# COMP 345: Data Mining
# Analysis of Large Graphs:
## Link Analysis, PageRank
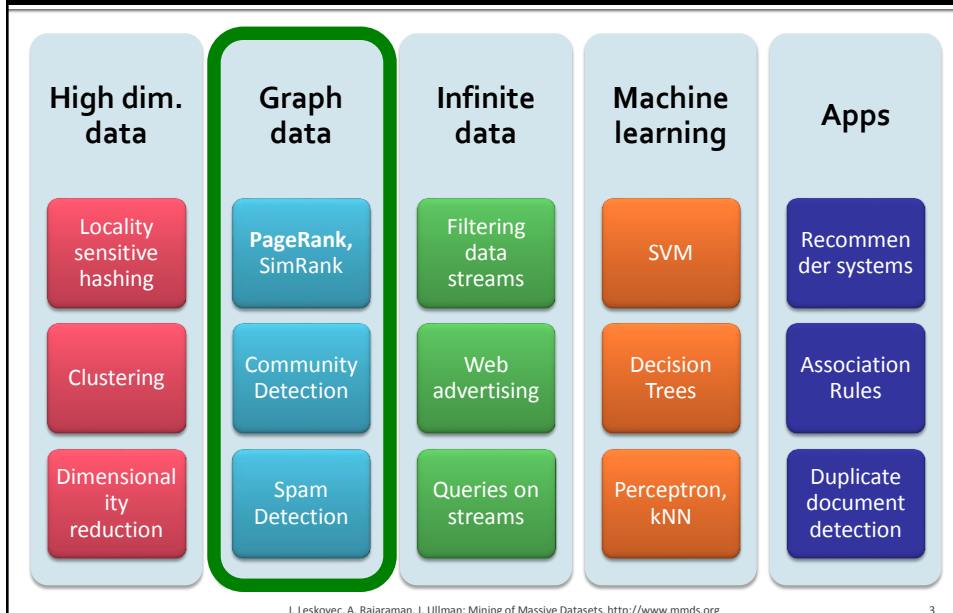
Slides Adapted From: www.mmds.org (Mining Massive Datasets)

Rhodes College

---

# Announcements

- For next time, watch the 3 video lectures on Moodle about MapReduce and take the online quiz.

2

# New Topic: Graph Data!

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

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org    3

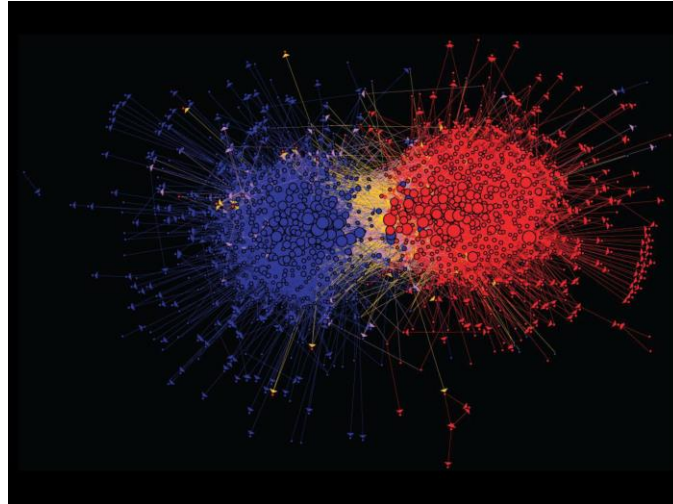# Graph Data: Social Networks



**Facebook social graph**
**4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]**

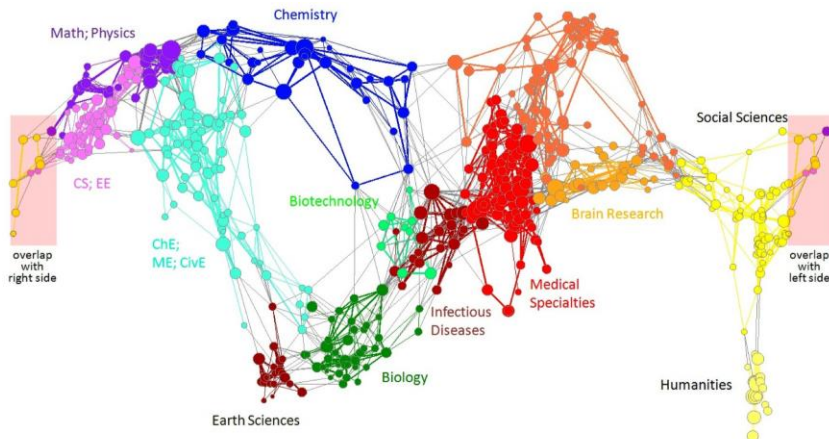J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org    4

2

# Graph Data: Media Networks



**Connections between political blogs**
Polarization of the network [Adamic-Glance, 2005]

# Graph Data: Information Nets



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

# Graph Data: Communication Nets
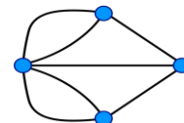


Internet

# Graph Data: Technological Networks



**Seven Bridges of Königsberg**
[Euler, 1735]
Return to the starting point by traveling each
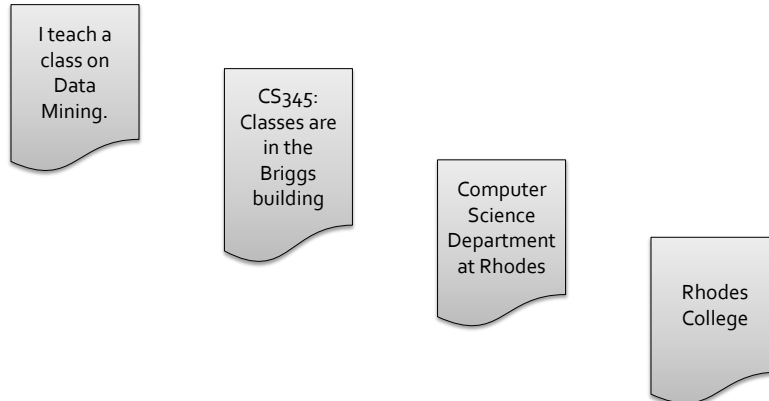link of the graph once and only once.

# Web as a Graph

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

I teach a class on Data Mining.

CS345: Classes are in the Briggs building

Computer Science Department at Rhodes

Rhodes College

# Web as a Graph

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

I teach a class on Data Mining.

CS345: Classes are in the Briggs building

Computer Science Department at Rhodes

Rhodes College

# Web as a Directed Graph



11

# 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.

12

## Web Search: 2 Challenges
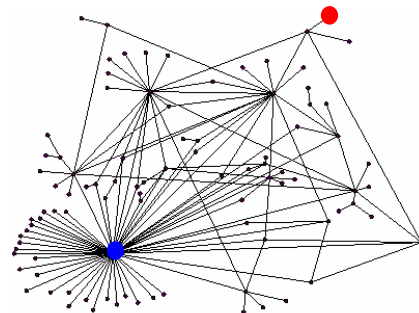
**2 challenges of web search:**
- **(1) Web contains many sources of information Who to "trust"?**
  - **Trick:** Trustworthy pages may point to each other!

- **(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

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org          13

## Ranking Nodes on the Graph

- **All web pages are not equally "important"**
  www.joe-schmoe.com vs. www.stanford.edu

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



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org          14

# Link Analysis Algorithms

- We will cover the following **Link Analysis approaches** for computing **importance** of nodes in a graph:
  - Page Rank
  - Topic-Specific (Personalized) Page Rank
  - Web Spam Detection Algorithms

15
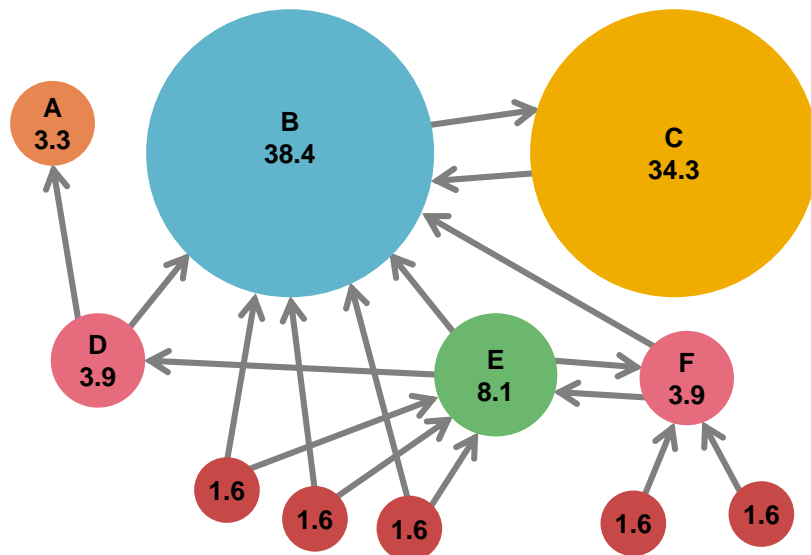
# PageRank:
# The "Flow" Formulation

# 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!

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org     17
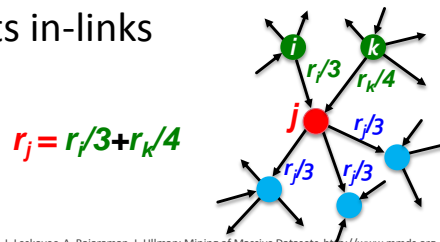
# Example: PageRank Scores



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org     18

# 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

$r_j = r_i/3 + r_k/4$



$r_i/3$    $r_k/4$

$r_j/3$

$r_j/3$   $r_j/3$

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org    19

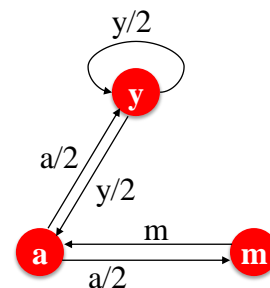# PageRank: The "Flow" Model

- **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



y/2

a/2   y/2

m   a/2

**"Flow" equations:**
$r_y = r_y/2 + r_a/2$
$r_a = r_y/2 + r_m$
$r_m = r_a/2$

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org    20

# 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:** $r_y = \frac{2}{5}, \ r_a = \frac{2}{5}, \ r_m = \frac{1}{5}$
- **Gaussian elimination method 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**

$$r = M \cdot r$$

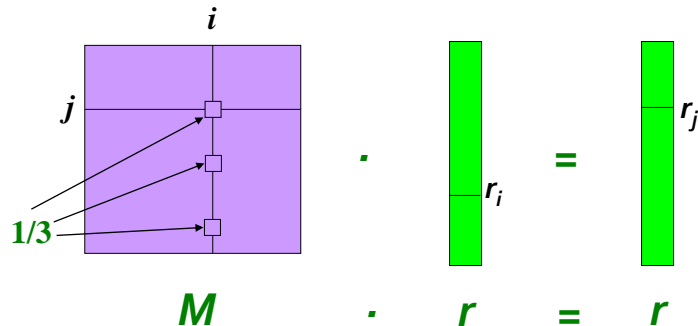$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

# Example

- **Remember the flow equation:** $r_j = \sum_{i \to j} \dfrac{r_i}{d_i}$
- **Flow equation in the matrix form**

$$M \cdot r = r$$

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



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

23

# 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**

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

    - Largest eigenvalue of **M** is **1** since **M** is column stochastic (with non-negative entries)
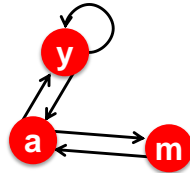      - *We know **r** is unit length and each column of **M** sums to one, so $Mr \le 1$*

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

24

12

# 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: $\mathbf{r}^{(0)} = [1/N,....,1/N]^T$
  - 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

    $|\mathbf{x}|_1 = \sum_{1 \leq i \leq N} |x_i|$ is the **L₁** norm
    Can use any other vector norm, e.g., Euclidean
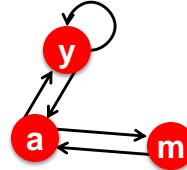
# PageRank: How to solve?

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

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 \\ 1/3 \\ 1/3 \end{matrix}$$

Iteration 0, 1, 2, …

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

$r_y = r_y/2 + r_a/2$
$r_a = r_y/2 + r_m$
$r_m = r_a/2$
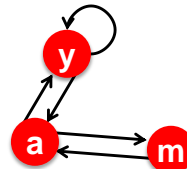
# PageRank: How to solve?

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

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

Iteration 0, 1, 2, …

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

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

# Random Walk Interpretation

- **Imagine a random web surfer:**
  - At any time $t$, surfer is on some page $i$
  - At time $t + 1$, the surfer follows an out-link from $i$ uniformly at random
  - Ends up on some page $j$ linked from $i$
  - Process repeats indefinitely
- **Let:**
  - $p(t)$ ... vector whose $i$th coordinate is the prob. that the surfer is at page $i$ at time $t$
  - So, $p(t)$ is a probability distribution over pages

$$r_j = \sum_{i \to j} \frac{r_i}{d_{out}(i)}$$

# The Stationary Distribution

- **Where is the surfer at time $t+1$?**
  - Follows a link uniformly at random
    $$p(t + 1) = M \cdot p(t)$$
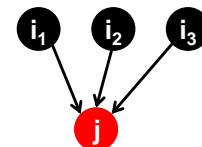
$p(t+1) = M \cdot p(t)$

- Suppose the random walk reaches a state
  $$p(t + 1) = M \cdot p(t) = p(t)$$
  then $p(t)$ is **stationary distribution** of a random walk
- **Our original rank vector $r$ satisfies $r = M \cdot r$**
  - **So, $r$ is a stationary distribution for the random walk**

# Existence and Uniqueness

- **A central result from the theory of random walks (a.k.a. Markov processes):**

For graphs that satisfy **certain conditions**, the **stationary distribution is unique** and eventually will be reached no matter what the initial probability distribution at time **t = 0**