# CS 422 Data Mining
# Lecture 11
# November 8, 2018

❑ Acknowledgment:

  ❑ This presentation is based on the book "Mining of Massive Datasets" by Anand Rajaraman and Jeff Ullman and the presentations by Jure Leskovec

❑ PageRank

❑ Page Ranks

❑ Mining Massive Datasets Jure Leskovec, Stanford UnivCS246: Mining Massive Datasets Jure Leskovec, Stanford University

  http://cs246.stanford.edu

# Social Network Graphs

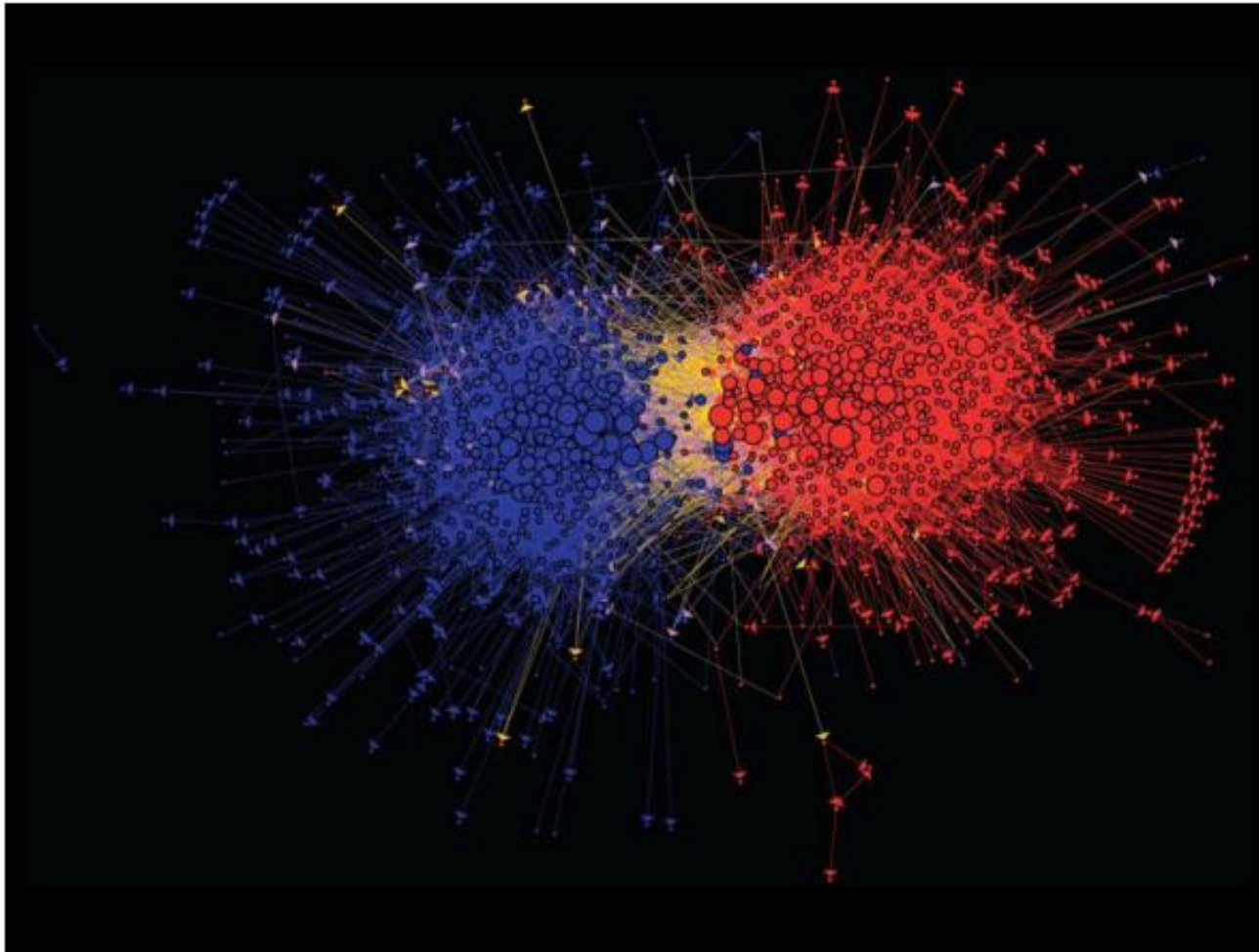

**Facebook social graph**
4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]

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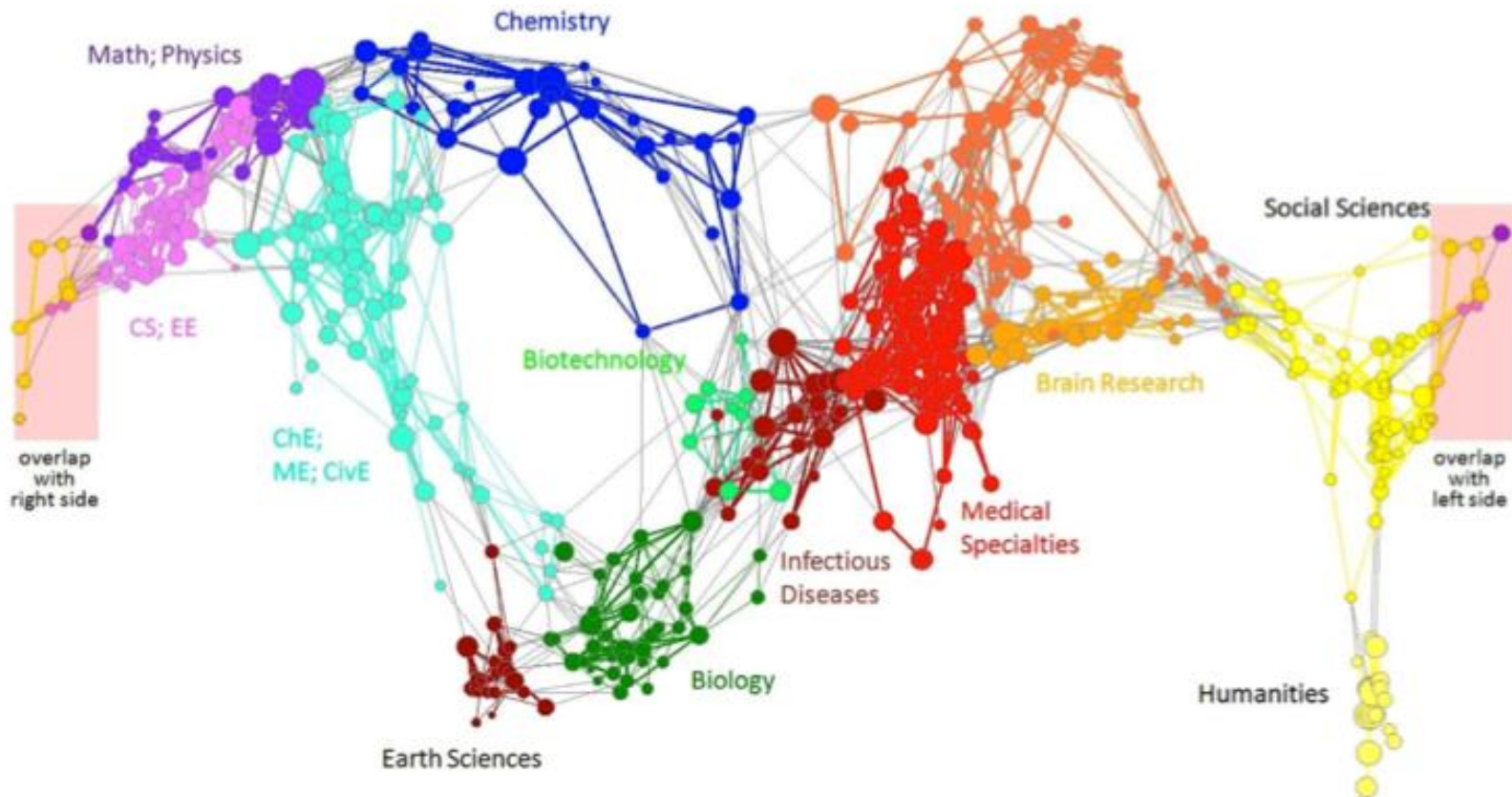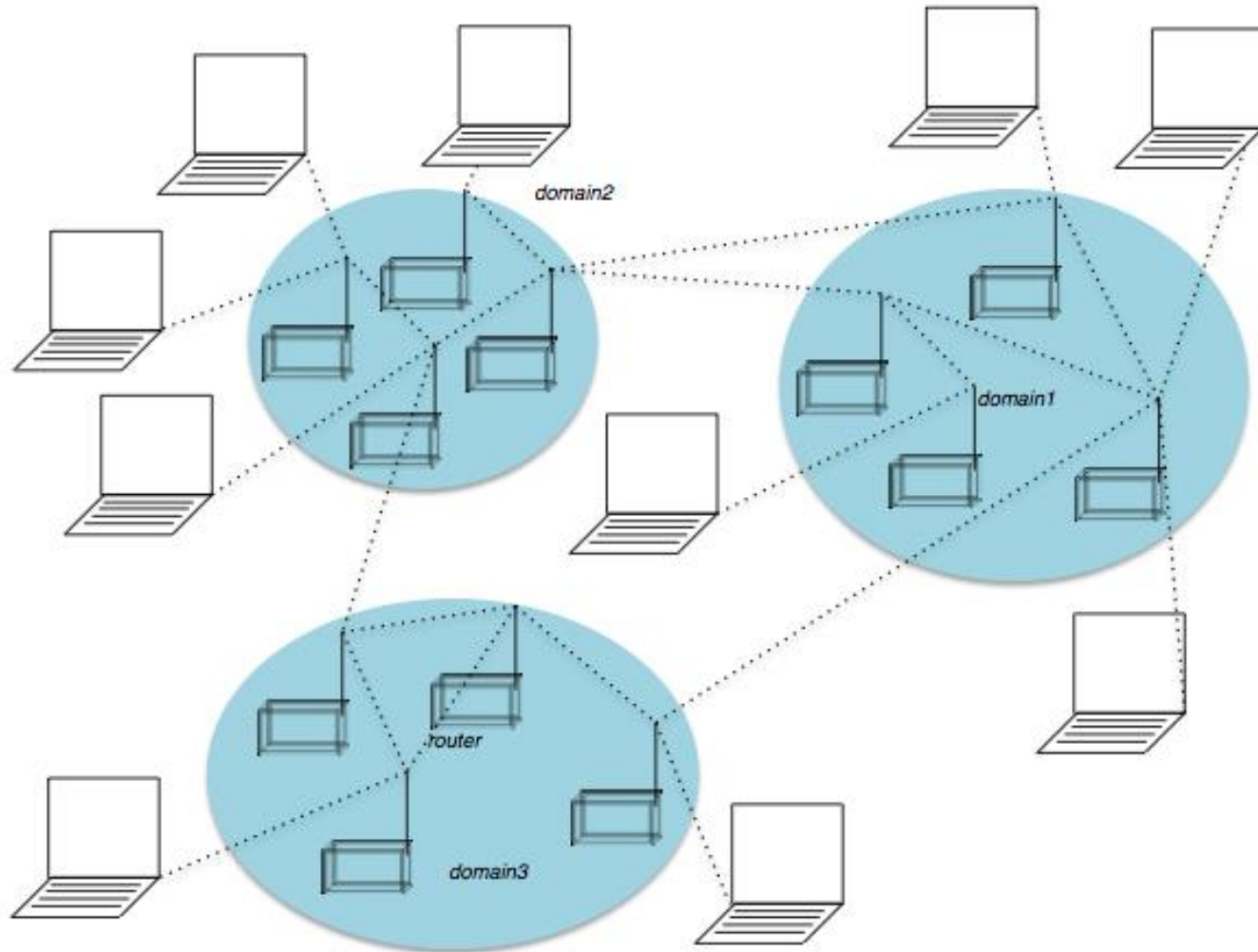Polarization of the network [Adamic-Glance, 2005]

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

domain2

domain1

router

domain3
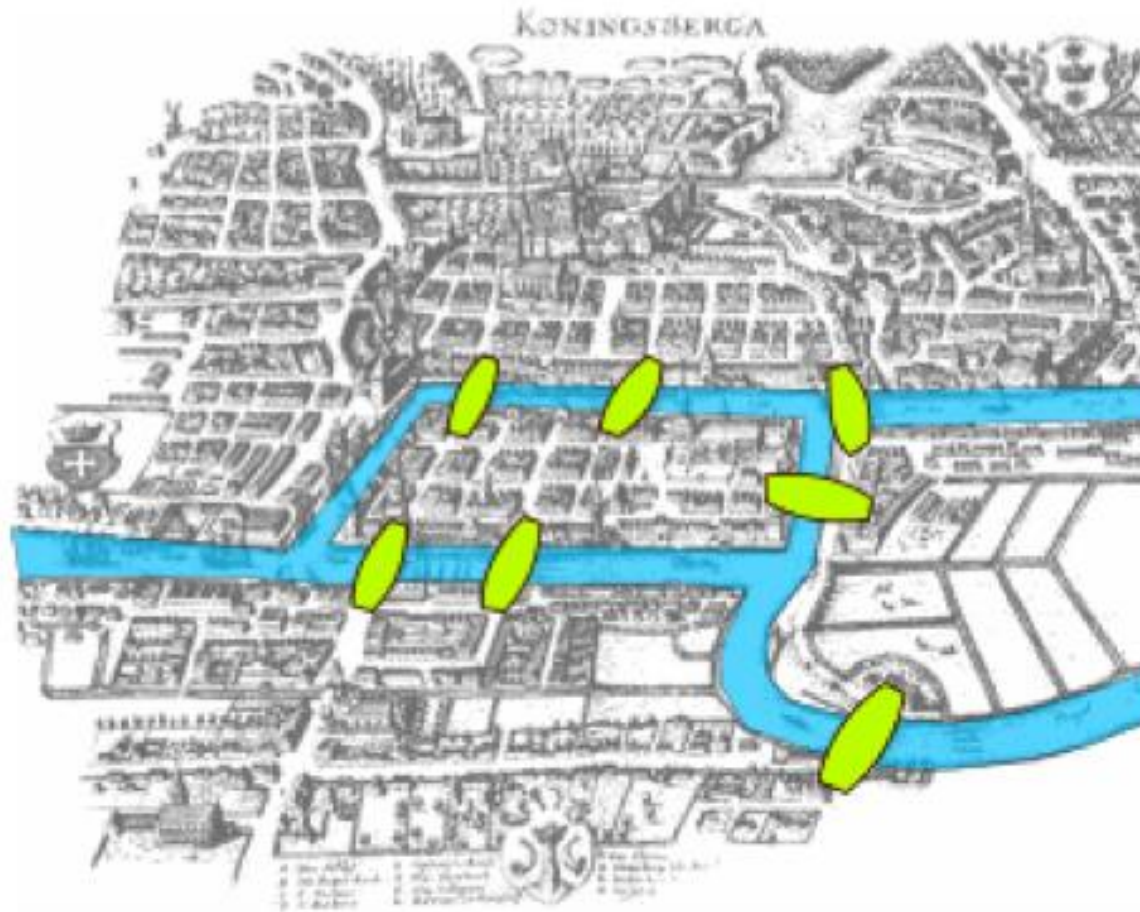
**Internet**

**Seven Bridges of Königsberg**

[Euler, 1735]

Return to the starting point by traveling each
link of the graph once and only once.

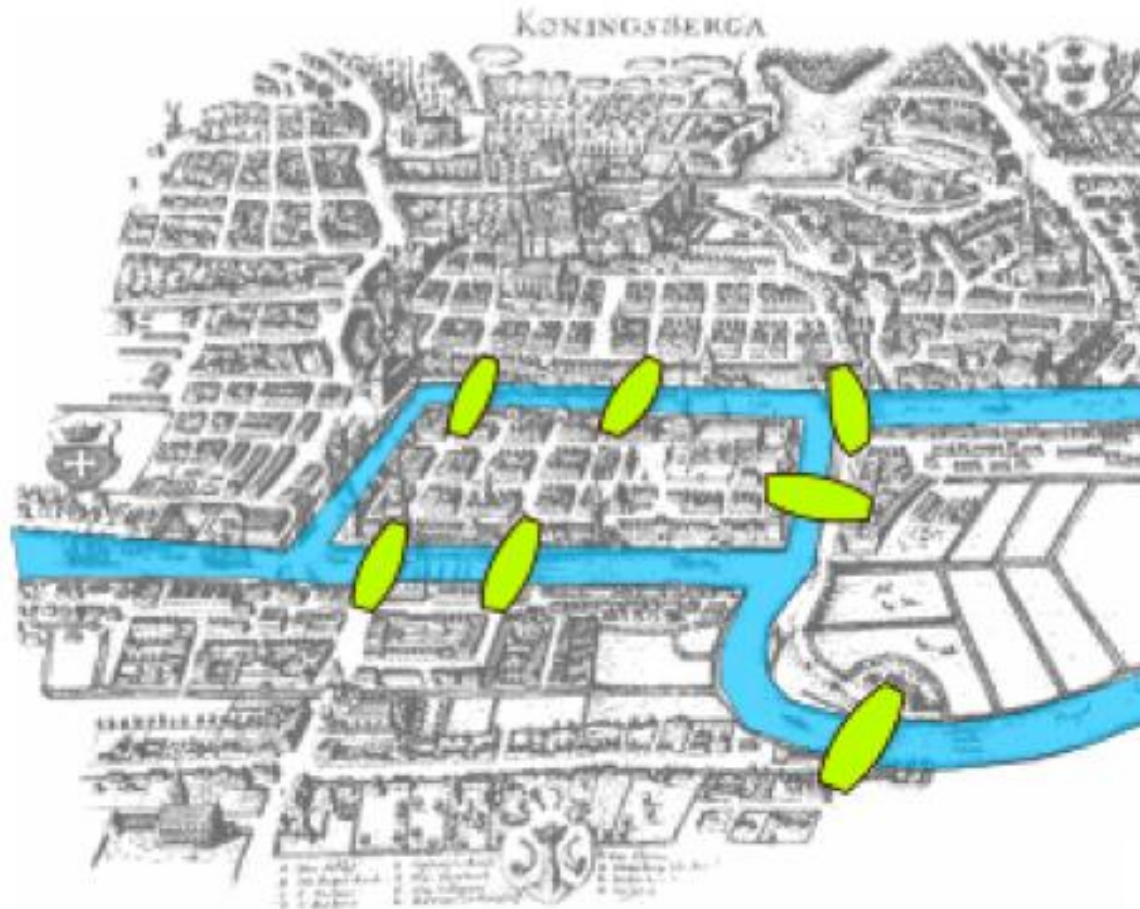**Seven Bridges of Königsberg**

[Euler, 1735]

Return to the starting point by traveling each
link of the graph once and only once.

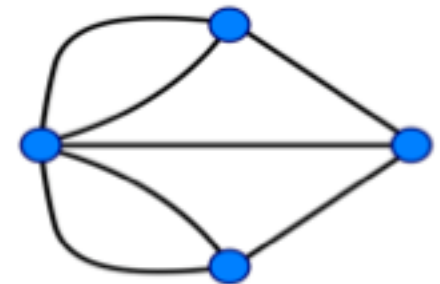- **Web as a directed graph:**
  - **Nodes: Webpages**
  - **Edges: Hyperlinks**

- **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.

**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

- **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!**

- **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!

❑ Rank nodes for a particular query

- ❑ Top k matches for "Random Walks" from Citeseer
- ❑ Who are the most likely co-authors of "Manuel Blum".
- ❑ Top k book recommendations  for Purna from Amazon
- ❑ Top k websites matching "Sound of Music"
- ❑ Top k friend recommendations for Purna when she joins "Facebook"

# Lecture Outline

- ❑ Basic definitions
    - ❑ Random walks
    - ❑ Stationary distributions
- ❑ Properties
    - ❑ Perron frobenius theorem
- ❑ Applications
    - ❑ Pagerank
        - ❑Power iteration
        - ❑Convergencce
    - ❑ Personalized pagerank
    - ❑ Rank stability

- ❑ nxn Adjacency matrix A.
  - ❑ A(i,j) = weight on edge from i to j
  - ❑ If the graph is undirected A(i,j)=A(j,i), i.e. A is symmetric

- ❑ nxn Transition matrix P.
  - ❑ P is row stochastic
  - ❑ P(i,j) = probability of stepping on node j from node i
  - ❑          = $A(i,j)/\sum_i A(i,j)$

- ❑ nxn Laplacian Matrix L.
  - ❑ $L(i,j)=\sum_i A(i,j)-A(i,j)$
  - ❑ Symmetric positive semi-definite for undirected graphs
  - ❑ Singular

$$\begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}$$

**Adjacency matrix A**

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

**Transition matrix P**

❑ Random Walk

t=0

1

1

1/2

1/2

t=0

1

1

1/2

1/2

t=1

1

1

1/2

1/2

# What is a random walk



t=0

1

1/2

1

1/2

t=1

1

1/2

1

1/2

t=2

1

1/2

1

1/2

# What is a random walk

❑ $x_{t(i)}$ = probability that the surfer is at node i at time t

❑ $x_{t+1}(i) = \sum_j$(Probability of being at node j)*Pr(j->i)
$= \sum_j xt_{(j)}*P(j,i)$

❑ $x_{t+1} = x_t*P = x_{t-1}*P*P = x_{t-2}*P*P*P = \ldots = x_0\ P^t$

❑ What happens when the surfer keeps walking for a long time?

# Stationary Distribution

❑ When the surfer keeps walking for a long time

❑ When the distribution does not change anymore

    ❑ i.e. $x_{T+1} = x_T$

❑ For "well-behaved" graphs this does not depend on the start distribution!!

❑ What is a stationary distribution?
Intuitively and Mathematically

- ❑ The stationary distribution at a node is related to the amount of time a random walker spends visiting that node.

❑ The stationary distribution at a node is related to the amount of time a random walker spends visiting that node.

❑ Remember that we can write the probability distribution at a node as

❑ $x_{t+1} = x_t P$

# Stationary Distribution

❑ The stationary distribution at a node is related to the amount of time a random walker spends visiting that node.

❑ Remember that we can write the probability distribution at a node as

  ❑ $x_{t+1} = x_t P$

❑ For the stationary distribution $v_0$ we have

  ❑ $v_0 = v_0 P$

# Stationary Distribution

- The stationary distribution at a node is related to the amount of time a random walker spends visiting that node.

- Remember that we can write the probability distribution at a node as
    - $x_{t+1} = x_t P$

- For the stationary distribution v0 we have
    - $v_0 = v_0 P$

- Whoa! that's just the left eigenvector of the transition matrix !

❑ Back to PageRank

# Example of Page Rank Scores

- 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

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

- 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:

$r_y = r_y/2 + r_a/2$

$r_a = r_y/2 + r_m$

$r_m = r_a/2$

# Solve the Flow Equation

**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!**

- **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**

$$r = M \cdot r$$

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

- **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 $j$**



$M \cdot r = r$

- **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

      - *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 λ if:

$$Ax = \lambda x$$

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

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

$$r = M \cdot r$$

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

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

- **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: $\mathbf{r}^{(0)} = [1/N,....,1/N]^{\mathsf{T}}$
    - Iterate: $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$
    - Stop when $|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}|_1 < \varepsilon$
        - $|\mathbf{x}|_1 = \sum_{1 \leq i \leq N} |x_i|$ is the L$_1$ norm

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

$d_i$ .... out-degree of node i

- **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 |

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

- **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:**



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

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

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

Iteration 0, 1, 2, …

❑ Why should this work?

❑ Write $x_0$ as a linear combination of the left eigenvectors $\{v_0, v_1, \ldots, v_{n-1}\}$ of P

❑ Remember that $v_0$ is the stationary distribution.

❑ $x_0 = c_0 v_0 + c_1 v_1 + c_2 v_2 + \ldots + c_{n-1} v_{n-1}$

$c_0 = 1$ .

$$X_0$$

$$v_0 \quad v_1 \quad v_2 \quad \dots\dots \quad v_{n-1}$$

$$1 \qquad c_1 \qquad\quad c_2 \qquad\qquad c_{n-1}$$

$$x_1 = x_0 \tilde{P}$$

$v_0 \quad v_1 \quad v_2 \quad \ldots\ldots \quad v_{n-1}$

$\sigma_0 \quad \sigma_1 c_1 \quad \sigma_2 c_2 \quad \sigma_{n-1} c_{n-1}$

$$x_2 = x_1 \tilde{P} = x_0 \tilde{P}^2$$

$v_0 \quad v_1 \quad v_2 \quad \ldots\ldots. \quad v_{n-1}$

$\sigma_0^2 \quad \sigma_1^2 c_1 \quad \sigma_2^2 c_2 \quad \sigma_{n-1}^2 c_{n-1}$

$$x_t = x_0 \tilde{P}^t$$

$$v_0 \quad v_1 \quad v_2 \quad \ldots\ldots \quad v_{n-1}$$

$$\sigma_0^\dagger \quad \sigma_1^\dagger c_1 \quad \sigma_2^\dagger c_2 \quad \sigma_{n-1}^\dagger$$

$$c_{n-1}$$

$$x_t = x_0 \tilde{P}^t$$

$$\sigma_0 = 1 > \sigma_1 \geq \ldots \geq \sigma_n$$

$v_0 \quad v_1 \quad v_2 \quad \ldots \ldots \quad v_{n-1}$

$1 \quad \sigma_1^t c_1 \quad \sigma_2^t c_2 \quad \sigma_{n-1}^t$

$c_{n-1}$

$$\mathbf{X}_\infty$$

$$\sigma_0 = 1 \; > \; \sigma_1 \geq \ldots \geq \sigma_n$$

$v_0 \quad v_1 \quad v_2 \quad \ldots\ldots. \quad v_{n-1}$

$1 \quad\;\; 0 \quad\quad 0 \quad\quad\quad\;\; 0$

# Convergence Issues

- Formally $||x_0 P^t - v_0|| \leq |\lambda|^t$
  - $\lambda$ is the eigenvalue with second largest magnitude

- The smaller the second largest eigenvalue (in magnitude), the faster the mixing.

- For $\lambda < 1$ there exists an unique stationary distribution, namely the first left eigenvector of the transition matrix.

- **Power iteration:**

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

  - $r^{(1)} = M \cdot r^{(0)}$
  - $r^{(2)} = M \cdot r^{(1)} = M(Mr^{(1)}) = M^2 \cdot r^{(0)}$
  - $r^{(3)} = M \cdot r^{(2)} = M(M^2 r^{(0)}) = M^3 \cdot r^{(0)}$

- **Claim:**

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

- ■ **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

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

- **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**

## 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)
  - Eventually spider traps absorb all importance

- **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 |

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

- **Example:**

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} =$$

| 1/3 | 2/6 | 3/12 | 5/24 | | 0 |
|---|---|---|---|---|---|
| 1/3 | 1/6 | 2/12 | 3/24 | … | 0 |
| 1/3 | 3/6 | 7/12 | 16/24 | | 1 |

Iteration 0, 1, 2, …

- **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**

- **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:**

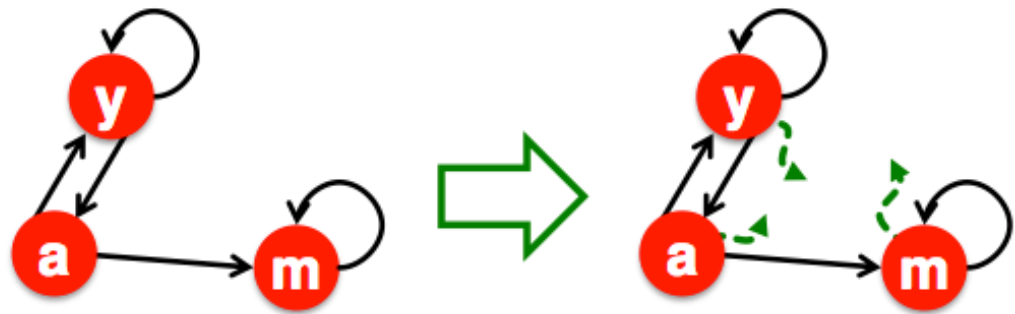$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & & 0 \end{matrix}$$

Iteration 0, 1, 2, …

- **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 |

$$r^{(t+1)} = Mr^{(t)}$$

## Markov chains

- Set of states $X$
- Transition matrix $P$ where $P_{ij} = P(X_t = i \mid X_{t-1} = j)$
- $\pi$ specifying the stationary probability of being at each state $x \in X$
- Goal is to find $\pi$ such that $\pi = P\,\pi$

- **Theory of Markov chains**

- **Fact:** For any start vector, the power method applied to a Markov transition matrix $P$ will converge to a unique positive stationary vector as long as $P$ is stochastic, irreducible and aperiodic.

- **Stochastic:** Every column sums to 1
- **A possible solution:** Add green links

$$A = M + a^T (\frac{1}{n} e)$$

- $a_i \ldots = 1$ if node $i$ has out deg 0, $=0$ else
- $e \ldots$ vector of all 1s

|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 1/3 |
| a | ½ | 0 | 1/3 |
| m | 0 | ½ | 1/3 |

$r_y = r_y/2 + r_a/2 + r_m/3$
$r_a = r_y/2 + r_m/3$
$r_m = r_a/2 + r_m/3$

- A chain is **periodic** if there exists $k > 1$ such that the interval between two visits to some state $s$ is always a multiple of $k$.
- **A possible solution:** Add green links

- From any state, there is a non-zero probability of going from any one state to any another
- **A possible solution:** Add green links

- **Google's solution that does it all:**
  - Makes *M* **stochastic, aperiodic, irreducible**
- 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.

- **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 $A$:**

$$A = \beta M + (1 - \beta) \frac{1}{n} e \cdot e^T$$

$e$…vector of all 1s

- **$A$ is stochastic, aperiodic and irreducible, so**

$$r^{(t+1)} = A \cdot r^{(t)}$$

- **What is $\beta$ ?**
  - In practice $\beta = 0.8, 0.9$ (make $5$ steps and jump)

$$
\begin{array}{cc}
\mathbf{M} & 1/n \cdot \mathbf{1} \cdot \mathbf{1}^T \\
0.8 \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} & +\ 0.2 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}
\end{array}
$$

$$
\begin{array}{c|ccc}
y & 7/15 & 7/15 & 1/15 \\
a & 7/15 & 1/15 & 1/15 \\
m & 1/15 & 7/15 & 13/15
\end{array}
$$

$$\mathbf{A}$$

$$
\begin{array}{ccccccc}
y & & 1/3 & 0.33 & 0.24 & 0.26 & & 7/33 \\
a & = & 1/3 & 0.20 & 0.20 & 0.18 & \cdots & 5/33 \\
m & & 1/3 & 0.46 & 0.52 & 0.56 & & 21/33
\end{array}
$$

- **Key step is matrix-vector multiplication**
  - $r^{new} = A \cdot r^{old}$
- Easy if we have enough main memory to hold **A**, $\mathbf{r}^{old}$, $\mathbf{r}^{new}$
- **Say N = 1 billion pages**
  - We need 4 bytes for each entry (say)
  - 2 billion entries for vectors, approx 8GB
  - Matrix **A** has $N^2$ entries
    - $10^{18}$ is a large number!

$$A = \beta \cdot M + (1-\beta)\,[1/N]_{NxN}$$

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

$$= \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix}$$

- Suppose there are **N** pages
- Consider page **j**, with **d**$_j$ out-links
- We have $M_{ij} = 1/|d_j|$ when $j \rightarrow i$
  and $M_{ij} = 0$ otherwise
- **The random teleport is equivalent to:**
  - Adding a **teleport link** from **j** to every other page and setting transition probability to **(1-β)/N**
  - Reducing the probability of following each out-link from $1/|d_j|$ to $\beta/|d_j|$
  - **Equivalent:** Tax each page a fraction **(1-β)** of its score and redistribute evenly

- $r = A \cdot r,$ where $A_{ij} = \beta \, M_{ij} + \frac{1-\beta}{N}$

- $r_i = \sum_{j=1}^{N} A_{ij} \cdot r_j$

- $r_i = \sum_{j=1}^{N} \left[ \beta \, M_{ij} + \frac{1-\beta}{N} \right] \cdot r_j$

  $= \sum_{j=1}^{N} \beta \, M_{ij} \cdot r_j + \frac{1-\beta}{N} \sum_{j=1}^{N} r_j$

  $= \sum_{j=1}^{N} \beta \, M_{ij} \cdot r_j + \frac{1-\beta}{N}$    since $\sum r_j = 1$

- So we get: $r = \beta \, M \cdot r + \left[ \frac{1-\beta}{N} \right]_N$

- We just rearranged the **PageRank equation**

$$r = \beta M \cdot r + \left[\frac{1 - \beta}{N}\right]_N$$

  - where $[(1-\beta)/N]_N$ is a vector with all $N$ entries $(1-\beta)/N$

- **$M$** is a **sparse matrix!** (with no dead-ends)
  - 10 links per node, approx 10N entries
- So in each iteration, we need to:
  - Compute $r^{new} = \beta M \cdot r^{old}$
  - Add a constant value $(1-\beta)/N$ to each entry in $r^{new}$
    - **Note if M contains dead-ends then $\sum_i r_i^{new} < 1$ and we also have to renormalize $r^{new}$ so that it sums to 1**

# PageRank

- **Input:** Graph $G$ and parameter $\beta$
  - Directed graph $G$ with **spider traps** and **dead ends**
  - Parameter $\beta$
- **Output: PageRank vector $r$**
  - **Set:** $r_j^{(0)} = \frac{1}{N}, \quad t = 1$
  - **do:**
    - $\forall j:\ r'^{(t)}_j = \sum_{i \to j} \beta \frac{r_i^{(t-1)}}{d_i}$
      $r'^{(t)}_j = 0$ if in-deg. of $j$ is **0**
    - **Now re-insert the leaked PageRank:**
      $\forall j:\ r_j^{(t)} = r'^{(t)}_j + \frac{1-S}{N} \quad$ where: $S = \sum_j r'^{(t)}_j$
    - $t = t + 1$
  - **while** $\sum_j \left| r_j^{(t)} - r_j^{(t-1)} \right| > \varepsilon$

- **Measures generic popularity of a page**
  - Biased against topic-specific authorities
  - **Solution:** Topic-Specific PageRank
- **Uses a single measure of importance**
  - Other models e.g., hubs-and-authorities
  - **Solution:** Hubs-and-Authorities
- **Susceptible to Link spam**
  - Artificial link topographies created in order to boost page rank
  - **Solution:** TrustRank

□ Hubs and Authorities

- **HITS (Hypertext-Induced Topic Selection)**
  - **Is a measure of importance of pages or documents, similar to PageRank**
  - Proposed at around same time as PageRank ('98)
- **Goal**: Say we want to find good newspapers
  - Don't just find newspapers. Find "experts" – people who link in a coordinated way to good newspapers
- **Idea: Links as votes**
  - **Page is more important if it has more links**
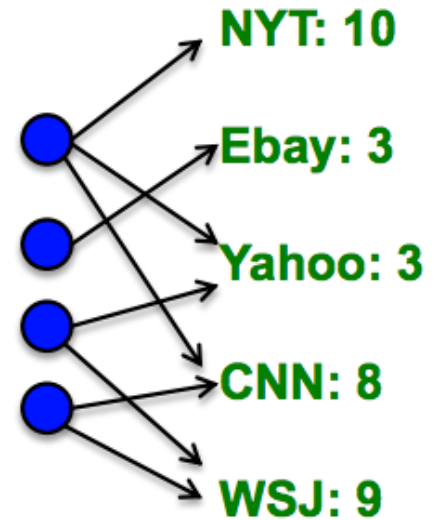    - In-coming links? Out-going links?

- **Hubs and Authorities**
  Each page has 2 scores:
  - **Quality as an expert (hub):**
    - Total sum of votes of authorities pointed to
  - **Quality as a content (authority):**
    - Total sum of votes coming from experts

NYT: 10

Ebay: 3

Yahoo: 3

CNN: 8

WSJ: 9

- **Principle of repeated improvement**

**Interesting pages fall into two classes:**

1. **Authorities** are pages containing useful information
   - Newspaper home pages
   - Course home pages
   - Home pages of auto manufacturers

2. **Hubs** are pages that link to authorities
   - List of newspapers
   - Course bulletin
   - List of US auto manufacturers

Each page starts with **hub** score 1. **Authorities** collect their votes

(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

SJ Merc News — 2 votes

Wall St. Journal — 2 votes

New York Times — 4 votes

Sum of **hub** scores of nodes pointing to NYT.

USA Today — 3 votes

Facebook — 1 vote

Each page starts with **hub** score 1. **Authorities** collect their votes

Yahoo! — 3 votes

Amazon — 3 votes

(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

Sum of authority scores of nodes that the node points to.

SJ Merc News — 2 votes

Wall St. Journal — 2 votes

New York Times — 4 votes

USA Today — 3 votes

Facebook — 1 vote

Yahoo! — 3 votes

Amazon — 3 votes

**Hubs** collect authority scores

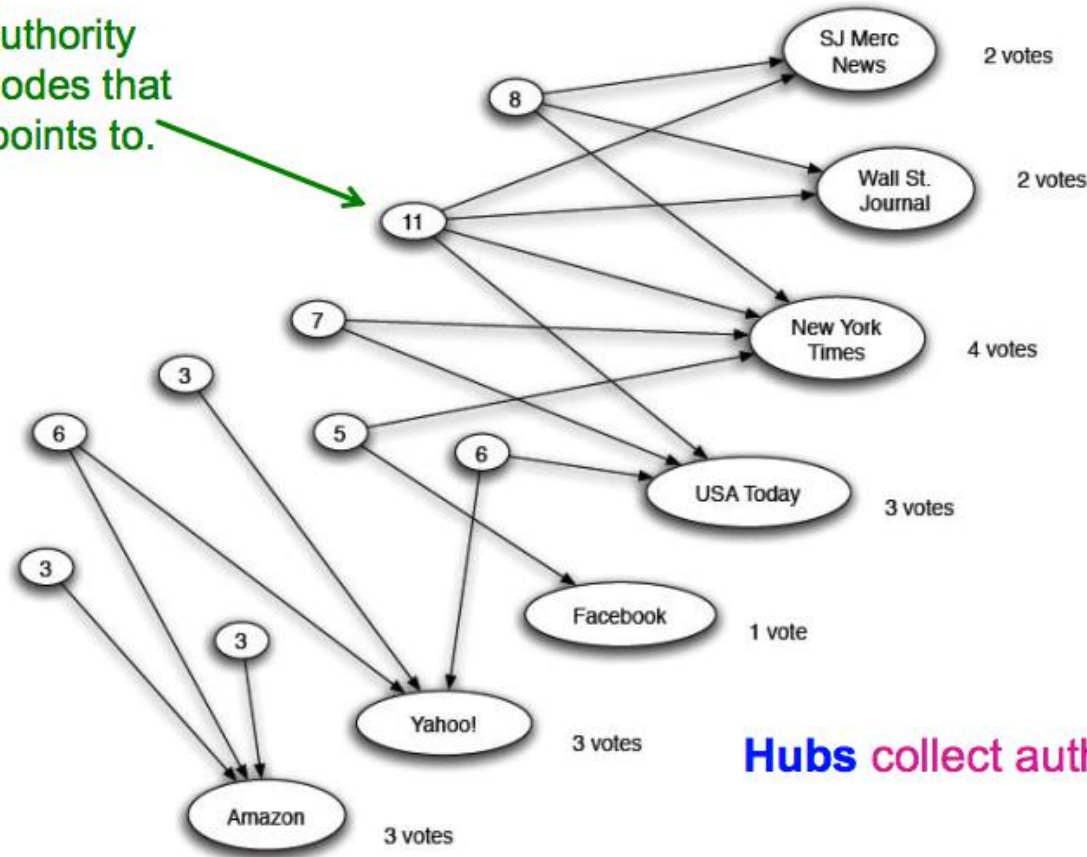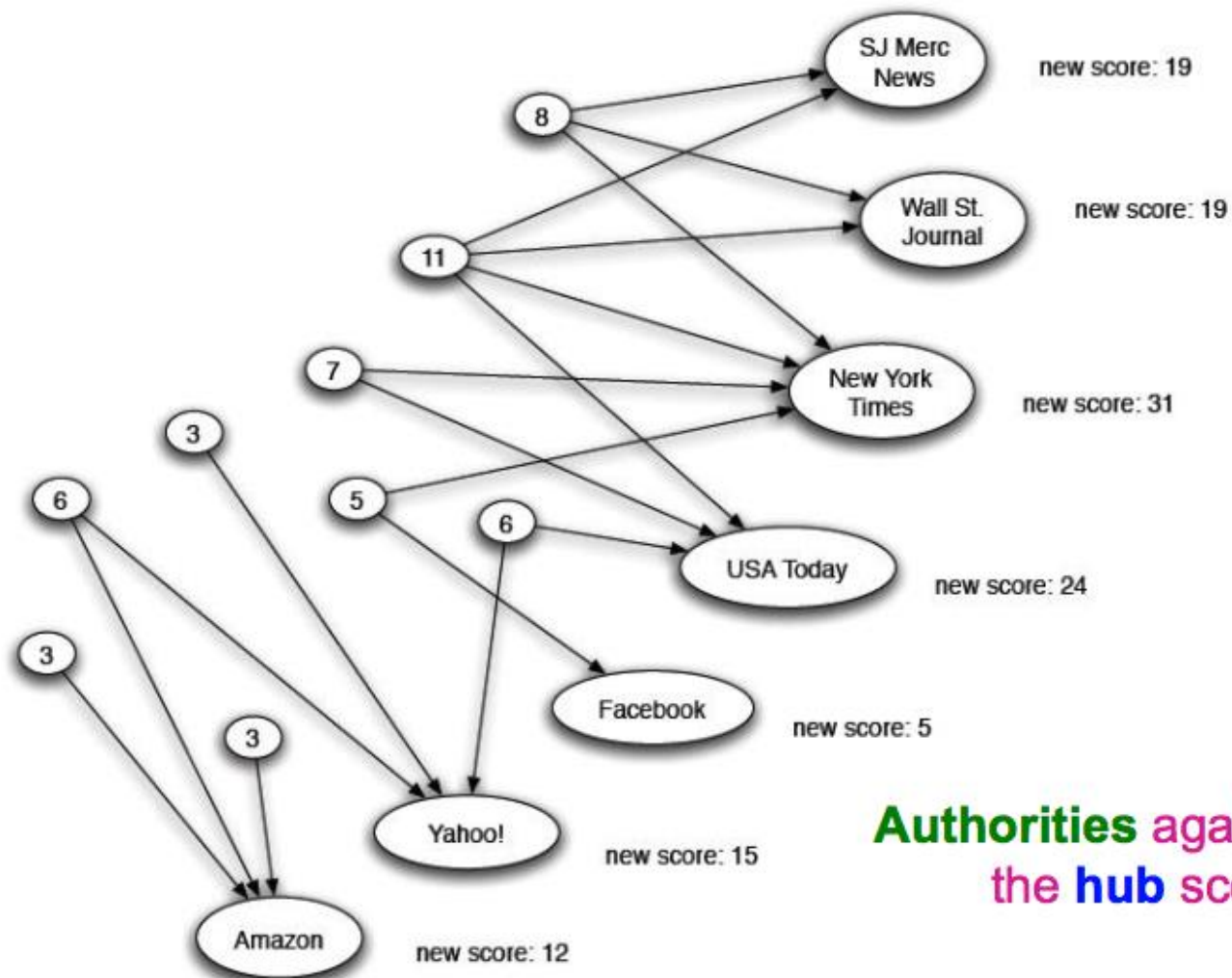(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

(Note this is idealized example. In reality graph is not bipartite
each page has both the hub and authority score)

- **A good hub links to many good authorities**

- **A good authority is linked from many good hubs**

- **Model using two scores for each node:**
  - **Hub** score and **Authority** score
  - Represented as vectors $h$ and $a$

- **Each page $i$ has 2 scores:**
  - Authority score: $a_i$
  - Hub score: $h_i$

**HITS algorithm:**

- Initialize: $a_i = 1/\sqrt{n}$, $h_i = 1/\sqrt{n}$
- Then keep iterating until convergence:
  - $\forall i$: **Authority**: $a_i = \sum_{j \to i} h_j$
  - $\forall i$: **Hub**: $h_i = \sum_{i \to j} a_j$
  - $\forall i$: **Normalize $a$, $h$ such that:**
    $\sum_i a_i^2 = 1$, $\sum_i h_i^2 = 1$

$a_i = \sum_{j \to i} h_j$

$h_i = \sum_{i \to j} a_j$

- **HITS converges to a single stable point**
- **Notation:**
    - Vector $a = (a_1 ..., a_n)$, $h = (h_1 ..., h_n)$
    - Adjacency matrix $A$ ($n$ x $n$): $A_{ij} = 1$ if $i \rightarrow j$
- **Then $h_i = \sum_{i \rightarrow j} a_j$**

    **can be rewritten as $h_i = \sum_j A_{ij} \cdot a_j$**

    **So: $h = A \cdot a$**
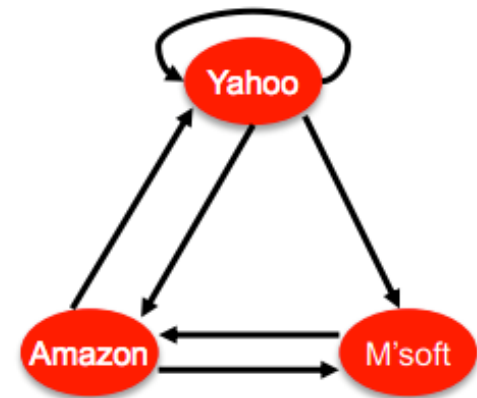- **Similarly, $a_i = \sum_{j \rightarrow i} h_j$**

    **can be rewritten as $a_i = \sum_j A_{ji} \cdot h_i = A^T \cdot h$**

$$A = \begin{vmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{vmatrix} \qquad A^T = \begin{vmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{vmatrix}$$



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| a(yahoo) | = | .58 | .80 | .80 | .79 | $\cdots$ | .788 |
| a(amazon) | = | .58 | .53 | .53 | .57 | $\cdots$ | .577 |
| a(m'soft) | = | .58 | .27 | .27 | .23 | $\cdots$ | .211 |
| | | | | | | | |
| h(yahoo) | = | .58 | .58 | .62 | .62 | $\cdots$ | .628 |
| h(amazon) | = | .58 | .58 | .49 | .49 | $\cdots$ | .459 |
| h(m'soft) | = | .58 | .58 | .62 | .62 | $\cdots$ | .628 |

- **HITS algorithm in vector notation:**

  - Set: $a_i = h_i = \dfrac{1}{\sqrt{n}}$

  **Repeat until convergence:**

  - $h = A \cdot a$

  - $a = A^T \cdot h$

  - Normalize $a$ and $h$

- **Then:** $a = A^T \cdot (\underbrace{\overbrace{A \cdot a}^{\text{new } h}}_{\text{new } a})$

- **Thus, in $2k$ steps:**

  $a = (A^T \cdot A)^k \cdot a$
  $h = (A \cdot A^T)^k \cdot h$

**Convergence criterion:**

$$\sum_i \left( h_i^{(t)} - h_i^{(t-1)} \right)^2 < \varepsilon$$

$$\sum_i \left( a_i^{(t)} - a_i^{(t-1)} \right)^2 < \varepsilon$$

$a$ is updated (in 2 steps):
$$a = A^T (A\,a) = (A^T A)\,a$$
$h$ is updated (in 2 steps):
$$h = A (A^T h) = (A\,A^T)\,h$$

**Repeated matrix powering**

- $\boldsymbol{h} = \lambda\, A\, \boldsymbol{a}$
- $\boldsymbol{a} = \mu\, A^T\, \boldsymbol{h}$
- $\boldsymbol{h} = \lambda\, \mu\, A\, A^T\, \boldsymbol{h}$
- $\boldsymbol{a} = \lambda\, \mu\, A^T\, A\, \boldsymbol{a}$

$$\lambda = 1/\textstyle\sum h_i$$
$$\mu = 1/\textstyle\sum a_i$$

- Under reasonable assumptions about **A**, HITS **converges to vectors $\boldsymbol{h}^*$ and $\boldsymbol{a}^*$**:
  - $\boldsymbol{h}^*$ is the **principal eigenvector** of matrix $\boldsymbol{A}\,\boldsymbol{A}^T$
  - $\boldsymbol{a}^*$ is the **principal eigenvector** of matrix $\boldsymbol{A}^T\boldsymbol{A}$