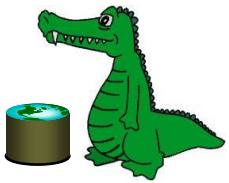# Logistics

- Lab 3 due Friday 11:55pm
  - TA office hour Thursday, possibly additional OH on Friday
- Sample midterm 1 questions
  - Review Lecture Monday
  - Solutions will be discussed and posted then
- Midterm 1 in class on Wed. (10/14)
  - Similar question topics/types as in sample
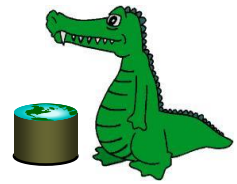  - 50 min, close book, paper based, no aids

# Midterm 1 topics

1. Data modeling and similarity metrics
   – Vectors, sets, strings, graphs
2. NLP techniques
3. Map-Reduce API
4. Exploratory Analysis using visualizations

# Review

- 3 more examples of M/R applications
- Web as a graph
  - Ways to search and rank the Web
- Link Analysis and Page Rank
  - Computing importance scores for each web page based on link structures
- Mathematics abstraction of the problem with flow equations and solutions
  - Gaussian elimination not scalable
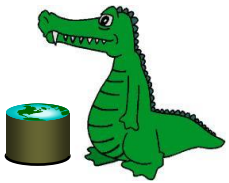- Matrix formulation of the problem

# Rank Vector r = Eigenvector of M

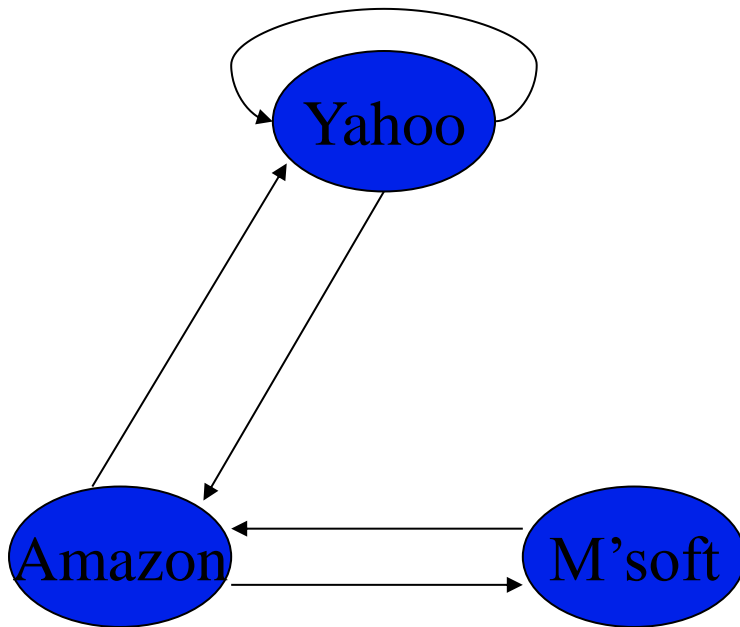- The flow equations can be written

$$\mathbf{r} = \mathbf{Mr}$$

- The rank vector **r** is an eigenvector of the stochastic web matrix **M**
  - with corresponding eigenvalue 1

- We can now efficiently solve for **r**!
  - The method is called Power iteration

# Power Iteration Example



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

**r = Mr**

$y = y/2 + a/2$

$a = y/2 + m$

$m = a/2$

$$\begin{bmatrix} y \\ a \\ m \end{bmatrix} = \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1 \\ 0 & 1/2 & 0 \end{bmatrix} \begin{bmatrix} y \\ a \\ m \end{bmatrix}$$
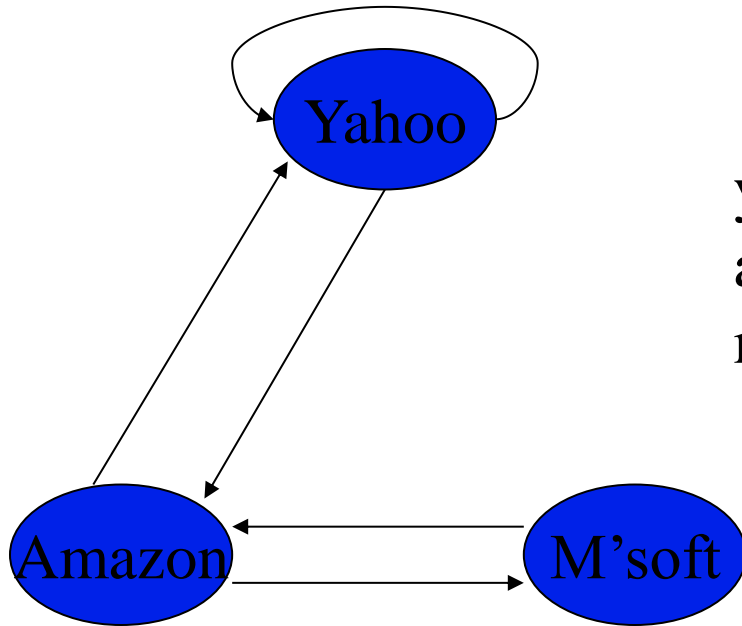
# Power Iteration method

- 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: $\mathbf{r}^0 = [1/N,....,1/N]^T$
  - Iterate: $\mathbf{r}^{k+1} = \mathbf{M}\mathbf{r}^k$
  - Stop when $|\mathbf{r}^{k+1} - \mathbf{r}^k|_1 < \varepsilon$

    $r^{(t+1)}_j = \sum_{i \to j} r^{(t)}_i / d_i$
    $d_i$ *is out-degree of node i*

    - $|\mathbf{x}|_1 = \sum_{1 \le i \le N} |x_i|$ is the $L_1$ norm
    - Can use any other vector norm e.g., Euclidean norm

# Power Iteration Example (cont.)

|   | y | a | m |
|---|---|---|---|
| y | 1/2 | 1/2 | 0 |
| a | 1/2 | 0 | 1 |
| m | 0 | 1/2 | 0 |

**Power Iteration:**
- Set $r_j = 1/N$
- **1:** $r'_j = \sum_{i \to j} r_i / d_i$
  
  **(i.e., r' = Mr)**
- **2:** $r = r'$
- Goto **1**

$$
\begin{matrix}
y \\
a \\
m
\end{matrix}
=
\begin{matrix}
1/3 & 1/3 & 5/12 & 3/8 & & 2/5 \\
1/3 & 1/2 & 1/3 & 11/24 & \dots & 2/5 \\
1/3 & 1/6 & 1/4 & 1/6 & & 1/5
\end{matrix}
$$

# Random Walk Interpretation

- Imagine a random web surfer
  - At any time t, surfer is on some page P
  - At time t+1, the surfer follows an outlink from P uniformly at random
  - Ends up on some page Q linked from P
  - Process repeats indefinitely
- Let $\mathbf{p}(t)$ be a vector whose $i^{th}$ component is the probability that the surfer is at page i at time t
  - $\mathbf{p}(t)$ is a probability distribution on pages
  - $\mathbf{p}(t)$ is a link-based importance rank t->∞
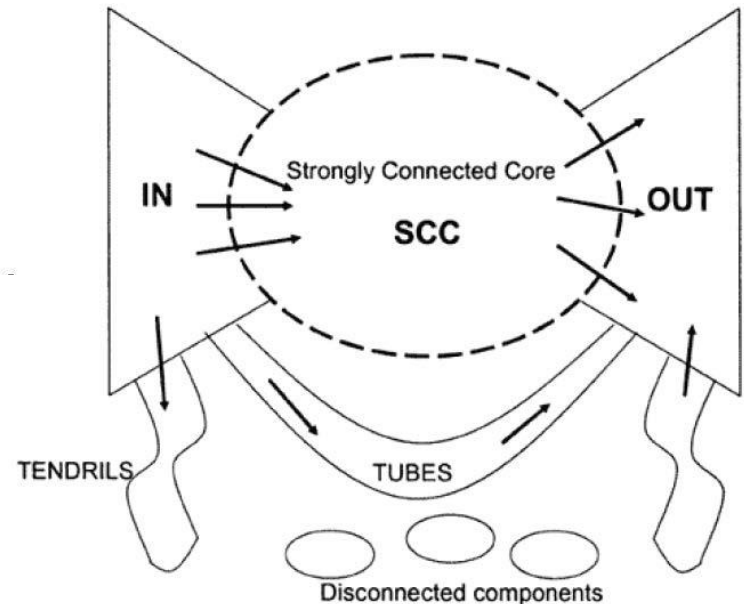
# PageRank: Problems

**2 problems:**

**(1) Spider traps** (all out-links are within the group)

- Eventually spider traps absorb all importance
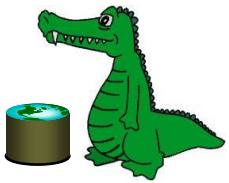
**(2)** Some pages are **dead ends** (have no out-links)

- Such pages cause importance to "leak out"

Strongly Connected Core
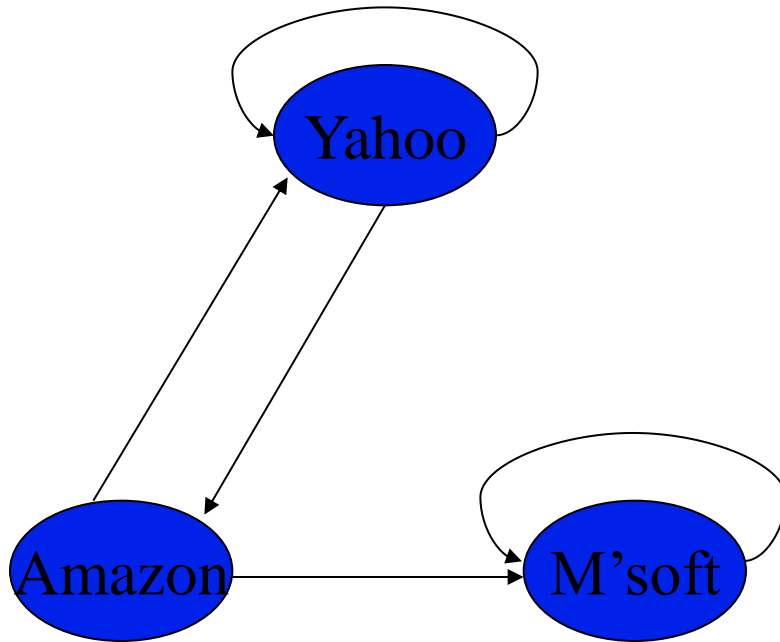IN — SCC — OUT
TENDRILS    TUBES
Disconnected components

# Spider traps

- A group of pages is a spider trap if there are no links from within the group to outside the group
  - Random surfer gets trapped

- Spider traps violate the conditions needed for the random walk theorem
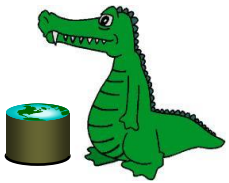
# Microsoft becomes a spider trap

|   | y | a | m |
|---|---|---|---|
| y | 1/2 | 1/2 | 0 |
| a | 1/2 | 0 | 0 |
| m | 0 | 1/2 | 1 |

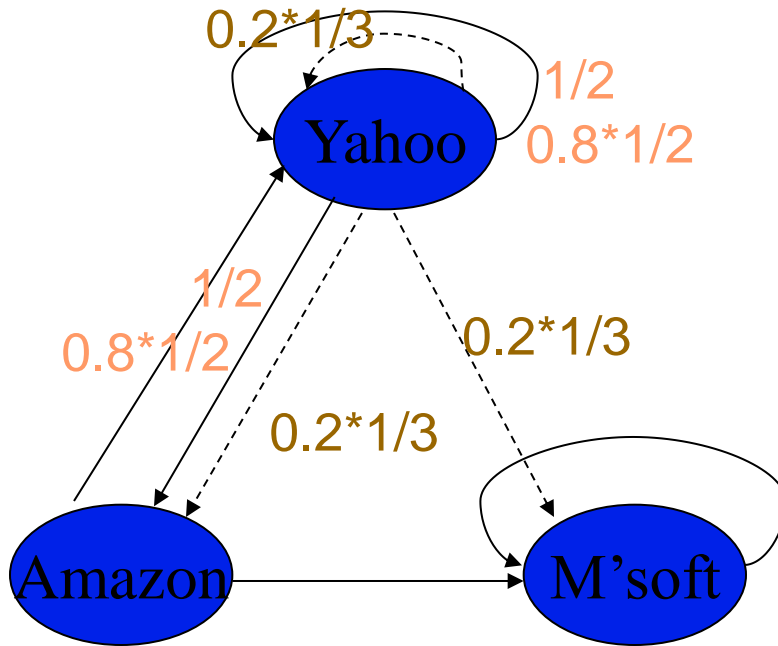| y |   | 1 | 1 | 3/4 | 5/8 |  | 0 |
|---|---|---|---|-----|-----|-----|---|
| a | = | 1 | 1/2 | 1/2 | 3/8 | . . . | 0 |
| m |   | 1 | 3/2 | 7/4 | 2 |  | 3 |

# Random teleports

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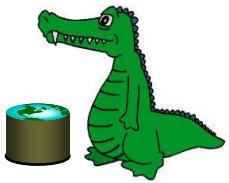- Results: Surfer will teleport out of spider trap within a few time steps
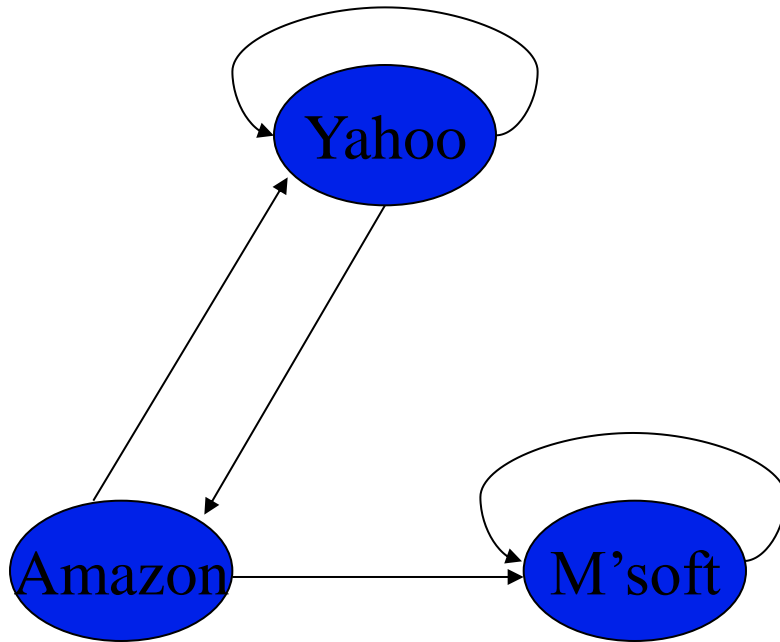
# Random teleports ($\beta = 0.8$)



0.2*1/3

1/2
0.8*1/2

1/2
0.8*1/2

0.2*1/3

0.2*1/3

Yahoo

Amazon          M'soft

|     | y   |
| --- | --- |
| y   | 1/2 |
| a   | 1/2 |
| m   | 0   |

$0.8 *$

|     |
| --- |
| 1/2 |
| 1/2 |
| 0   |

$+ 0.2 *$

|     |
| --- |
| 1/3 |
| 1/3 |
| 1/3 |

$0.8$

| 1/2 | 1/2 | 0 |
| --- | --- | --- |
| 1/2 | 0   | 0 |
| 0   | 1/2 | 1 |

$+ 0.2$

| 1/3 | 1/3 | 1/3 |
| --- | --- | --- |
| 1/3 | 1/3 | 1/3 |
| 1/3 | 1/3 | 1/3 |

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

# 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 |

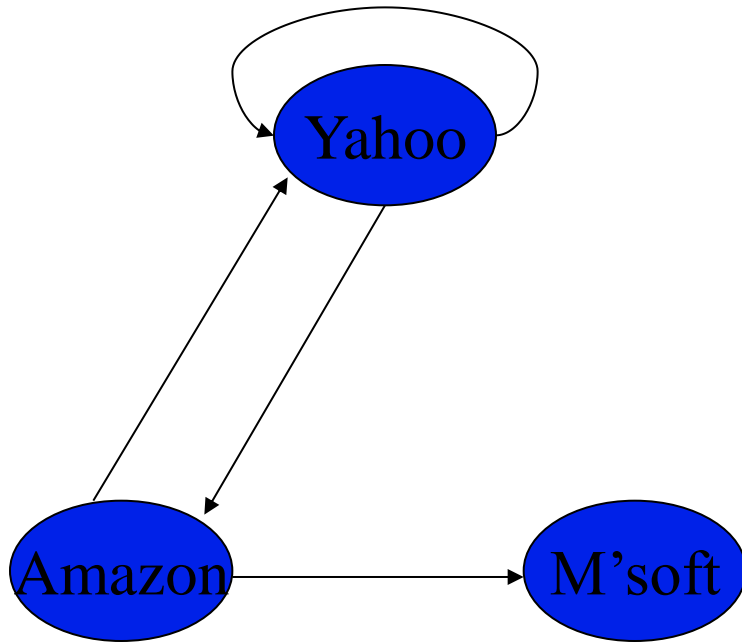| y |   | 1 | 1.00 | 0.84 | 0.776 |     | 7/11  |
|---|---|---|------|------|-------|-----|-------|
| a | = | 1 | 0.60 | 0.60 | 0.536 | . . . | 5/11  |
| m |   | 1 | 1.40 | 1.56 | 1.688 |     | 21/11 |

# Dead ends

- Pages with no outlinks are "dead ends" for the random surfer
  - Nowhere to go on next step
  - All importance becomes zero

# Microsoft becomes a dead end

Yahoo

Amazon → M'soft

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

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

$$\begin{matrix} y \\ a & = \\ m \end{matrix} \begin{matrix} 1 & 1 & 0.787 & 0.648 & & 0 \\ 1 & 0.6 & 0.547 & 0.430 & \ldots & 0 \\ 1 & 0.6 & 0.387 & 0.333 & & 0 \end{matrix}$$

# Dealing with dead-ends

- Teleport
  - Follow random teleport links with probability 1.0 from dead-ends
  - Adjust matrix accordingly

- Prune and propagate
  - Preprocess the graph to eliminate dead-ends
  - Might require multiple passes
  - Compute page rank on reduced graph
  - Approximate values for deadends by propagating values from reduced graph