

CS6220: DATA MINING TECHNIQUES

Mining Graph/Network Data

Instructor: Yizhou Sun

yzsun@ccs.neu.edu

March 16, 2016

Methods to Learn

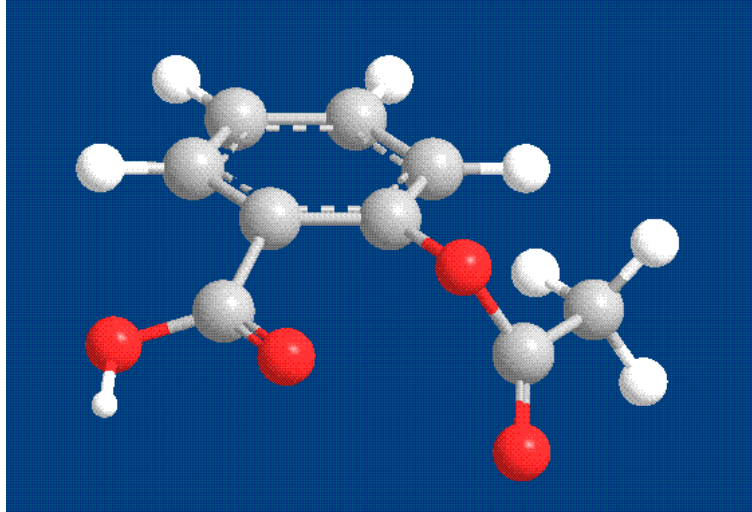
	Matrix Data	Text Data	Set Data	Sequence Data	Time Series	Graph & Network	Images
Classification	Decision Tree; Naïve Bayes; Logistic Regression SVM; kNN			HMM		Label Propagation*	Neural Network
Clustering	K-means; hierarchical clustering; DBSCAN; Mixture Models; kernel k-means*	PLSA				SCAN*; Spectral Clustering	
Frequent Pattern Mining			Apriori; FP-growth	GSP; PrefixSpan			
Prediction	Linear Regression				Autoregression	Recommendation	
Similarity Search					DTW	P-PageRank	
Ranking						PageRank	

Mining Graph/Network Data

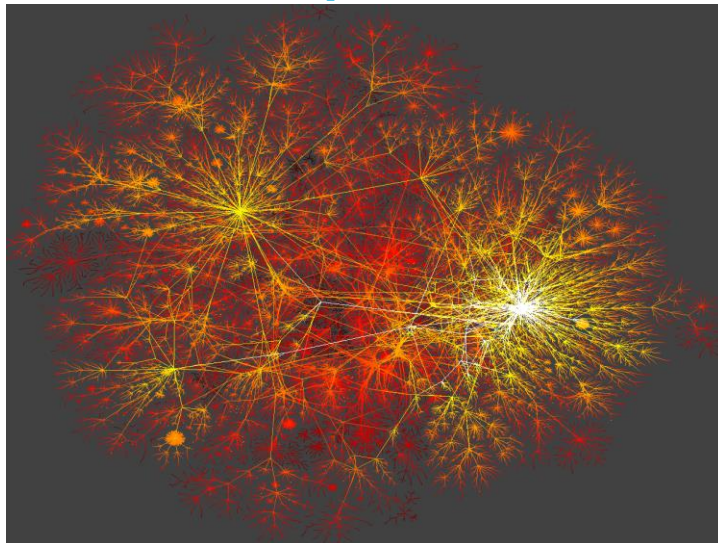
- Introduction to Graph/Network Data
- PageRank
- Proximity Definition in Graphs
- Clustering
- Summary



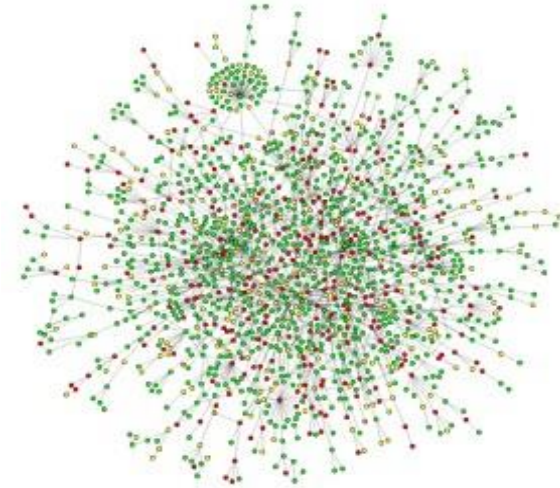
Graph, Graph, Everywhere



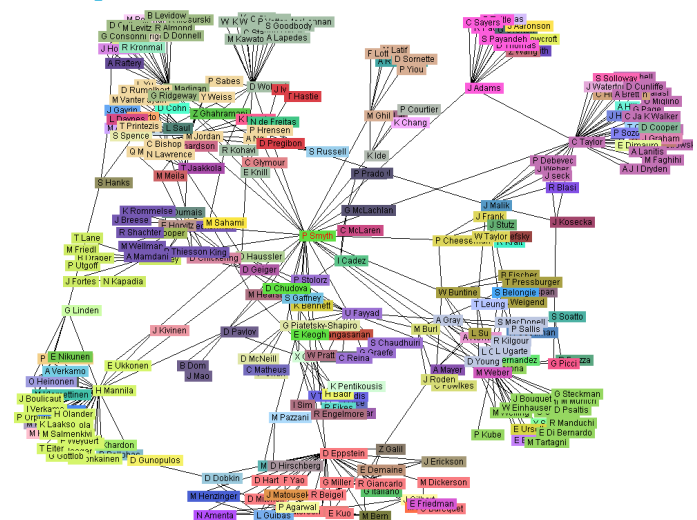
Aspirin



Internet



Yeast protein interaction network



Co-author network

from H. Jeong et al Nature 411, 41 (2001)

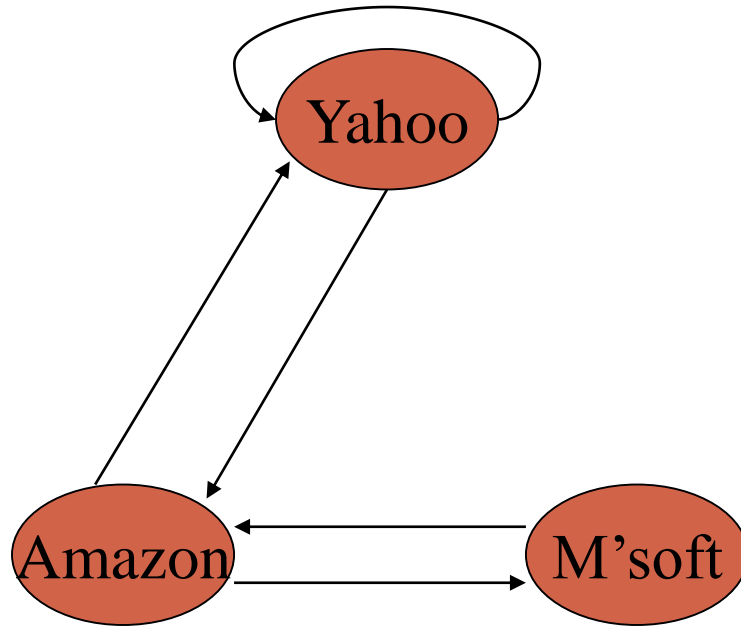
Why Graph Mining?

- Graphs are ubiquitous
 - Chemical compounds (Cheminformatics)
 - Protein structures, biological pathways/networks (Bioinformatics)
 - Program control flow, traffic flow, and workflow analysis
 - XML databases, Web, and social network analysis
- Graph is a general model
 - Trees, lattices, sequences, and items are degenerated graphs
- Diversity of graphs
 - Directed vs. undirected, labeled vs. unlabeled (edges & vertices), weighted, with angles & geometry (topological vs. 2-D/3-D)
- Complexity of algorithms: many problems are of high complexity

Representation of a Graph

- $G = \langle V, E \rangle$
 - $V = \{u_1, \dots, u_n\}$: node set
 - $E \subseteq V \times V$: edge set
- Adjacency matrix
 - $A = \{a_{ij}\}, i, j = 1, \dots, N$
 - $a_{ij} = 1, \text{ if } \langle u_i, u_j \rangle \in E$
 - $a_{ij} = 0, \text{ if } \langle u_i, u_j \rangle \notin E$
 - Undirected graph vs. Directed graph
 - $A = A^T$ vs. $A \neq A^T$
 - Weighted graph
 - Use W instead of A , where w_{ij} represents the weight of edge $\langle u_i, u_j \rangle$


Example



	y	a	m
y	1	1	0
a	1	0	1
m	0	1	0

Adjacency matrix A

Mining Graph/Network Data

- Introduction to Graph/Network Data
- PageRank 
- Personalized PageRank
- Summary

The History of PageRank

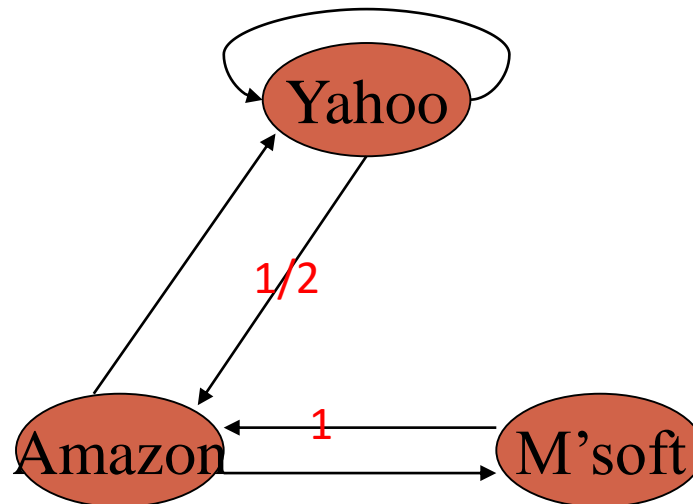
- PageRank was developed by Larry Page (hence the name *Page-Rank*) and Sergey Brin.
- It is first as part of a research project about a new kind of search engine. That project started in 1995 and led to a functional prototype in 1998.
- Shortly after, Page and Brin founded Google.

Ranking web pages

- Web pages are not equally “important”
 - www.cnn.com vs. a personal webpage
- Inlinks as votes
 - The more inlinks, the more important
- Are all inlinks equal?
 - Higher ranked inlink should play a more important role
 - Recursive question!

Simple recursive formulation

- Each link's vote is proportional to the **importance** of its source page
- If page **P** with importance **x** has **n** outlinks, each link gets x/n votes
- Page **P**'s own importance is the sum of the votes on its inlinks



Matrix formulation

- Matrix **M** has one row and one column for each web page

- Suppose page j has n outlinks

- If $j \rightarrow i$, then $M_{ij} = 1/n$
- Else $M_{ij} = 0$

	y	a	m	
y	1	1	0	→ 1/2, 0, 1
a	1	0	1	
m	0	1	0	

- **M** is a **column stochastic matrix**

- Columns sum to 1

- Suppose **r** is a vector with one entry per web page

- r_i is the importance score of page i
- Call it the **rank vector**
- $|\mathbf{r}| = 1$ (i.e., $r_1 + r_2 + \dots + r_N = 1$)

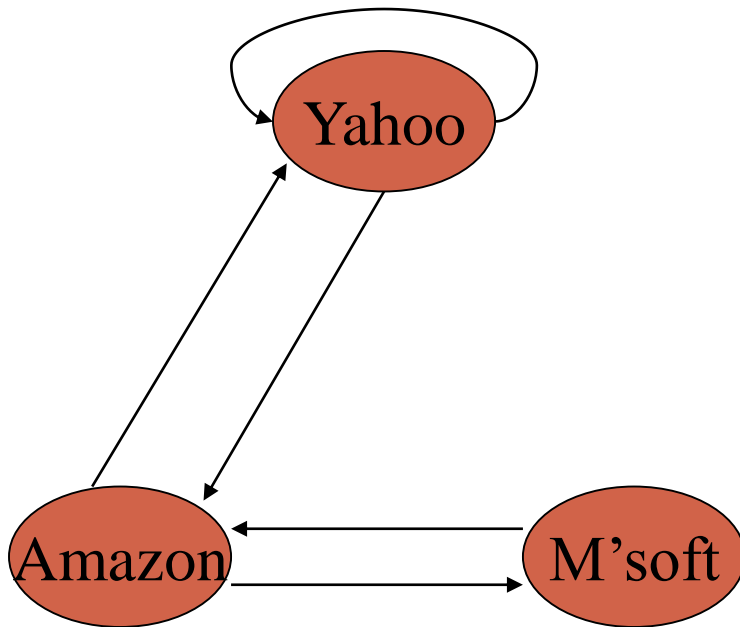
Eigenvector formulation

- The flow equations can be written

$$r = Mr$$

- So the rank vector is an eigenvector of the stochastic web matrix
 - In fact, its first or principal eigenvector, with corresponding eigenvalue 1

Example



$$y = y/2 + a/2$$

$$a = y/2 + m$$

$$m = a/2$$

	y	a	m
y	1/2	1/2	0
a	1/2	0	1
m	0	1/2	0

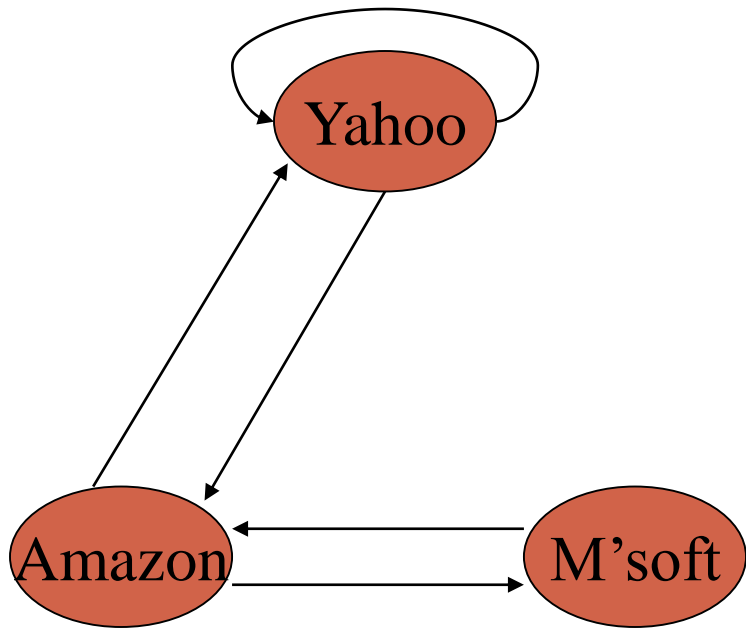
$$\mathbf{r} = \mathbf{M} * \mathbf{r}$$

y	1/2	1/2	0	y
a	1/2	0	1	a
m	0	1/2	0	m

Power Iteration method

- Simple iterative scheme
- Suppose there are N web pages
 - Initialize: $\mathbf{r}^0 = [1/N, \dots, 1/N]^T$
 - Iterate: $\mathbf{r}^{k+1} = \mathbf{M}\mathbf{r}^k$
 - Stop when $\|\mathbf{r}^{k+1} - \mathbf{r}^k\|_1 < \varepsilon$
 - $\|\mathbf{x}\|_1 = \sum_{1 \leq i \leq N} |x_i|$ is the L_1 norm
 - Can use any other vector norm e.g., Euclidean

Power Iteration Example



	y	a	m
y	1/2	1/2	0
a	1/2	0	1
m	0	1/2	0

y	=	1/3	1/3	5/12	3/8		2/5
a		1/3	1/2	1/3	11/24	...	2/5
m		1/3	1/6	1/4	1/6		1/5
		r_0	r_1	r_2	r_3	...	r^*

Random Walk Interpretation

- Imagine a **random web surfer**
 - At any time t , surfer is on some page P
 - At time $t+1$, the surfer follows an outlink from P uniformly at random
 - Ends up on some page Q linked from P
 - Process repeats indefinitely
- Let $\mathbf{p}(t)$ be a vector whose i^{th} component is the probability that the surfer is at page i at time t
 - $\mathbf{p}(t)$ is a probability distribution on pages

The stationary distribution

- Where is the surfer at time $t+1$?
 - Follows a link uniformly at random
 - $\mathbf{p}(t+1) = \mathbf{M}\mathbf{p}(t)$
- Suppose the random walk reaches a state such that $\mathbf{p}(t+1) = \mathbf{M}\mathbf{p}(t) = \mathbf{p}(t)$
 - Then $\mathbf{p}(t)$ is called a **stationary distribution** for the random walk
- Our rank vector \mathbf{r} satisfies $\mathbf{r} = \mathbf{M}\mathbf{r}$
 - So it is a stationary distribution for the random surfer

Existence and Uniqueness

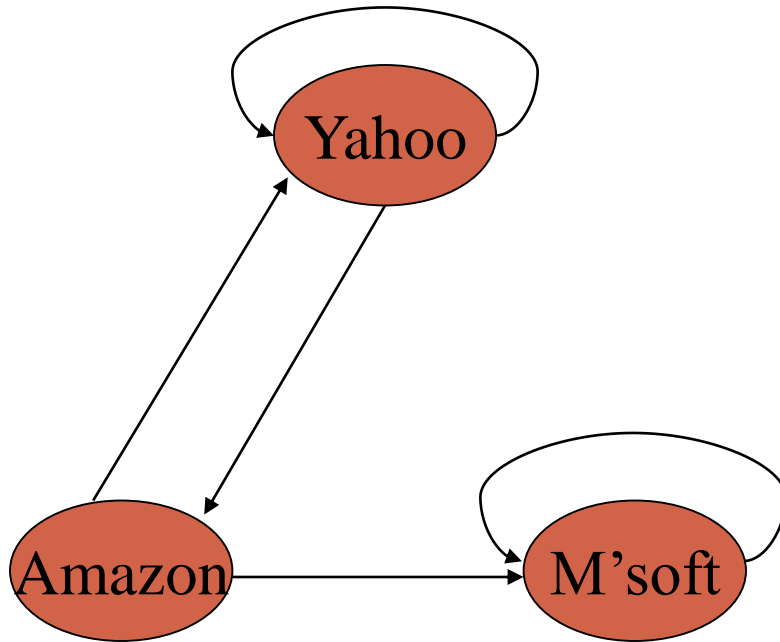
A central result from the theory of random walks (aka Markov processes):

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 $t = 0$.

Spider traps

- A group of pages is a **spider trap** if there are no links from within the group to outside the group
 - Random surfer gets trapped
- Spider traps violate the conditions needed for the random walk theorem

Microsoft becomes a spider trap



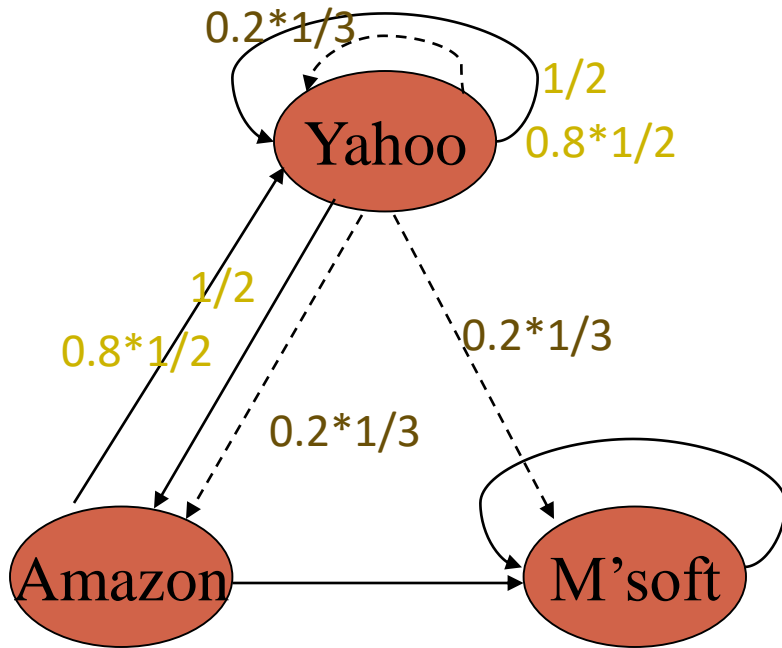
	y	a	m
y	1/2	1/2	0
a	1/2	0	0
m	0	1/2	1

y	=	1/3	1/3	1/4	5/24		0
a		1/3	1/6	1/6	1/8	...	0
m		1/3	1/2	7/12	2/3		1

Random teleports

- The Google solution for spider traps
- At each time step, the random surfer has two options:
 - With probability β , follow a link at random
 - With probability $1-\beta$, jump to some page uniformly at random
 - Common values for β are in the range 0.8 to 0.9
- Surfer will teleport out of spider trap within a few time steps

Random teleports ($\beta = 0.8$)



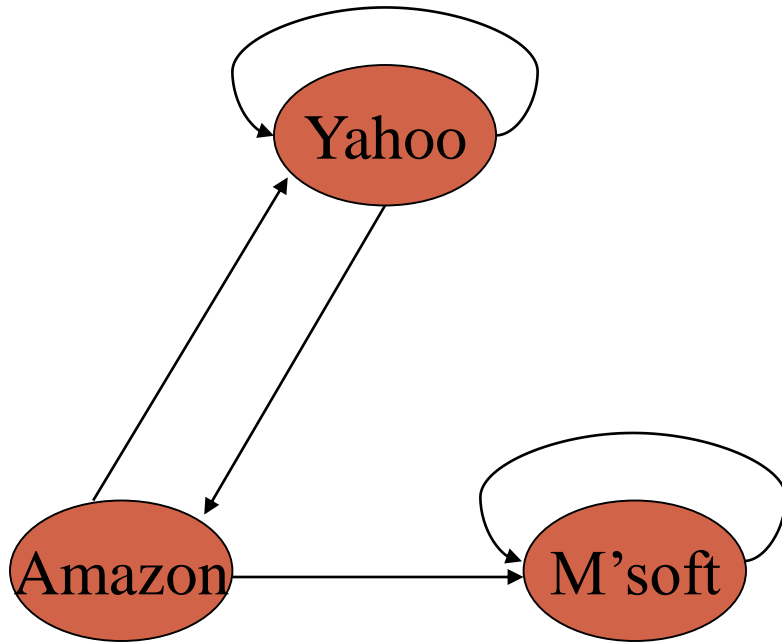
$$\begin{array}{c} y \\ a \\ m \end{array} \begin{array}{c} y \\ 1/2 \\ 1/2 \\ 0 \end{array} \quad 0.8 * \begin{array}{c} y \\ 1/2 \\ 1/2 \\ 0 \end{array} \quad + \quad 0.2 * \begin{array}{c} y \\ 1/3 \\ 1/3 \\ 1/3 \end{array}$$

$$0.8 \begin{array}{ccc} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{array} \quad + \quad 0.2 \begin{array}{ccc} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{array}$$

-----> : teleport links from "Yahoo"

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

Random teleports ($\beta = 0.8$)



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

$$\begin{matrix} y \\ a \\ m \end{matrix} \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix}$$

$$\begin{matrix} y \\ a \\ m \end{matrix} = \begin{bmatrix} 0.333 \\ 0.333 \\ 0.333 \end{bmatrix} \begin{bmatrix} 0.333 \\ 0.200 \\ 0.467 \end{bmatrix} \begin{bmatrix} 0.280 \\ 0.200 \\ 0.520 \end{bmatrix} \begin{bmatrix} 0.259 \\ 0.179 \\ 0.563 \end{bmatrix} \dots \begin{bmatrix} 7/33 \\ 5/33 \\ 21/33 \end{bmatrix}$$

Matrix formulation

- Suppose there are N pages
 - Consider a page j , with set of outlinks $O(j)$
 - We have $M_{ij} = 1/|O(j)|$ 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 with probability $(1-\beta)/N$
 - reducing the probability of following each outlink from $1/|O(j)|$ to $\beta/|O(j)|$
 - Equivalent: tax each page a fraction $(1-\beta)$ of its score and redistribute evenly

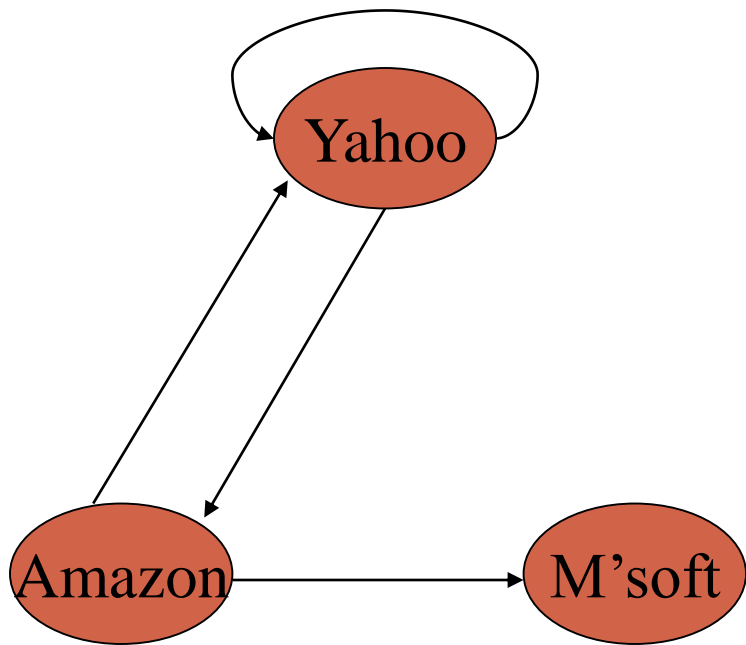
PageRank

- Construct the N-by-N matrix **A** as follows
 - $A_{ij} = \beta M_{ij} + (1-\beta)/N$
- Verify that **A** is a stochastic matrix
- The **page rank vector** **r** is the principal eigenvector of this matrix
 - satisfying $\mathbf{r} = \mathbf{A}\mathbf{r}$
- Equivalently, **r** is the stationary distribution of the random walk with teleports

Dead ends

- Pages with no outlinks are “dead ends” for the random surfer
 - Nowhere to go on next step

Microsoft becomes a dead end



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

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

$$\begin{matrix} y \\ a \\ m \end{matrix} \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 1/15 \end{bmatrix}$$

$$\begin{matrix} y \\ a \\ m \end{matrix} = \begin{bmatrix} 1/3 & 1/3 \\ 1/3 & 0.2 \\ 1/3 & 0.2 \end{bmatrix}$$

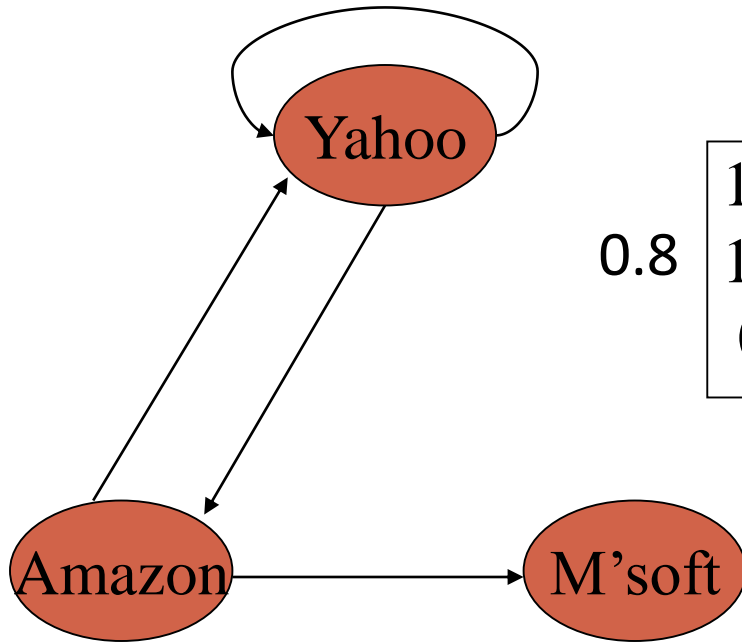
$$\begin{matrix} \dots \\ \dots \\ \dots \end{matrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

↓
Non-stochastic!

Dealing with dead-ends

- **Teleport**
 - Follow random teleport links with probability 1.0 from dead-ends
 - Adjust matrix accordingly
- **Prune and propagate**
 - Preprocess the graph to eliminate dead-ends
 - Might require multiple passes
 - Compute page rank on reduced graph
 - Approximate values for deadends by propagating values from reduced graph

Dealing dead end: teleport

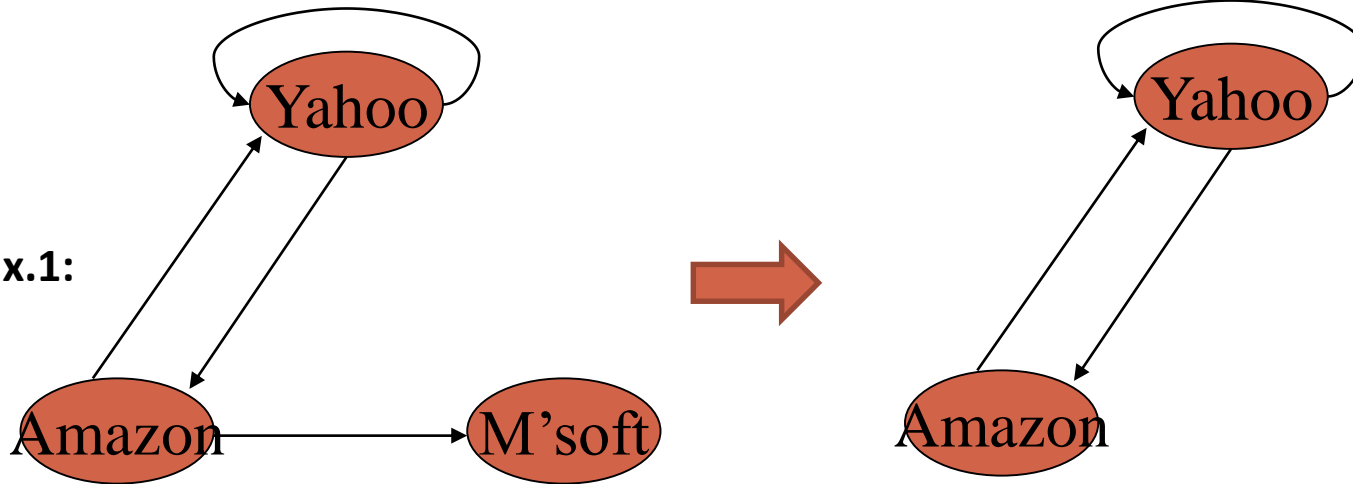


$$0.8 \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 0 \end{bmatrix} + \begin{bmatrix} 0.2*1/3 & 0.2*1/3 & 1*1/3 \\ 0.2*1/3 & 0.2*1/3 & 1*1/3 \\ 0.2*1/3 & 0.2*1/3 & 1*1/3 \end{bmatrix}$$

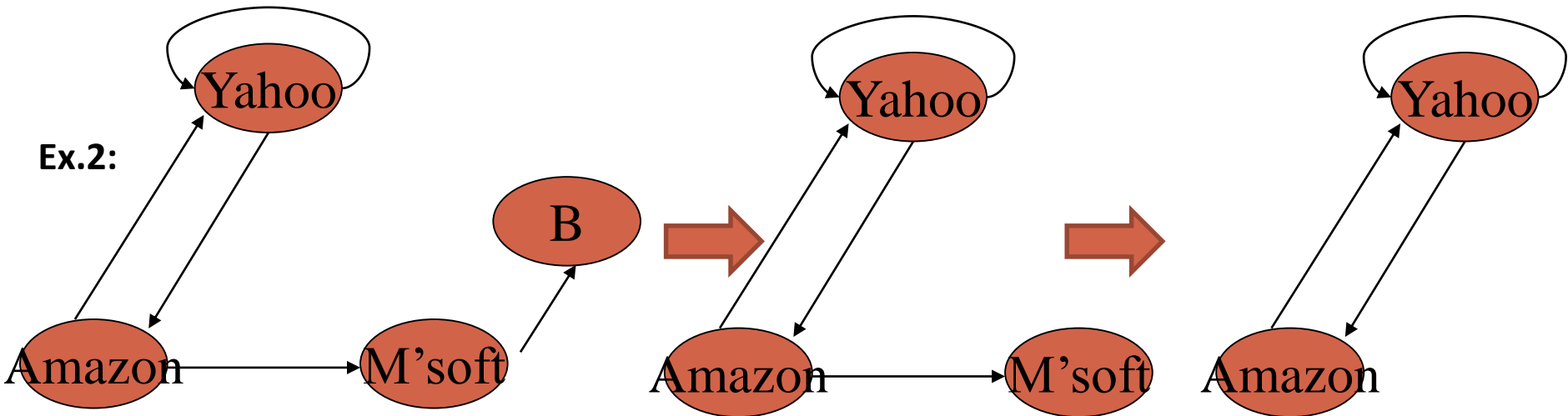
$$\begin{matrix} y \\ a \\ m \end{matrix} \begin{bmatrix} 7/15 & 7/15 & 1/3 \\ 7/15 & 1/15 & 1/3 \\ 1/15 & 7/15 & 1/3 \end{bmatrix}$$

Dealing dead end: reduce graph

Ex.1:



Ex.2:



Computing PageRank

- Key step is matrix-vector multiplication
 - $\mathbf{r}^{\text{new}} = \mathbf{A}\mathbf{r}^{\text{old}}$
- Easy if we have enough main memory to hold \mathbf{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 \mathbf{A} has N^2 entries
 - 10^{18} is a large number!

Rearranging the equation

$\mathbf{r} = \mathbf{A}\mathbf{r}$, where

$$A_{ij} = \beta M_{ij} + (1-\beta)/N$$

$$r_i = \sum_{1 \leq j \leq N} A_{ij} r_j$$

$$r_i = \sum_{1 \leq j \leq N} [\beta M_{ij} + (1-\beta)/N] r_j$$

$$= \beta \sum_{1 \leq j \leq N} M_{ij} r_j + (1-\beta)/N \sum_{1 \leq j \leq N} r_j$$

$$= \beta \sum_{1 \leq j \leq N} M_{ij} r_j + (1-\beta)/N, \text{ since } |\mathbf{r}| = 1$$

$$\mathbf{r} = \beta \mathbf{M}\mathbf{r} + [(1-\beta)/N]_N$$

where $[x]_N$ is an N-vector with all entries x

Sparse matrix formulation

- We can rearrange the page rank equation:
 - $\mathbf{r} = \beta \mathbf{M} \mathbf{r} + [(1-\beta)/N]_N$
 - $[(1-\beta)/N]_N$ is an N -vector with all entries $(1-\beta)/N$
- \mathbf{M} is a sparse matrix!
 - 10 links per node, approx $10N$ entries
- So in each iteration, we need to:
 - Compute $\mathbf{r}^{\text{new}} = \beta \mathbf{M} \mathbf{r}^{\text{old}}$
 - Add a constant value $(1-\beta)/N$ to each entry in \mathbf{r}^{new}

Sparse matrix encoding

- Encode sparse matrix using only nonzero entries
 - Space proportional roughly to number of links
 - say $10N$, or $4 * 10 * 1$ billion = 40GB
 - still won't fit in memory, but will fit on disk

source node	degree	destination nodes
0	3	1, 5, 7
1	5	17, 64, 113, 117, 245
2	2	13, 23


Basic Algorithm

- Assume we have enough RAM to fit \mathbf{r}^{new} , plus some working memory
 - Store \mathbf{r}^{old} and matrix \mathbf{M} on disk

Basic Algorithm:

- Initialize: $\mathbf{r}^{\text{old}} = [1/N]_N$
- Iterate:
 - **Update:** Perform a sequential scan of \mathbf{M} and \mathbf{r}^{old} to update \mathbf{r}^{new}
 - Write out \mathbf{r}^{new} to disk as \mathbf{r}^{old} for next iteration
 - Every few iterations, compute $|\mathbf{r}^{\text{new}} - \mathbf{r}^{\text{old}}|$ and stop if it is below threshold
 - Need to read in both vectors into memory

Mining Graph/Network Data

- Introduction to Graph/Network Data
- PageRank
- Proximity Definition in Graphs 
- Clustering
- Summary

Personalized PageRank

- Query-dependent Ranking
 - For a query webpage u , which webpages are most important to u ?
 - We need a measure $s(u,v)$
 - The relative important webpages to different queries would be different

Calculation of P-PageRank

- Recall PageRank calculation:

- $\mathbf{r} = \beta \mathbf{M} \mathbf{r} + [(1-\beta)/N] \mathbf{1}_N$ or

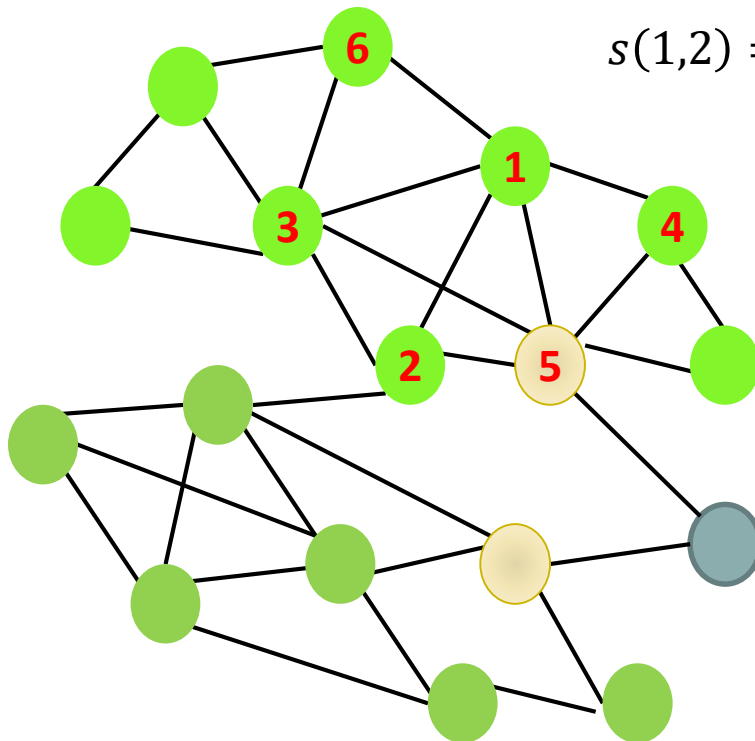
- $\mathbf{r} = \beta \mathbf{M} \mathbf{r} + (1-\beta) \mathbf{r}_0$, where $\mathbf{r}_0 = \begin{pmatrix} 1/N \\ 1/N \\ \dots \\ 1/N \end{pmatrix}$

- For P-PageRank, $s(u,v) = r(v)$

by replacing \mathbf{r}_0 with $\mathbf{r}_0 = \begin{pmatrix} 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \end{pmatrix}$ ← uth webpage

Common Neighbors

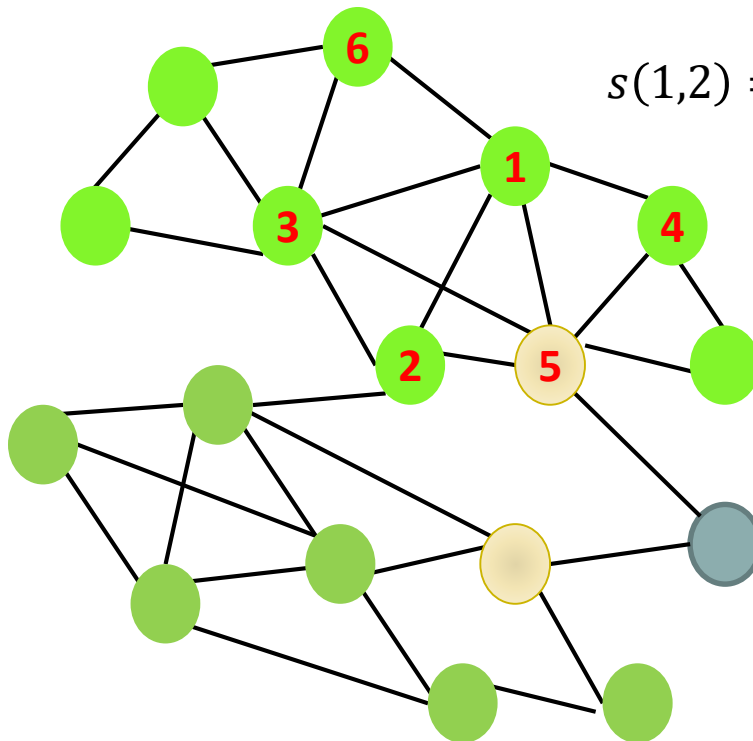
- $s(u, v) = |\Gamma(u) \cap \Gamma(v)|$,
where $\Gamma(u)$ denotes the neighbors of u



$$s(1,2) = |\{4, 5, 2, 3, 6\} \cap \{1, 3, 5\}| = |\{3, 5\}| = 2$$

Jaccard's Coefficient

- $s(u, v) = \frac{|\Gamma(u) \cap \Gamma(v)|}{|\Gamma(u) \cup \Gamma(v)|}$

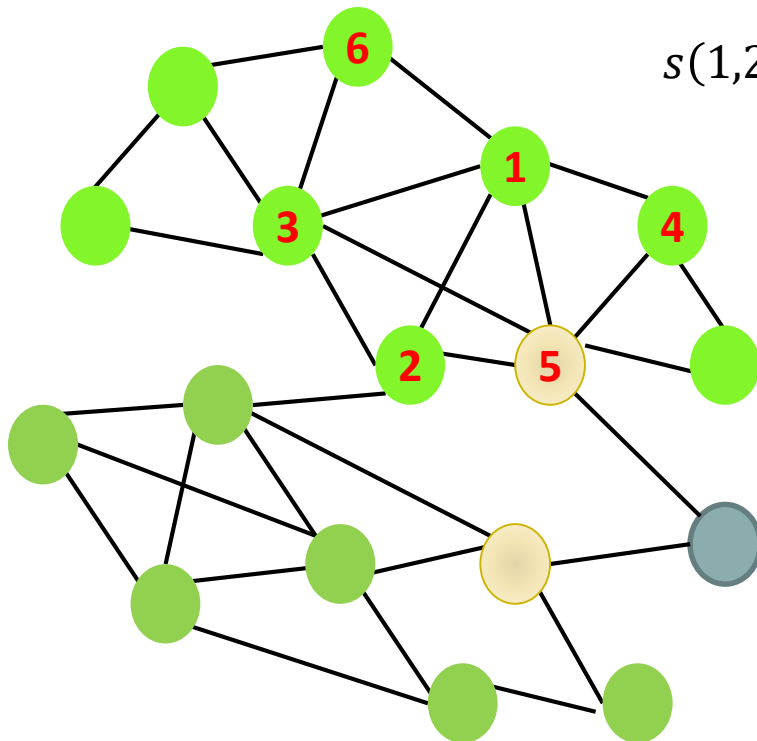


$$s(1,2) = \frac{|\{4, 5, 2, 3, 6\} \cap \{1, 3, 5\}|}{|\{4, 5, 2, 3, 6\} \cup \{1, 3, 5\}|} = \frac{2}{6} = \frac{1}{3}$$

Adamic/Adar


- $s(u, v) = \sum_{w \in \Gamma(u) \cap \Gamma(v)} \frac{1}{\log |\Gamma(w)|}$

- A more connected node will be punished



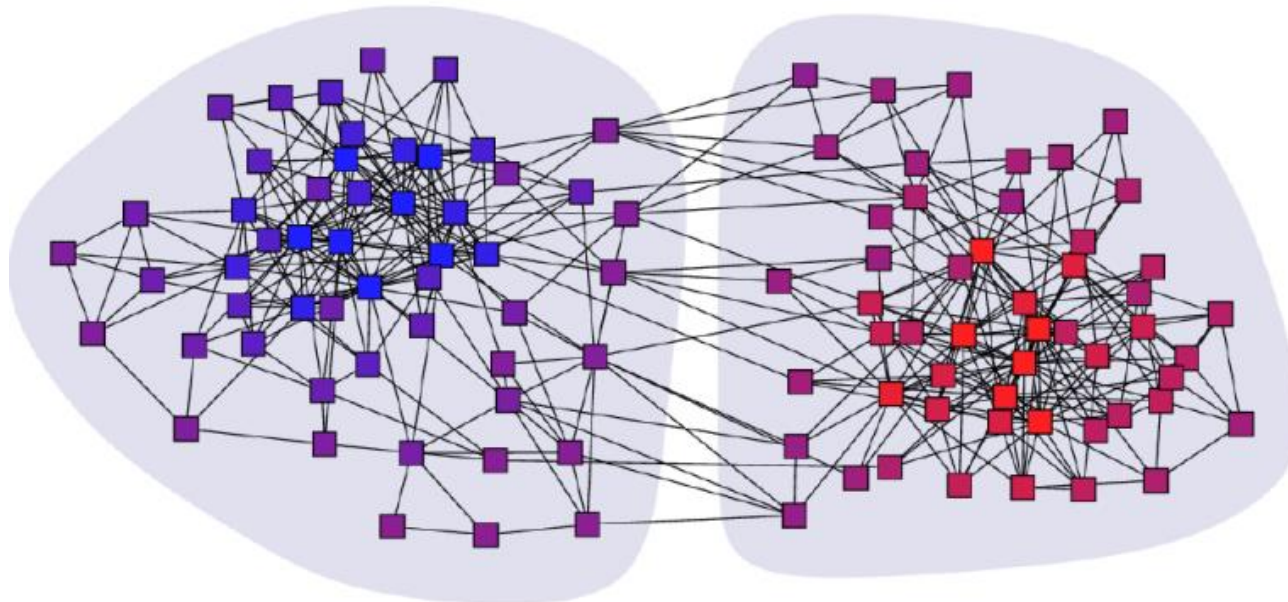
$$s(1,2) = \frac{1}{\log |\Gamma(3)|} + \frac{1}{\log |\Gamma(5)|} = \frac{1}{\log 6} + \frac{1}{\log 6} \\ = 1.12 \text{ (in the original paper, take } e \text{ as base)}$$

Mining Graph/Network Data

- Introduction to Graph/Network Data
- PageRank
- Proximity Definition in Graphs
- Clustering 
- Summary

Clustering Graphs and Network Data

- Applications
 - Bi-partite graphs, e.g., customers and products, authors and conferences
 - Web search engines, e.g., click through graphs and Web graphs
 - Social networks, friendship/coauthor graphs



Clustering books about politics [Newman, 2006]

Spectral Clustering

- Reference: ICDM'09 Tutorial by Chris Ding
- Example:
 - Clustering supreme court justices according to

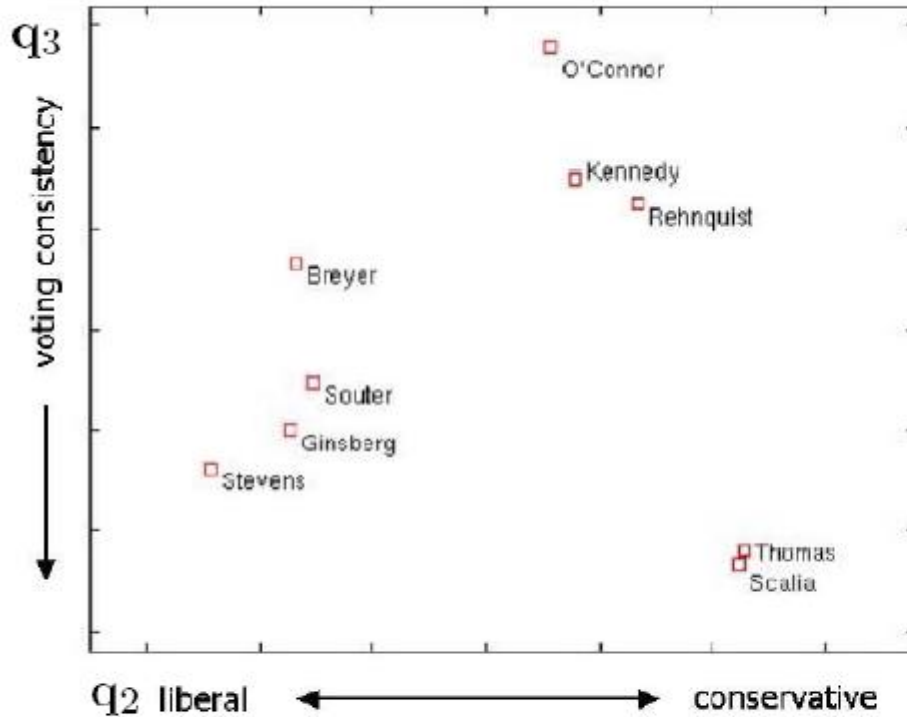
Number of times (%) two Justices voted in agreement

	Ste	Bre	Gin	Sou	O'Co	Ken	Reh	Scal	Tho
Stevens	–	62	66	63	33	36	25	14	15
Breyer	62	–	72	71	55	47	43	25	24
Ginsberg	66	72	–	78	47	49	43	28	26
Souter	63	71	78	–	55	50	44	31	29
O'Connor	33	55	47	55	–	67	71	54	54
Kennedy	36	47	49	50	67	–	77	58	59
Rehnquist	25	43	43	44	71	77	–	66	68
Scalia	14	25	28	31	54	58	66	–	79
Thomas	15	24	26	29	54	59	68	79	–

$W =$

Table 1: From the voting record of Justices 1995 Term – 2004 Term, the number of times two justices voted in agreement (in percentage). (Data source: from July 2, 2005 *New York Times*. Originally from *Legal Affairs; Harvard Law Review*)

Example: Continue



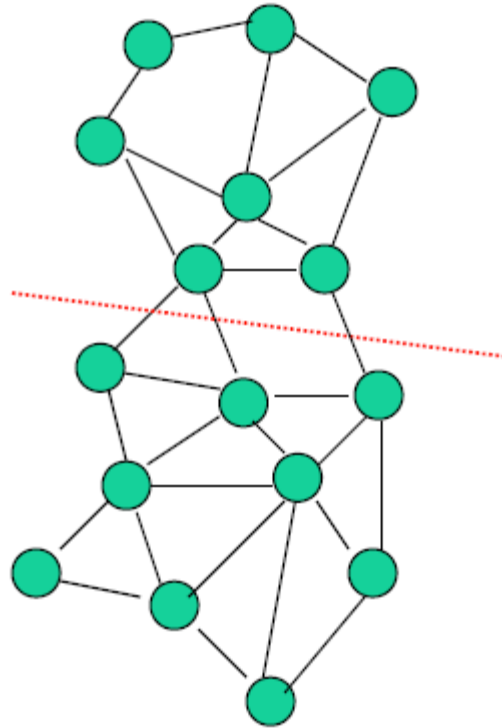
$$C = q_2 q_2^T + q_3 q_3^T$$

	Stevens	Breyer	Ginsberg	Souter	O'Connor	Kennedy	Rehnquist	Scalia	Thomas
Stevens	Green	Green	Green	Green	Red	Red	Red	Red	Red
Breyer	Green	Green	Green	Green	Green	Red	Red	Red	Red
Ginsberg	Green	Green	Green	Green	Red	Red	Red	Red	Red
Souter	Green	Green	Green	Green	Red	Red	Red	Red	Red
O'Connor	Red	Green	Red	Red	Green	Green	Green	Red	Red
Kennedy	Red	Red	Red	Red	Green	Green	Green	Red	Red
Rehnquist	Red	Red	Red	Red	Green	Green	Green	Red	Red
Scalia	Red	Red	Red	Red	Red	Red	Red	Green	Green
Thomas	Red	Red	Red	Red	Red	Red	Red	Green	Green

- Three groups in the Supreme Court:
 - Left leaning group, center-right group, right leaning group.

Spectral Graph Partition

- Min-Cut
 - Minimize the # of cut of edges



Objective Function

2-way Spectral Graph Partitioning

Partition membership indicator: $q_i = \begin{cases} 1 & \text{if } i \in A \\ -1 & \text{if } i \in B \end{cases}$

$$\begin{aligned} J = \text{CutSize} &= \frac{1}{4} \sum_{i,j} w_{ij} [q_i - q_j]^2 \\ &= \frac{1}{4} \sum_{i,j} w_{ij} [q_i^2 + q_j^2 - 2q_i q_j] = \frac{1}{2} \sum_{i,j} q_i [d_i \delta_{ij} - w_{ij}] q_j \\ &= \frac{1}{2} q^T (D - W) q \end{aligned}$$

Relax indicators q_i from discrete values to continuous values, the solution for $\min J(q)$ is given by the eigenvectors of

$$(D - W)q = \lambda q$$

(Fiedler, 1973, 1975)

(Pothen, Simon, Liou, 1990)

Algorithm

- Step 1:
 - Calculate Laplacian matrix: $L = D - W$
- Step 2:
 - Calculate the second eigenvector q
- Step 3:
 - Bisect q (e.g., 0) to get two clusters

$$(D - W)q = \lambda q$$

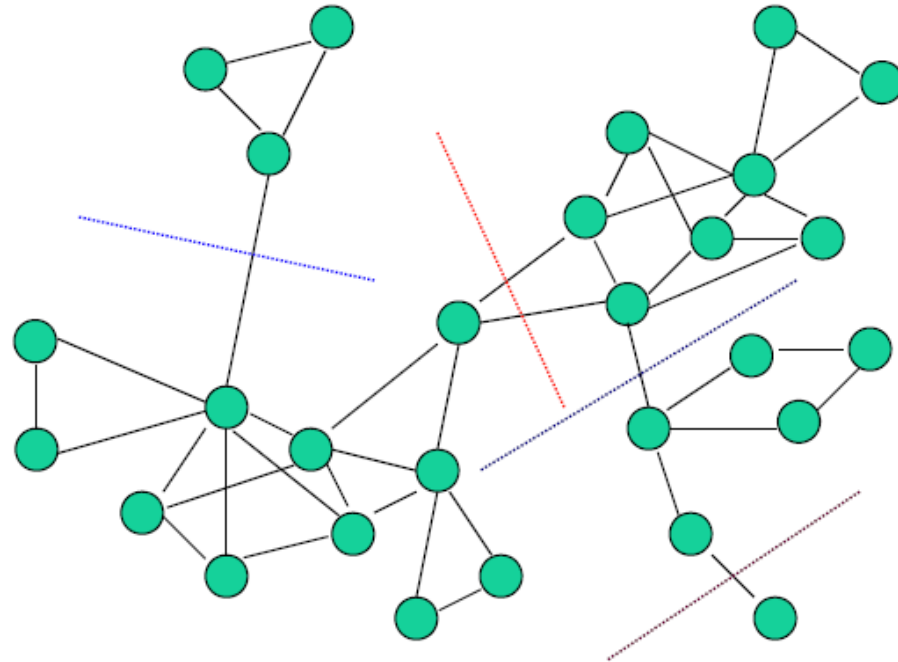
(Fiedler, 1973, 1975)

(Pothen, Simon, Liou, 1990)

*Minimum Cut with Constraints

minimize cutsizes without explicit size constraints

But where to cut ?



Need to balance sizes

*New Objective Functions

- Ratio Cut (Hangen & Kahng, 1992)

$$s(A,B) = \sum_{i \in A} \sum_{j \in B} w_{ij}$$

$$J_{Rcut}(A,B) = \frac{s(A,B)}{|A|} + \frac{s(A,B)}{|B|}$$

- Normalized Cut (Shi & Malik, 2000)

$$d_A = \sum_{i \in A} d_i$$

$$J_{Ncut}(A,B) = \frac{s(A,B)}{d_A} + \frac{s(A,B)}{d_B}$$

$$= \frac{s(A,B)}{s(A,A) + s(A,B)} + \frac{s(A,B)}{s(B,B) + s(A,B)}$$

- Min-Max-Cut (Ding et al, 2001)


$$J_{MMC}(A,B) = \frac{s(A,B)}{s(A,A)} + \frac{s(A,B)}{s(B,B)}$$

Other References

- A Tutorial on Spectral Clustering by U. Luxburg

http://www.kyb.mpg.de/fileadmin/user_upload/files/publications/attachments/Luxburg07_tutorial_4488%5B0%5D.pdf

Mining Graph/Network Data

- Introduction to Graph/Network Data
- PageRank
- Proximity Definition in Graphs
- Clustering
- Summary 

Summary

- Ranking on Graph / Network
 - PageRank
- Proximities
 - Personalized PageRank, common neighbors, Jaccard's coefficient, Adamic/Adar
- Clustering
 - Spectral clustering