# COMP4434 Big Data Analytics

## Lecture 11 PageRank

HUANG Xiao

xiaohuang@comp.polyu.edu.hk

**Big Data Analytics**

Machine learning

- Supervised learning
  - Logistic regression classifier
  - Support vector machine
- Unsupervised learning
  - Dimensionality reduction (autoencoder, SVD)
  - Clustering: K-means
  - Overfitting & cross validation
  - F1 score, precision, recall (Evaluation metrics)
- Deep learning
  - Backpropagation
  - Multilayer perceptron (autoencoder)
  - Recurrent neural network (text, time-series)
  - Convolutional neural network (image)

Large-scale data analytics systems
- MapReduce
- Hadoop

Characteristics of big data
- Volume
- Velocity
- Variety (tabular, text, time-series, image, graph)
- Veracity

Applications: AI with big data
- ChatGPT
- AlphaGo
- AlphaFold 2
- Autonomous driving
- Facial recognition
- Recommender system
  - Content-based Recommendation
  - Collaborative filtering
  - Factorization (SVD)
- Web search
  - PageRank + MapReduce

Basic statistical analysis
- Linear regression
- Gradient descent
- Singular value decomposition (SVD)

Graph analytics
- Adjacency matrix
  - User-item interaction network
- PageRank

# PageRank Motivation: How to organize the Web?

- **First try:** Human curated **Web directories**
  - Yahoo, DMOZ, LookSmart
- **Second try: Web search**
  - **Information Retrieval** investigates: Finding relevant docs in a small and trusted set
    - Newspaper articles, Patents, etc.
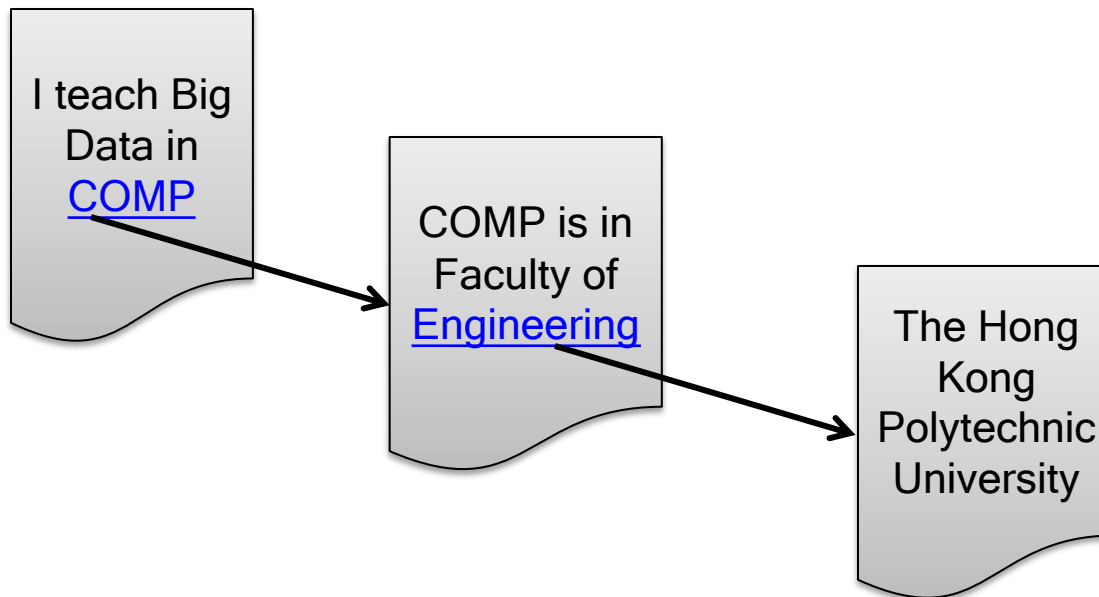  - **But:** Web is **huge**, full of untrusted documents, random things, web spam, etc.
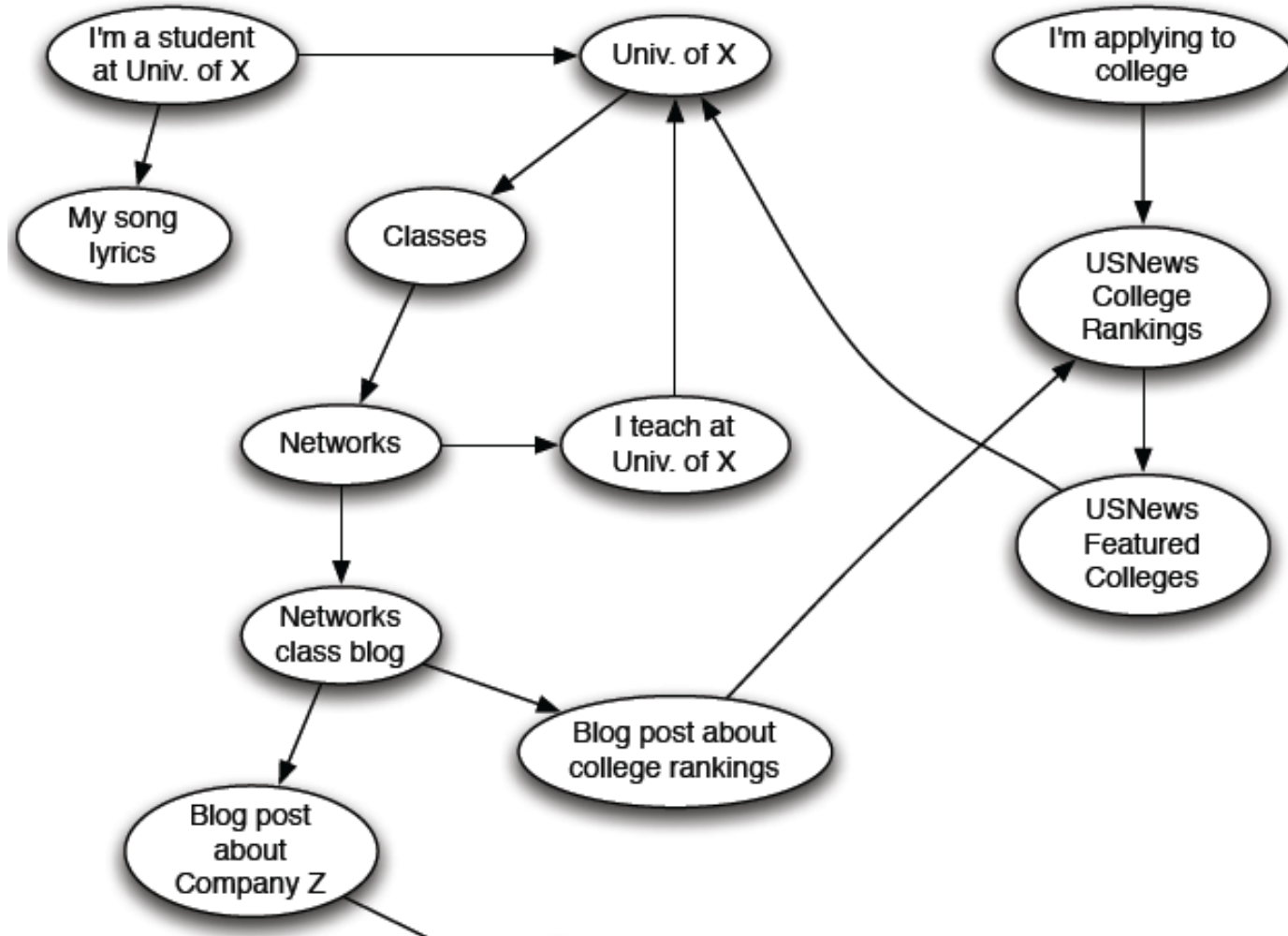
# Challenges in Web Search

- **(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

# Hint: Web as a Directed Graph

- **Nodes: Webpages**
- **Edges: Hyperlinks**

# Web as a Directed Graph

# Ranking Nodes on the Graph

- **All web pages are not equally "important"**

  https://xhuang31.github.io vs.
  https://www.polyu.edu.hk

- There is large diversity in the web-graph node connectivity.
  **Let's rank the pages by the link structure!**

# Example of Node Ranking

- Page Ranking

- Social Ranking

- Paper Ranking

- Scholar Ranking

- ......

# 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
    - https://xhuang31.github.io has 0 in-link

- **Are all in-links are equal?**
    - **Links from important pages count more**
    - Recursive question!

# Google PageRank

- **In-coming links!** Out-going links?
- A page with high PageRank value
  - Many pages pointing to it, or
  - There are some pages that point to it and have high PageRank values
- Example:
  - Page C has a higher PageRank than Page E, even though it has fewer links to it
  - The link it has is of a much higher value

# Is Page == "Webpage"?

- Born in March 26, 1973

- Found Google at September 4, 1998

- As of Nov 2024, own an estimated net worth of $163 billion (No.15 Richest)

- Begins from "Larry Page and Sergey Brin developed PageRank at Stanford University in 1996" as part of a research project about a new kind of search engine.



Larry Page
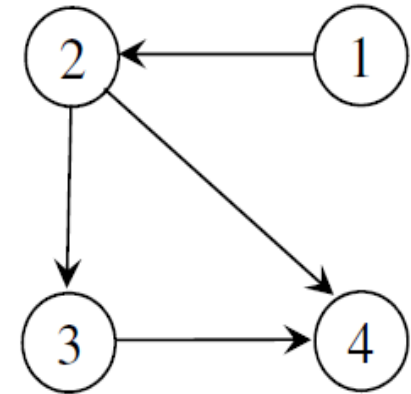Co-founder of Google

# Simple Recursive Formulation

- Each link's vote is proportional to the **importance** of its source page

- If page *j* with importance *PR(j)* has *n* out-links, each link gets *PR(j) / n* votes

- Page *j*'s own importance is the sum of the votes on its in-links

$$PR(j) = PR(i)/3 + PR(k)/4$$

# How to Represent a Graph

- **Graph model** $G = (V, E)$
  - $V$ is a set of pages
  - $E$ is a set of edges
  - Each edge $(u, v) \in E$ represents that page $u$ points/references to page $v$

- **Adjacent List**
  - A data structure for a graph
  - $Adj[u] = \{v: (u, v) \in E\}$ contains each vertex $v$ being adjacent to $u$
  - Example: $Adj[2] = \{3, 4\}$



1: $\boxed{2}$

2: $\boxed{3 \mid 4}$

3: $\boxed{4}$

4: $\boxed{\phantom{x}}$

# 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 $j$**

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

$d_i$ … **out-degree of node $i$**

y/2

a/2

y/2

m

a/2

y

a

m

"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**
  - No unique solution
  - All solutions equivalent modulo the scale factor
- **Additional constraint forces uniqueness:**
  - $r_y + r_a + r_m = 1$
  - **Solution:** $r_y = \frac{2}{5}, \; r_a = \frac{2}{5}, \; r_m = \frac{1}{5}$
- **But, we need a better method for large web-size graphs**

# PageRank: Matrix Formulation

- **Stochastic adjacency matrix $M$**
  - Let page $i$ has $d_i$ out-links
  - If $i \rightarrow j$, then $M_{ji} = \dfrac{1}{d_i}$ 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**

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

$$\boldsymbol{r = M \cdot r}$$

# Example



- **Remember the flow equation:**
- **Flow equation in the matrix form**

$$M \cdot r = r$$

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



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

*M* · *r* = *r*

# Example: Flow Equations & M



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

$$r = M \cdot r$$

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

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

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

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

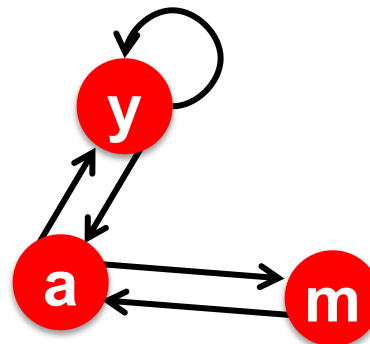$$r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i}$$

$d_i$ …. out-degree of node i

$|\mathbf{x}|_1 = \sum_{1 \le i \le N} |x_i|$ is the **L1** norm
Can use any other vector norm, e.g., Euclidean

# PageRank: How to solve?

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



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

$r_y = r_y/2 + r_a/2$

$r_a = r_y/2 + r_m$

$r_m = r_a/2$

- **Example:**

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

Iteration 0, 1, 2, …

# Why Power Iteration works? (1)

- **Power iteration:**

  A method for finding dominant eigenvector (the vector corresponding to the largest eigenvalue)

  - $r^{(1)} = M \cdot r^{(0)}$

  - $r^{(2)} = M \cdot r^{(1)} = M(Mr^{(0)}) = M^2 \cdot r^{(0)}$

  - $r^{(3)} = M \cdot r^{(2)} = M(M^2 r^{(0)}) = M^3 \cdot r^{(0)}$

- **Claim:**

  Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \ldots M^k \cdot r^{(0)}, \ldots$ approaches the dominant eigenvector of $M$ ($M$ is stochastic/Markov matrix)

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

$$Mx = \lambda x$$

  **Optimal r is the first or principal eigenvector of M, with corresponding eigenvalue 1**

# Why Power Iteration works? (2)




- **Claim:** Sequence $\boldsymbol{M} \cdot \boldsymbol{r}^{(0)}, \boldsymbol{M}^{2} \cdot \boldsymbol{r}^{(0)}, \ldots \boldsymbol{M}^{k} \cdot \boldsymbol{r}^{(0)}, \ldots$ approaches the dominant eigenvector of $\boldsymbol{M}$

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

- **Proof:**
  - Assume ***M*** has ***n*** linearly independent eigenvectors, $x_1, x_2, \ldots, x_n$ with corresponding eigenvalues $\lambda_1, \lambda_2, \ldots, \lambda_n$, where $\lambda_1 > \lambda_2 > \cdots > \lambda_n$
  - Vectors $x_1, x_2, \ldots, x_n$ form a basis and thus we can write: $r^{(0)} = c_1\, x_1 + c_2\, x_2 + \cdots + c_n\, x_n$
  - $\boldsymbol{M}\boldsymbol{r}^{(0)} = \boldsymbol{M}(\boldsymbol{c}_1\, \boldsymbol{x}_1 + \boldsymbol{c}_2\, \boldsymbol{x}_2 + \cdots + \boldsymbol{c}_n\, \boldsymbol{x}_n)$

    $= c_1(Mx_1) + c_2(Mx_2) + \cdots + c_n(Mx_n)$

    $= c_1(\lambda_1 x_1) + c_2(\lambda_2 x_2) + \cdots + c_n(\lambda_n x_n)$
  - **Repeated multiplication on both sides produces**

$$M^k r^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \cdots + c_n(\lambda_n^k x_n)$$

# Why Power Iteration works? (3)

- **Claim:** Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \ldots M^k \cdot r^{(0)}, \ldots$ approaches the dominant eigenvector of $M$

- **Proof (continued):**

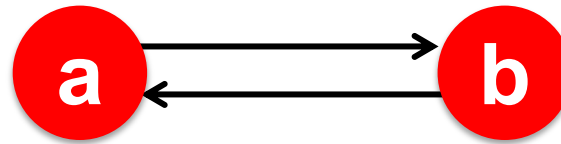  - Repeated multiplication on both sides produces

  $$M^k r^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \cdots + c_n(\lambda_n^k x_n)$$

  - $M^k r^{(0)} = \lambda_1^k \left[ c_1 x_1 + c_2 \left(\frac{\lambda_2}{\lambda_1}\right)^k x_2 + \cdots + c_n \left(\frac{\lambda_n}{\lambda_1}\right)^k x_n \right]$

  - Since $\lambda_1 > \lambda_2$ then fractions $\frac{\lambda_2}{\lambda_1}, \frac{\lambda_3}{\lambda_1} \ldots < 1$

    and so $\left(\frac{\lambda_i}{\lambda_1}\right)^k = 0$ as $k \to \infty$ (for all $i = 2 \ldots n$).

  - **Thus: $M^k r^{(0)} \approx c_1\left(\lambda_1^k x_1\right)$**

    - Note if $c_1 = 0$ then the method won't converge

  - The largest eigenvalue of a stochastic matrix is always 1.

# PageRank: Three Questions

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

- **Does this converge?**

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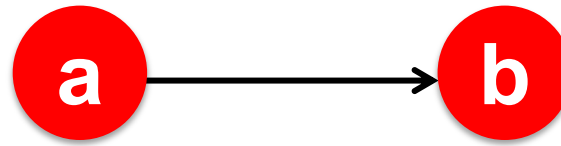
- **Are results reasonable?**

# Does this converge?



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

- **Example:**

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

Iteration 0, 1, 2, …

# Does it converge to what we want?



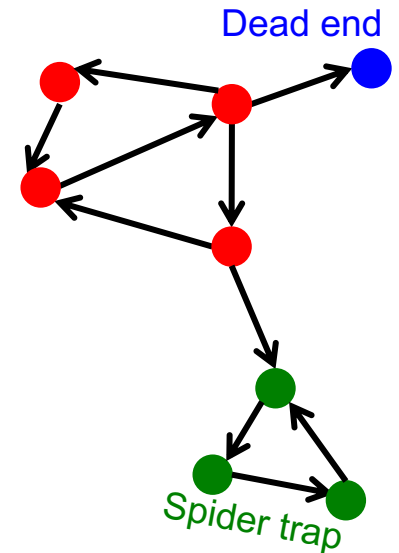$$r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i}$$

- **Example:**

$$\begin{matrix} r_a \\ r_b \end{matrix} \quad = \quad \begin{matrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{matrix}$$

Iteration 0, 1, 2, …
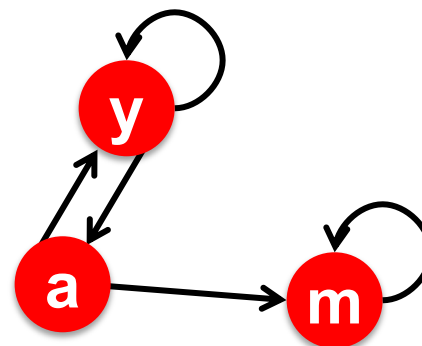
# PageRank: Problems

**2 problems:**

- **(1)** Some pages are
  **dead ends** (have no out-links)
    - "Vote" has "nowhere" to go to
    - Such pages cause importance to "leak out"

- **(2)** **Spider traps:**
  (all out-links are within the group)
    - "Vote" gets "stuck" in a trap
    - And eventually spider traps absorb all importance

# Problem: Spider Traps

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



m is a spider trap

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

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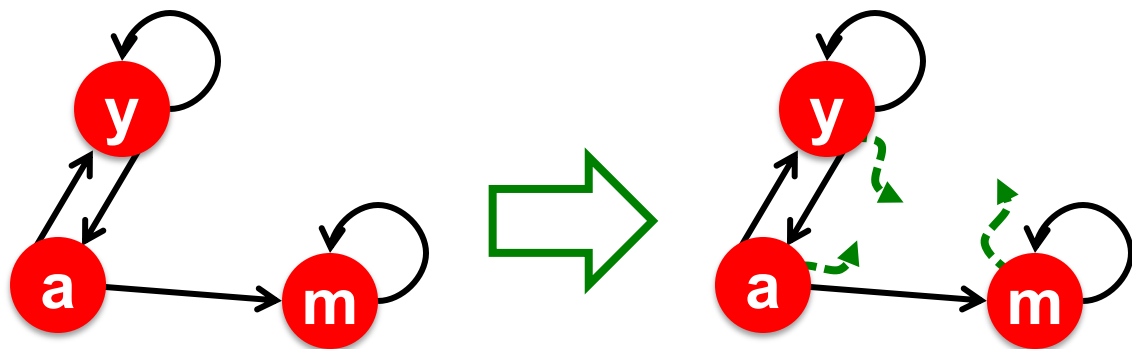
$r_a = r_y/2$

$r_m = r_a/2 + r_m$

- **Example:**

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

Iteration 0, 1, 2, …

All the PageRank score gets "trapped" in node m.
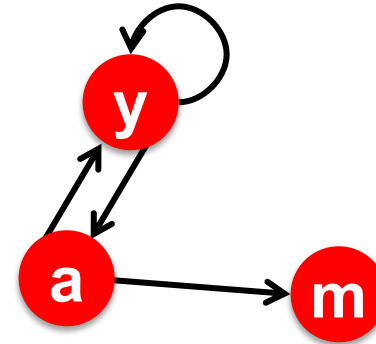
28

# Solution: Teleports!

- **The Google solution for spider traps: At each time step, the "vote" 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

- **"Vote" will teleport out of spider trap within a few time steps**

# 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 |

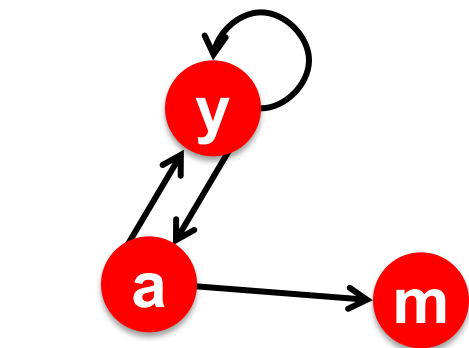$r_y = r_y/2 + r_a/2$

$r_a = r_y/2$

$r_m = r_a/2$

- **Example:**

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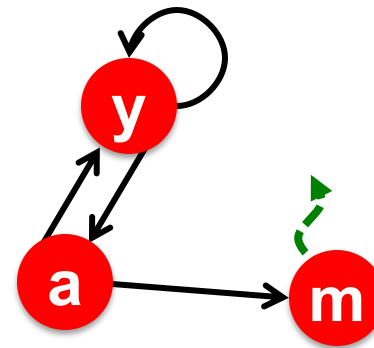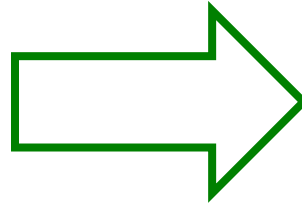
Iteration 0, 1, 2, …

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



|   | 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 |

# Why 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** [Larry Page and Sergey Brin 1998]

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

$d_i$ ... 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 $A$:**

$$A = \beta \, M + (1 - \beta) \left[ \frac{1}{N} \right]_{N \times N}$$
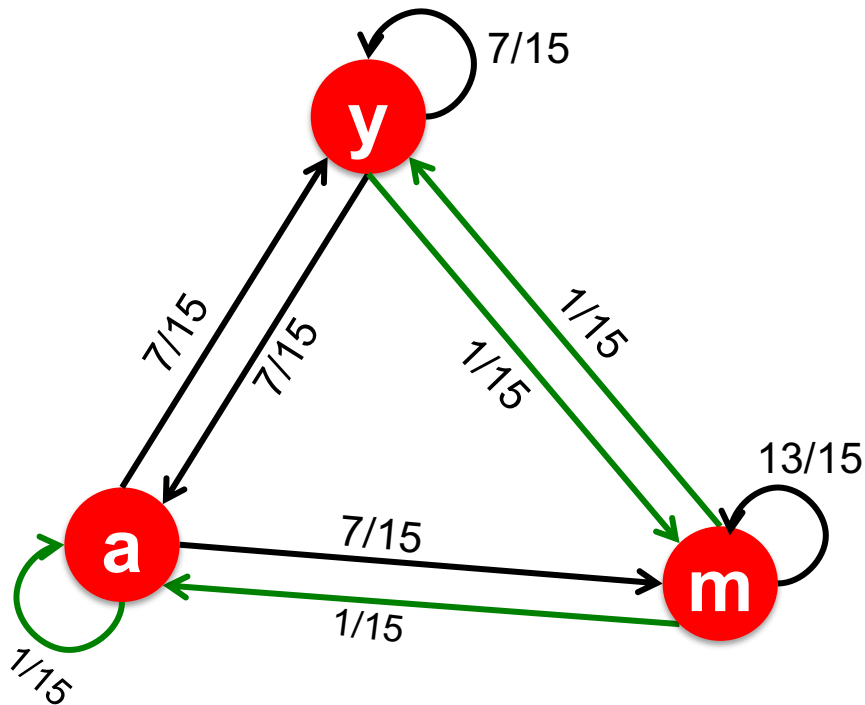
$[1/N]_{NxN}$…N by N matrix where all entries are 1/N

- **We have a recursive problem: $r = A \cdot r$**
  **And the Power method still works!**

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

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

# Random Teleports ($\beta$ = 0.8)



**M**

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

**[1/N]$_{N \times N}$**

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

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

# MapReduce Program for PageRank

Map(*key*, *value*) {
    // *key*:           a page,
    // *value*:       page rank of the page

    For each *page* in Adj[*key*]
        **emit(*page*, PR(*key*)/sizeof(Adj[*key*]));**
}

Reduce(*key*, *values*) {
    // *key*:           a page,
    // *values*:     a list of page ranks from all its incoming pages

    PR(*key*)=1-$\beta$;
    For each *pagerank* in *values*
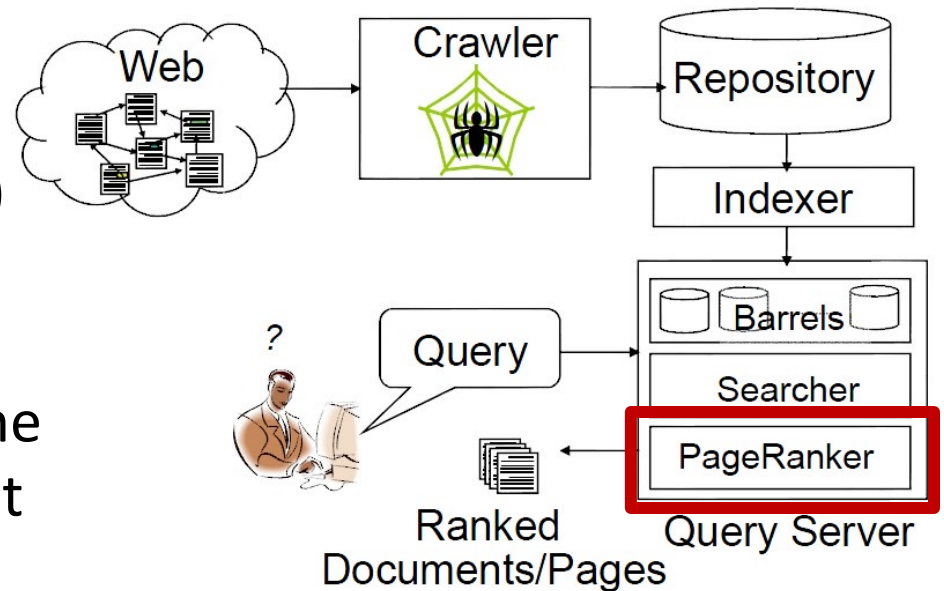        PR(*key*) = PR(*key*) + $\beta$**pagerank*;
    **emit(*key*, PR(*key*));**
}

# MapReduce Program for PageRank

```
A B C D
B A D
C A B
D B C
```
Links.txt

```
A 1
B 1
C 1
D 1
```
Initial PR

Map →

```
B 1/3
C 1/3
D 1/3
```

```
A 1/2
D 1/2
```

```
A 1/2
B 1/2
```

```
B 1/2
C 1/2
```

Reduce →

```
A 1/2
A 1/2
```

```
B 1/3
B 1/2
B 1/2
```

```
C 1/3
C 1/2
```

```
D 1/3
D 1/2
```

# Web Search Engines

- Indexer

    - Process the retrieved pages/documents and represents them in efficient search data structures (inverted files)

- Query Server

    - Accept the query from the user and return the result pages by consulting the search data structure