# Link Analysis & Prediction Heuristics (PageRank)
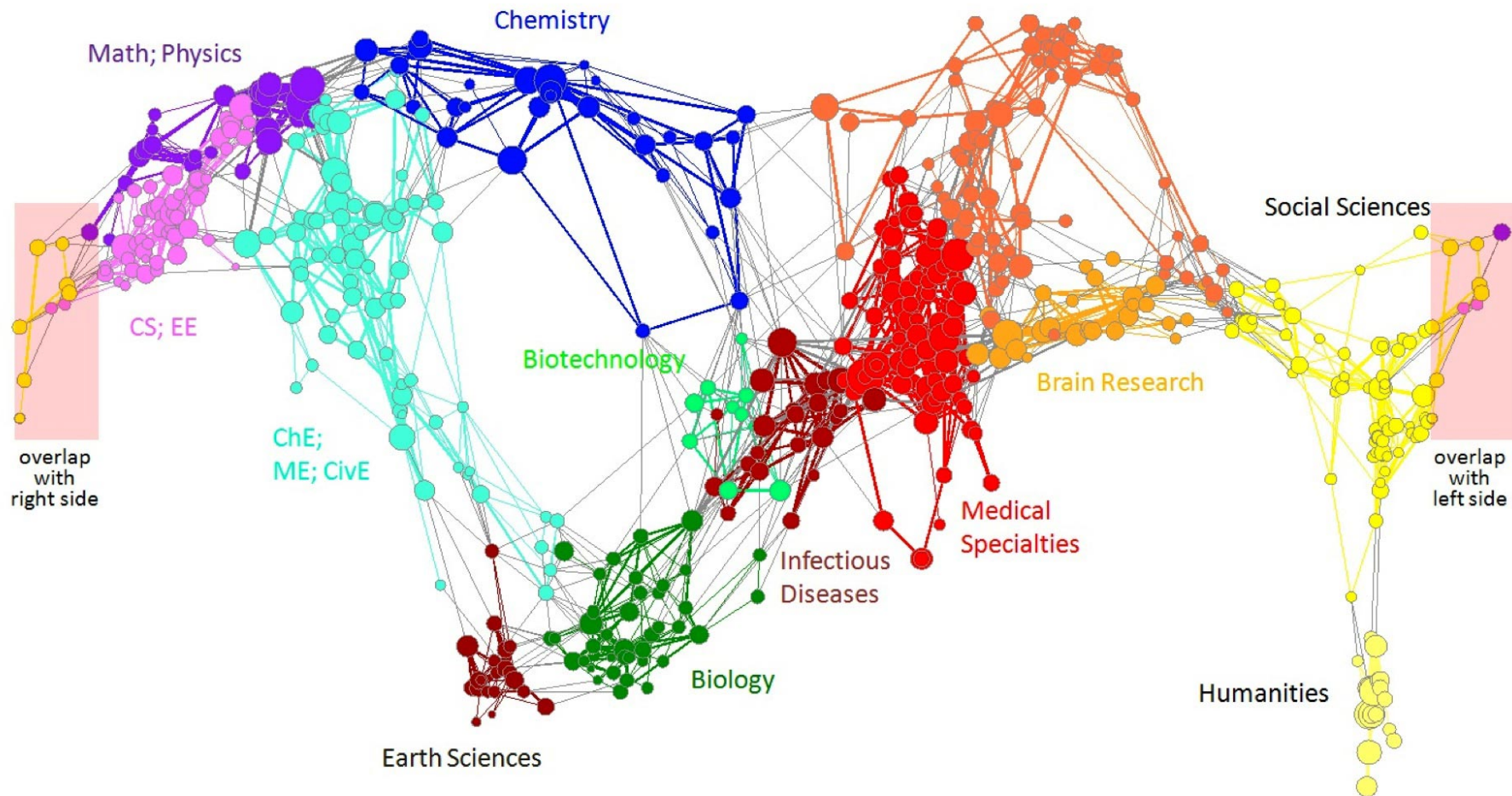
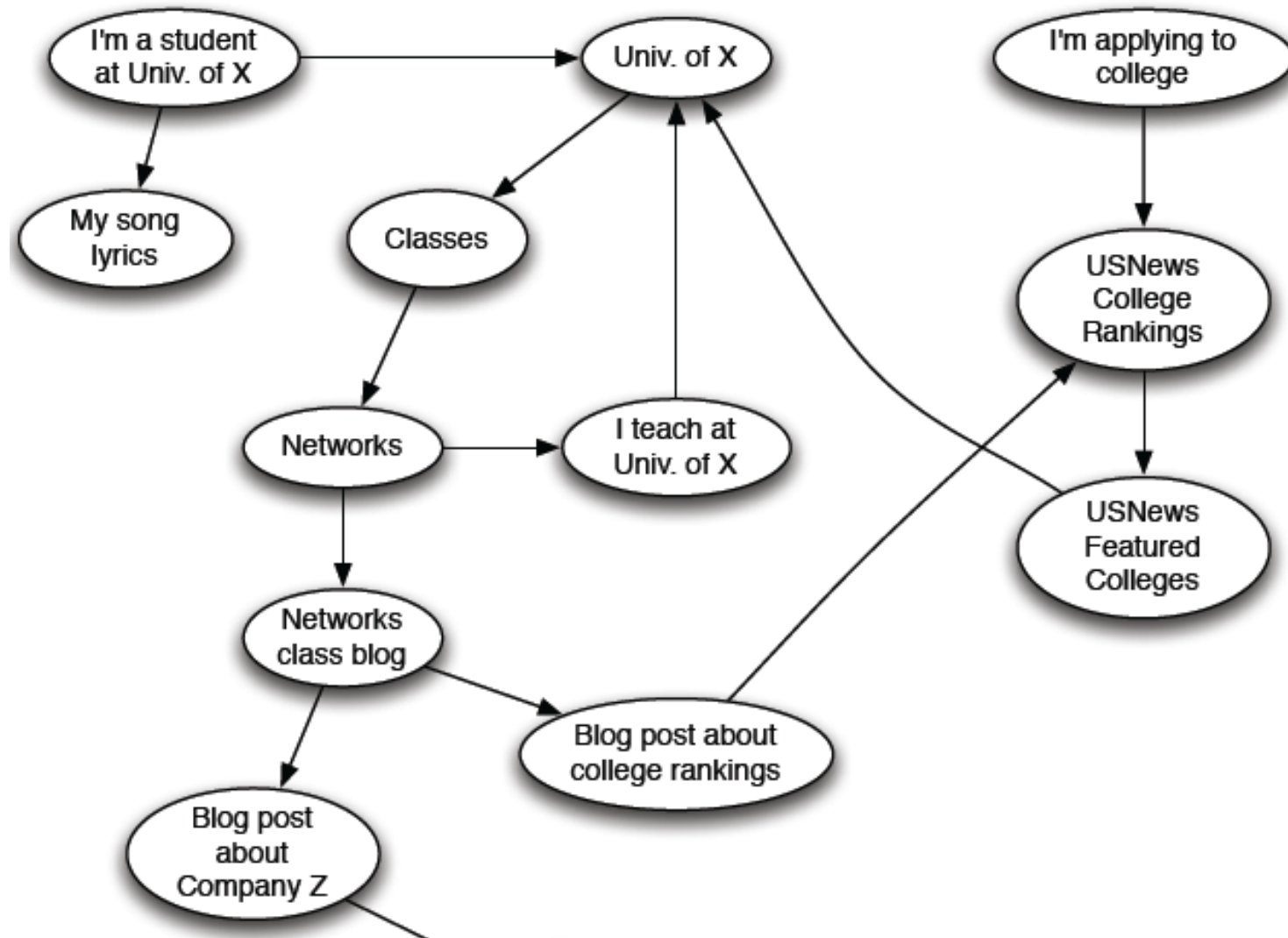## CS57300 Data Mining
## Spring 2016

## Instructor: Bruno Ribeiro

# Graph Data: Information Nets



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

# 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!
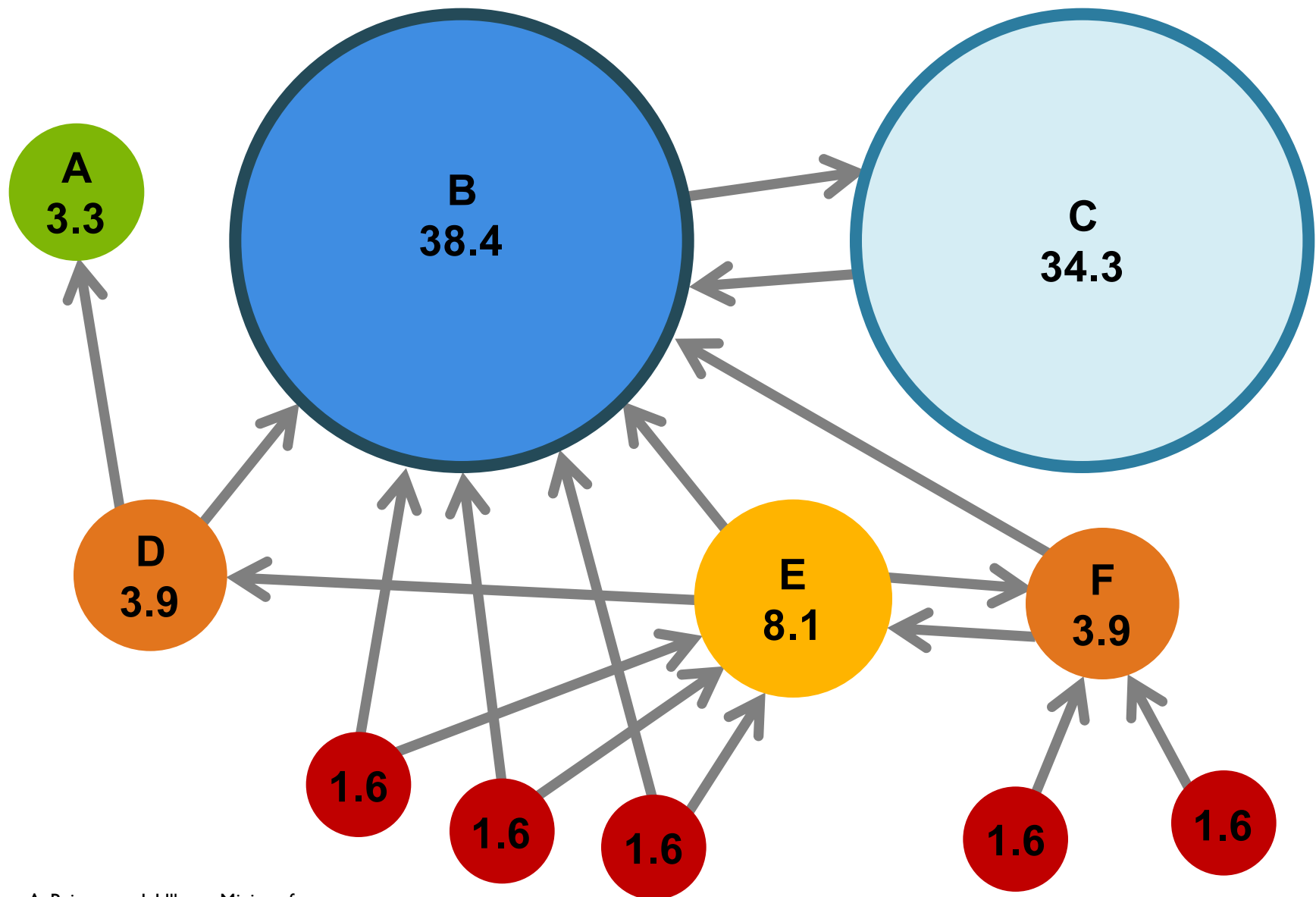
(2) What is the "best" answer to query "newspaper"?

- ◦ No single right answer
- ◦ Trick: Pages that actually know about newspapers might all be pointing to many newspapers
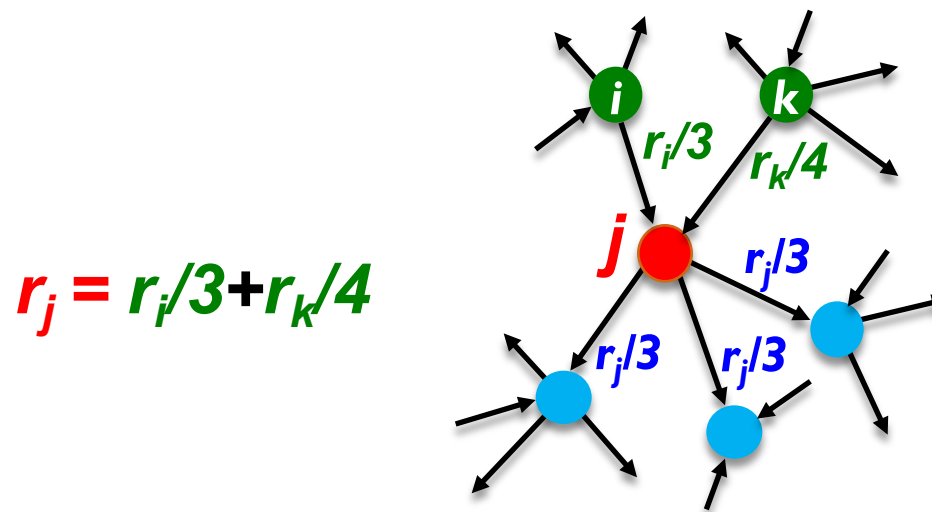
# PageRank: Flow Formulation

# 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.purdue.edu has **1,910** in-links
  - ◦ www.fake-school-name.com has **1** in-link

- ▸ 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

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

# Solving the Flow Equations

**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, no constants
- No unique solution
- All solutions equivalent modulo the scale factor

▸ Additional constraint forces uniqueness:
- $r_y + r_a + r_m = 1$
- Solution: $r_y = \frac{2}{5}$, $r_a = \frac{2}{5}$, $r_m = \frac{1}{5}$

▸ 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 $i$
  - ◦ $\sum_i r_i = 1$
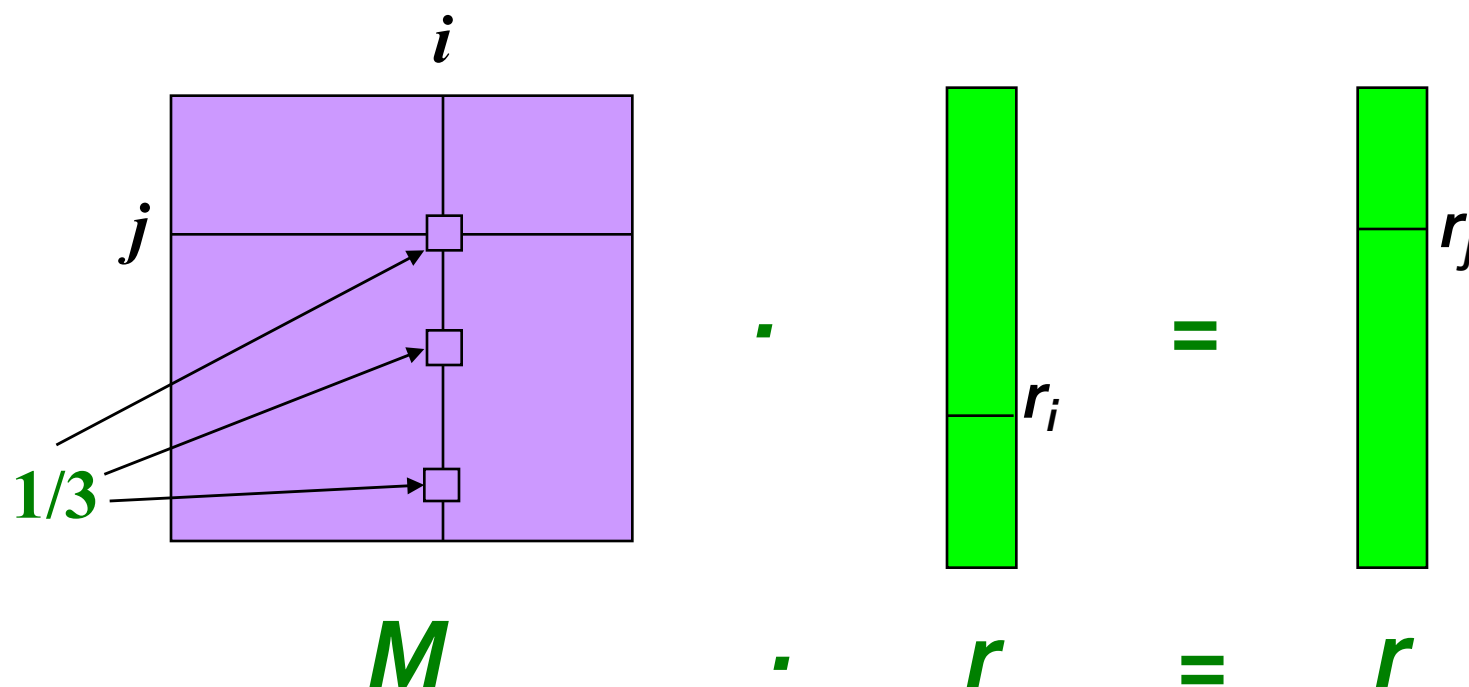- ▸ The flow equations can be written

$$r = M \cdot r$$

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

# Example

▸ **Remember the flow equation:** $\quad r_j = \sum_{i \to j} \dfrac{r_i}{d_i}$

▸ **Flow equation in the matrix form**

$$M \cdot r = r$$

◦ **Suppose page _i_ links to 3 pages, including _j_**



$$M \quad \cdot \quad r \quad = \quad r$$

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

# 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 $M$

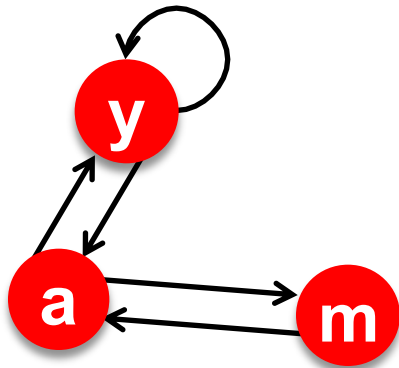  ◦ In fact, its first or principal eigenvector, with corresponding eigenvalue $1$

    • 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 $Mr \leq 1$*

**NOTE: $x$** is an eigenvector with the corresponding eigenvalue $\boldsymbol{\lambda}$ if:
$$Ax = \lambda x$$

▸ We can now efficiently solve for $r$ *via Power iteration*

# Example: Flow Equations & M



|   | y | a | m |
|---|---|---|---|
| **y** | ½ | ½ | 0 |
| **a** | ½ | 0 | 1 |
| **m** | 0 | ½ | 0 |

$$r = M \cdot r$$

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

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

# PageRank: Three Questions

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

▸ Does this converge?

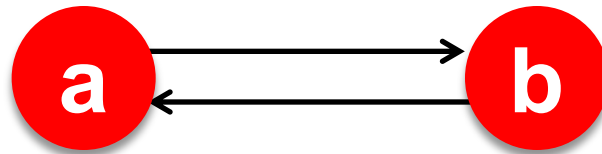▸ Does it converge to what we want?

▸ Are results reasonable?
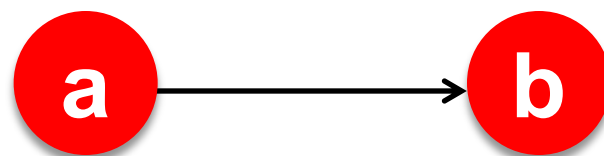
# Does this converge?

▸ **Example:**



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

| $r_a$ | | 1 | 0 | 1 | 0 |
|---|---|---|---|---|---|
| $r_b$ | = | 0 | 1 | 0 | 1 |

Iteration 0, 1, 2, …

# Does it converge to what we want?

▸ **Example:**



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

$r_a$      1     0     0     0

$r_b$      0     1     0     0

=

Iteration 0, 1, 2, …

# PageRank: Problems

**2 problems:**
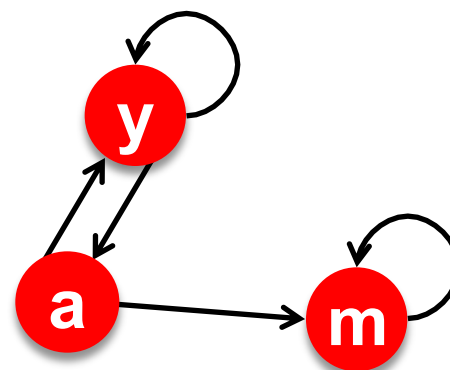
▸ (1) Some pages are
dead ends (have no out-links)

- 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



Dead end

Spider trap

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

# Problem: Spider Traps

▸ Power Iteration:
  ◦ Set $r_j = 1$
  ◦ $r_j = \sum_{i \to j} \frac{r_i}{d_i}$
    • And iterate

|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 0 |
| a | ½ | 0 | 0 |
| m | 0 | ½ | 1 |

m is a spider trap

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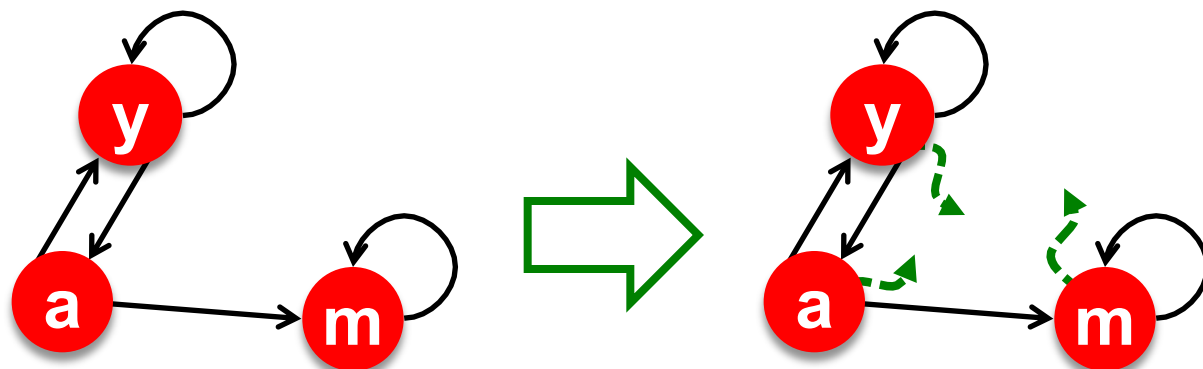$$r_a = r_y /2$$
$$r_m = r_a /2 + r_m$$

▸ **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!

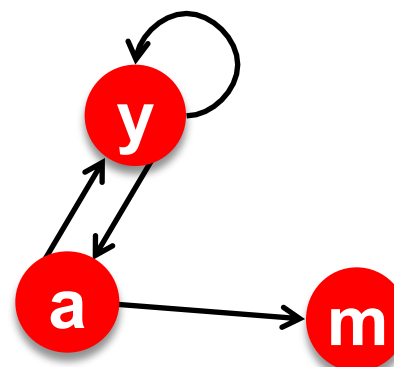- The Google solution for spider traps: At each time step, the random surfer has two options
  - With prob. $\beta$, follow a link at random
  - With prob. $1-\beta$, jump to some random page
  - Common values for $\beta$ are in the range 0.8 to 0.9
- Surfer will teleport out of spider trap within a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

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

# Problem: Dead Ends

- ▶ Power Iteration:
  - ◦ Set $r_j = 1$
  - ◦ $r_j = \sum_{i \to j} \frac{r_i}{d_i}$
    - • And iterate

|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 0 |
| a | ½ | 0 | 0 |
| m | 0 | ½ | 0 |

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

- ▶ **Example:**
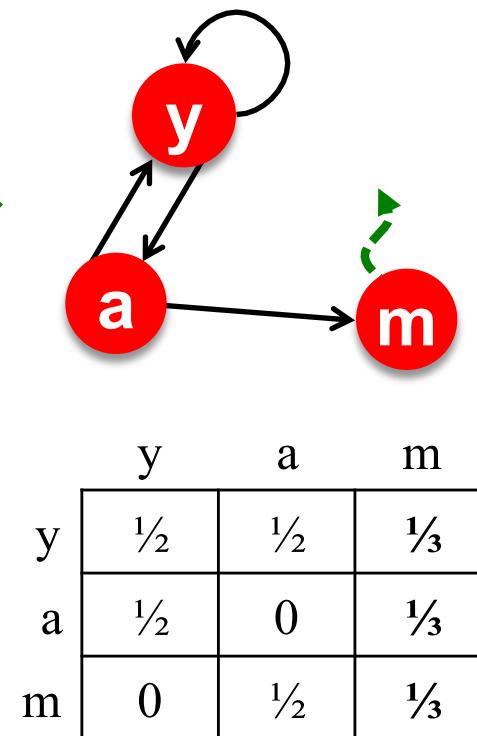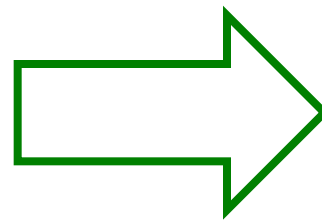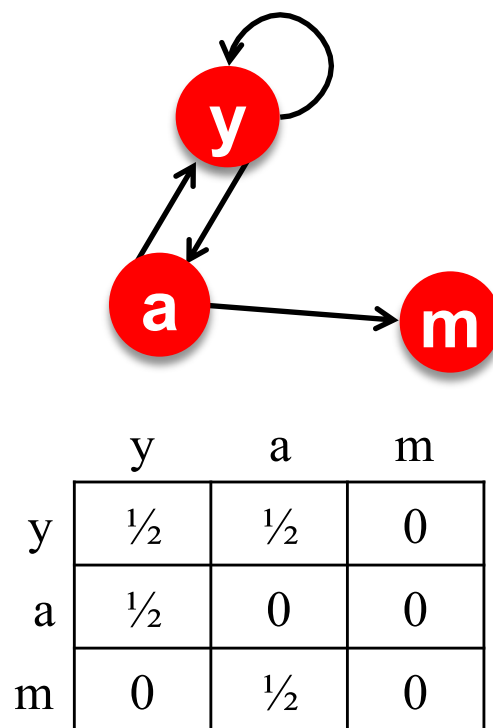
$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 \\ 1/3 & 1/6 & 2/12 & 3/24 \\ 1/3 & 1/6 & 1/12 & 2/24 \end{matrix} \quad \dots \quad \begin{matrix} 0 \\ 0 \\ 0 \end{matrix}$$

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



|   | y   | a   | m   |
|---|-----|-----|-----|
| y | ½   | ½   | 0   |
| a | ½   | 0   | 0   |
| m | 0   | ½   | 0   |

|   | y   | a   | m   |
|---|-----|-----|-----|
| y | ½   | ½   | ⅓   |
| a | ½   | 0   | ⅓   |
| m | 0   | ½   | ⅓   |

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

# 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]

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

$d_i$ ... 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]

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

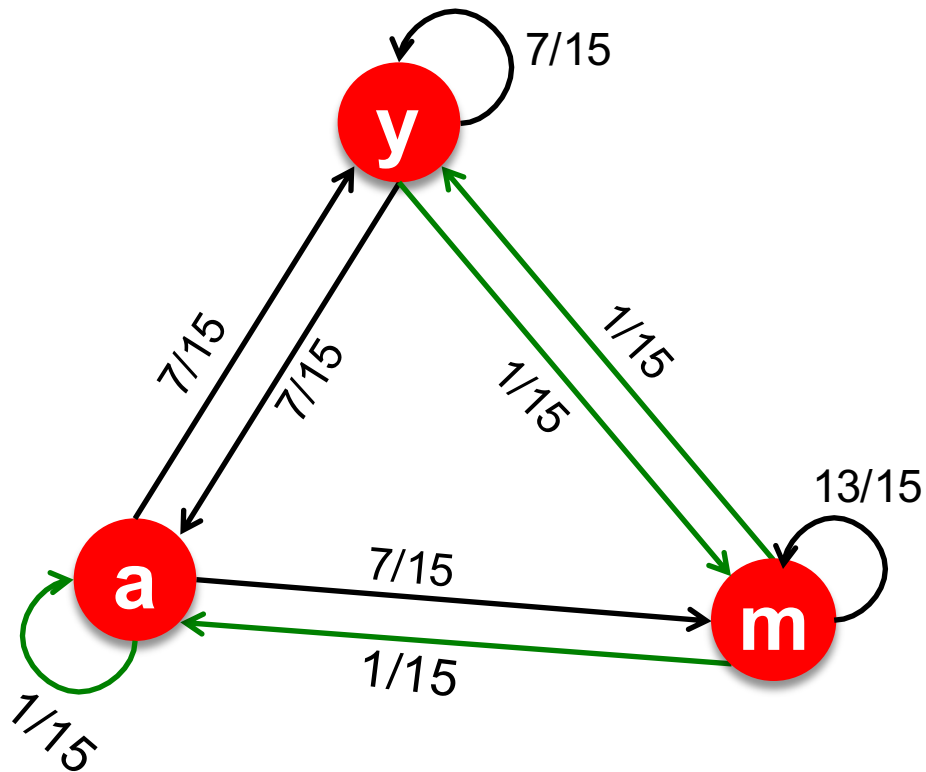- The Google Matrix $A$:

$$A = \beta\, M + (1 - \beta) \left[\frac{1}{N}\right]_{N \times N}$$

$[1/N]_{NxN}$…N by N matrix where all entries are 1/N

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

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

# Random Teleports (β = 0.8)



**M**

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

**[1/N]**$_{NxN}$

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

|   |        |        |        |
|---|--------|--------|--------|
| 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 |

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