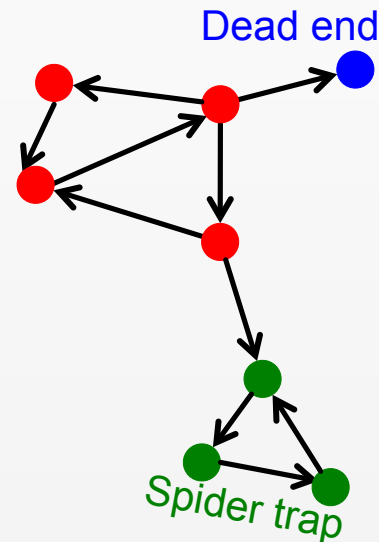
$$p^{t+1} = \tilde{M}p^t$$

$$\begin{matrix} \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} & \dots & \begin{bmatrix} .033 \\ .491 \\ .475 \end{bmatrix} & \begin{bmatrix} .033 \\ .491 \\ .475 \end{bmatrix} & \begin{bmatrix} .033 \\ .491 \\ .475 \end{bmatrix} \\ p^0 & & p^{100} & p^{101} & p^{102} \end{matrix}$$





# PageRank: Problems



Not irreducible

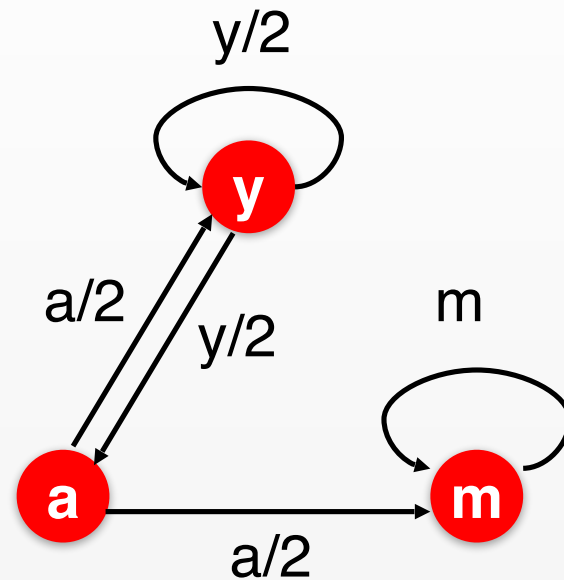
## 1. *Dead Ends*

- Nodes with no outgoing links.
- Where do surfers go next?

## 2. *Spider Traps*

- Subgraph with no outgoing links to wider graph
- Surfers are “trapped” with no way out.

# Power Iteration: Spider Traps



$$\mathbf{p}^t = M\mathbf{p}^{t-1} = M^t\mathbf{p}^0$$

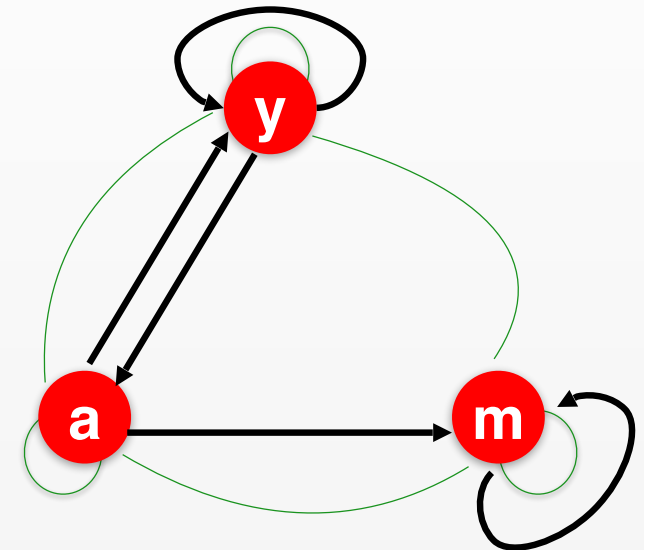
$$\mathbf{p}^0 = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad M = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}_{p^0} \quad \begin{bmatrix} 2/6 \\ 1/6 \\ 3/6 \end{bmatrix}_{p^1} \quad \begin{bmatrix} 3/12 \\ 2/12 \\ 7/12 \end{bmatrix}_{p^2} \quad \begin{bmatrix} 5/24 \\ 3/24 \\ 16/24 \end{bmatrix}_{p^3} \quad \dots \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Probability accumulates in traps (surfers get stuck)

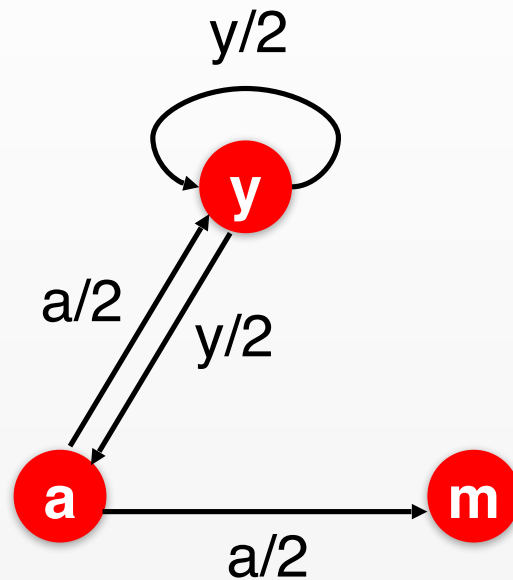
# Power Iteration: Teleports

$$\tilde{M} = 4/5 \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} & 1 \end{bmatrix} + 1/5 \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix} = \begin{bmatrix} \frac{7}{15} & \frac{7}{15} & \frac{1}{15} \\ \frac{7}{15} & \frac{1}{15} & \frac{7}{15} \\ \frac{1}{15} & \frac{7}{15} & \frac{1}{15} \end{bmatrix}$$



$$\begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}_{p^0} \quad \begin{bmatrix} 0.33 \\ 0.20 \\ 0.46 \end{bmatrix}_{p^1} \quad \begin{bmatrix} 0.24 \\ 0.20 \\ 0.56 \end{bmatrix}_{p^2} \quad \dots \quad \begin{bmatrix} 7/33 \\ 5/33 \\ 21/33 \end{bmatrix}$$

# Power Iteration: Dead Ends



$$\mathbf{p}^t = M\mathbf{p}^{t-1} = M^t\mathbf{p}^0$$

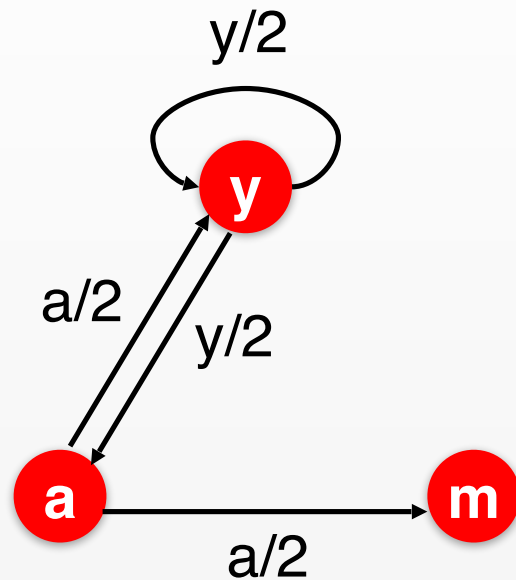
$$\mathbf{p}^0 = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad M = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} & 0 \end{bmatrix}$$

$$\mathbf{p}^t = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad \begin{bmatrix} 2/6 \\ 1/6 \\ 1/6 \end{bmatrix} \quad \begin{bmatrix} 3/12 \\ 1/12 \\ 1/12 \end{bmatrix} \quad \dots \quad \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$p^0 \qquad p^1 \qquad p^2$

Probability not conserved

# Power Iteration: Dead Ends



$$\mathbf{p}^t = M\mathbf{p}^{t-1} = M^t\mathbf{p}^0$$

$$\mathbf{p}^0 = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad M = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & 0 & \frac{1}{3} \\ 0 & \frac{1}{2} & \frac{1}{3} \end{bmatrix}$$

(teleport at dead ends)

$$\mathbf{p}^t = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \quad \begin{bmatrix} 8/18 \\ 5/18 \\ 5/18 \end{bmatrix} \quad \begin{bmatrix} 49/108 \\ 34/108 \\ 35/108 \end{bmatrix} \quad \dots$$

Fixes “probability sink” issue

# Power Iteration: Dead Ends

$M$ : the original probability transition matrix

$\ddot{M}$ : Constructed from  $M$  by replacing any column corresponding to a dead end (contains only 0's) by a column containing all entries equal to  $1/N$ .

$$\tilde{M} = \beta \ddot{M} + (1 - \beta) \frac{1}{N} \mathbf{1} \mathbf{1}^T$$

$$\mathbf{1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

# Intermezzo: Markov Chains

Markov Property

$$p(z_{t+1} | z_t) = p(z_{t+1} | z_t, z_{t-1}, \dots, z_0)$$

Irreducibility

$$\forall i, j \exists n_{ij} > 0 : p(z_{n_{ij}} = i | z_0 = j) > 0$$

Ergodicity

$$\exists N : N \geq \max_{i,j} n_{ij}$$

Aperiodicity:

$$\exists i, \text{GCD}(T_i) = 1$$

$$T_i = \{t \geq 1 : p(z_t = i | z_0 = i) > 0\}$$

There exists  $\mathbf{r}$  with all entries strictly positive

$$\mathbf{r} = \mathbf{M}\mathbf{r}$$

Power iteration converges to  $\mathbf{r}$ , starting with any probability vector  $\mathbf{p}^0$

$$\lim_{t \rightarrow \infty} \mathbf{p}^t = \mathbf{r}$$

where

$$\mathbf{p}^t = \mathbf{M}\mathbf{p}^{t-1}$$