# Analysis of Large Graphs: Link Analysis, PageRank
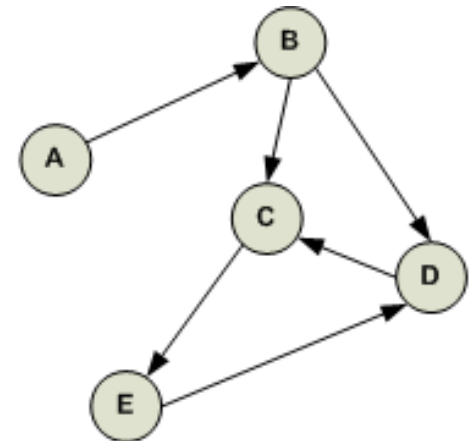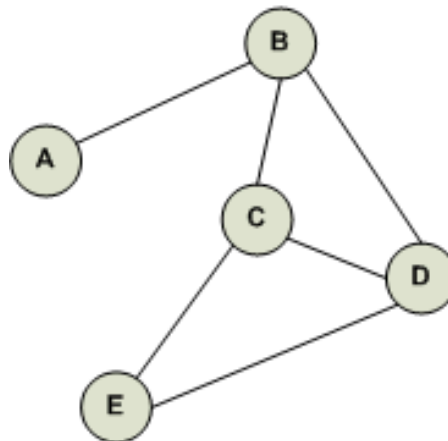
Big Data Analytics CSCI 4030

# New Topic: Graph Data!

| High dim. data | Graph data | Infinite data | Machine learning | Apps |
|---|---|---|---|---|
| Locality sensitive hashing | **PageRank,** SimRank | Filtering data streams | SVM | Recommen der systems |
| Clustering | Community Detection | Web advertising | Decision Trees | Association Rules |
| Dimensional ity reduction | Spam Detection | Queries on streams | Perceptron, kNN | Duplicate document detection |

Big Data Analytics CSCI 4030

# Graphs

- A **graph** is a representation of a set of objects (e.g., users, computers, …) where some pairs of objects are connected by links
- Objects are called nodes or vertices
- Links are called edges
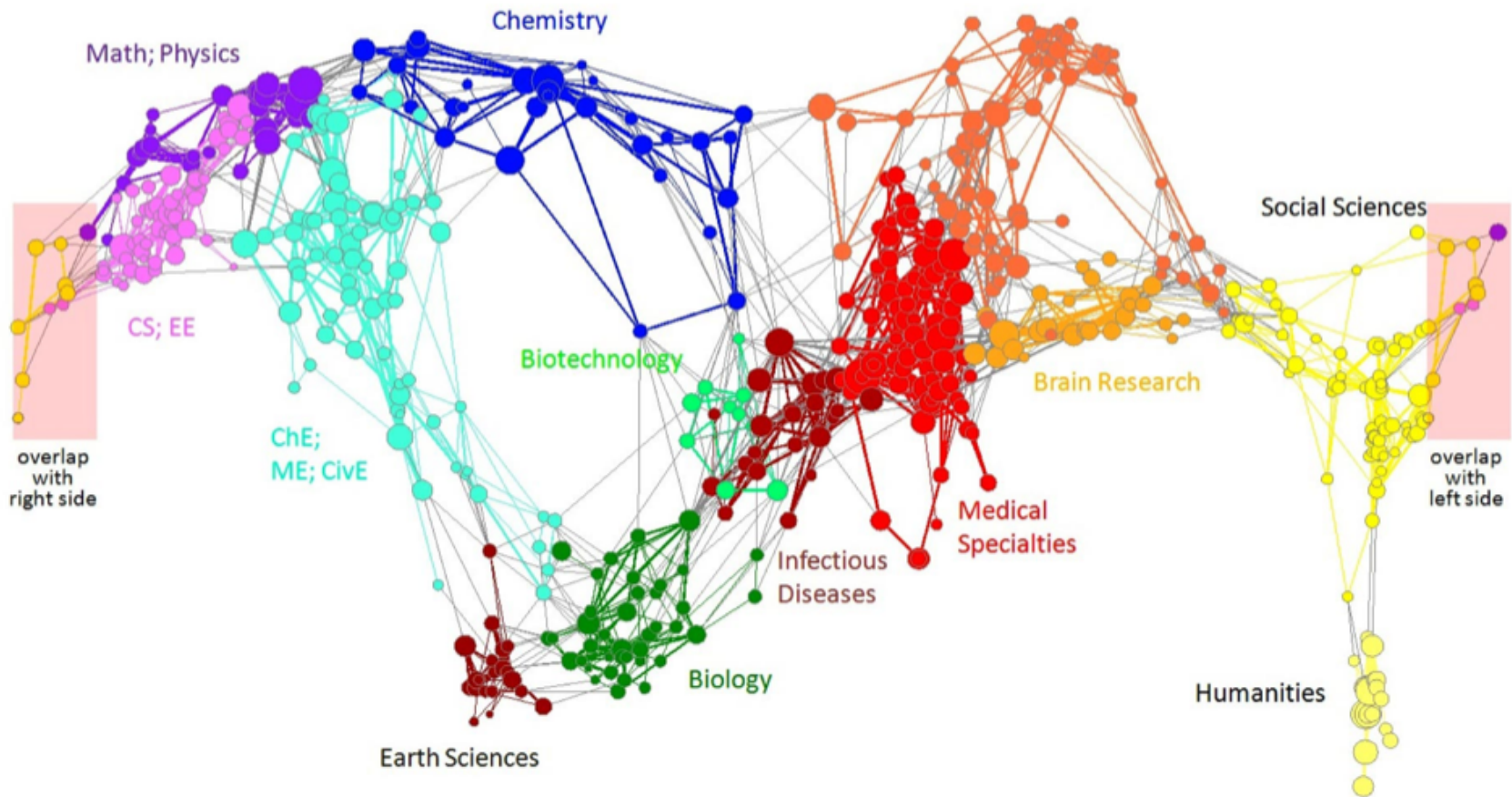- The edges may be **directed** or **undirected**

# Graph Data: Social Networks



**Facebook social graph**

4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]
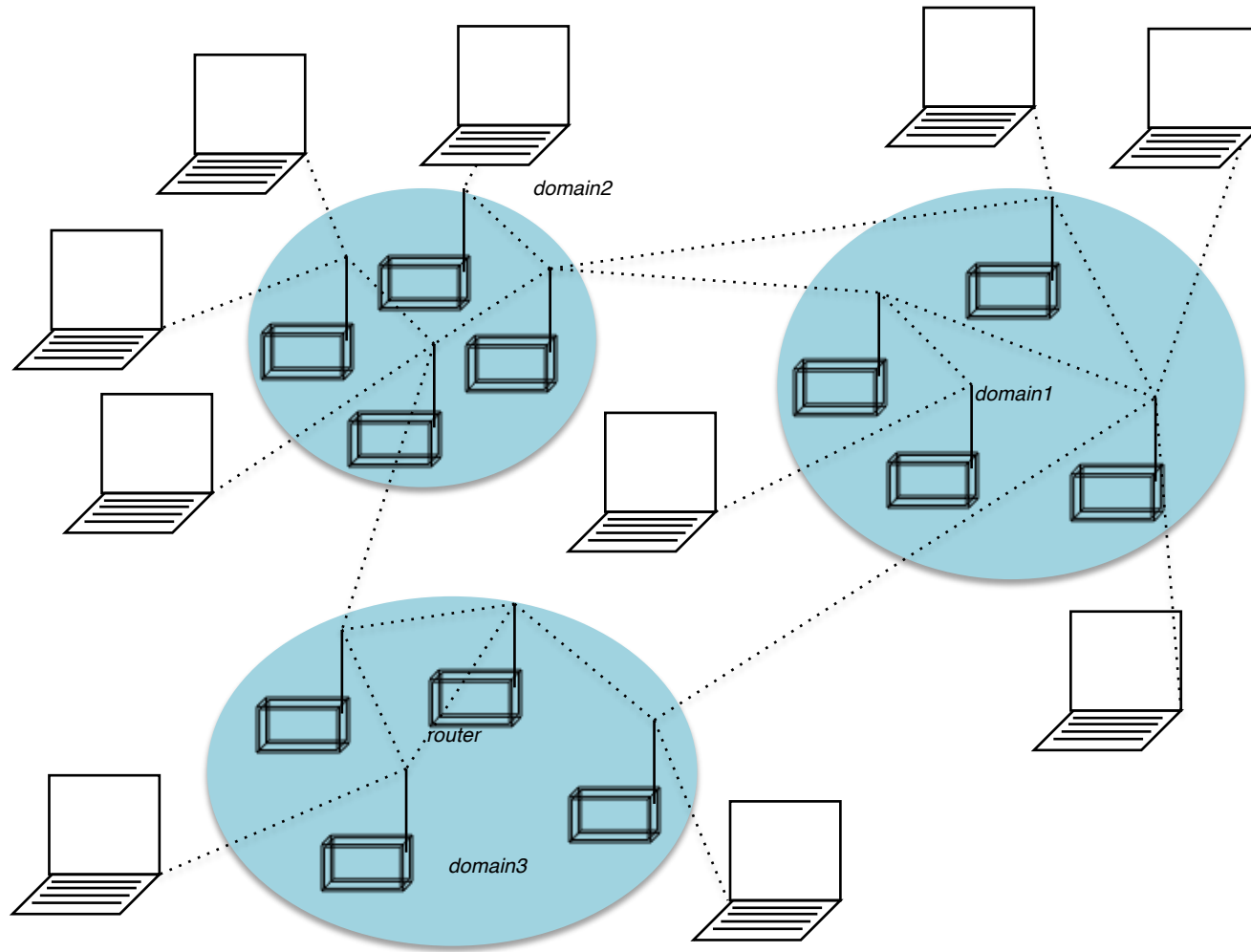
Big Data Analytics CSCI 4030

**Citation networks and Maps of science**

[Börner et al., 2012]

*domain2*

*domain1*

*router*

*domain3*

# Internet

# Web as a Graph

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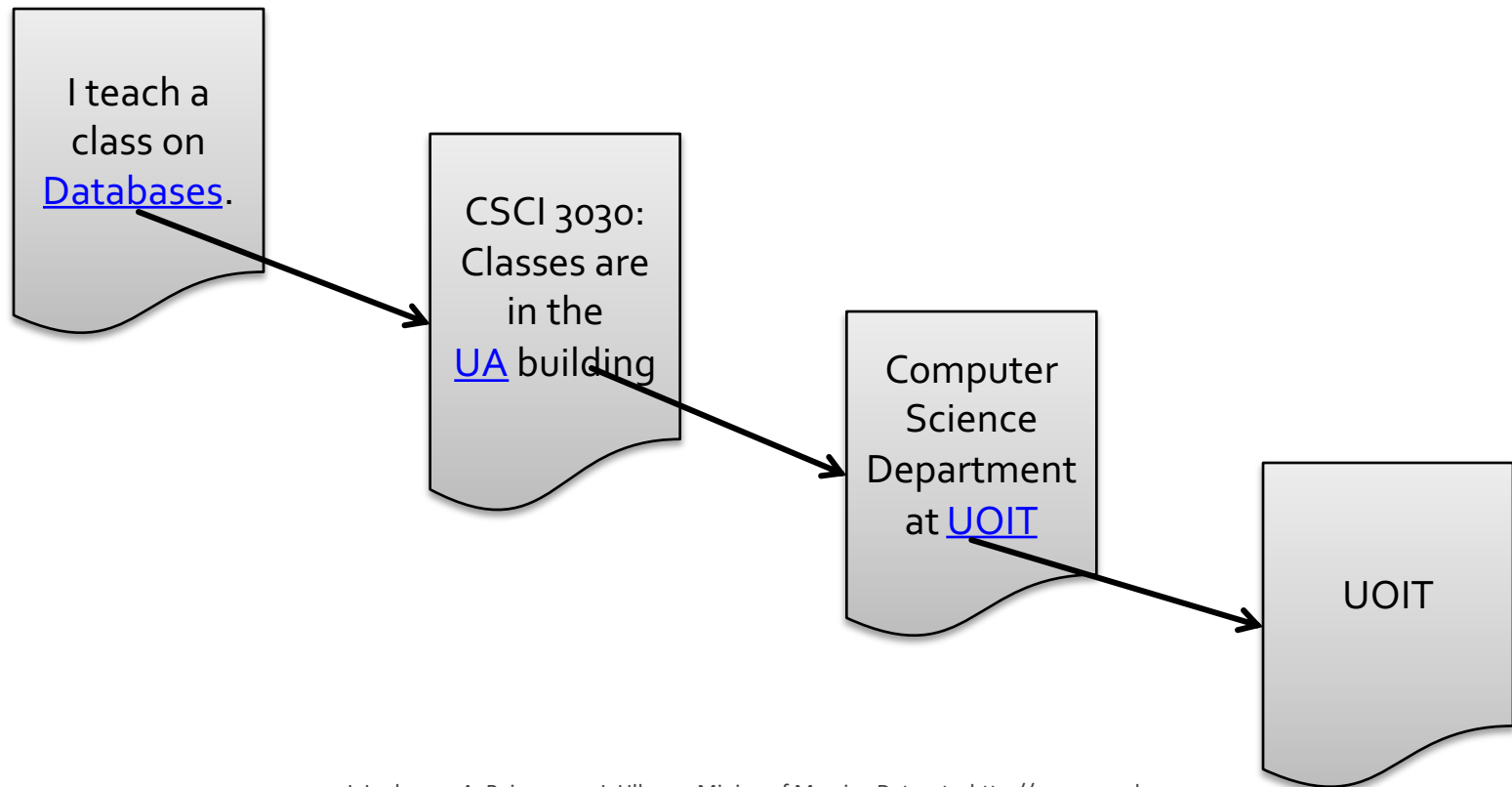  - **Edges: Hyperlinks**

I teach a class on Databases.

CSCI3030: Classes are in the UA building

Computer Science Department at UOIT

UOIT

# Web as a Graph

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

I teach a class on Databases.

CSCI 3030: Classes are in the UA building

Computer Science Department at UOIT

UOIT

# Broad Question

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

  **But:** Web is **huge**, full of **untrusted** documents, random things, web spam, etc.
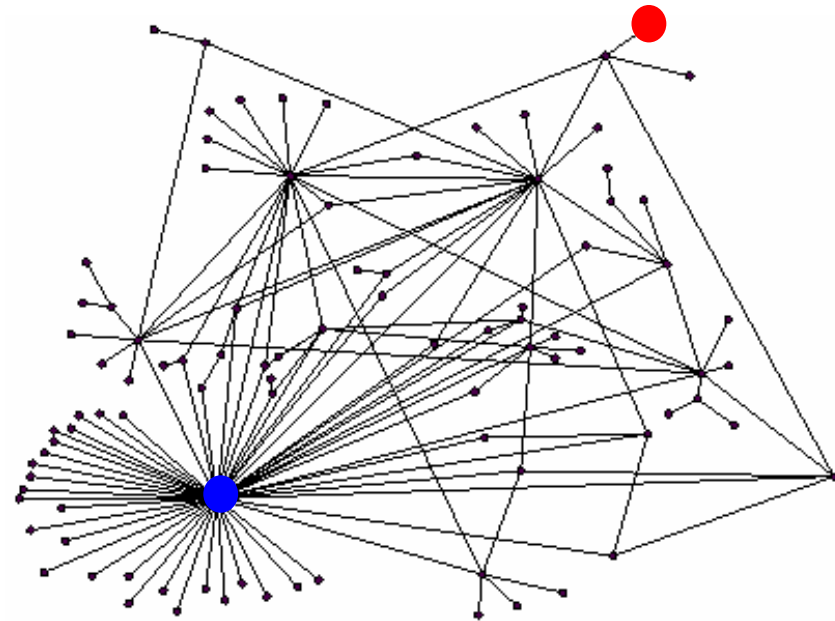
# Web Search: 2 Challenges

**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"**

  www.joe-schmoe.com vs. www.stanford.edu

- **Let's rank the pages by the link structure!**
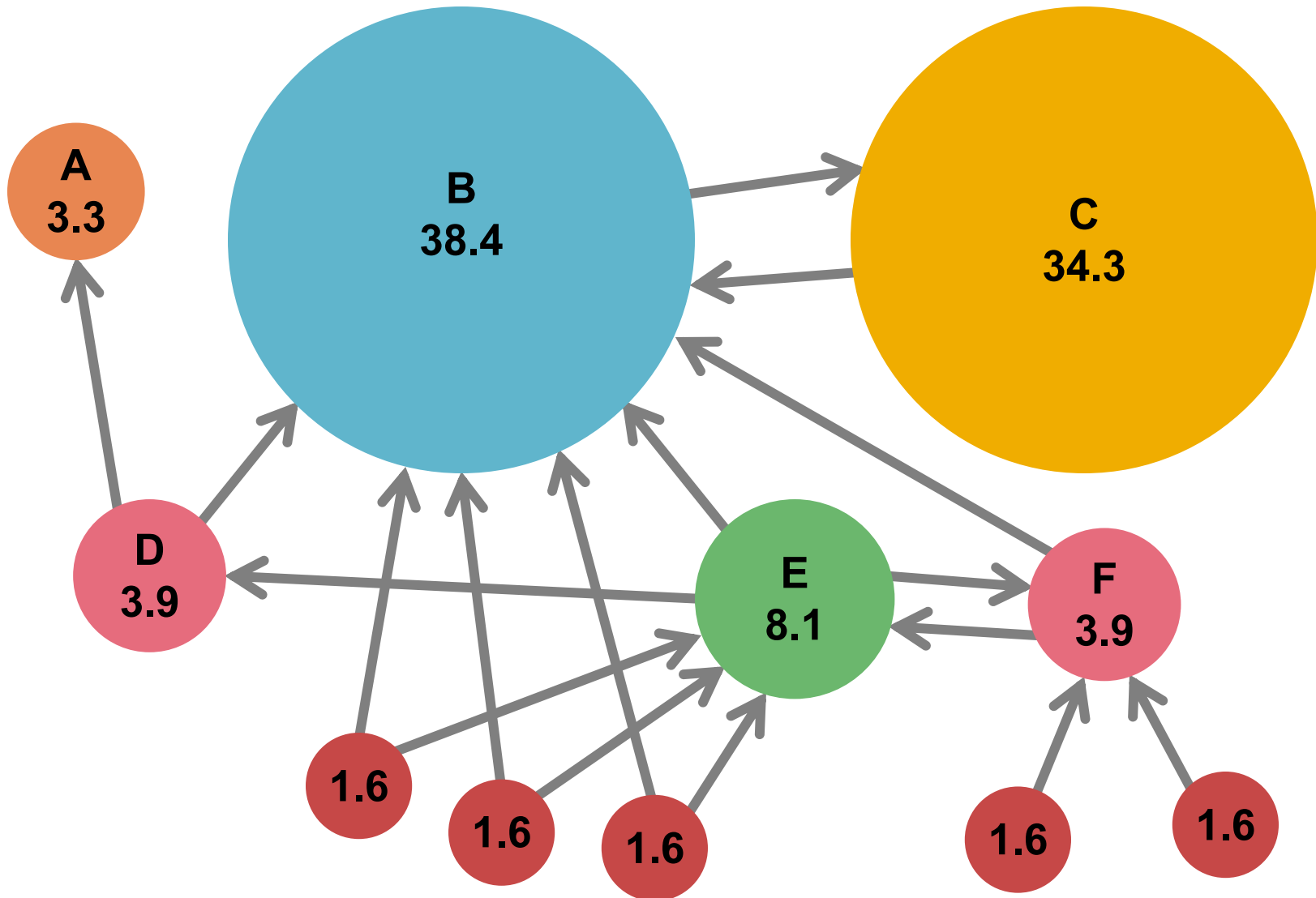
# Link Analysis Algorithms

- 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 has 23,400 in-links

  - www.joe-schmoe.com has 1 in-link

- ## Are all in-links are equal?

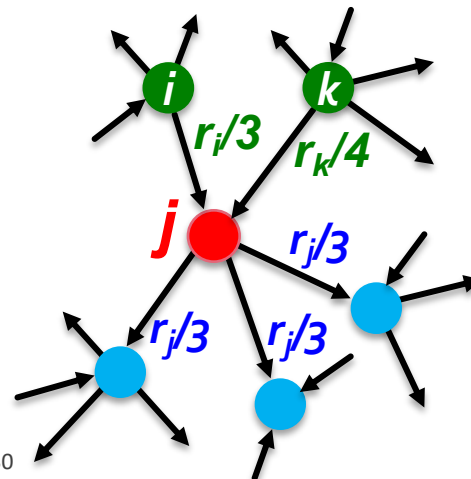  - ### Links from important pages count more

  - Recursive question!

# 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

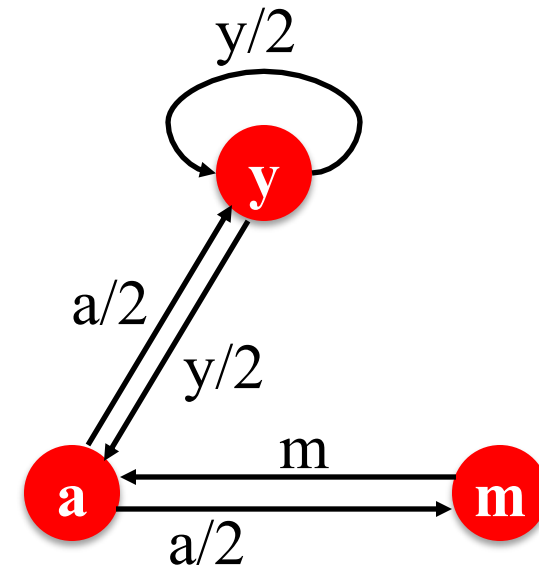- Page $j$'s own importance is the sum of the votes on its in-links

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

$r_i/3$   $r_k/4$

$r_j/3$   $r_j/3$   $r_j/3$

# 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 $i$**

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

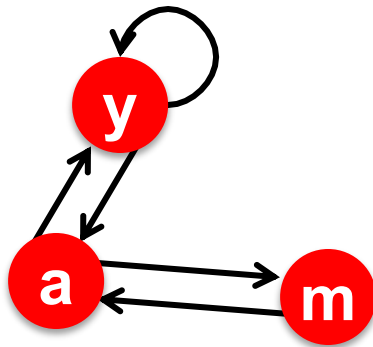- **3 equations, 3 unknowns, no constants**
  - No unique solution

- **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** (an algorithm for solving linear equations) **works for small examples, but we need a better method for large web-size graphs**
- **We need a new formulation!**

**Flow equations:**

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

# 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**

$$r = M \cdot r$$

# Example: Flow Equations & M



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

$$r = M \cdot r$$

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

$$\begin{bmatrix} y \\ a \\ m \end{bmatrix} = \begin{bmatrix} ½ & ½ & 0 \\ ½ & 0 & 1 \\ 0 & ½ & 0 \end{bmatrix} \begin{bmatrix} y \\ a \\ m \end{bmatrix}$$

# 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: $\mathbf{r}^{(0)} = [1/N, \ldots, 1/N]^T$
  - Iterate: $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$
  - Stop when $|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}| < \varepsilon$

# PageRank: How to solve?

- **Power Iteration:**
  - $r^{(0)} = [1/N,....,1/N]^T$
  - $r^{(t+1)} = M \cdot r^{(t)}$
  - $|r^{(t+1)} - r^{(t)}|_1 < \varepsilon$



|   | 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:**

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

Iteration 0, 1, 2, …

# Why Power Iteration works?

- **Power iteration:**
  - $r^{(1)} = M \cdot r^{(0)}$
  - $r^{(2)} = M \cdot r^{(1)} = M(Mr^{(1)}) = M^2 \cdot r^{(0)}$
  - $r^{(3)} = M \cdot r^{(2)} = M(M^2 r^{(0)}) = M^3 \cdot r^{(0)}$
- **Claim:**

  Sequence $M \cdot r^{(0)}, M^2 \cdot r^{(0)}, \dots M^k \cdot r^{(0)}, \dots$ will converge (under which **conditions**?!)

# PageRank:
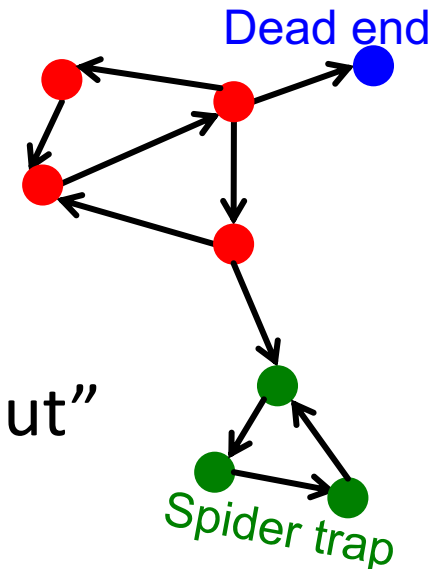# The Google Formulation

# PageRank: Three Questions

$$r = Mr$$

- **Does this converge?**

- **Does it converge to what we want?**

- **Are results reasonable?**
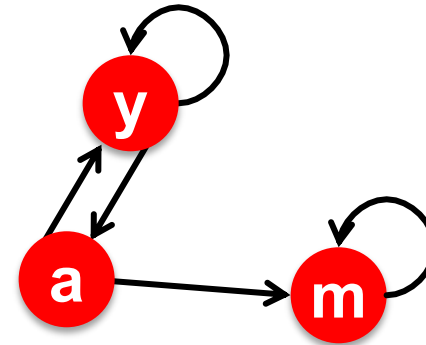
# PageRank: Problems

**2 problems:**

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

- **(2) Spider traps:**
  (all out-links are within the group)
  - And eventually spider traps absorb all importance

Dead end

Spider trap

# Problem: Spider Traps

- ## **Power Iteration:**



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

m is a spider trap

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

- ## **Example:**

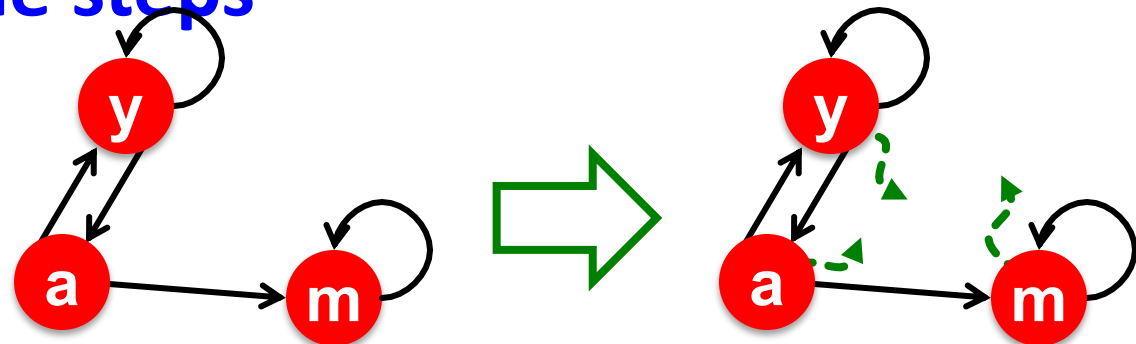$$r_a = r_y/2$$

$$r_m = r_a/2 + r_m$$

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

Iteration 0, 1, 2, …
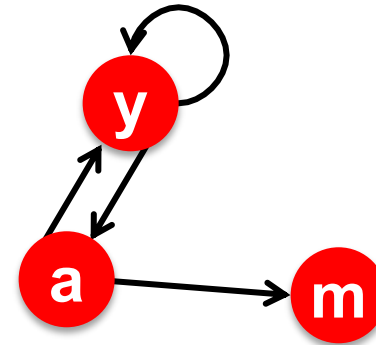
All the PageRank score gets "trapped" in node m.

# Solution: Teleports!

- **The Google solution for spider traps: At each time step, the random surfer has two options**

  - With probability $\beta$, follow a link at random

  - With probability **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:**

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

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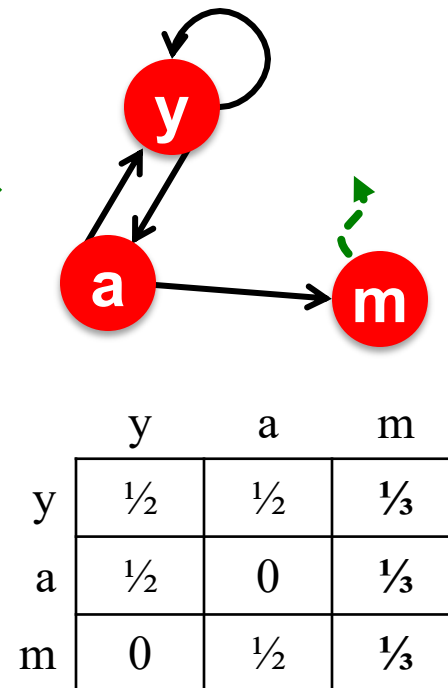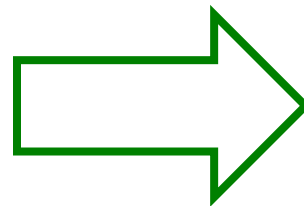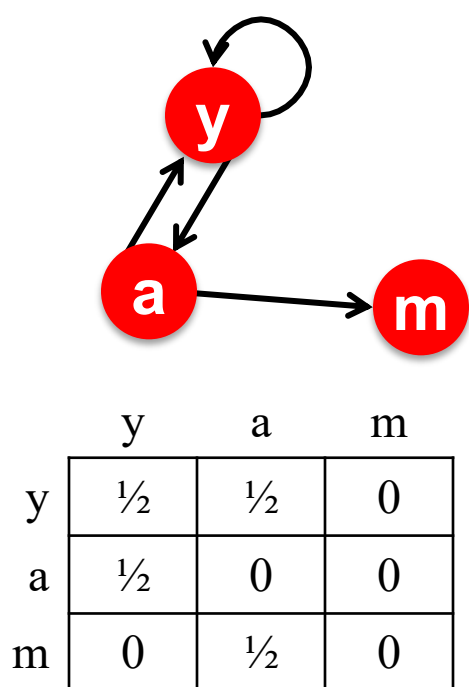$$r_a = r_y/2$$
$$r_m = r_a/2$$

- ## **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 |

# Solution: Random Teleports

- **Google's solution that does it all:**
  At each step, random surfer has two options:

  - With probability $\boldsymbol{\beta}$, follow a link at random
  - With probability $\boldsymbol{1-\beta}$, jump to some random page

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

# The Google Matrix

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

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

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

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

- **We have a recursive equation:**

  **And the Power method still works!**

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

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

**M**

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

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

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

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

# Some Problems with Page Rank

- **Measures generic popularity of a page**
  - Biased against topic-specific authorities
  - **Solution:** Topic-Specific PageRank (**next**)

- **Sensitive to Link spam**
  - Artificial link topographies created in order to boost page rank
  - **Solution:** TrustRank

# Topic-Specific PageRank

# Topic-Specific PageRank

- **Instead of generic popularity, can we measure popularity within a topic?**
- **Example:** Query "Trojan" wants different pages depending on whether you are interested in sports, history and computer security
- **Goal:** Evaluate Web pages not just according to their popularity, but by how close they are to a particular topic, e.g. "sports" or "history"
- **Allows search queries to be answered based on interests of the user**

# Topic-Specific PageRank

- Random walker has a small probability of teleporting at any step
- **Teleport can go to:**
  - **Standard PageRank: Any page with equal probability**
    - To avoid dead-end and spider-trap problems
  - **Topic Specific PageRank: A topic-specific set of "relevant" pages (teleport set)**
- **Idea: Bias the random walk**
  - When walker teleports, she pick a page from a set *S*
  - *S* contains only pages that are relevant to the topic
    - E.g., Open Directory (DMOZ) pages for a given topic/query

# Matrix Formulation

- **To make this work all we need is to update the teleportation part of the PageRank formulation:**

$$A_{ij} = \begin{cases} \beta\, M_{ij} + (1-\beta)/|S| & \text{if } i \in S \\ \beta\, M_{ij} + 0 & \text{otherwise} \end{cases}$$

  - *A* is stochastic!

- We weighted all pages in the teleport set **S** equally

- **Compute as for regular PageRank:**

  - Multiply by **M**, then add a vector

## Suppose $S = \{1\}$, $\beta = 0.8$

| Node | Iteration | | | | |
|------|------|------|------|------|--------|
|      | **0** | **1** | **2** | **…** | **stable** |
| 1 | 0.25 | 0.4 | 0.28 | | 0.294 |
| 2 | 0.25 | 0.1 | 0.16 | | 0.118 |
| 3 | 0.25 | 0.3 | 0.32 | | 0.327 |
| 4 | 0.25 | 0.2 | 0.24 | | 0.261 |

**S={1,2,3,4}, β=0.8:**
r=[0.13, 0.10, 0.39, 0.36]

**S={1,2,3} , β=0.8:**
r=[0.17, 0.13, 0.38, 0.30]

**S={1,2} , β=0.8:**
r=[0.26, 0.20, 0.29, 0.23]

**S={1} , β=0.8:**
r=[0.29, 0.11, 0.32, 0.26]

**S={1}, β=0.90:**
r=[0.17, 0.07, 0.40, 0.36]

**S={1} , β=0.8:**
r=[0.29, 0.11, 0.32, 0.26]

**S={1}, β=0.70:**
r=[0.39, 0.14, 0.27, 0.19]

# TrustRank:
# Combating the Web Spam

# What is Web Spam?

- **Spamming:**
  - Any deliberate action to boost a web page's position in search engine results, incommensurate with page's real value
- **Spam:**
  - Web pages that are the result of spamming
- This is a very broad definition
  - **SEO** industry might disagree!
  - SEO = search engine optimization

- Approximately **10-15%** of web pages are spam

# Web Search

- **Early search engines:**
  - Crawl the Web
  - Index pages by the **words** they contained
  - Respond to search queries (lists of words) with the pages **containing those words**
- **Early page ranking:**
  - Attempt to order **pages matching a search query** by "importance"
  - **First search engines considered:**
    - **(1)** Number of times query words appeared
    - **(2)** Prominence of word position, e.g. title, header

# First Spammers

- As people began to use search engines, those with commercial interests tried to **exploit search engines**

- **Example:**
  - Shirt-seller might pretend to be about "movies"
- **Techniques for achieving high relevance/importance for a web page**

# First Spammers: Term Spam

- **How do you make your page appear to be about movies?**
  - **(1)** Add the word movie 1,000 times to your page
  - Set text color to the background color, so only search engines would see it
  - **(2)** Or, run the query "movie" on your target search engine
  - See what page came first in the listings
  - Copy it into your page, make it "invisible"
- **These and similar techniques are term spam**

# Google Bomb
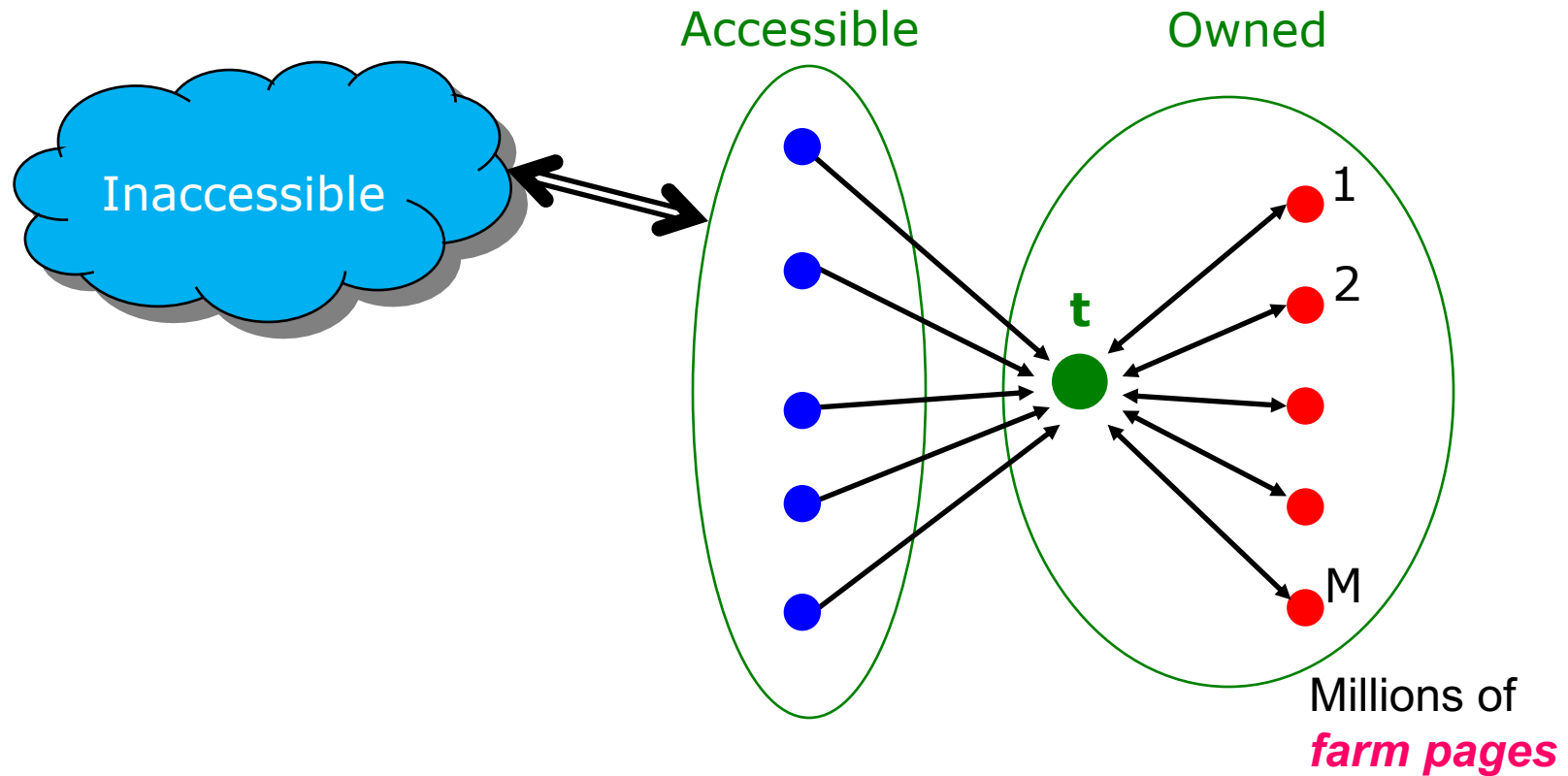
SPAM FARMING

# Google vs. Spammers: Round 2!

- Once Google became the dominant search engine, spammers began to work out ways to fool Google

- **Spam farms** were developed to concentrate PageRank on a single page

- **Link spam:**
  - Creating link structures that boost PageRank of a particular page

# Link Spamming

- **Three kinds of web pages from a spammer's point of view**

  - **Inaccessible pages**

  - **Accessible pages**

    - e.g., blog comments pages

    - spammer can post links to his pages

  - **Owned pages**

    - Completely controlled by spammer

    - May span multiple domain names

# Link Farms



Accessible  Owned

Inaccessible

t

1
2

M

Millions of
*farm pages*

## Get as many links from accessible pages as possible to target page t

# TrustRank:
# Combating the Web Spam

# Combating Spam

- **Combating link spam**

  - **Detection and blacklisting of structures that look like spam farms**

    - Leads to another war – hiding and detecting spam farms

  - **TrustRank**

# TrustRank: Idea

- **Basic principle: Approximate isolation**
  - It is rare for a "good" page to point to a "bad" (spam) page

- Choose a set of seed pages that are identified as **good** the **trusted pages**

- Perform a topic-sensitive PageRank with **teleport set = trusted pages**
  - **Propagate trust through links**
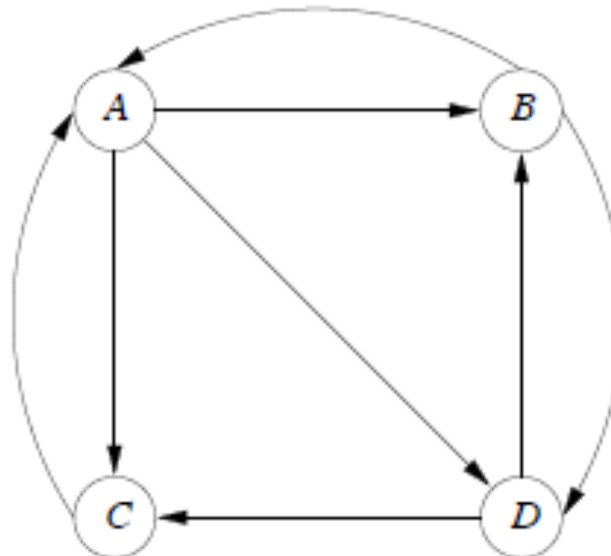
# Picking the Seed Set

- **Two conflicting considerations:**

  - Human has to inspect each seed page, so seed set must be as small as possible

  - Must ensure every **good page** gets adequate trust rank, so we need to make all good pages reachable from seed set by short paths

# Approaches to Picking Seed Set

- Suppose we want to pick a seed set of *k* pages
- **How to do that?**
- **(1) PageRank:**
  - Pick the top *k* pages by PageRank
  - Theory is that you can't get a bad page's rank really high
- **(2) Use trusted domains** whose membership is controlled, like .edu, .mil, .gov

# Quiz: Page Rank

- Given an example of the Web
  - Compute transition matrix
  - What is the meaning of the transition matrix?
  - Compute Page Rank.

# Quiz: Page Rank

- ## Answer
  - The transition matrix for the presented Web is

$$
M = \begin{bmatrix}
0 & 1/2 & 1 & 0 \\
1/3 & 0 & 0 & 1/2 \\
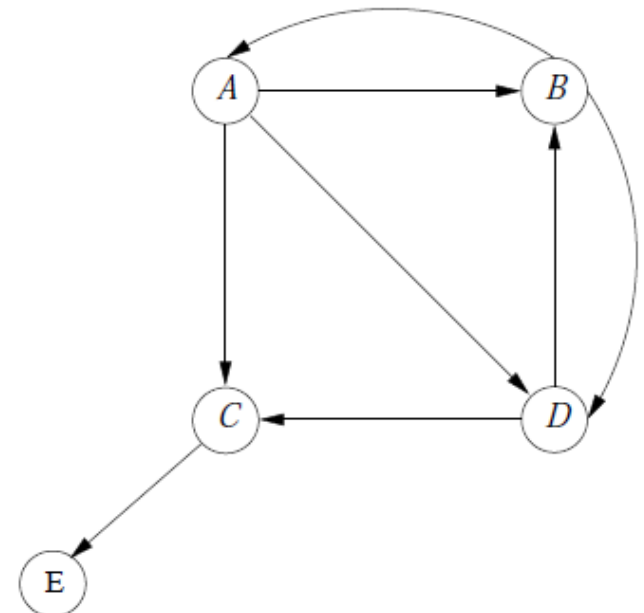1/3 & 0 & 0 & 1/2 \\
1/3 & 1/2 & 0 & 0
\end{bmatrix}
$$

- The sequence of approximations to the limit that we get by multiplying at each step by M is:
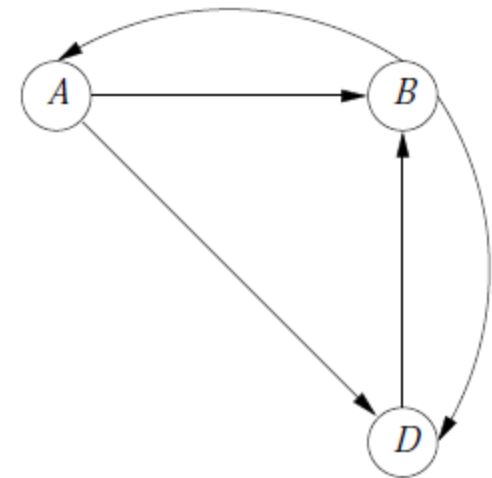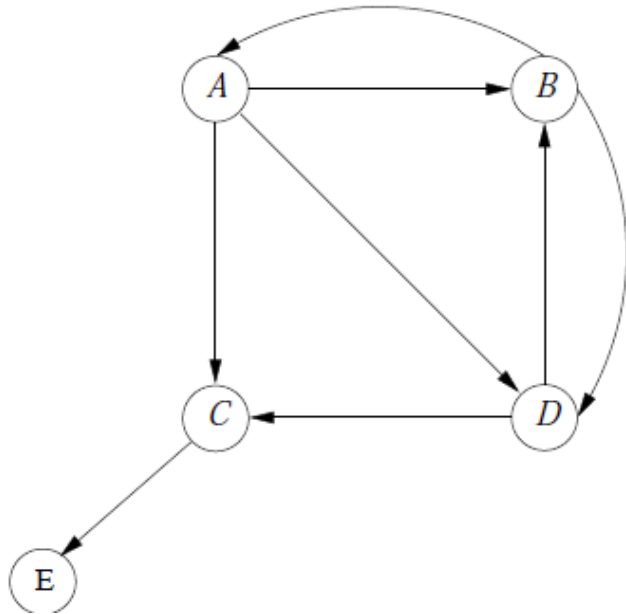
$$\begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \end{bmatrix}, \begin{bmatrix} 9/24 \\ 5/24 \\ 5/24 \\ 5/24 \end{bmatrix}, \begin{bmatrix} 15/48 \\ 11/48 \\ 11/48 \\ 11/48 \end{bmatrix}, \begin{bmatrix} 11/32 \\ 7/32 \\ 7/32 \\ 7/32 \end{bmatrix}, \ldots, \begin{bmatrix} 3/9 \\ 2/9 \\ 2/9 \\ 2/9 \end{bmatrix}$$

# Quiz: Dead Ends

- One way to deal with dead ends is to drop the dead ends from the graph, and also drop their incoming arcs. Doing so may create more dead ends, which have to be dropped, recursively.
- Apply this method on the dead ends of the following graph and draw the new graph.
- Compute the transition matrix of the new graph.

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

# Actions

- Review slides!

- Read Chapter 5 (Link Analysis) from course book.

    - You can find electronic version of the book on Blackboard.