

# Web Search Engines

- SI650 / EECS 549  
Information Retrieval
- September 25, 2024

# Today's lecture plan

- Gain a better understanding of how web search works
- Learn how web search isn't just IR ranking—or is it?
- Rank using non-textual information

# Web Search is Hard

Google digital camere Search Advanced Search

Web Hide options

All results Images Videos

olympus digital camera digital video camera canon powershot s90 digital camera digital camera best buy canon digital camera

nikon digital camera digital camera compare sony digital camera

Did you mean: [digital camera](#)

Digital Camera Reviews and News: Digital Photography Review ... ✓ Mar 12, 2010 ... Digital Photography Review: All the latest digital camera reviews and digital imaging news. Lively discussion forums. Canon - Reviews - Nikon - Most popular cameras www.dpreview.com/ - Cached - Similar

Camera Labs: Digital Camera, DSLR, and Lens Reviews. Recommended ... ✓ Camera Labs: Digital camera, digital SLR and lens reviews, workshops, news. www.cameralabs.com/ - Cached - Similar

Unbiased Digital Camera Reviews and News | Digital Camera Resource ... ✓ The Digital Camera Resource Page has been providing unbiased digital camera reviews, news, discussion forums, buyers guides, and frequently asked questions ... www.dcresource.com/ Cached Similar

News results for [digital camere](#) ✓

Panasonic's Lumix DMC-ZS5 Digital Camera Features Manual Controls ... ✓ - 14 hours ago by PCWorld Videos | 0 Comments Recommends Panasonic's Lumix DMC-ZS5 is a 12.1 megapixel digital camera with an ultrawide lens that has excellent picture ... PC World - 6 related articles »

Canon U.S.A. Joins Creative Forces with Fashion Brand LeSportsac ✓ - MarketWatch (press release) - 29 related articles »

Image results for [digital camere](#) ✓ - Report images

Nikon Coolpix P90 Digital Camera ✓ 7 min - Apr 7, 2009 www.youtube.com

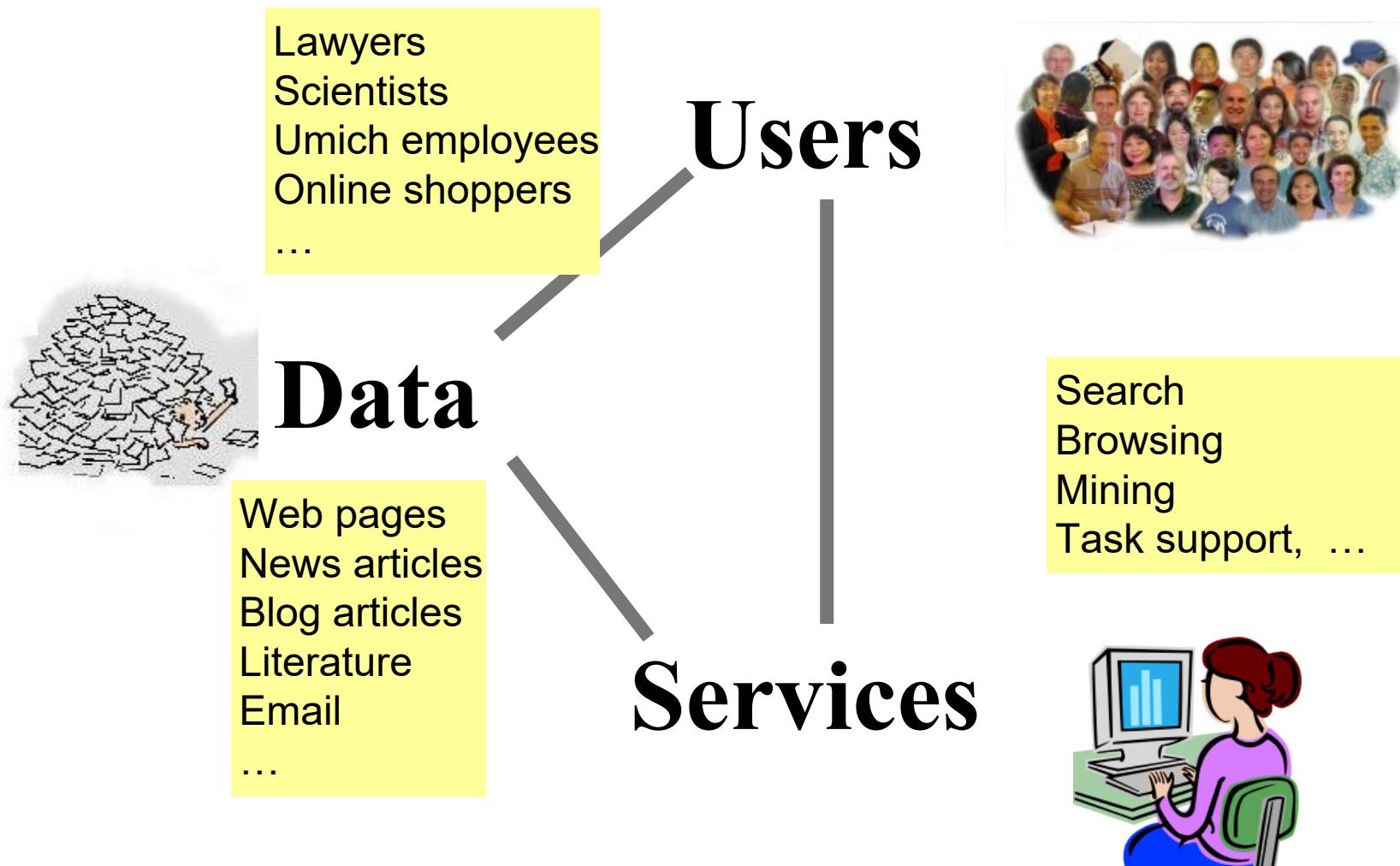
FUJIFILM FinePix S2000HD Digital Camera ✓ 1 min 46 sec - Oct 2, 2008 www.youtube.com

Sponsored Links

Digital Camere ✓ Compare Deals on Electronics & More Low Prices On Digital Camere shopping.yahoo.com See your ad here »

Results 1 - 10 of about 76,700,000 for [digital camere](#). (0.27 seconds)

# The Data-User-Service (DUS) Triangle



- Slide from ChengXiang Zhai

# Challenges in Web Search

- Content relevance isn't the only objective
  - Quality, diversity, novelty, ...
- Queries are vague, poorly formulated
- Relevance is a highly personalized notion
- The web is messy
  - 30% content on the web are near-duplicate
- Spam! Spam!! Spam!!!
- So, really quality! How do we get at quality?
  - ...without looking at text!

# Search as a Network Problem

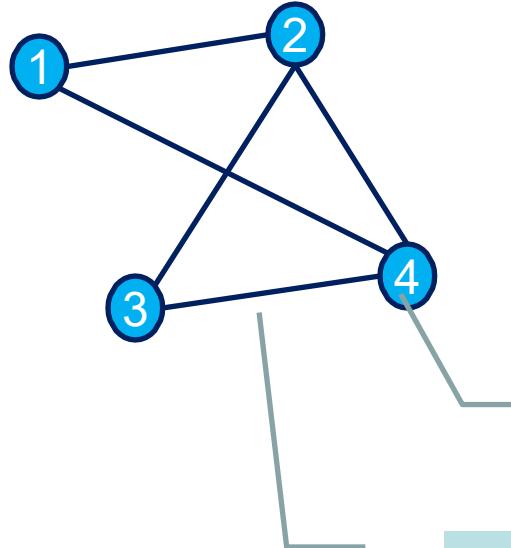
# The Web Network



Source - <http://www.research.att.com/~yifanhu/GIF/www.html>

# Representation of Networks

- Networks are collections of points joined by lines.



Mathematically (Topologically):  
“Network”  $\equiv$  “Graph”

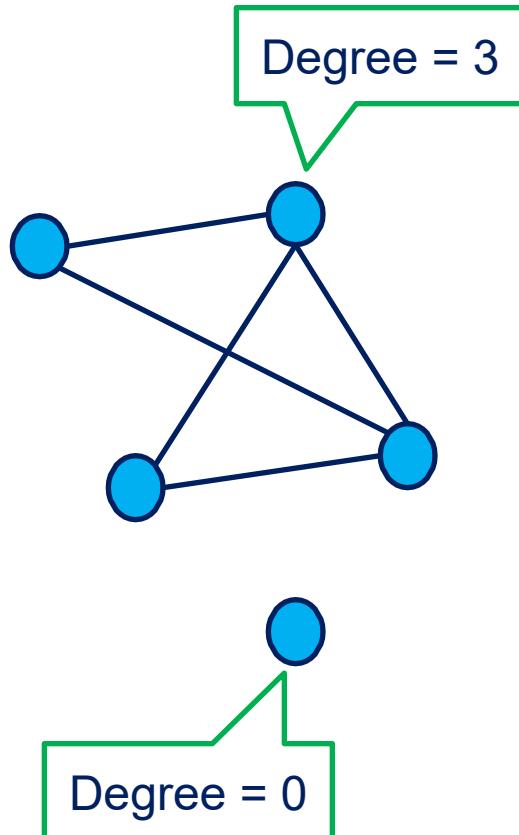
$$\begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

Vertex, node, site, actor, ...

Edge, arc, link, bond, tie, relation ...

But a real network is much more than a topological structure

# Network Measures: Characterize Networks



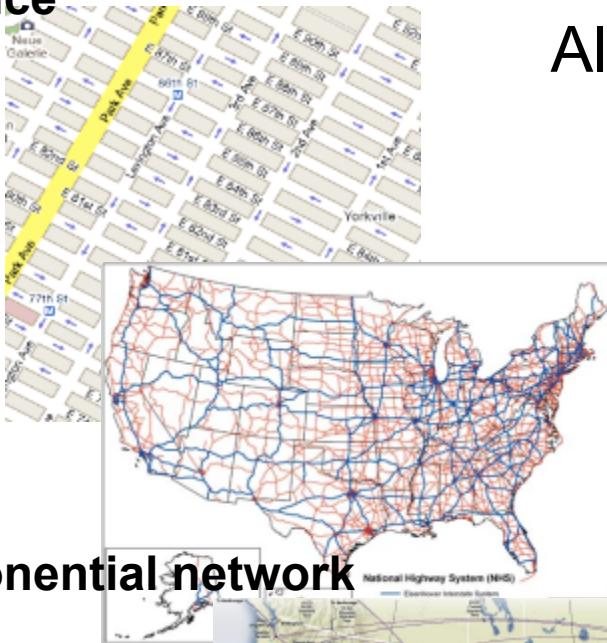
- Degree: how many friends do I have?
- Weights: how strong are the ties?
- Path: how far am I from another vertex?
- Connectivity: can I reach all other vertices?
- Diameter: how dense are they?
- Centrality (e.g., betweenness, closeness): Am I in the center of everyone?
- Density: are the vertices well connected?
- Clustering Coefficient: Are my neighbors well connected?

# The Web as a Network

- The web graph
  - degree distribution
  - clustering
  - motif profile
  - communities

# Properties: Degree Distribution

Lattice



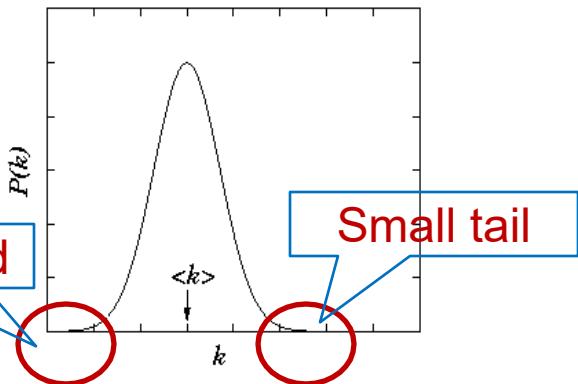
All vertices have the same degree

Exponential network

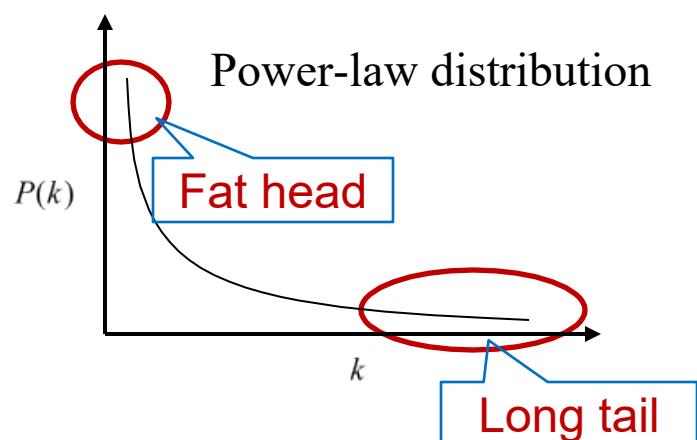


scale-free network

Poisson distribution



Power-law distribution



# The Web is Scale Free

- The web is a scale free network
- Has power law degree distribution. The probability of observing an item of size (degree) 'x' is given by

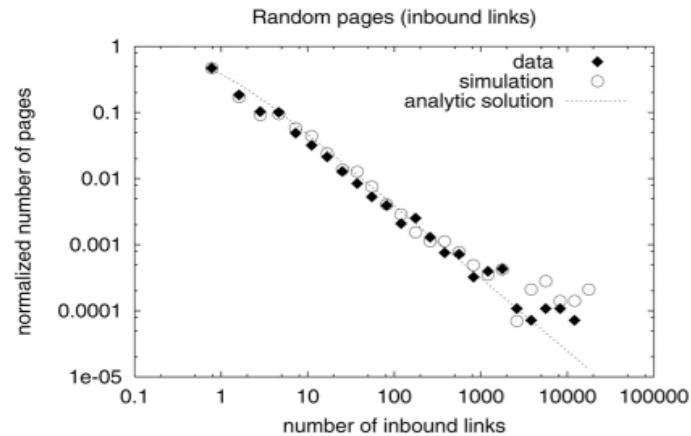
$$p(x) = Cx^{-\alpha}$$

normalization constant  
(probabilities over all  $x$  must sum to 1)

$\alpha$  : scaling exponent,  
or power law exponent

- Straight line on a log-log plot

$$\ln(p(x)) = c - \alpha \ln(x)$$



indegree,  $\alpha \sim 2.1$   
outdegree,  $\alpha \sim 2.4$

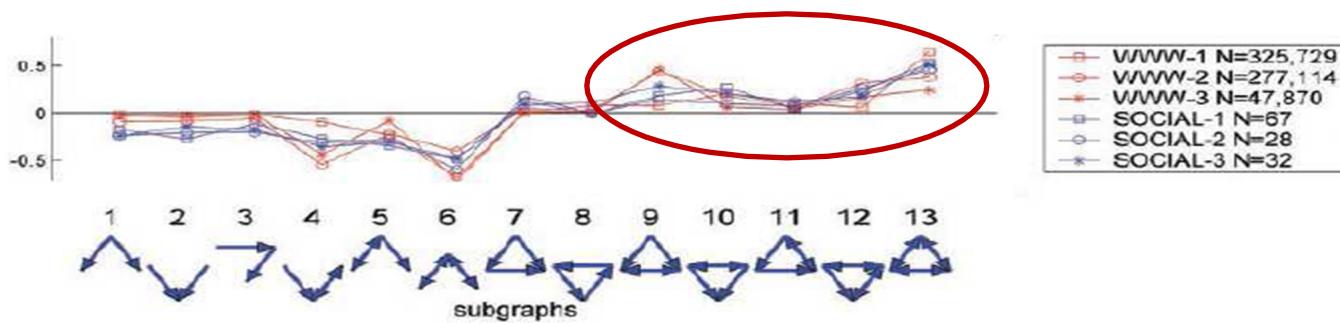
source: Pennock et al.: Winners don't take all:  
Characterizing the competition for links on  
the web  
PNAS April 16, 2002 vol. 99 no. 8 5207-5211

度为  $x$  的节点出现的概率随着  $x$  的增大而快速下降，但仍然存在一些极少数度数特别大的节点。

12

# Clustering & Motifs

- Web network has high clustering compared to random networks
- Pages linked to from a given domain tend to link to each other as well
- Clustering coefficient  $\sim 0.11$  (at the site level)



Source: Milo et al., "Superfamilies of evolved and designed networks", *Science* 303 (5663), p. 1538-1542, 2004.

## 1. Clustering (聚类/聚集系数)

- 概念：在网络中，某个节点的邻居之间是否也彼此相连。
  - 举例：如果网页 A 链接到网页 B 和 C，那么 B 和 C 之间是否也互相链接？
- Web 的特点：
  - 相比随机网络，Web 的聚类系数更高。
  - 例如：同一个域名下的网页，往往会互相链接（如某大学网站的不同院系页面）。
- 数据：在站点级别，Web 的聚类系数大约为 0.11 (Milo et al., Science 2004)

 SI650-Week-05-WebSearch

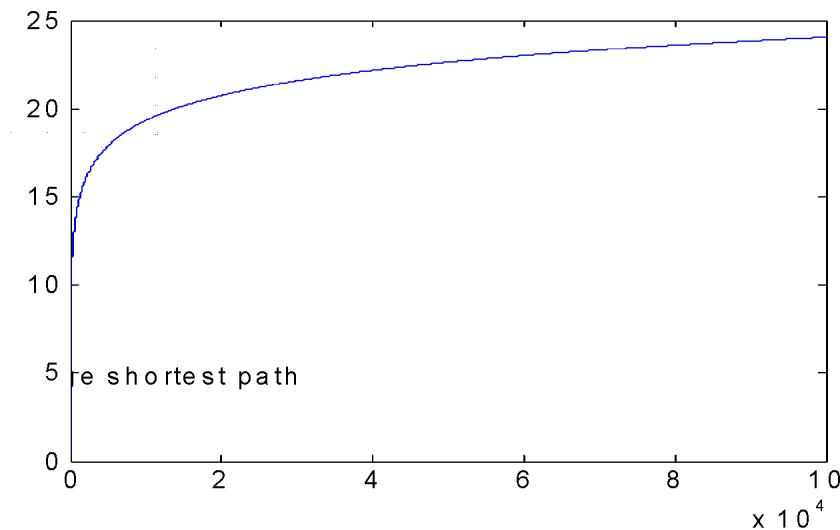
这说明：Web 网络里存在大量“小圈子”，局部紧密连接。

## 2. Motifs (结构基元/模式)

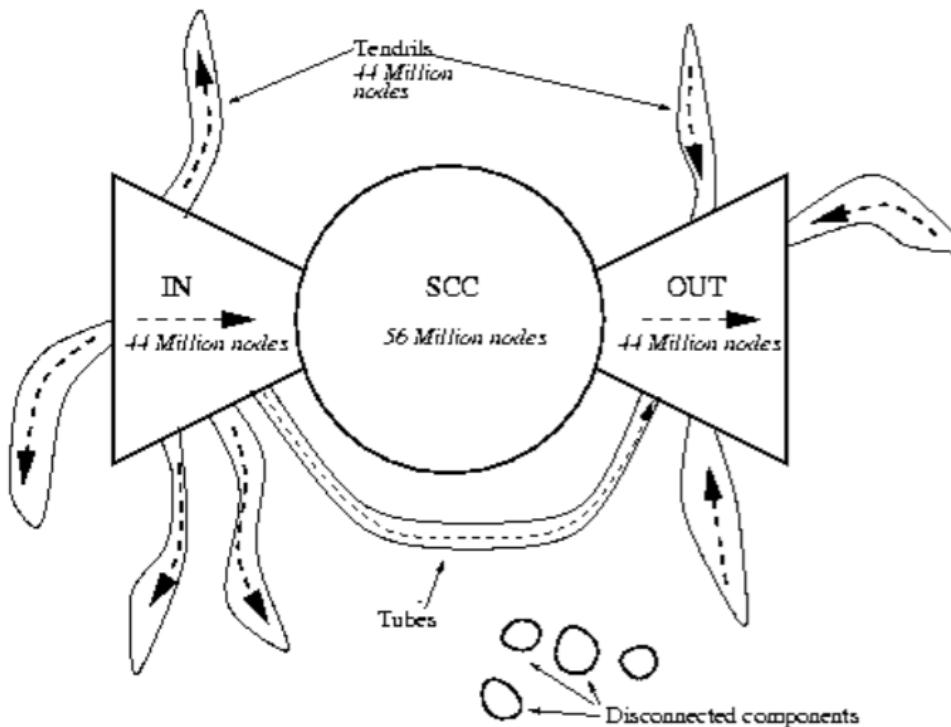
- 定义：在网络中反复出现的小规模连接模式。
  - 例如：三角形结构 ( $A \rightarrow B, B \rightarrow C, C \rightarrow A$ )，双向结构 ( $A \leftrightarrow B$ )，星型结构（一个节点指向多个节点）。
- 意义：Motifs 被认为是网络的“功能模块”，不同类型的网络（如社交、Web、基因调控网络）都有其典型 motif 分布。
- 在 Web 上：
  - 高聚类意味着许多 三角形结构（三个页面相互链接）。
  - 同一类网站（如新闻站、学术站点）可能有特定的 motif 特征。

# Shortest Paths

- On average web pages are within a small number of hops from each other
- $\langle d \rangle = 0.35 + 2.06 \log(N)$
- prediction:  $\langle d \rangle = 17.5$  for 200 million nodes
- actual:  $\langle d \rangle = 16$  for reachable pairs



# Components: the Bowtie Model



Outdated: from 1999 AltaVista data

## 2. Bowtie 模型的主要组成部分

“蝴蝶结”之所以得名，是因为整个网络结构看起来像一个蝴蝶结，中间有核心，两边有入口和出口：

### 1. SCC (Strongly Connected Component, 强连通核心)

- 位于中心。
- 从其中任何一个网页出发，都可以沿超链接路径返回到另一个网页。
- 这是 Web 的“核心”区域，大型网站和高度互链的页面通常在这里。

### 2. IN 组件

- 能进入 SCC，但 无法从 SCC 返回。
- 就像“入口”，比如一些只指向核心的网站，但核心不会再链接回它们。

### 3. OUT 组件

- 可以从 SCC 到达，但 无法返回到 SCC。
- 就像“出口”，例如一些终端网站（如个人主页、文件下载页）。

### 4. Tendrils (触须)

- 连接在 IN 和 OUT 的边缘，但既不在 SCC，也不属于纯 IN 或 OUT。
- 它们可能能连到 IN/OUT，但和核心没有直接联系。

### 5. Disconnected Components (孤立部分)

- 完全与主结构不连通的网页或小型网络。

# That was the Web Graph Overall

Web Images Videos Maps News Shopping Gmail more ▾

**Google** network  Search Advanced Search

Web [Show options...](#) Results 1 - 10 of about 853,000,000 for **network [definition]**. (0.14 seconds)

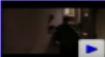
**Network (1976)**  
Directed by Sidney Lumet. With Faye Dunaway, William Holden, Peter Finch. A TV **network** cynically exploits a deranged ex-TV anchor's ravings and revelations ...  
[Full cast and crew](#) - [Memorable quotes](#) - [Awards](#) - [Plot Summary](#)  
[www.imdb.com/title/tt0074958/](http://www.imdb.com/title/tt0074958/) - Cached - Similar

**Network - Wikipedia, the free encyclopedia**  
Look up **network** or networking in Wiktionary, the free dictionary. ... **Network** (mathematics), a type of digraph in graph theory; **Network** theory, ...  
[en.wikipedia.org/wiki/Network](http://en.wikipedia.org/wiki/Network) - Cached - Similar

**Computer network - Wikipedia, the free encyclopedia**  
A computer **network** is a group of computers that are connected to each other for the purpose of communication. **Networks** may be classified according to a wide ...  
[en.wikipedia.org/wiki/Computer\\_network](http://en.wikipedia.org/wiki/Computer_network) - Cached - Similar

[Show more results from en.wikipedia.org](#)

**Video results for network**

 **Network: I'm Mad as Hell**  
4 min 12 sec - Jul 22, 2006  
[www.youtube.com](http://www.youtube.com)

 **Network**  
3 min 59 sec - Aug 30, 2006  
[www.youtube.com](http://www.youtube.com)

**Cisco Systems, Inc.**  
Borderless Networks. Bringing interactions closer to your customers. Learn More - Cisco Nexus 4000 Series Blade Switches ... [Show stock quote for CSCO](#)  
[www.cisco.com/](http://www.cisco.com/) - Cached - Similar

**NETWORK - A National Catholic Social Justice Lobby**  
Catholic organization involved in education, lobbying, and organizing around economic and social justice issues in the United States.  
[www.networklobby.org/](http://www.networklobby.org/) - Cached - Similar

**Domain Names, Web Hosting and Online Marketing Services | Network ...**  
Find domain names, web hosting and online marketing for your website -- all in one place. **Network** Solutions helps businesses get online and grow online with ...  
[www.networksolutions.com/](http://www.networksolutions.com/) - Cached - Similar

**Cartoon Network | Free Games and Online Video from Ben 10, Star ...**

Sponsored Links

**Network Monitoring**  
See what is really happening with DVR-like **network** recorder  
[www.SoleraNetworks.com](http://www.SoleraNetworks.com)

**Network Training Courses**  
Get the same IT training used by 1000s of Schools & Companies on DVD  
[www.edulearn.com](http://www.edulearn.com)

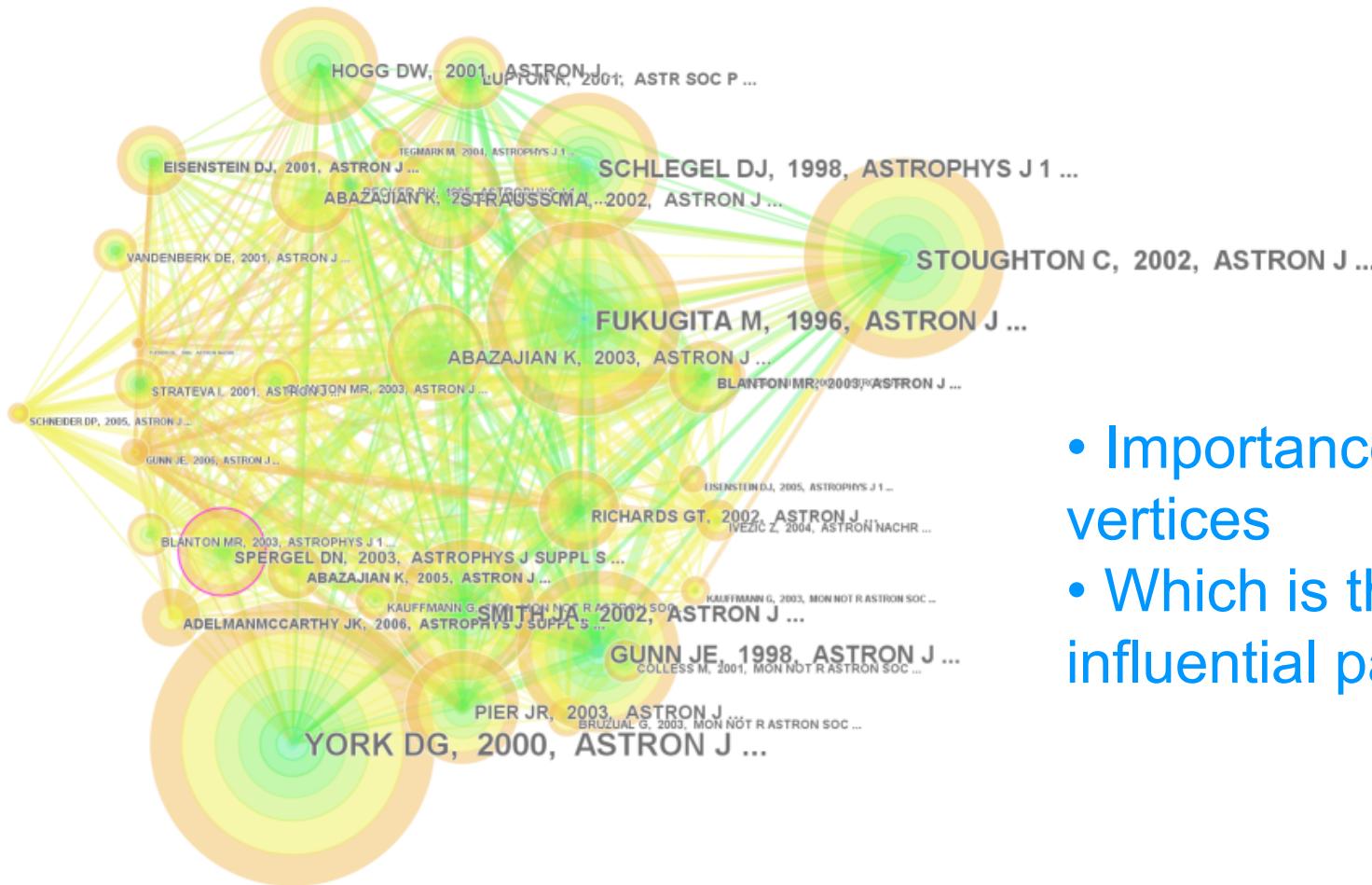
**Network**  
Find Answers & Useful Links On **Network**  
[Network.blurtit.com](http://Network.blurtit.com)  
Michigan

**Canon S630 Network Ink**  
80% off Ink & Laser Cartridges for the Canon S630 **Network** Printer!  
[www.ink-bates.com](http://www.ink-bates.com)

[See your ad here »](#)

**Problem:**  
**How do we know which individual pages are more important?**

# Ranking in Networks



- Importance of vertices
- Which is the most influential paper?

Co-citation network of Sloan Digital Sky Survey  
- <http://nevac.ischool.drexel.edu/~james/infovis09/FP-tree-visual.html>

# First Try: Network Centrality

- A good ranking measure:
  - Making sense;
  - Robust to spamming;
  - Easy to compute;
- Can we simply go with degree centrality (number of connections)?
- Who is more important?
  - A page that links to many other pages
  - Or a page that is linked by many other pages
- Not stable: one can easily spam the ranking

常见的中心性指标：

- **Degree Centrality (度中心性)**
  - 入度：多少页面指向它？（被多少人推荐 → 像“权威”）
  - 出度：它指向多少页面？（链接很多 → 像“导航”）
- **Betweenness Centrality (介数中心性)**
  - 一个页面在多少最短路径上出现？
  - 类似“交通枢纽”，影响信息流通。
- **Closeness Centrality (接近度中心性)**
  - 一个页面到所有其他页面的平均距离有多近？
  - 越近说明位置越核心。

# Betweenness/Closeness Centrality

- Betweenness: how many pairs of vertices have to go through you on their shortest path?

$$C_B(i) = \sum_{j, k \neq i, j < k} g_{jk}(i) / g_{jk}$$

Sum over all pairs of  
nodes  
excluding  $i$  itself

Where:

- $g_{jk}$  = the number of shortest paths connecting  $j$  and  $k$
- $g_{jk}(i)$  = the number of SPs that vertex  $i$  is on.

- Closeness: average distance from everyone else to you

$$C_C(i) = \left[ \sum_{j=1, j \neq i}^N d(i, j) \right]^{-1}$$

All vertices but  $i$

The smaller the length, the  
larger the closeness

# Betweenness/Closeness Centrality

- Can we use them to rank webpages?
- Pros:
  - Makes sense (can be explained by how information propagates)
  - More robust than in-degree (adding many spam in-links won't change much on the average distance/shortest path of all nodes)
- But no one really tried this on ranking web pages
  - Hard to compute: need to compute all pair shortest paths, which takes  $O(n^3)$ .

# Ranking by Network Centrality

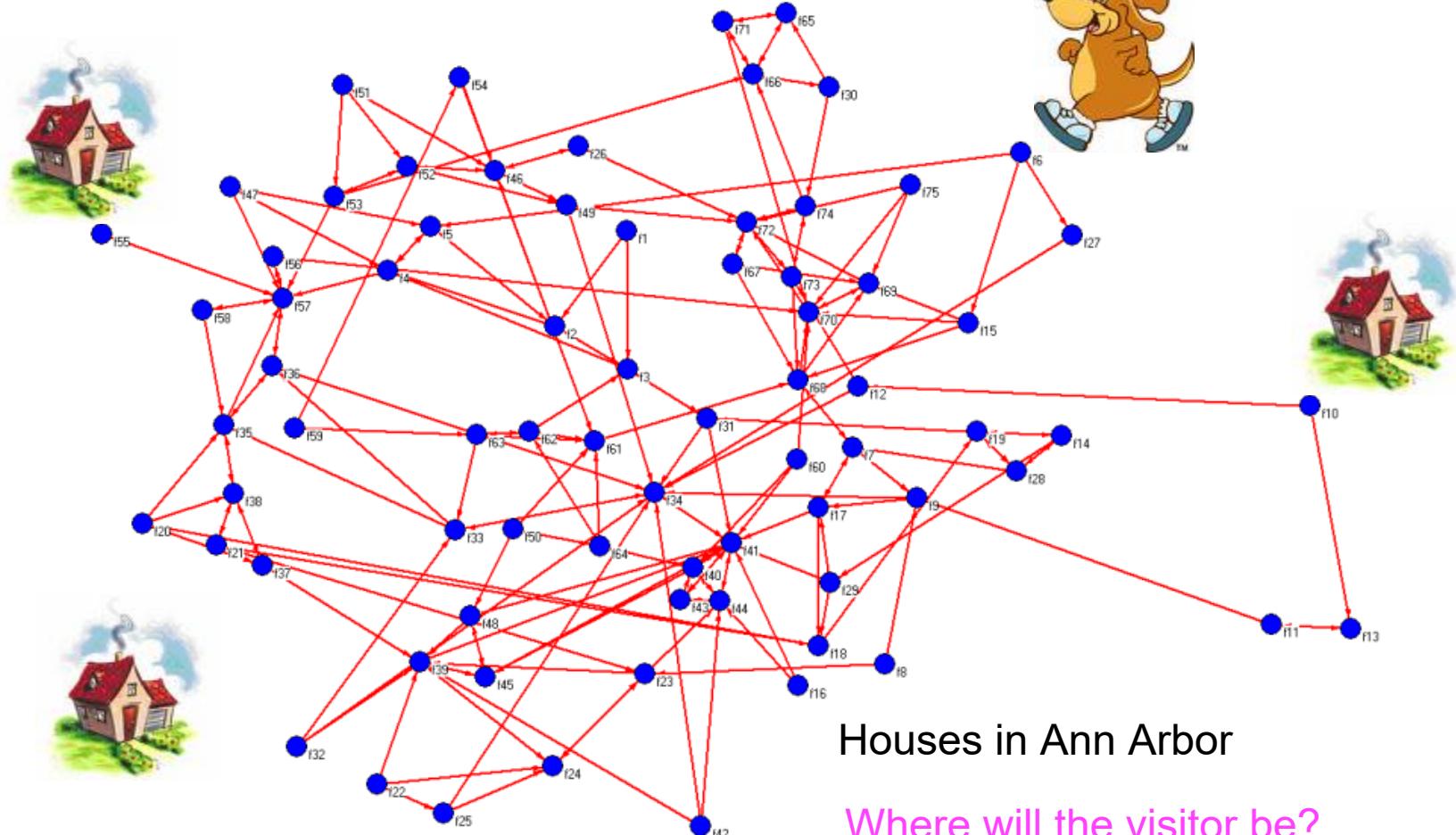
- Degree
- Betweenness
- Closeness
- ...
- Vertices that are linked to by many important vertices are important – recursive definition!
- Question: how to model this in a network?

23

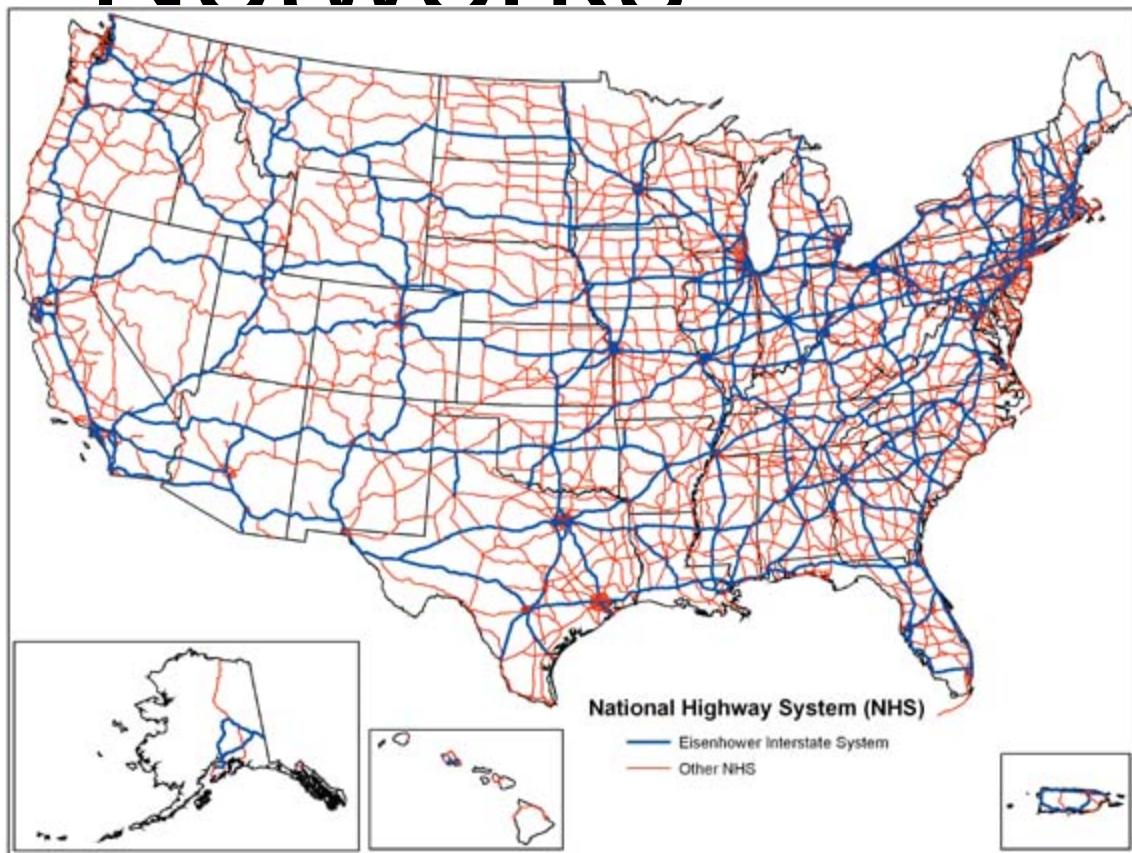
# Two methods we will cover today

- PageRank
  - To understand PageRank we will first cover **Random Walk**
- HITS

# Random Walk on Networks



# Random Walk on Networks

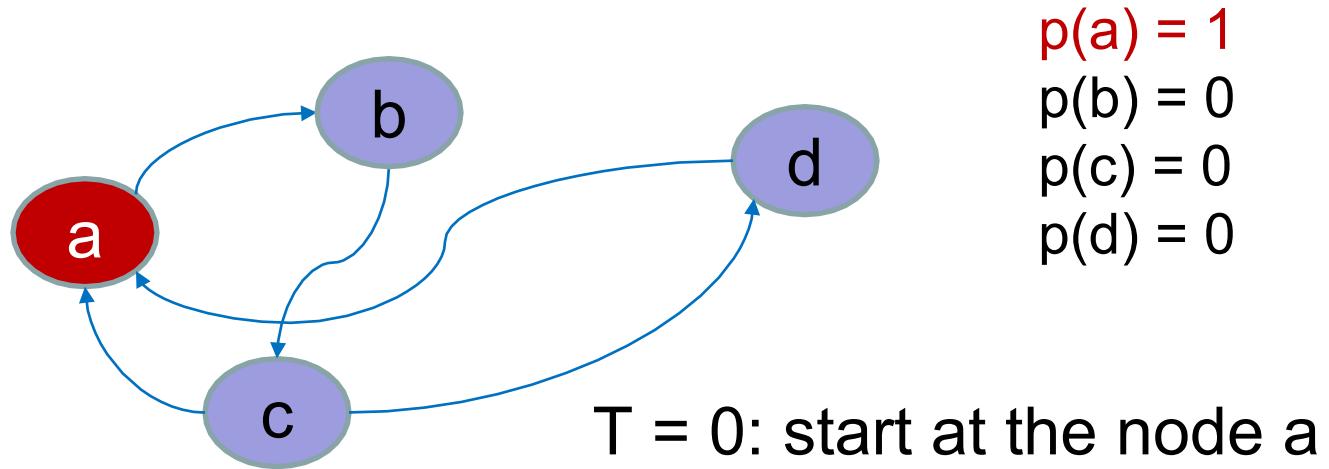


- Track the trace of Forrest Gump:
  - Where to find Forrest in a month from now?
  - Where will Forrest spend most of his time?

# What is a Random Walk?

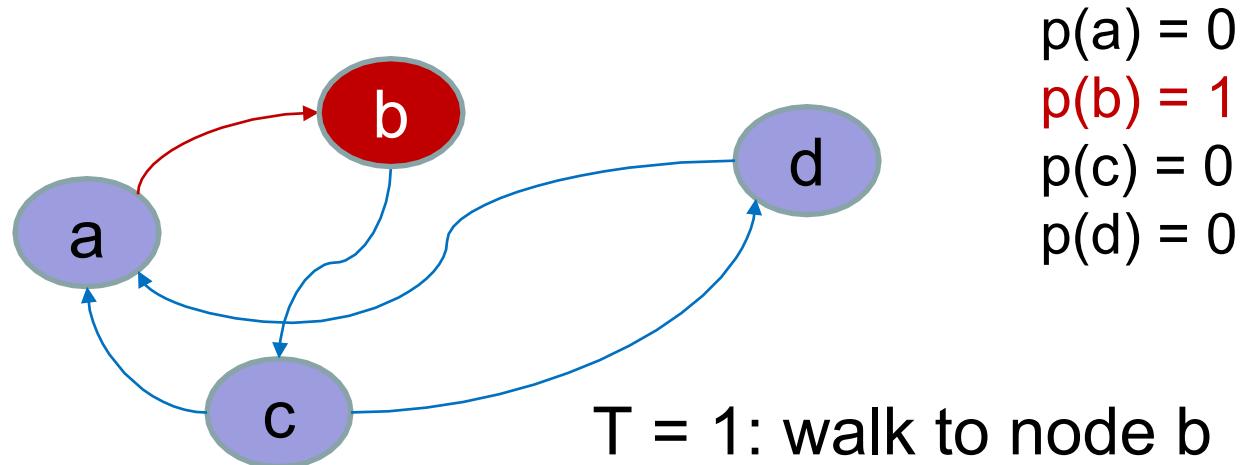
- A “stochastic process” of walking on a network
- At every step, “walk” from the current position to one of its (directed) neighbors
- If more than one neighbor, “randomly” choose one neighbor to walk to.
- $p(X_t = i)$ : the probability that at time  $t$ , the walk is at vertex  $i$ . (sometimes abbreviated as  $p_t(i)$ , or  $p(i)$  )
- For all vertices  $i$ , at any given time  $t$ , the sum of  $p(X_t = i)$  is 1.

# Random Walk: Illustration



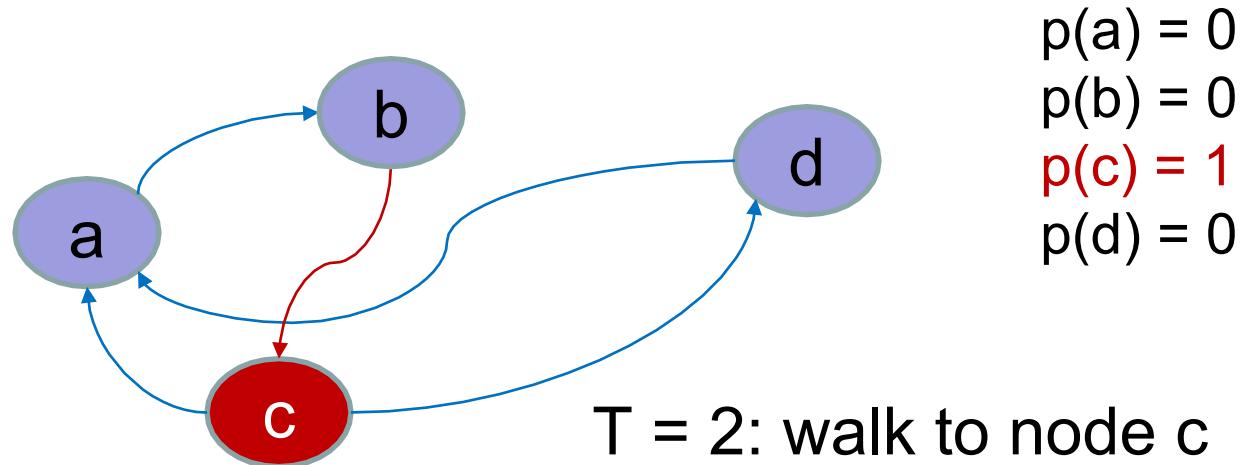
28

# Random Walk: Illustration



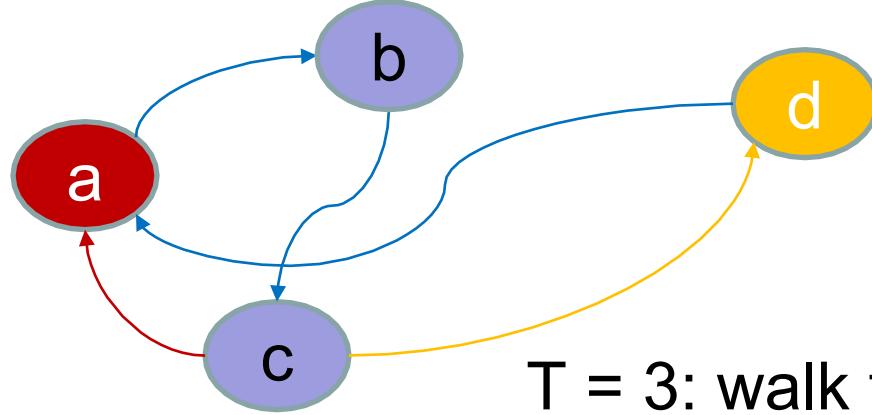
29

# Random Walk: Illustration



30

# Random Walk: Illustration

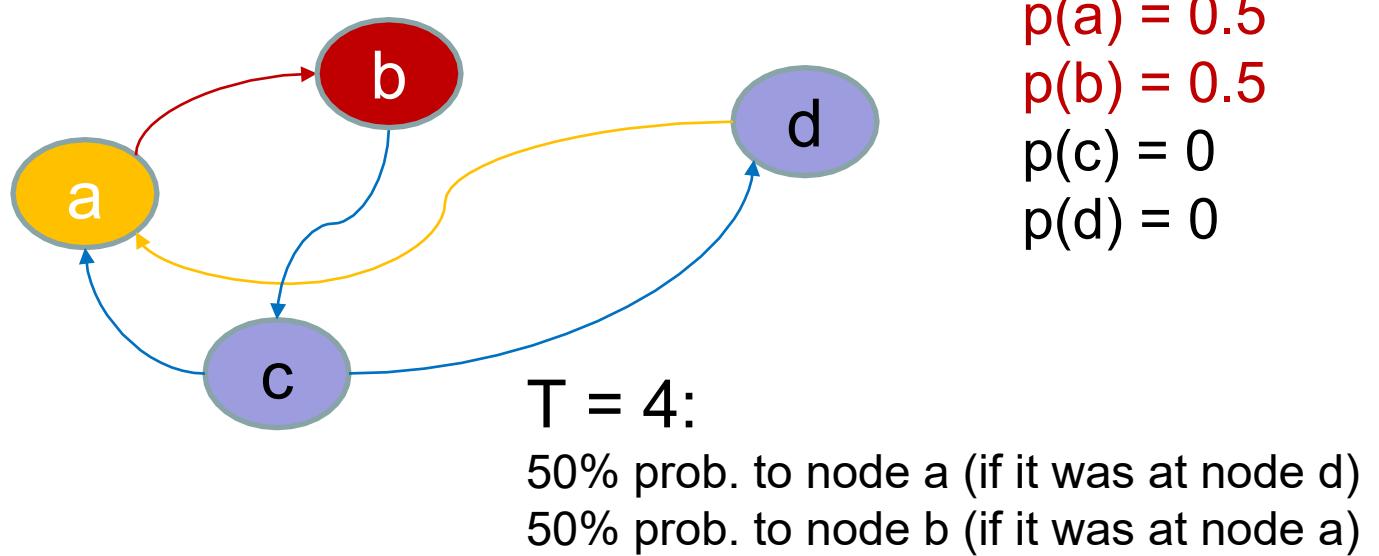


$$\begin{aligned} p(a) &= 0.5 \\ p(b) &= 0 \\ p(c) &= 0 \\ p(d) &= 0.5 \end{aligned}$$

$T = 3$ : walk to either a or d

Question: where will you be next?

# Random Walk: Illustration



Where will you be  
next?

32

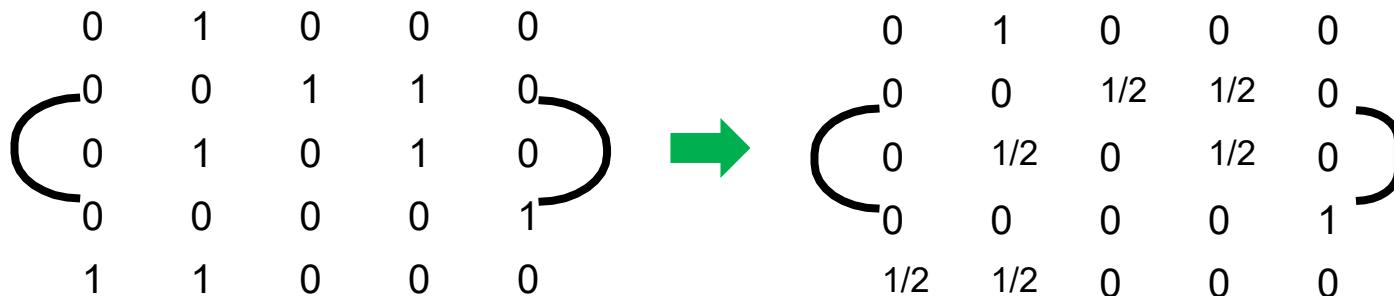
# Network as a Markov Chain

- With such a random walk process, the network is known as a **Markov chain**
  - Named after Andrey Markov
- Basic property:
  - Directed graph (an undirected graph can be viewed as a directed graph with mutual links)
  - Every vertex is a *state*
  - The probability of visiting each vertex at time t is known as the *state distribution* at time t.
  - The state to visit next only depends on the current state

# Transition Probability

- Transition Probability:  $p(X_{t+1} = j | X_t = i)$
- **转移概率** 表示: 如果现在在节点  $i$ , 下一步走到节点  $j$  的概率是多少。
  - Usually abbreviated as  $p(j|i)$ , or  $p(i \rightarrow j)$ , or  $p_{ij}$
  - If the current state is  $i$ , how likely the walk will visit state  $j$  next?
  - $P(j|i) = 1/\text{out-degree}(i)$
  - Or  $w(i, j)/\text{out-degree}(i)$  if the edges are weighted.
- Adjacency matrix → transition matrix:

$$p(i \rightarrow j) = \frac{1}{\text{outdegree}(i)}$$



34

# Random Walk in a Markov Chain

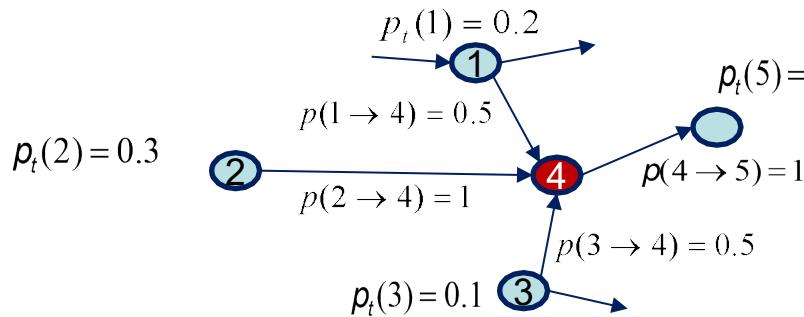
- Given
  - Probabilities:  $p_t(i)$
  - Transition Probabilities:  $p(i \rightarrow j)$
- The distribution at time  $t+1$  can be computed as:

$$p_{t+1}(j) = \sum_{i \in V} p_t(i) \cdot p(i \rightarrow j)$$

How likely the walk will be at  $j$  next

How likely the walk is currently at  $i$

How likely the walk is going from  $i$  to  $j$



Poll: What is  $p_{t+1}(4)$ ?

Respond at [pollev.com/cerenbudak421](https://pollev.com/cerenbudak421)

$$\begin{aligned}p_{t+1}(4) &= p_t(1) \cdot p(1 \rightarrow 4) + p_t(2) \cdot p(2 \rightarrow 4) + p_t(3) \cdot p(3 \rightarrow 4) \\&= 0.3 * 1 + 0.2 * 0.5 + 0.1 * 0.5 = 0.45\end{aligned}$$

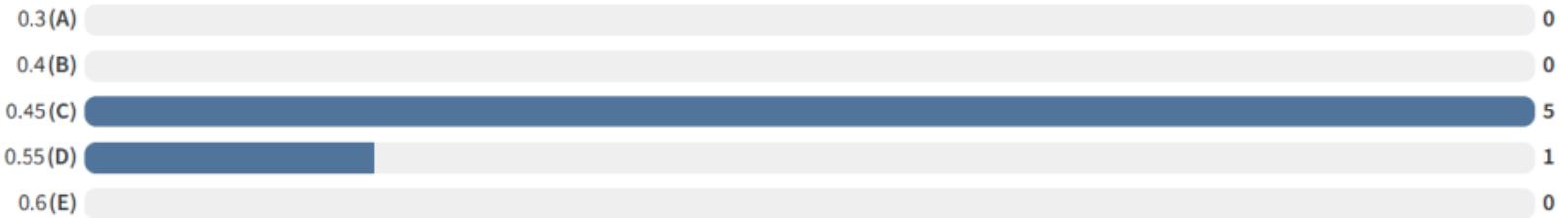


< Back



Join by Web [PollEv.com/cerenbudak421](https://PollEv.com/cerenbudak421)

Join by Text Send **cerenbudak421** to **37607**



Powered by Poll Everywhere

# Stationary Distribution

- Two basic questions:
  - After walking for a sufficiently long time, how much time did you spend on each vertex? (How frequently did you visit each vertex?)
  - After walking for a sufficiently long time, which vertex will you be at (or visit next)?
- They can be explained by the “stationary distribution” of a random walk
- $\pi$ : stationary distribution
  - The “steady” state distribution
  - $\pi(i)$ : the probability that the walk is visiting vertex  $i$  after a sufficiently long time.
  - This is independent from the initial state distribution (i.e., the start point).

37

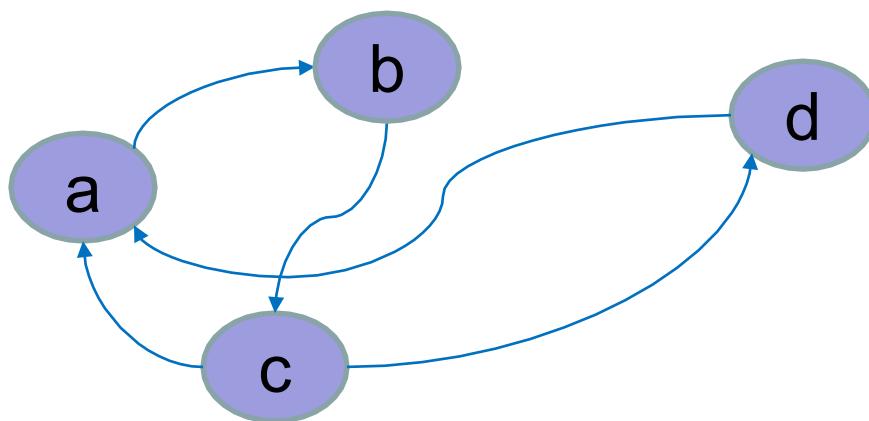
- **Stationary Distribution** 指的是：

当游走的步数趋向无穷大时，在各个状态上的概率分布趋于稳定，不再变化。

换句话说：

| 长期来看，随机游走者在某个节点出现的频率 = 该节点的平稳分布值。

# Example of a Stationary Distribution



Stationary distribution:

$$\pi(a) = 2/7$$

$$\pi(b) = 2/7$$

$$\pi(c) = 2/7$$

$$\pi(d) = 1/7$$

No matter where it starts, after a sufficiently large number of steps, the walk visited a, b, c twice as frequently as d; on the next step the walk is twice likely to be at a, b, or c than d.

# How Long is Sufficiently Long?

- “Stationary” distribution: the “steady” state distribution when the random walk reaches a stationary status
- When  $t$  is sufficiently large, the distribution is steady so that the distribution at time  $t$  equals the distribution at  $t+1$

$$p_t(j) = p_{t+1}(j) = \pi(j)$$

$$\pi(j) = \sum_{i \in V} \pi(i) \cdot p(i \rightarrow j)$$

- When this condition holds for every vertex  $j$ , we have reached the steady status - stationary distribution.
- Again, this is independent from the start point.

# Will this Eventually Happen?

- For some networks, yes
  - Irreducible: possible to reach any node from any node
    - One strongly connected component
  - Positive recurrent: from any node, the walk can return to the same node in finite time.
- For some networks, no
  - Networks with unconnected components;
    - Walks starting in one component cannot reach the others
  - Network has nodes (or groups of nodes) with zero out-links;
    - Walks reaching those nodes will be trapped.
  - ...

# How to Compute the Stationary Distribution?

$$\pi(j) = \sum_{i \in V} \pi(i) \cdot p(i \rightarrow j)$$

- Let  $P$  be the transition matrix ( $P_{ij} = p(i \rightarrow j)$ ), we have

$$\vec{\pi} = \vec{\pi} \cdot P$$

意思是：如果当前分布是  $\pi$ ，随机游走一步之后，分布仍然是  $\pi$ 。  
因此  $\pi$  是 不变的、长期稳定的概率分布。

- In other words, the stationary distribution  $\pi$  is a left eigenvector of the transition matrix  $P$  (with eigenvalue 1)
- $\pi$  can be found by solving a linear system
  - Finds the exact solution of  $\pi$ , but this is usually time consuming

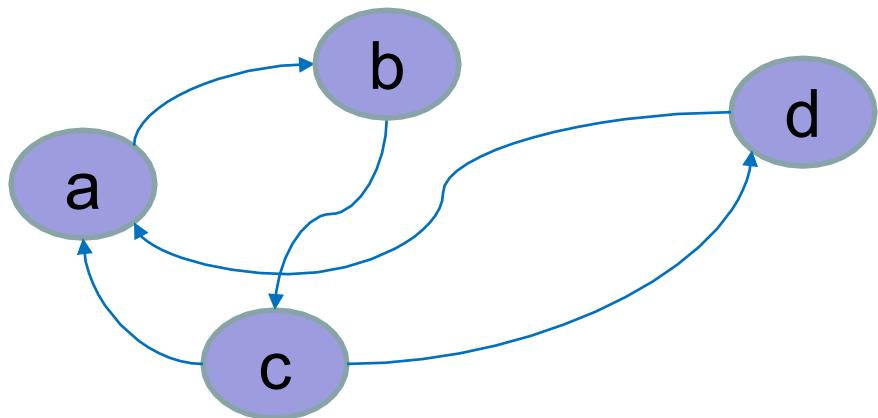
# How to Compute the Stationary Distribution?

- Alternative strategy:
  - Make use of the characteristics of  $\pi$ :
  - $\pi$  is independent from the starting state distribution
  - Choose a random start, and iteratively run the following random walk until convergence.

$$p_{t+1}(j) = \sum_{i \in V} p_t(i) \cdot p(i \rightarrow j)$$

- $\pi$  is found when for every vertex,  $p_t(j) = p_{t+1}(j)$ .

# Example of Finding $\pi$



$$\begin{aligned} T = 0 \\ p(a) &= 1 \\ p(b) &= 0 \\ p(c) &= 0 \\ p(d) &= 0 \end{aligned}$$



$$\begin{aligned} T = 5 \\ p(a) &= 0 \\ p(b) &= 0.5 \\ p(c) &= 0.5 \\ p(d) &= 0 \end{aligned}$$



$$\begin{aligned} T = 10 \\ p(a) &= 0.375 \\ p(b) &= 0.125 \\ p(c) &= 0.25 \\ p(d) &= 0.25 \end{aligned}$$



$$\begin{aligned} T \rightarrow \dots \\ \pi(a) &= 2/7 \\ \pi(b) &= 2/7 \\ \pi(c) &= 2/7 \\ \pi(d) &= 1/7 \\ &\quad \uparrow \\ T \sim 100 \\ p(a) &= 2/7 \\ p(b) &= 2/7 \\ p(c) &= 2/7 \\ p(d) &= 1/7 \end{aligned}$$

44

# Back to ranking the Web

The screenshot shows a Google search results page for the query "network". The search bar at the top contains the word "network". Below the search bar, there are tabs for "All", "Images", "News", "Videos", "Books", and "More". To the right of these tabs are "Settings" and "Tools" buttons. A message indicates "About 6,080,000,000 results (0.74 seconds)". The first result is a link to the Wikipedia page for "Network (1976 film)". The page summary mentions it's a 1976 film directed by Sidney Lumet, written by Paddy Chayefsky, and produced by Howard Gottfried and Fred C. Caruso. It was released by Metro-Goldwyn-Mayer with a box office of \$23.7 million. Below the summary are links for "Plot", "Production", "Release", "Awards and honors", and a "People also ask" section. The "People also ask" section includes questions like "What network means?", "What are the 4 types of networks?", "What is an example of a network?", and "What is the best definition of the word network?". At the bottom of the search results, there is a "Feedback" link.

www.merriam-webster.com › dictionary › network

[Network | Definition of Network by Merriam-Webster](#)

Network definition is - a fabric or structure of cords or wires that cross at regular intervals and are knotted or secured at the crossings. How to use network in a ...

www.imdb.com › title

[Network \(1976\) - IMDb](#)

Directed by Sidney Lumet. With Faye Dunaway, William Holden, Peter Finch, Robert Duvall. A television network cynically exploits a deranged former anchor's ...

Director: Sidney Lumet      Gross USA: \$23,689,877  
Plot Keywords: media | television | television n...      Stars: Faye Dunaway, William Holden, Peter ...

**Problem:**  
**How do we know which individual pages are more important?**

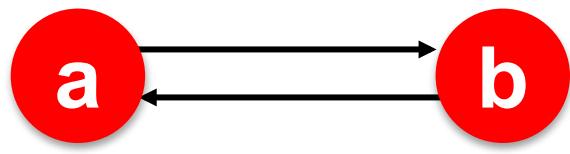
# PageRank: Bringing Order to the Web

- PageRank: Sergey Brin and Larry Page(判断哪个网页更重要)
- Intuition 1: links to a web page can be interpreted as endorsements or recommendations
  - the more links a page receives, the more likely it is to be a good/entertaining/provocative/authoritative/interesting information source
- Intuition 2: but not all link sources are equal
  - a link from a respected information source (e.g., NSF website)
  - a link from a page created by a spammer
- The autonomy of the web; bringing order to the web
- What's behind the primitive success of Google

# Behind PageRank

- "PageRank relies on the uniquely democratic nature of the web by using its vast link structure as an indicator of an individual page's value. In essence, Google interprets a link from page A to page B as a vote, by page A, for page B. But, Google looks at more than the sheer volume of votes, or links a page receives; it also analyzes the page that casts the vote. Votes cast by pages that are themselves "important" weigh more heavily and help to make other pages "important"." -- Google Technology
- What's the secret behind PageRank?
- Essentially a **probability distribution** of how likely a person visits a page by randomly surfing on the web
- Woa, isn't this simply the stationary distribution of a random walk on the web graph?
- Question: In what way is the web different?

# Does this converge?



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

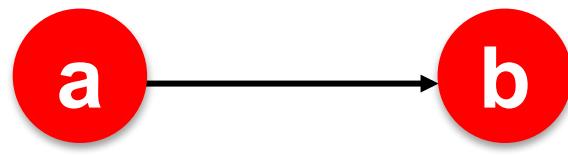
- **Example:**

$$r_a = 1$$

$$r_b = 0$$

Iteration 0, 1, 2, ...

# Does it converge to what we want?



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

- **Example:**

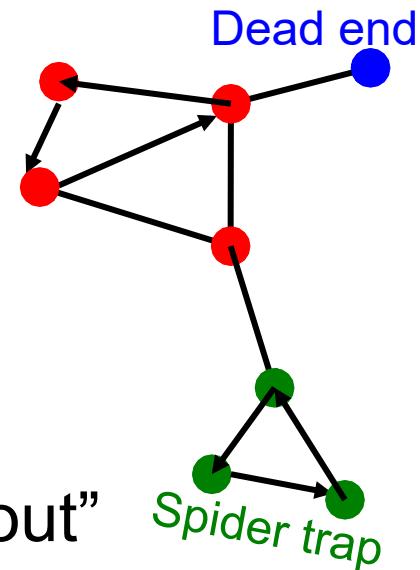
$$\begin{array}{ll} r_a & = \\ r_b & = \end{array} \quad \begin{array}{l} 1 \\ 0 \end{array}$$

Iteration 0, 1, 2, ...

# PageRank: Problems

## 2 problems:

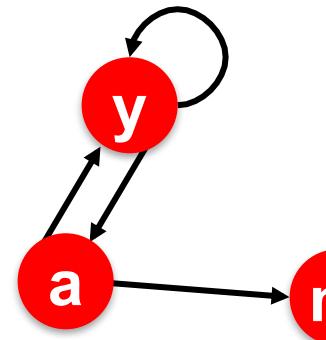
- (1) Some pages are **dead ends** (have no out-links)
  - Random walk has “nowhere” to go to
  - Such pages cause importance to “leak out”
- (2) **Spider traps:**  
(all out-links are within the group)
  - Random walked gets “stuck” in a trap
  - And eventually spider traps absorb all importance



# Problem: Spider Traps

## ■ Power Iteration:

- Set  $r_j = 1$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- And iterate



	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	0
a	$\frac{1}{2}$	0	0
m	0	$\frac{1}{2}$	1

m is a spider trap

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

$$r_a = r_y/2$$

$$r_m = r_a/2 + r_m$$

## ■ Example:

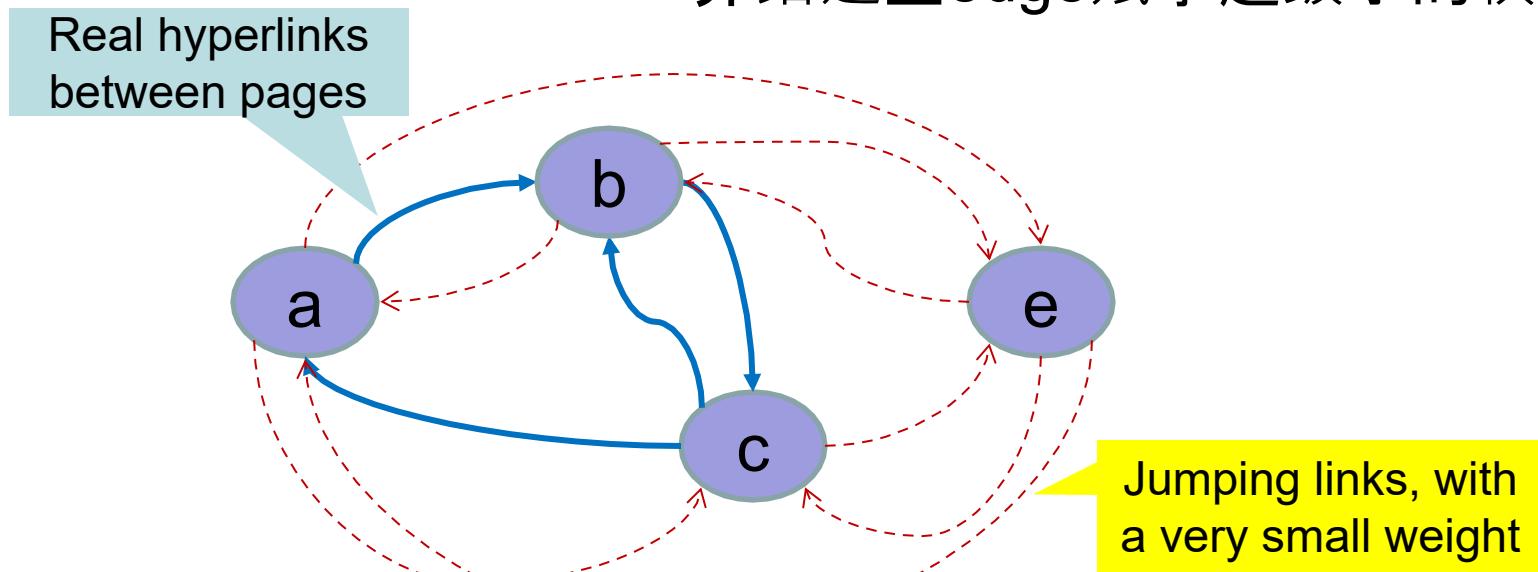
$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 3/6 & 7/12 & 16/24 & & 1 \end{matrix}$$

Iteration 0, 1, 2, ...

All the PageRank score gets “trapped” in node

# Adding Random Jumps to the Web Network

加edge把graph变成complete graph  
并给这些edge赋予超级小的权重



PageRank is essentially computing the stationary distribution of the random walk on this manipulated network

52

# Why do We Need Random Jumps?

- In PageRank, this is explained by the behavior of surfing without following the hyperlinks (surfing with jumping to a random web page).
- But essentially, they need this for a reason...
- Is the web graph connected at all? No!
- Is the web graph strongly connected? No!
  - Think about the bowtie model.
- Are there pages with no out-links? Yes!
- The truth is: there won't be a stationary distribution for a random walk on the web network without random jumps

# Computing PageRank

- d: damping factor which controls the jumping probability
- N: total number of vertices (web pages)
- D(i): the number of out-links of web page I

$$PR(j) = \frac{1-d}{N} + d \cdot \sum_{i:i \rightarrow j} \frac{PR(i)}{D(i)}$$

Visited by randomly jumping from other pages

Visited by surfing from one of its in-links

- Start with a random PageRank value for every vertex, iteratively updating the PageRank until converging

# PageRank: The Complete Algorithm

- **Input:** Graph  $G$  and parameter  $\beta$

- Directed graph  $G$  with **spider traps** and **dead ends**
  - Parameter  $\beta$

- **Output:** PageRank vector  $r$

- **Set:**  $r_j^{(0)} = \frac{1}{N}$ ,  $t = 1$
  - **do:**

- $\forall j: r_j^{(t)} = \sum_{i \rightarrow j} \beta \frac{r_i^{(t-1)}}{d_i}$
    - $r_j^{(t)} = \mathbf{0}$  if in-degree of  $j$  is 0

- Now re-insert the leaked PageRank:

- $\forall j: r_j^{(t)} = r_j^{(t)} + \frac{1-\beta}{N}$

- $t = t + 1$

- **while**  $\sum_j |r_j^{(t)} - r_j^{(t-1)}| > \varepsilon$       **where:**  $S = \sum_j r_j^{(t)}$

If the graph has no dead-ends then the amount of leaked PageRank is  $1-\beta$ . But since we have dead-ends the amount of leaked PageRank may be larger. We have to explicitly account for it by computing  $\mathbf{S}$ .

# Alternate Solution: HITS

- HITS algorithm (developed by Jon Kleinberg):
  - start with a set of pages matching a query
  - expand the set by following forward and back links
  - Keep two scores for each web page:
    - Authority score a;
    - Hub score h;
  - compute the authority scores a, and hub scores h through an iterative approach:

$$a(j) = \sum_{i: i \rightarrow j} h(i)$$

i: i → j

每个网页有两种角色：

$$h(j) = \sum_{i: j \rightarrow i} a(i)$$

i: j → i

Sum over the hub score from all in-links

## 1. Authority (权威)

- 网页本身内容有价值，被很多好页面指向。
- 例子：一篇学术论文首页。

## 2. Hub (枢纽)

- 网页本身可能没多少内容，但它指向很多好的页面。
- 例子：一个论文索引页面，或导航目录。

👉 好 Hub 指向好 Authority; 好 Authority 被好 Hub 指向。

这是一个互相强化 (mutual reinforcement) 的关系。 chigan

Sum over the authority score from all out-links

# Alternate Solution: HITS

- A page has good content (*authority*) thus is pointed by others; a page also points to other pages (*hub*)
- Vertex pointing to many *good authorities* is a *good hub*;
- Vertex pointed to by many *good hubs* is a *good authority*
  - *Page pointing to only a lot of junk pages is junk*
  - *Page only pointed by a lot of junk pages is junk*
- Keeps two scores:  $A(i)$  and  $H(i)$ 
  - $A(i) = \text{sum of } H(j) \text{ for all } j \text{ who points to } i;$
  - $H(i) = \text{sum of } A(j) \text{ for all } j \text{ pointed by } i.$
  - Iteratively updating
  - These values converge in the limit

# Example query: newspaper

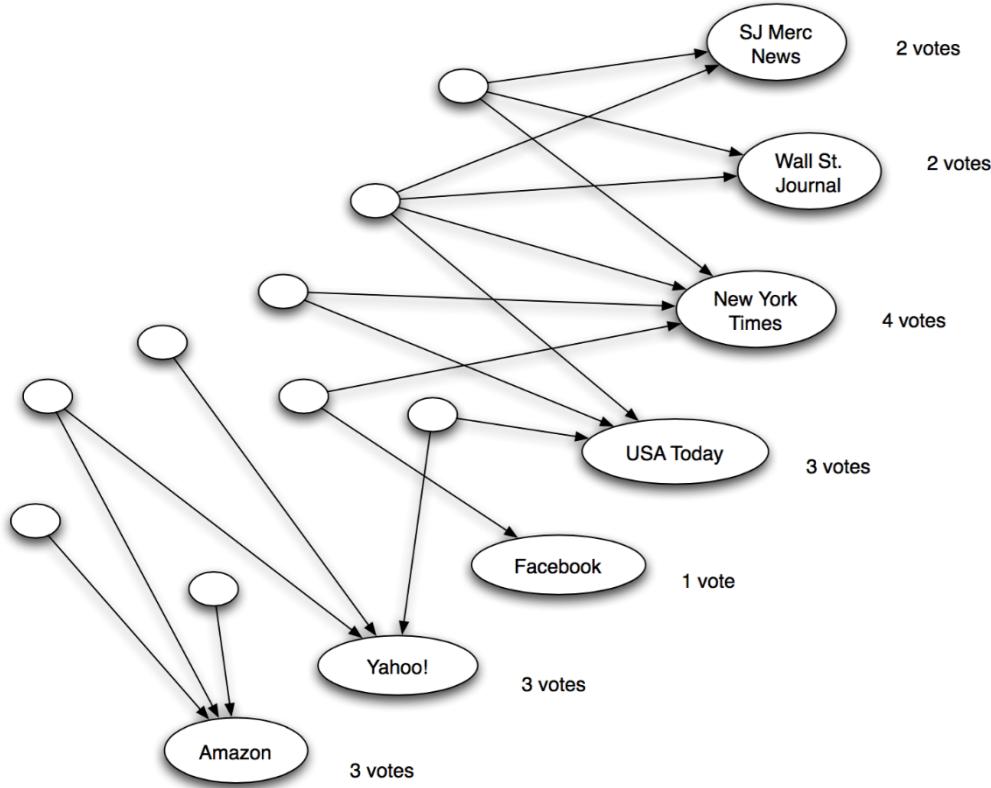


Figure 14.1: Counting in-links to pages for the query “newspapers.”

# Example query: newspaper

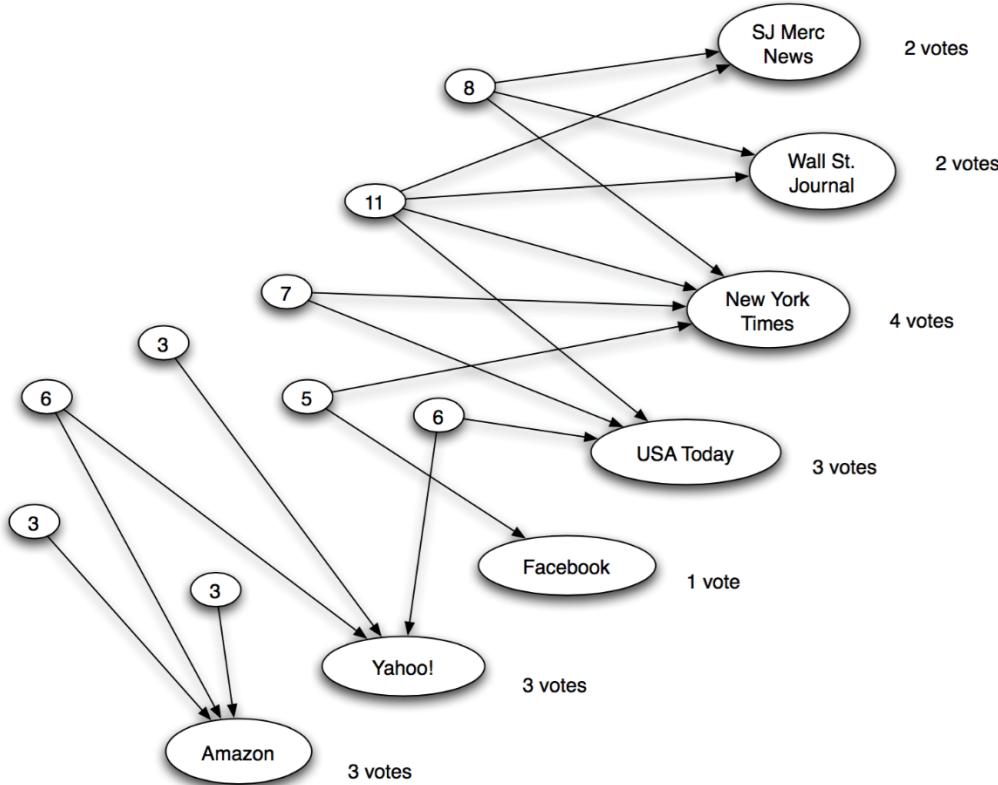


Figure 14.2: Finding good lists for the query “newspapers”: each page’s value as a list is written as a number inside it.

# Example query: newspaper

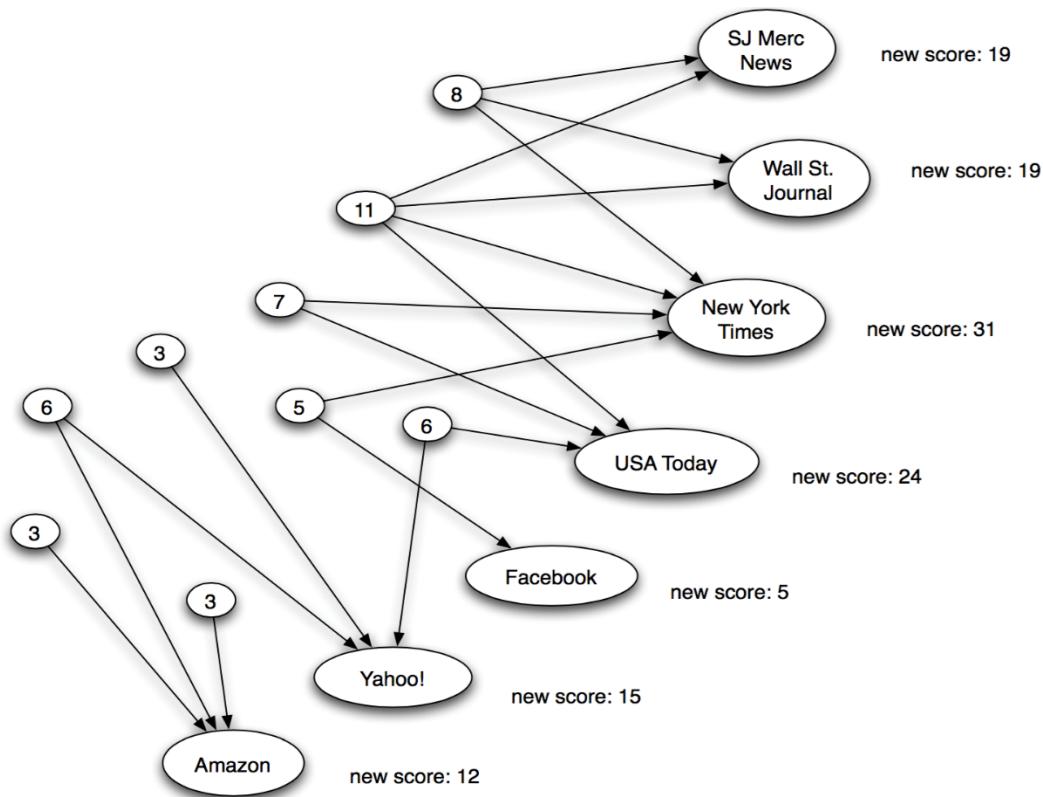


Figure 14.3: Re-weighting votes for the query “newspapers”: each of the labeled page’s new score is equal to the sum of the values of all lists that point to it.

# Task: Finding the best webpage for the query: newspaper

- In the limit:

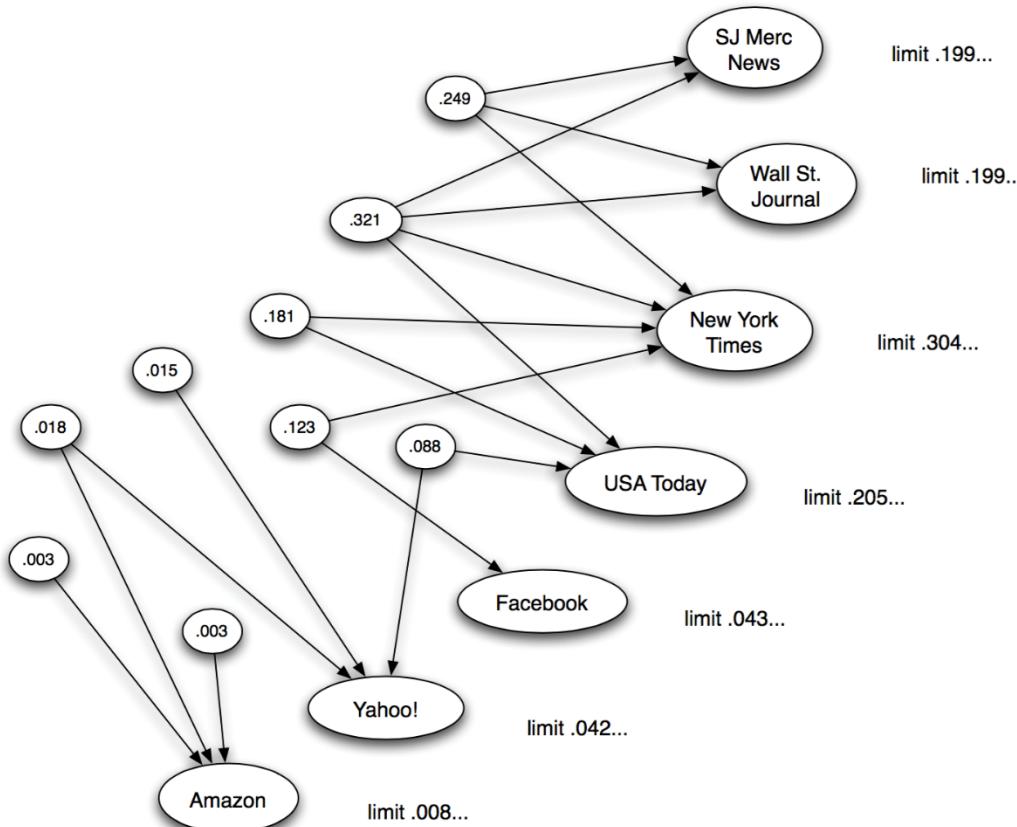


Figure 14.5: Limiting hub and authority values for the query “newspapers.”

### 3. 算法流程

#### 1. 构建子图

- 给定查询词（例如“newspaper”）。
- 先找到一批相关页面（候选集），再扩展它们的入链和出链页面。

#### 2. 初始化分数

- 每个页面的 **authority** 分数  $a(i)$  和 **hub** 分数  $h(i)$  初始设为 1。

#### 3. 迭代更新

- Authority 更新：

$$a(i) = \sum_{j:j \rightarrow i} h(j)$$

即：权威分数 = 所有指向它的 hub 的分数之和。

- Hub 更新：

$$h(i) = \sum_{j:i \rightarrow j} a(j)$$

即：枢纽分数 = 它所指向的所有 authority 的分数之和。

#### 4. 归一化

- 为防止数值无限增长，每次迭代后对向量归一化。

#### 5. 收敛

- 不断迭代，直到分数稳定。



# Implementation

- Given a query, first construct a focused subgraph, and then:

Iterate( $G, k$ )

$G$ : a collection of  $n$  linked pages

$k$ : a natural number

Let  $z$  denote the vector  $(1, 1, 1, \dots, 1) \in \mathbf{R}^n$ .

Set  $x_0 := z$ .

Set  $y_0 := z$ .

For  $i = 1, 2, \dots, k$

Apply the  $\mathcal{I}$  operation to  $(x_{i-1}, y_{i-1})$ , obtaining new  $x$ -weights  $x'_i$ .

Update authority scores

Apply the  $\mathcal{O}$  operation to  $(x'_i, y_{i-1})$ , obtaining new  $y$ -weights  $y'_i$ .

Update hub scores

Normalize  $x'_i$ , obtaining  $x_i$ .

Normalize (important for convergence)

Normalize  $y'_i$ , obtaining  $y_i$ .

End

Return  $(x_k, y_k)$ .

## Initialization

# Topic-Specific PageRank

# Topic-Specific PageRank

随机跳转不再是“均匀分布到所有页面”，而是 限制在一组与主题相关的页面集合  $S$ 。

- Instead of generic popularity, can we measure popularity within a topic?
- 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
  - Example: Query “Trojan” wants different pages depending on whether you are interested in sports, history and computer security

# 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
  - For each teleport set S, we get a different vector  $r_S$

# Matrix Formulation

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

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

- $A$  is stochastic!
- We weighted all pages in the teleport set  $S$  equally
  - Could also assign different weights to pages!
- Compute as for regular PageRank:
  - Multiply by  $M$ , then add a vector
  - Maintains sparseness

# Discovering the Topic Vector $S$

- **Create different PageRanks for different topics**
  - The 16 DMOZ top-level categories:
    - arts, business, sports,...
- **Which topic ranking to use?**
  - User can pick from a menu
  - Classify query into a topic
  - Can use the **context** of the query
    - E.g., query is launched from a web page talking about a known topic
    - History of queries e.g., “basketball” followed by “Jordan”
  - User context, e.g., user’s bookmarks, ...

# Diversity & Page Rank

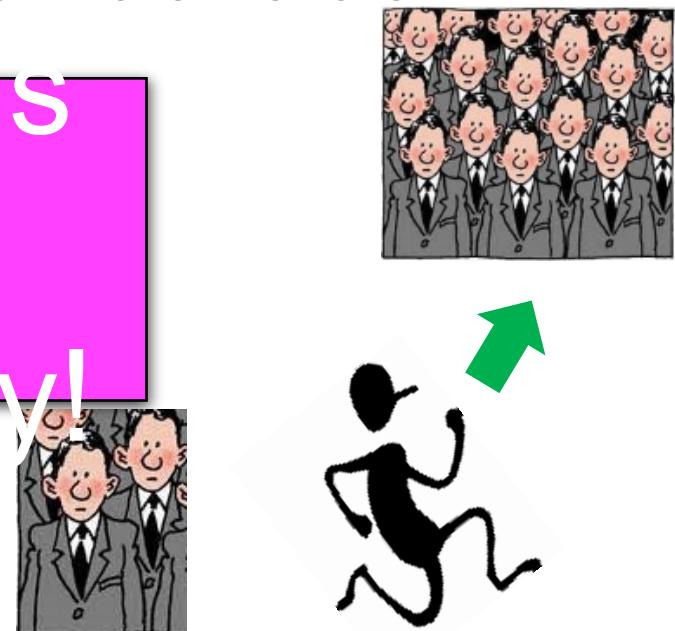
Why is there an issue and what to do

# Reinforcements in Random Walks

- Random walks are not random - rich gets richer;
  - e.g., civilization/immigration – big cities attract larger population
  - Tourism – busy restaurants attract more visitors



Source - <http://www.resettlementagency.co.uk/modern-world-migration/>



## 1. 背景：随机游走的“富者愈富”效应

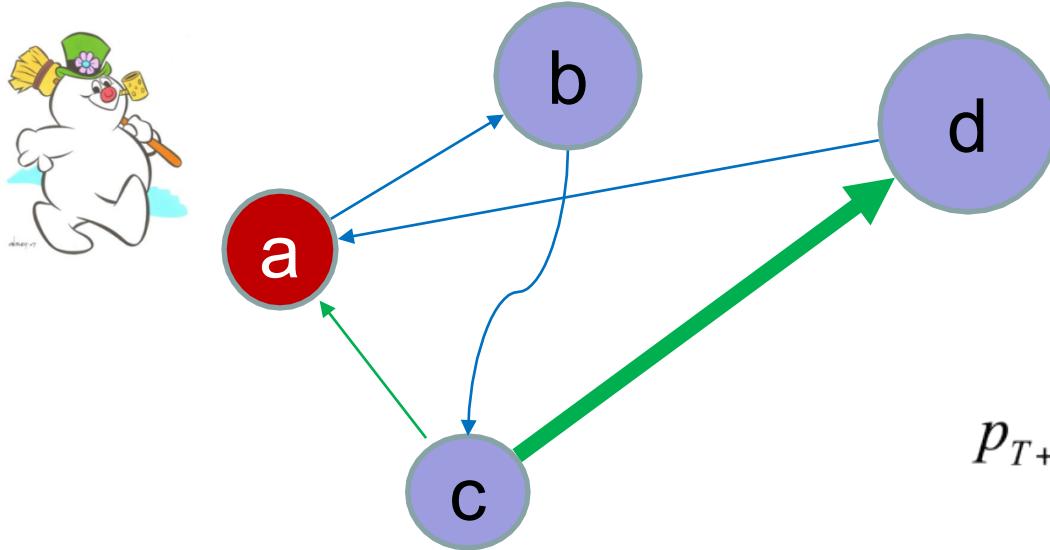
- 在普通的随机游走里，节点的访问概率会逐渐稳定在某个分布（比如 PageRank 的平稳分布）。
- 但现实中常见的现象是：
  - **文明/移民**：大城市越大，吸引的人越多。
  - **旅游/餐馆**：热门景点或餐馆因为人多而更受欢迎。  
👉 这是典型的“rich gets richer”（马太效应）。

## 2. 强化随机游走 (Reinforced Random Walk)

- 在这种模型里，转移概率会随着访问次数发生变化。
- 即：你走过某条边/访问过某个节点的次数越多，未来再次走这条边/访问这个节点的概率就更大。
- 这种机制会导致 收敛到少数“热门”节点，产生更强的集中效应。

# Vertex-Reinforced (顶点强化) Random Walk

(Pemantle 92)



transition probabilities  
change over time

$$p_{T+1}(v) = \sum_{(u,v) \in E} p_T(u,v) p_T(u)$$

Reinforced random walk: transition probability is reinforced by the weight (number of visits) of the target state

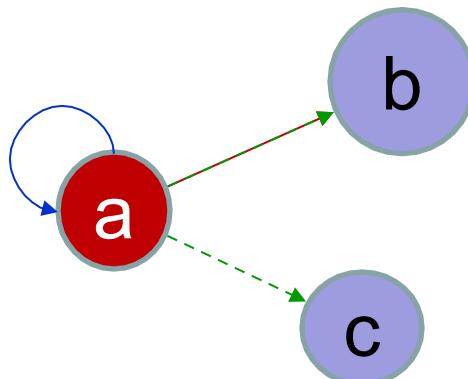
$$p_T(u,v) \propto N_T(v)$$

### 3. 区别总结

特性	强化随机游走 (RRW)	顶点强化随机游走 (VRRW)
强化对象	节点或边 (更一般)	只针对节点
转移概率依赖	边的权重 or 节点的权重	邻居节点的访问次数
表现	可以模拟路径依赖 (ERRW) 或节点吸引力 (VRRW)	更专注“节点富者愈富”现象
范围关系	更广泛的概念	RRW 的一个特例

# DivRank (Mei et al. 2010)

- A smoothed version of Vertex-reinforced Random Walk



$$p_T(u, v) = (1 - \lambda)p^*(v) + \lambda \cdot \frac{p_0(u, v)N_T(v)}{D_T(u)}$$

Random jump, could  
be personalized

“organic” transition  
probability

Number of times  $v$  has  
been visited up to time  $T$

- Adding self-links;
- Efficient approximations: use  $E[N_T(v)]$  to approximate  $N_T(v)$

Cumulative DivRank:

$$E[N_T(v)] \propto \sum_{t=0}^T p_t(v)$$

Pointwise DivRank:

$$E[N_T(v)] \propto p_T(v)$$

# PageRank vs. DivRank



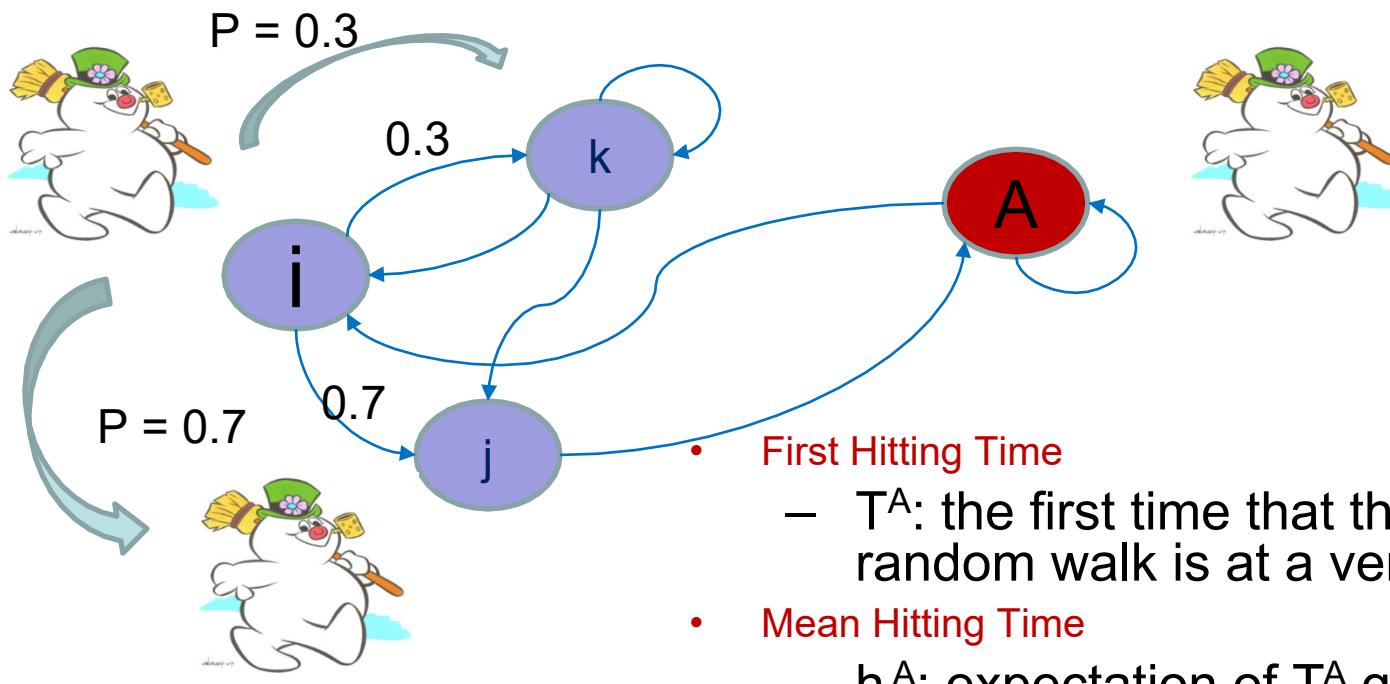
# Another application of random walks:

## Query Recommendation!

# Other Interesting Concepts of Random Walks

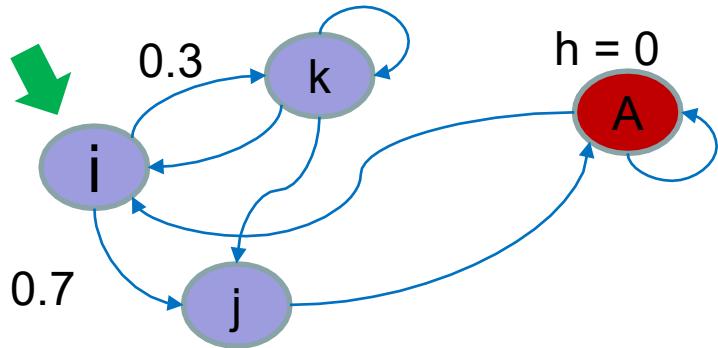
- **Hitting time** –  $h$  is the expected time in a random walk to reach vertex  $v$  starting from vertex  $u$
- **Absorbing random walks** – If the random walk starts from vertex  $i$ , how likely is it to be eventually absorbed by vertex  $s$ ?
- **Commute time** – The expected time to start from  $u$ , reach  $v$ , and then return to  $u$
- **Cover time** – The expected time to start from  $u$  and visit every other vertex at least once
- **Meeting time** – The expected time of two random walks to start from  $u$  and  $v$  and meet at some vertex in between
- ...

# Random Walk and Hitting Time



# Computing Hitting Time

$$h_i^A = 0.7 h_j^A + 0.3 h_k^A + 1$$



$T^A$ : the **first** time that the random walk is at a vertex in  $A$

$$T^A = \min\{ t : X_t \in A, t \geq 0 \}$$

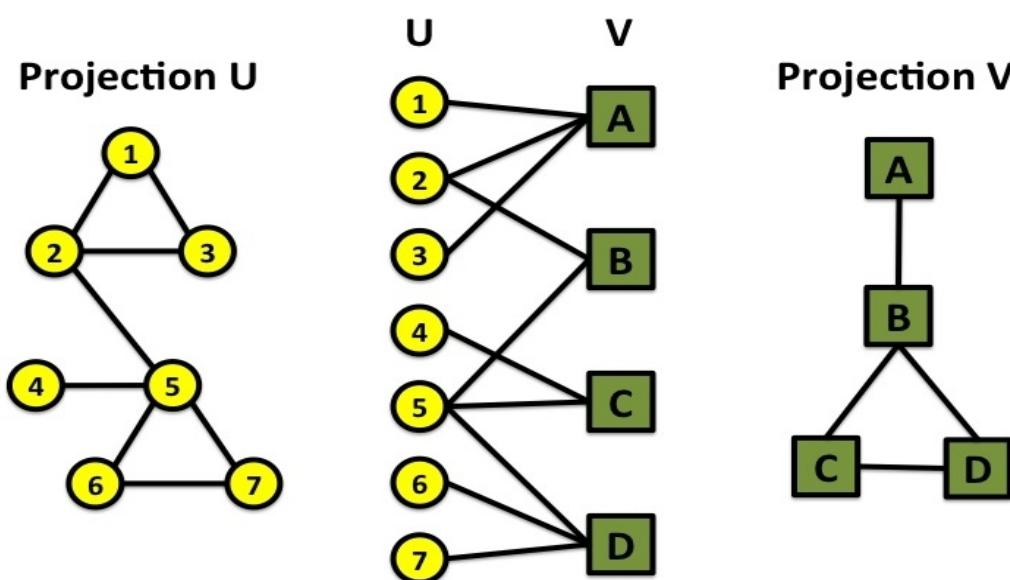
$h_i^A$ : **expectation** of  $T^A$  given that the walk starting from vertex  $i$

$$h_i^A = \begin{cases} \sum_{j \in V} p(i \rightarrow j) h_j^A + 1, & \text{for } i \notin A \\ 0, & \text{for } i \in A \end{cases}$$

Iterative  
Computation

# Quick digression: Bipartite Graphs

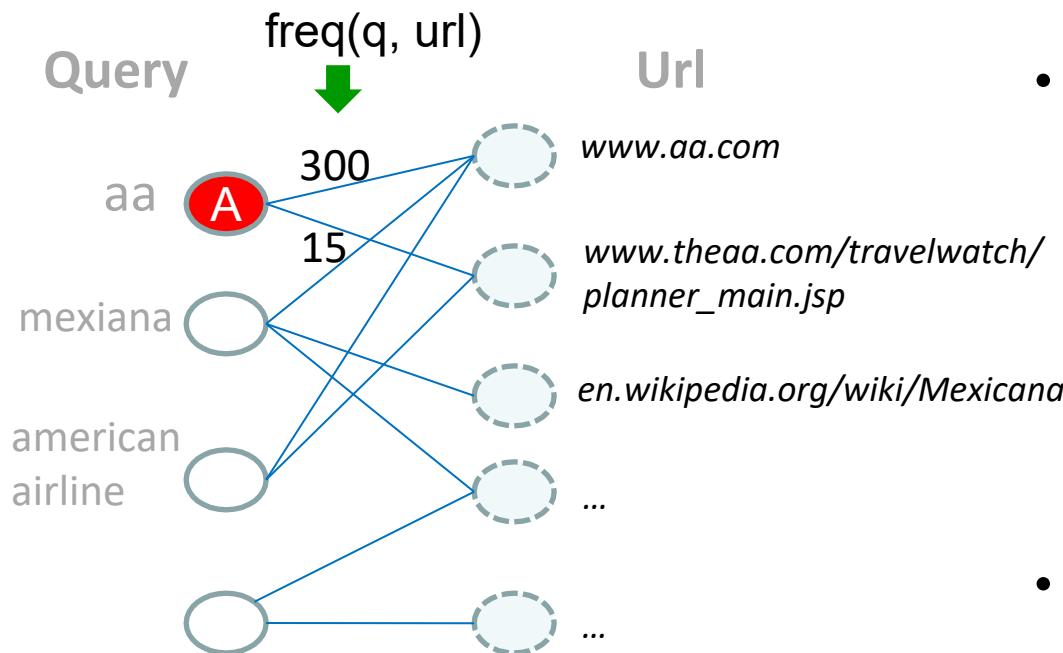
- Bipartite graph (or bigraph) is a [graph](#) whose nodes can be divided into two [disjoint sets](#)  $U$  and  $V$  such that every link connects a node in  $U$  to one in  $V$ ; that is,  $U$  and  $V$  are [independent sets](#).



Examples:

Hollywood actor network  
Collaboration networks  
Disease network (diseasome)  
**Query networks (queries leading to pages)**

# Using Hitting Time for Recommendation



- Construct a (kNN) subgraph from the query log data (of a predefined number of queries/urls)
  - Given a query, construct a subgraph using DFS from that node until certain number of queries is reached.
- Compute transition probabilities  $p(i \rightarrow j)$
- Compute hitting time  $h_i^A$
- Rank candidate queries using  $h_i^A$

Why not consider all <Query, Url> pairs?

- Mei et al. 2008. Query suggestion using hitting time

# Next Generation Search Engine (2010)?

- Better support for query formulation
  - Allow querying from any task context
  - Query by examples
  - Automatic query generation (recommendation)
- Better search accuracy
  - More accurate information need understanding (more personalization and context modeling)
  - More accurate document content understanding (more powerful content analysis)
- More complex retrieval criteria
  - Consider multiple utility aspects of information items (e.g., readability, quality, communication cost)
  - Consider collective value of information items (context-sensitive ranking)
- Better result presentation
  - Better organization of search results to facilitate navigation
  - Better summarization
- More effective and robust retrieval models
  - Automatic parameter tuning
- More scalable retrieval architecture
  - P2P

# Next Generation Search Engines?

- More specialized/customized
  - Special group of users (community engines, e.g., Citeseer)
  - Personalized (better understanding of users)
  - Special genre/domain (better understanding of documents)
- Learning over time (evolving)
- Integration of search, navigation, and recommendation/filtering (full-fledged information management)
- Beyond search to support tasks (e.g., shopping)

# What do You Think?

- Perspective 1: Google Rules
  - No matter what you call as “the next generation search engine”, it will be from Google
- Perspective 2: Bing (or Baidu, Yandex...) has a chance
  - Vertical search engines have their advantage that will beat Google one day
- Perspective 3: None of them
  - The next generation search engine will grow in a totally different soil (Meta? Tiktok? Twitter? Quora? Somewhere else? ).

# You Should Know

- There are much more challenge in a Web search engine than in a standard IR system
  - More data, more complexity, more opportunities
- How PageRank and HITS work
  - What are issues with PageRank and associated extensions
- Research issues in Web search engines