

Web大数据挖掘

Mining Social-Network Graphs

Slides adapted from <http://www.mmnds.org>

Agenda

High dim. data

Locality sensitive hashing

Clustering

Dimensionality reduction

Graph data

PageRank, SimRank

Community Detection

Spam Detection

Infinite data

Filtering data streams

Web advertising

Queries on streams

Machine learning

SVM

Decision Trees

Perceptron, kNN

Apps

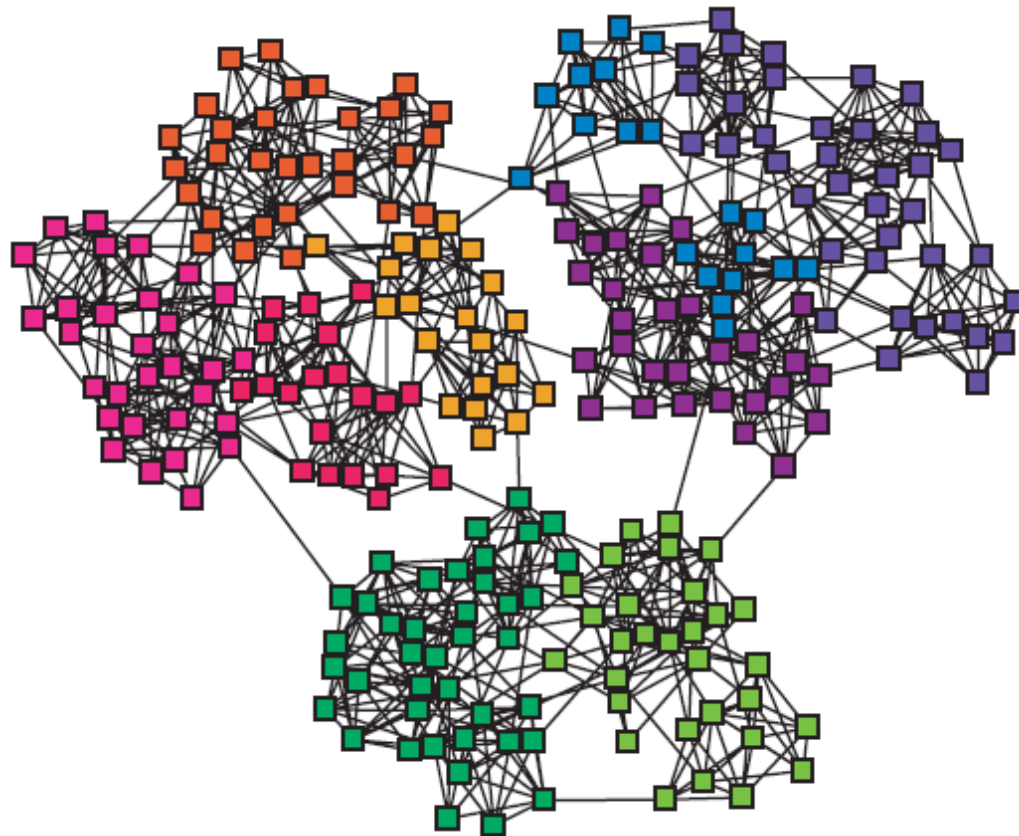
Recommender systems

Association Rules

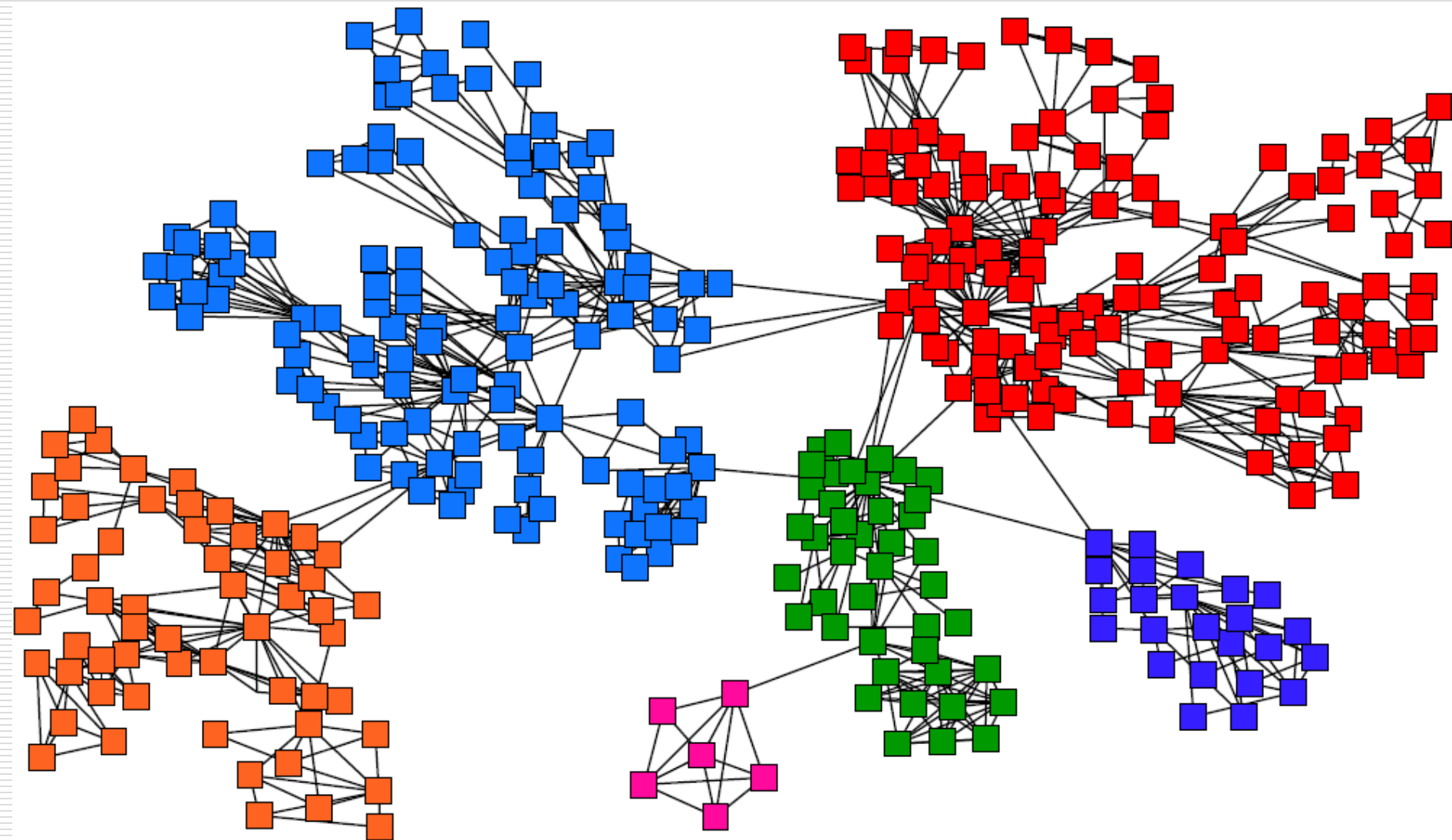
Duplicate document detection

Networks & Communities

- We often think of networks being organized into **modules, cluster, communities**:

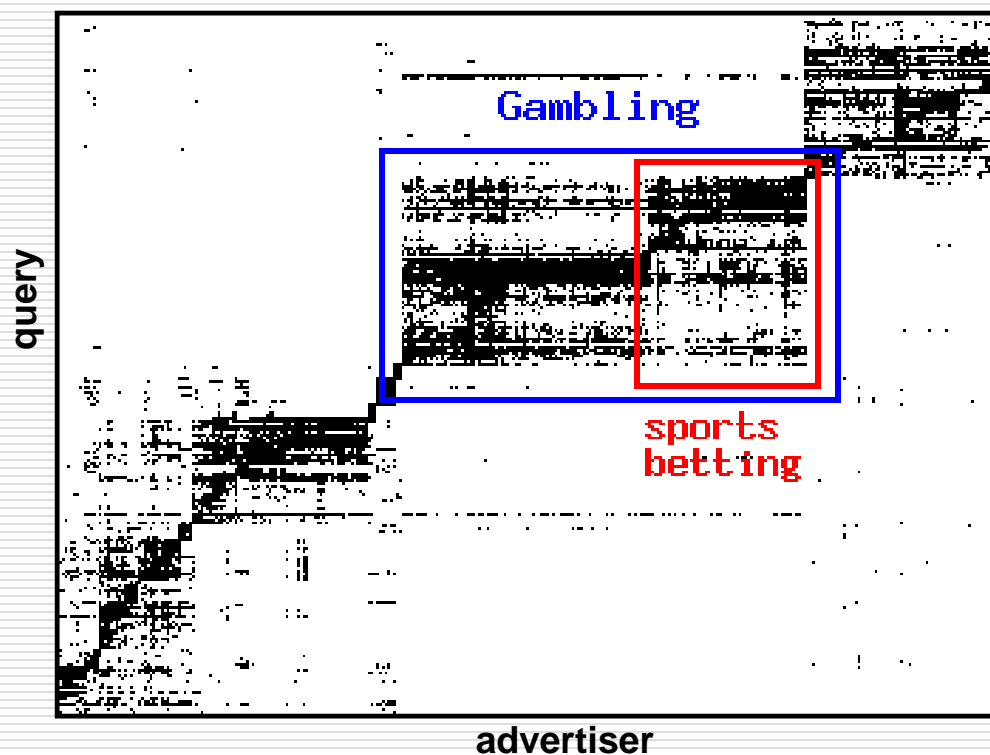


Goal: Find Densely Linked Clusters



Micro-Markets in Sponsored Search

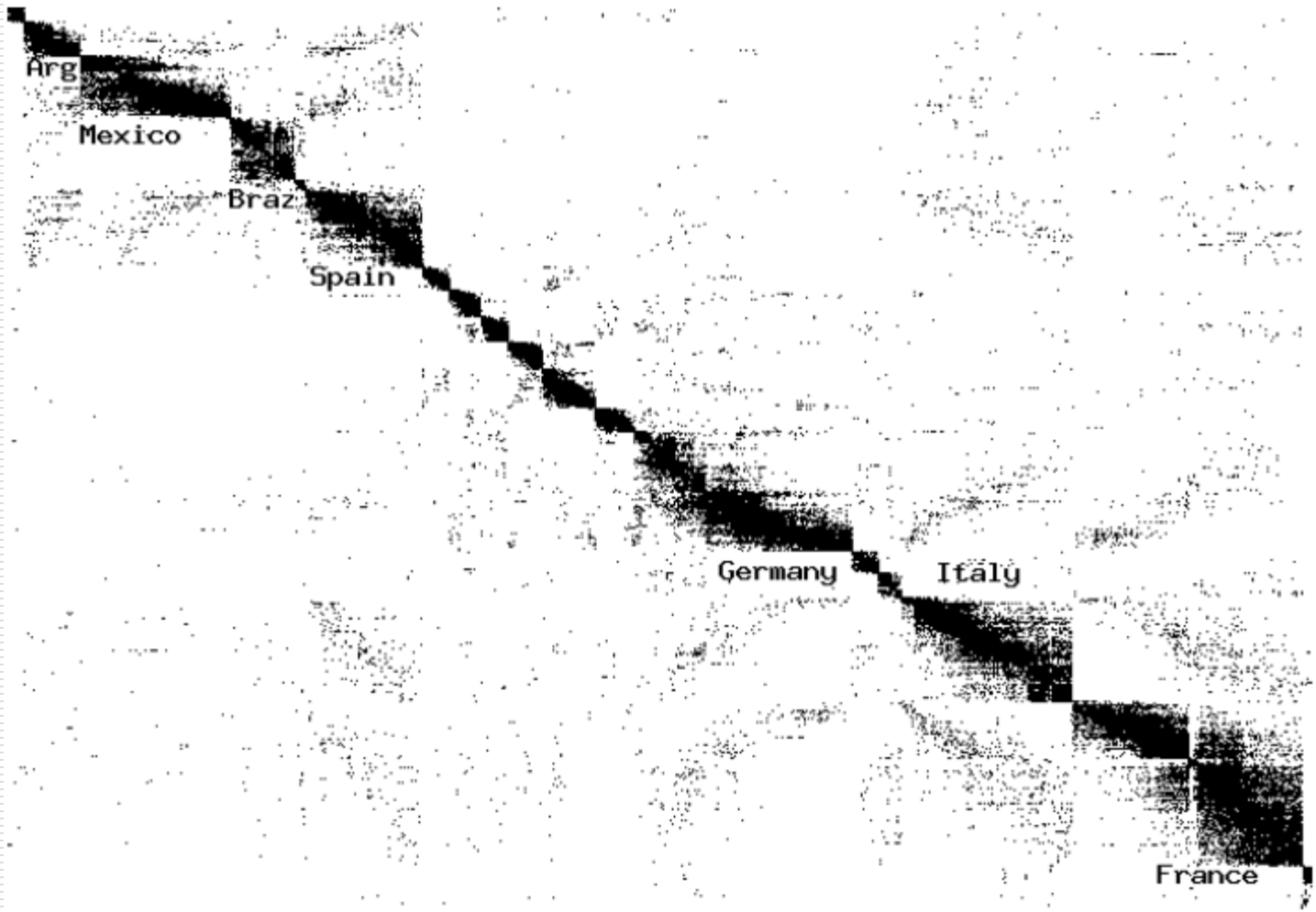
- Find micro-markets by partitioning the query-to-advertiser graph:



[Andersen, Lang: Communities from seed sets, 2006]

Movies and Actors

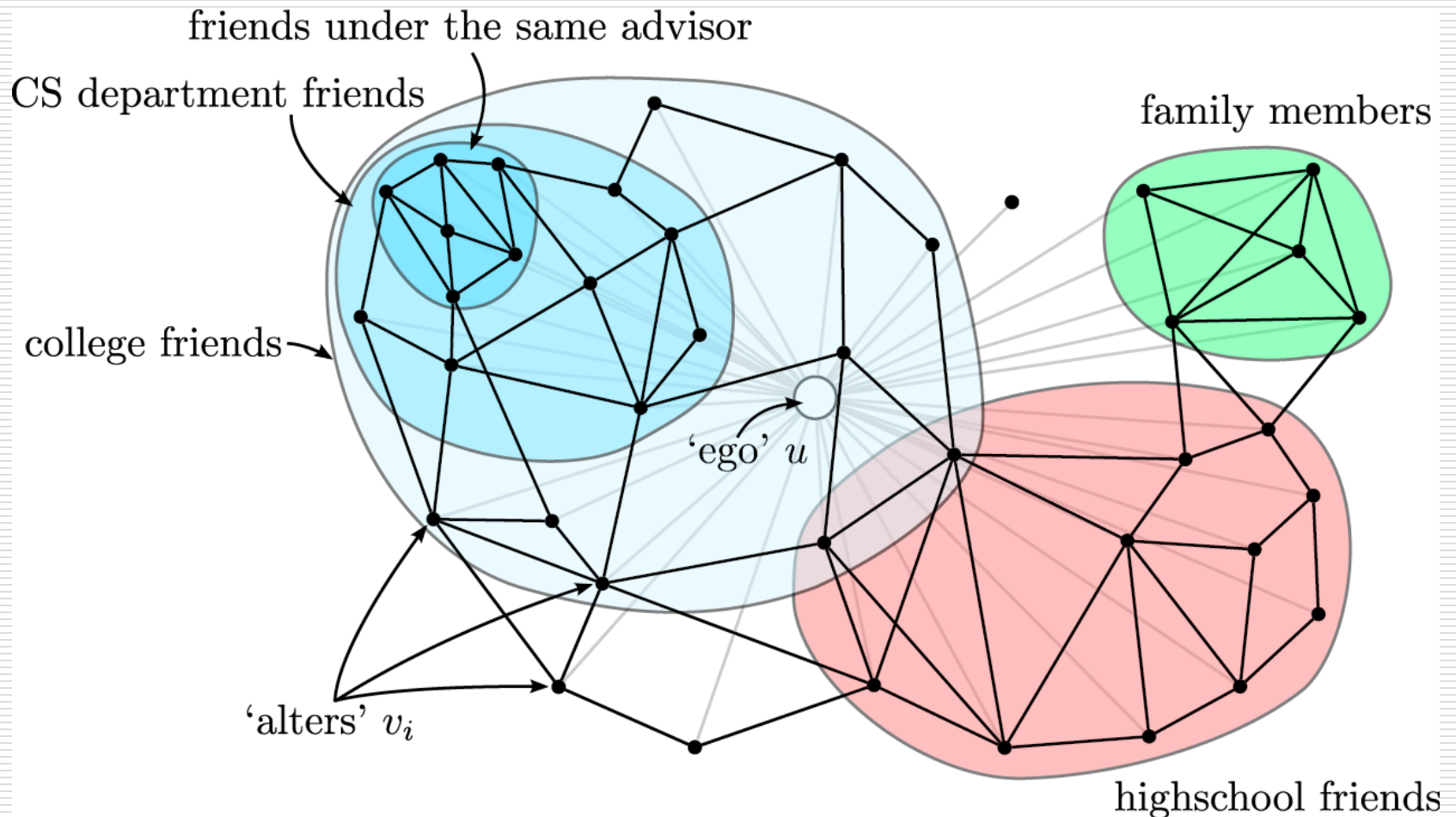
❑ Clusters in Movies-to-Actors graph:



[Andersen, Lang: Communities from seed sets, 2006]

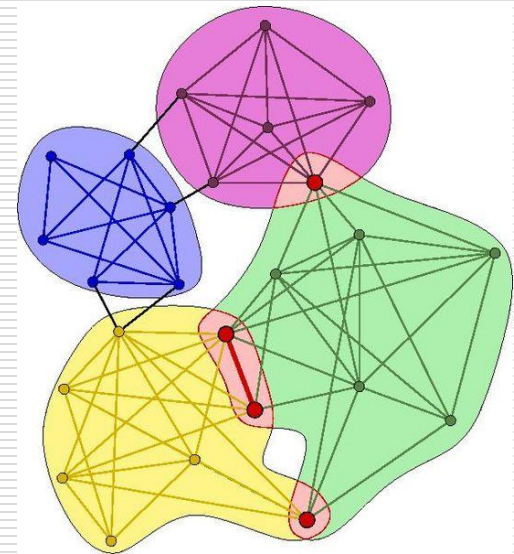
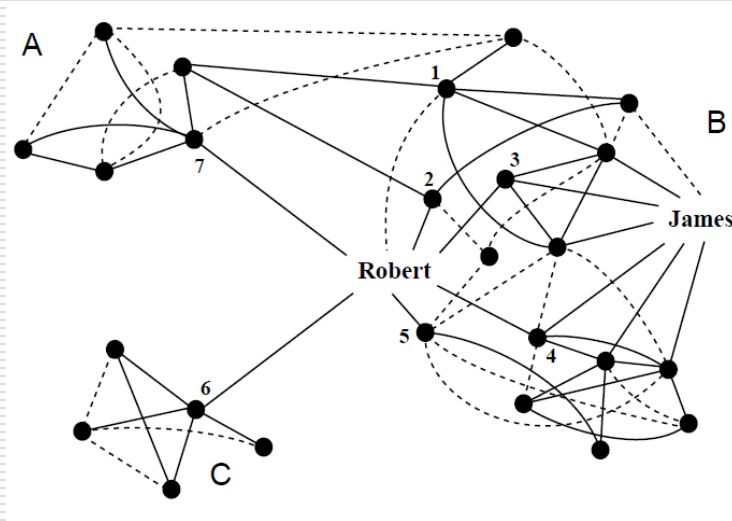
Twitter & Facebook

□ Discovering social circles, circles of trust:



[McAuley, Leskovec: Discovering social circles in ego networks, 2012]

How to find communities?

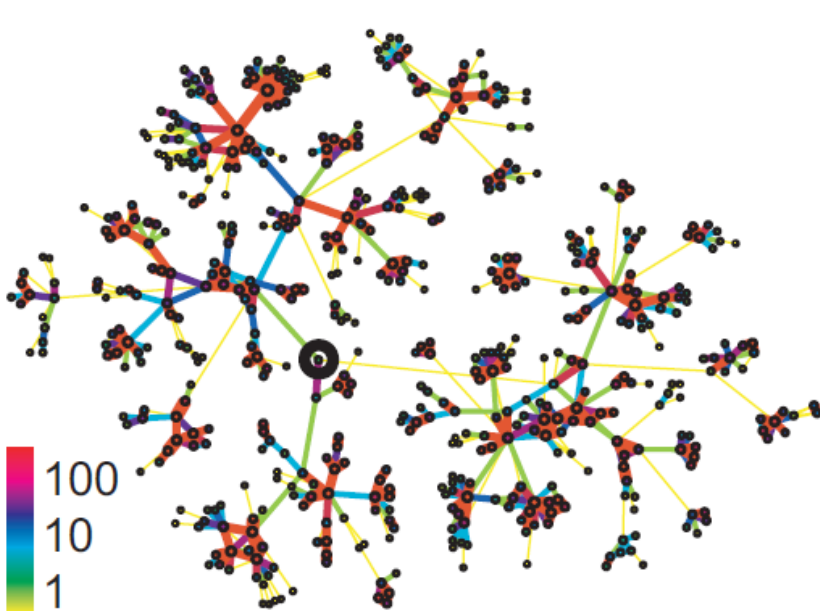
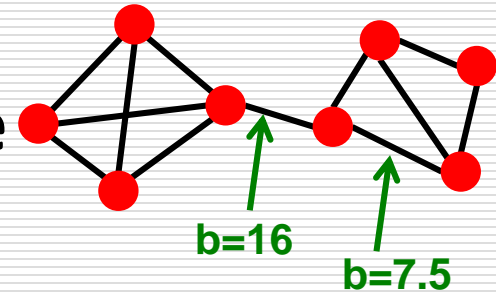


We will work with undirected (unweighted) networks

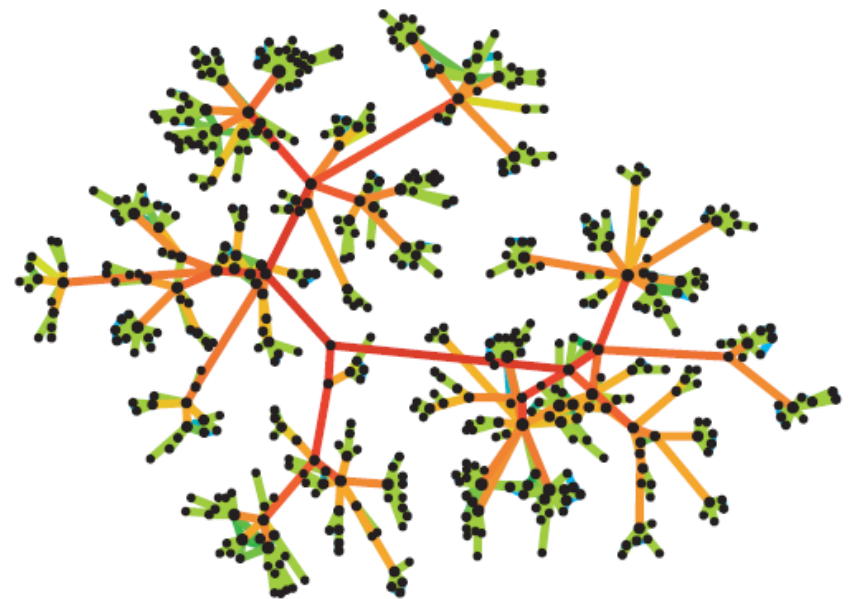
Method 1: Strength of Weak Ties

□ **Edge betweenness:** Number of shortest paths passing over the edge

□ **Intuition:**



Edge strengths (call volume)
in a real network



Edge betweenness
in a real network

Method 1: Girvan-Newman

- Divisive hierarchical clustering based on the notion of edge **betweenness**:

Number of shortest paths passing through the edge

- **Girvan-Newman Algorithm:**

- Undirected unweighted networks

- **Repeat until no edges are left:**

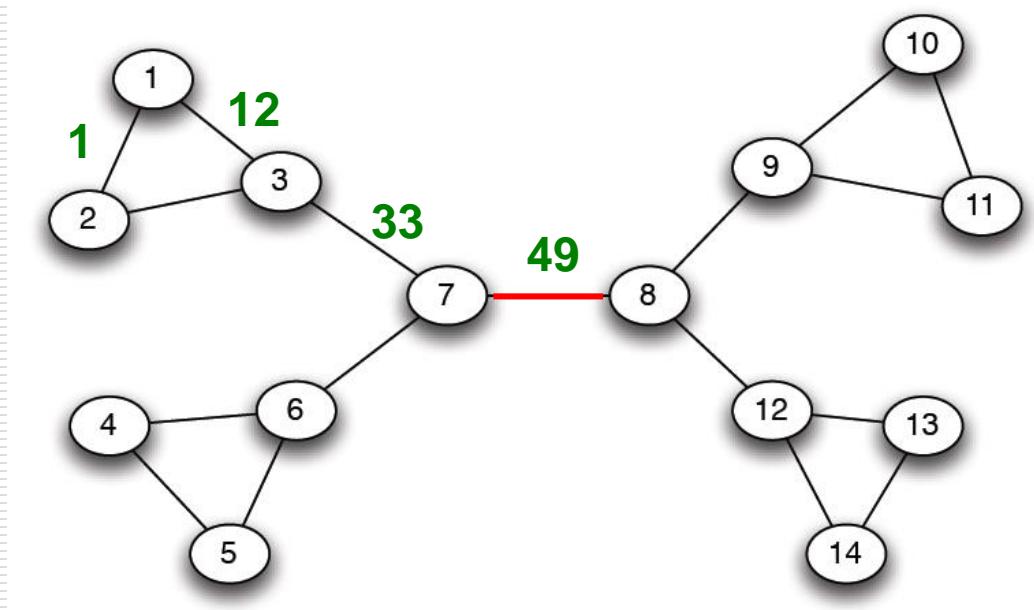
- Calculate betweenness of edges

- Remove edges with highest betweenness

- Connected components are communities

- Gives a hierarchical decomposition of the network

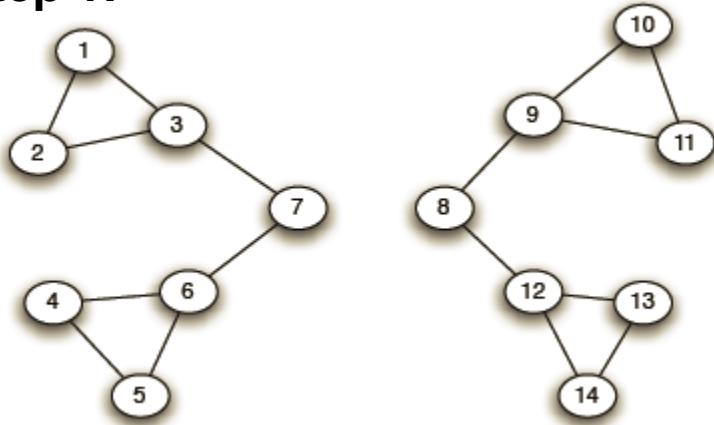
Girvan-Newman: Example



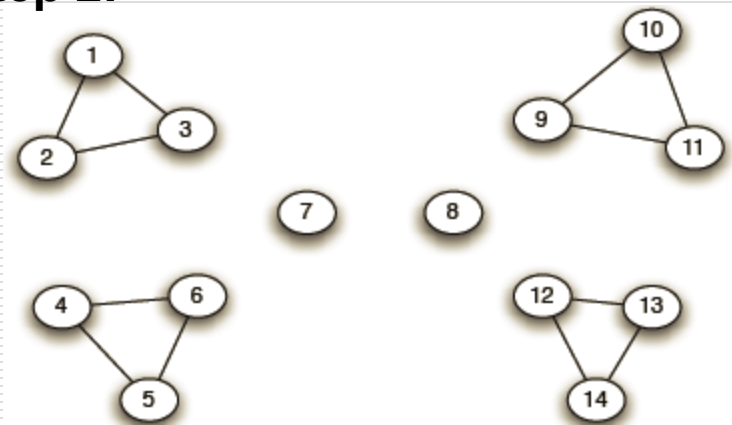
**Need to re-compute
betweenness at every
step**

Girvan-Newman: Example

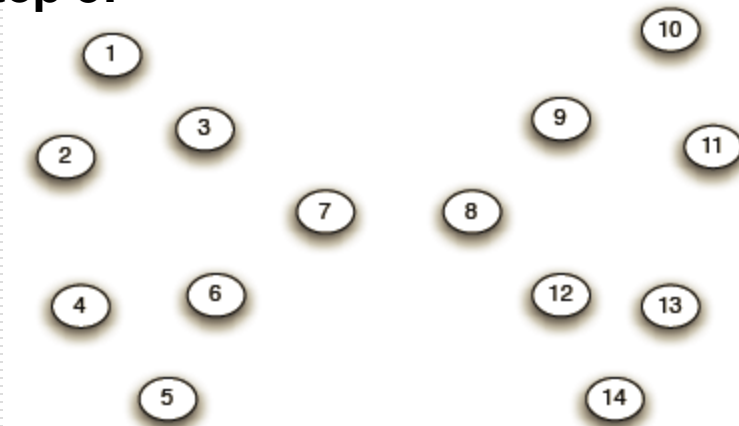
Step 1:



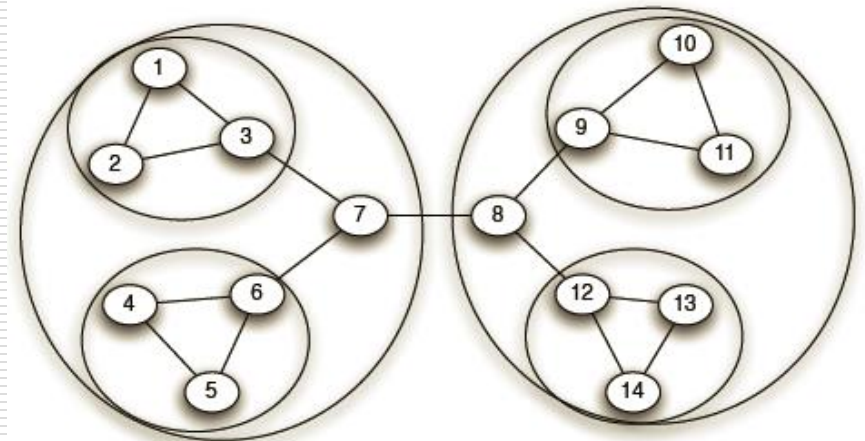
Step 2:



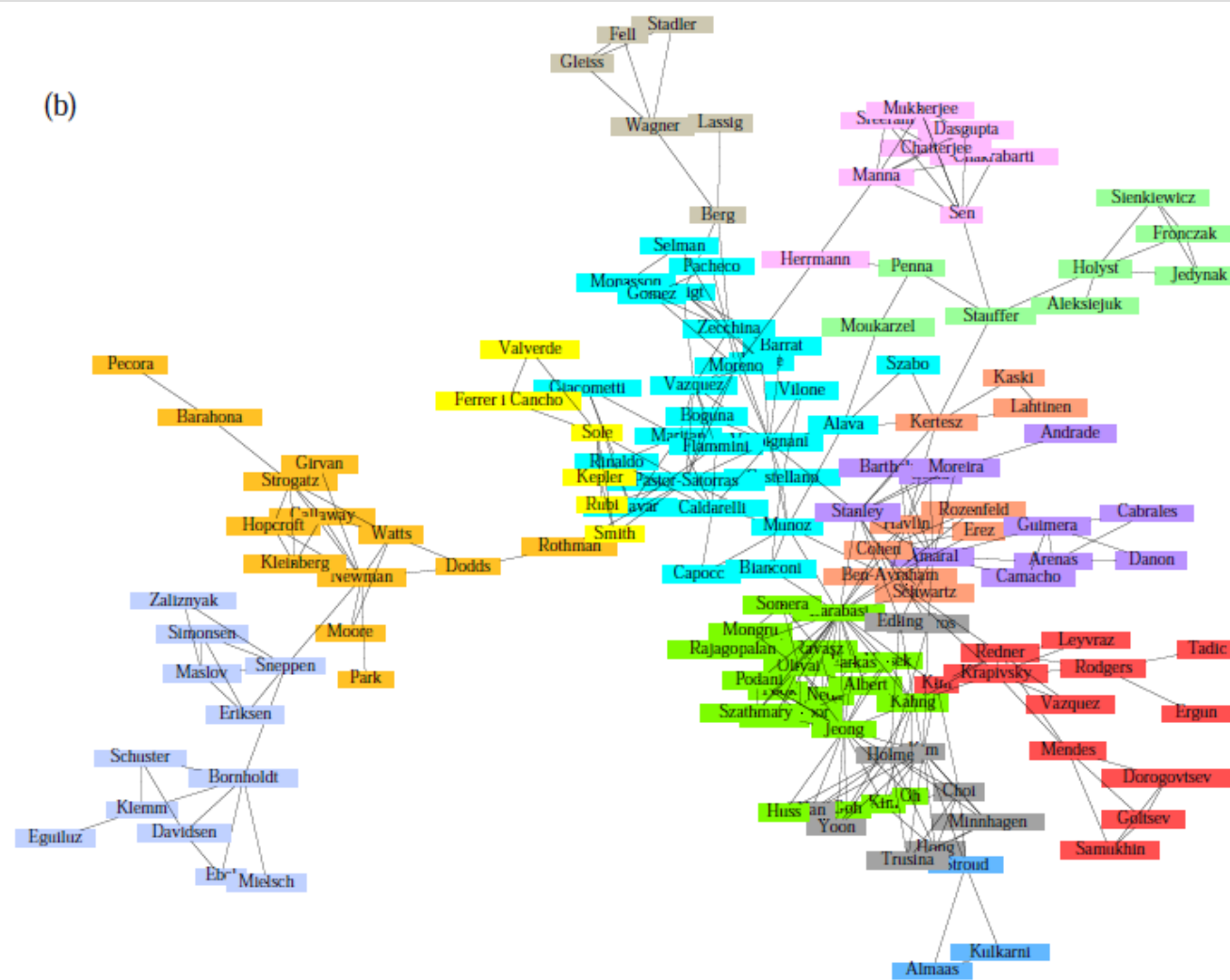
Step 3:



Hierarchical network decomposition:



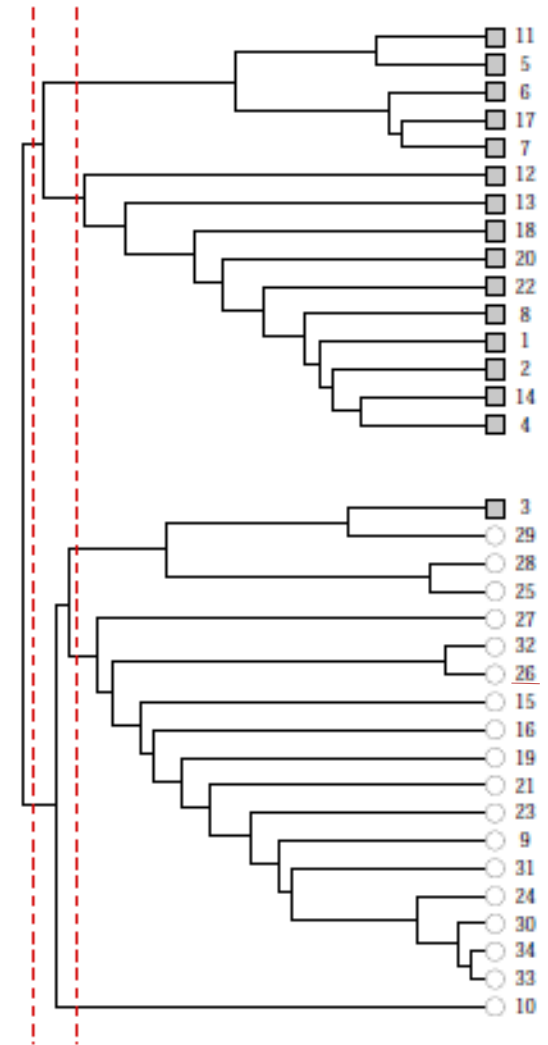
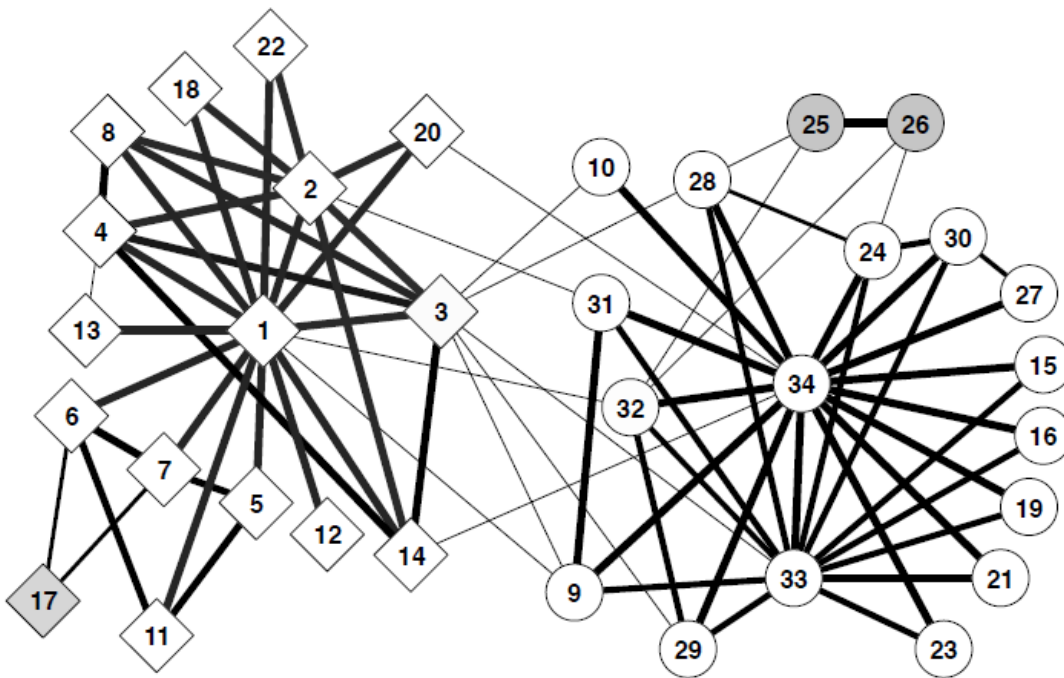
Girvan-Newman: Results



Communities in physics collaborations

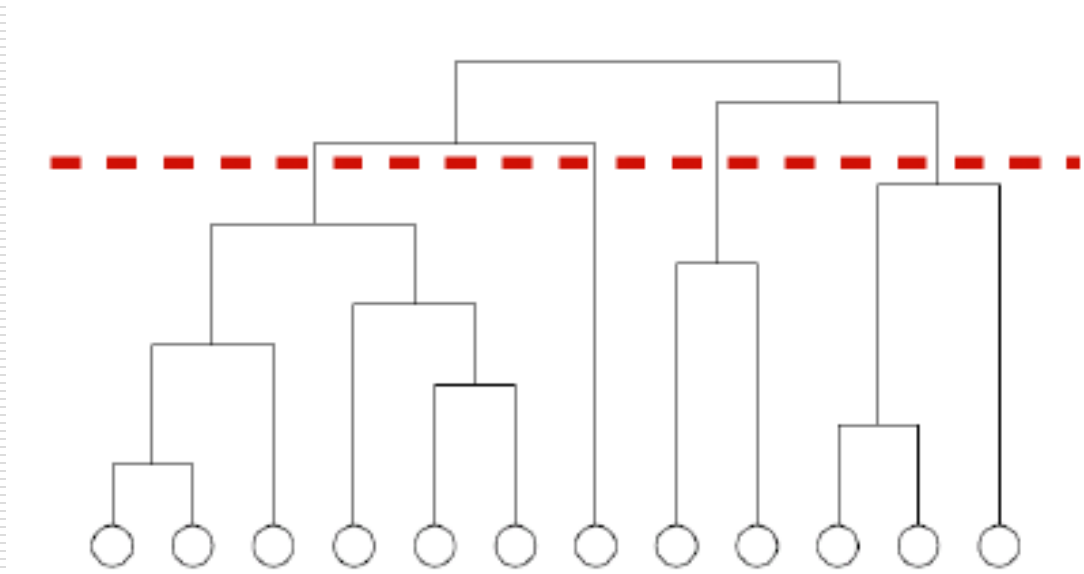
Girvan-Newman: Results

□ Zachary's Karate club: Hierarchical decomposition



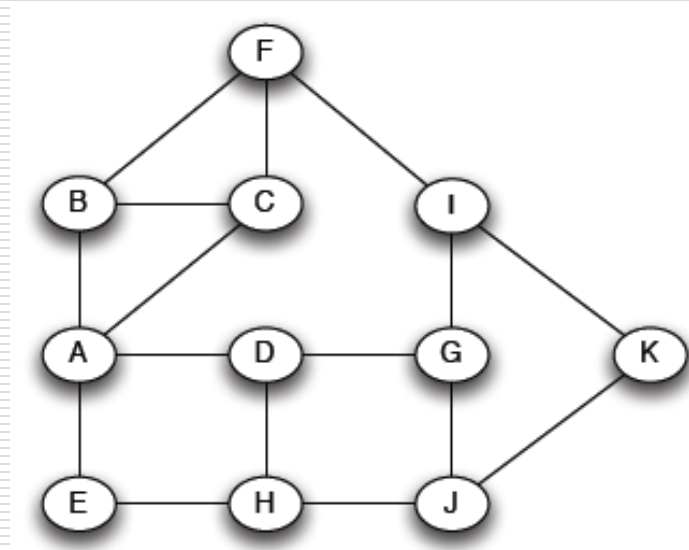
We need to resolve 2 questions

1. How to compute betweenness?
2. How to select the number of clusters?

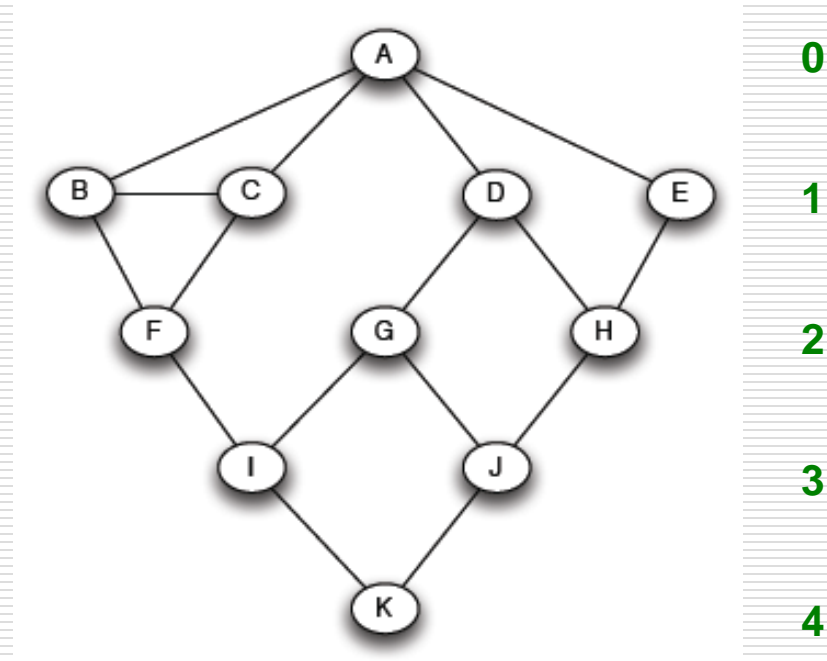


How to Compute Betweenness?

□ Want to compute betweenness of paths starting at node *A*

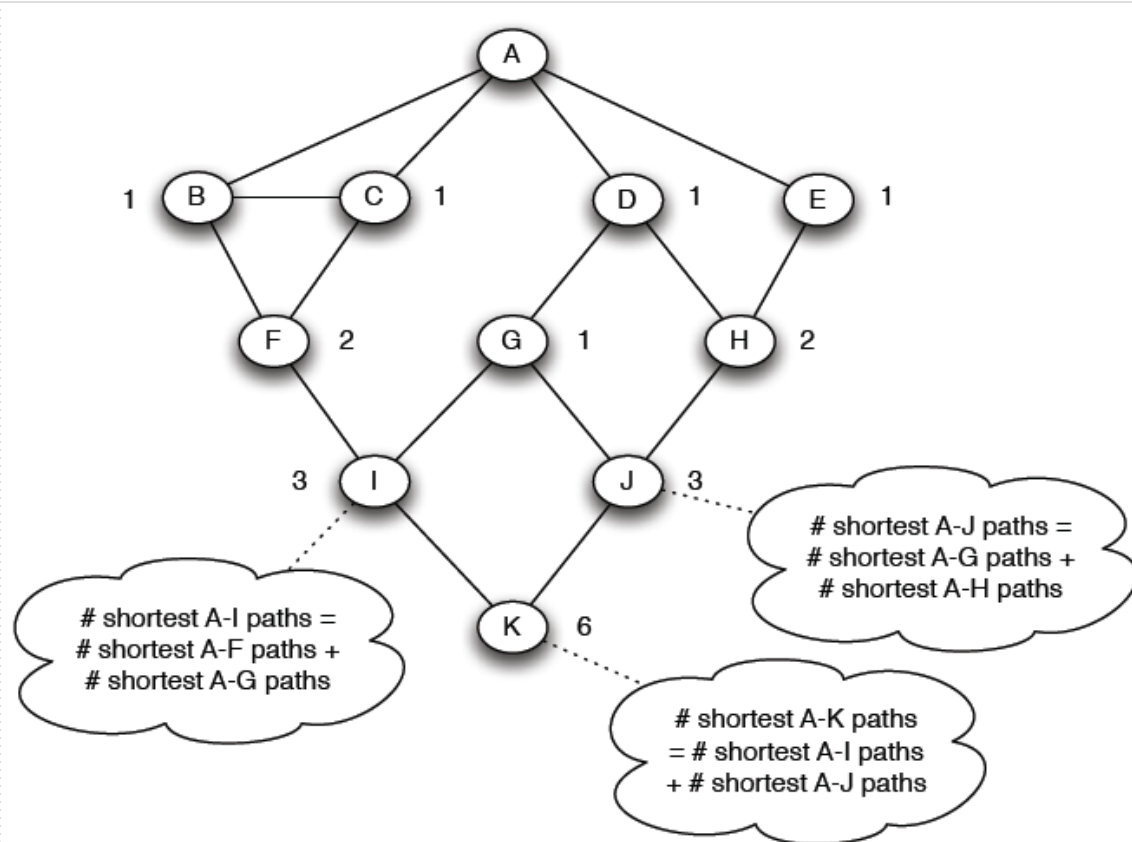


□ Breath first search starting from *A*:



How to Compute Betweenness?

- Count the number of shortest paths from *A* to all other nodes of the network:



How to Compute Betweenness?

□ Compute betweenness by working up the tree:

If there are multiple paths count them fractionally

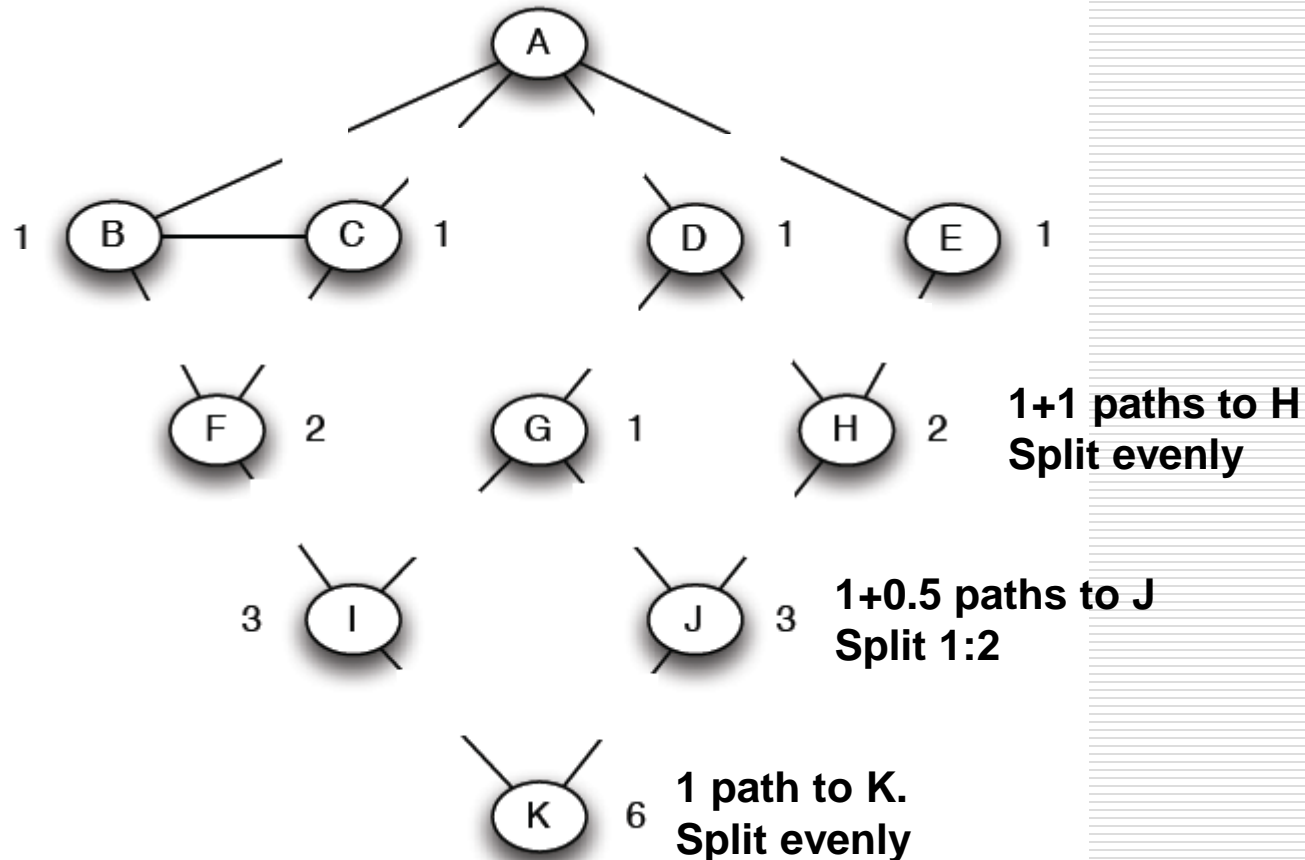
The algorithm:

- Add edge flows:

- node flow = $1 + \sum \text{child edges}$

- split the flow up based on the parent value

- Repeat the BFS procedure for each starting node U



How to Compute Betweenness

□ Compute betweenness by working up the tree:

If there are multiple paths count them fractionally

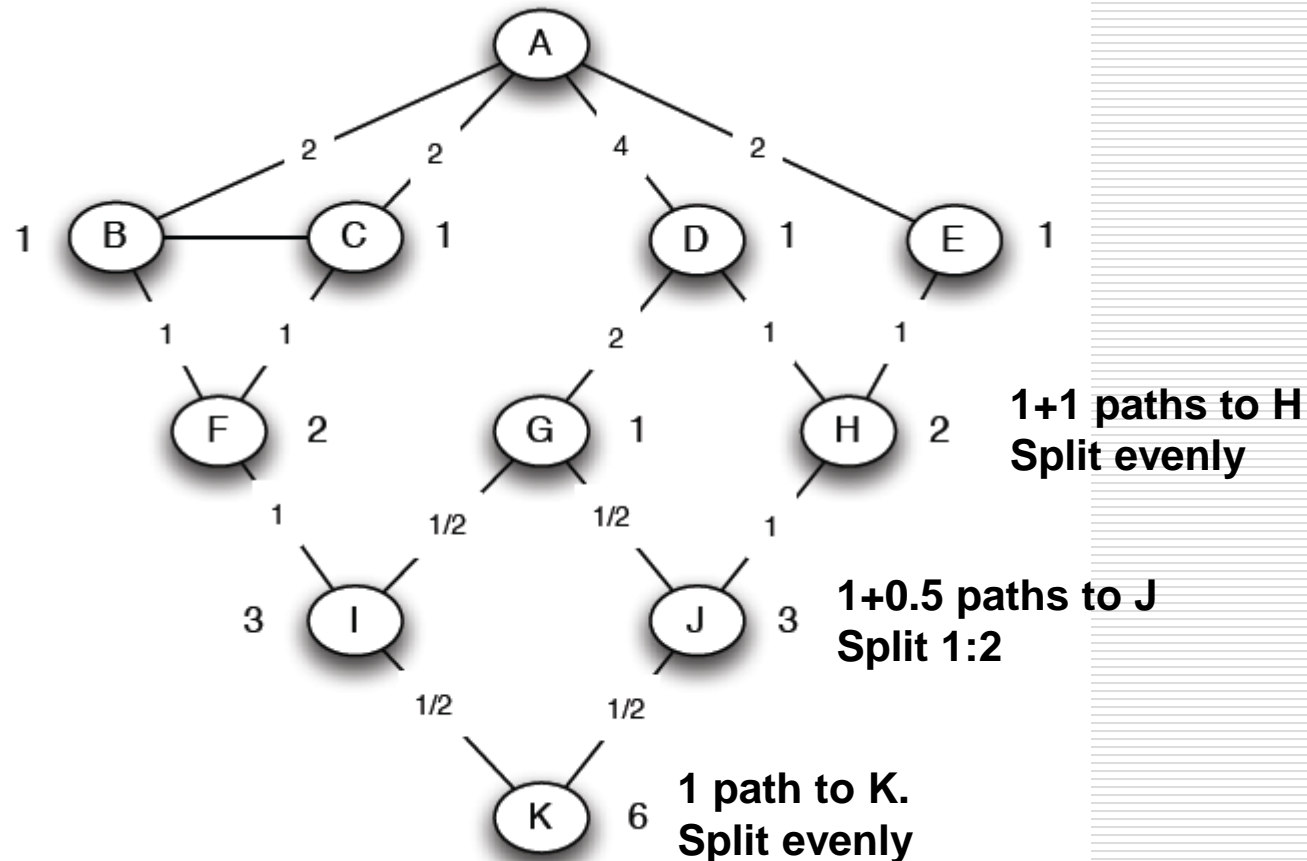
The algorithm:

- Add edge flows:

- node flow = $1 + \sum \text{child edges}$

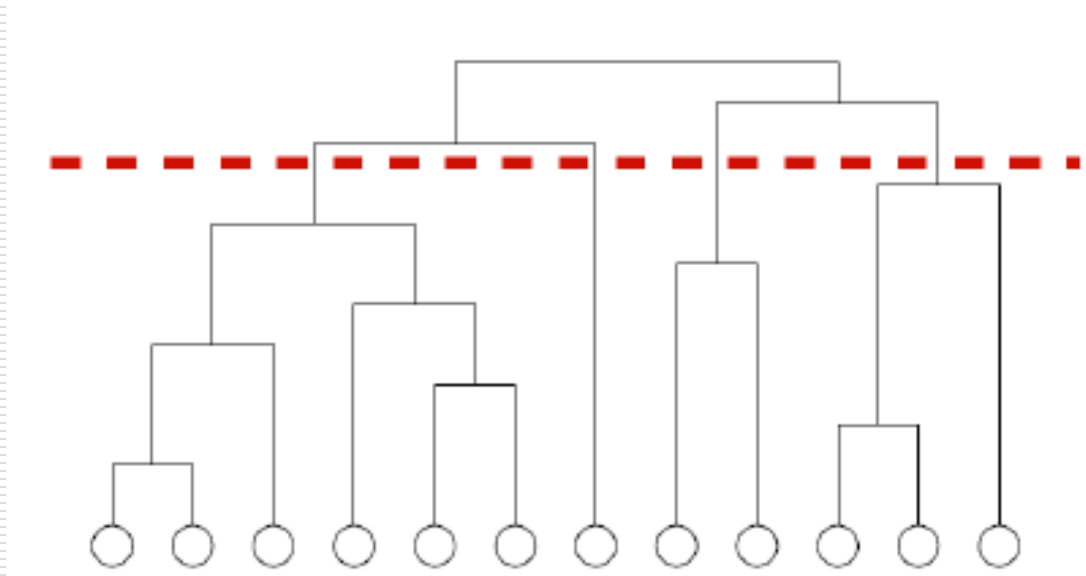
- split the flow up based on the parent value

- Repeat the BFS procedure for each starting node U



We need to resolve 2 questions

1. How to compute betweenness?
2. How to select the number of clusters?



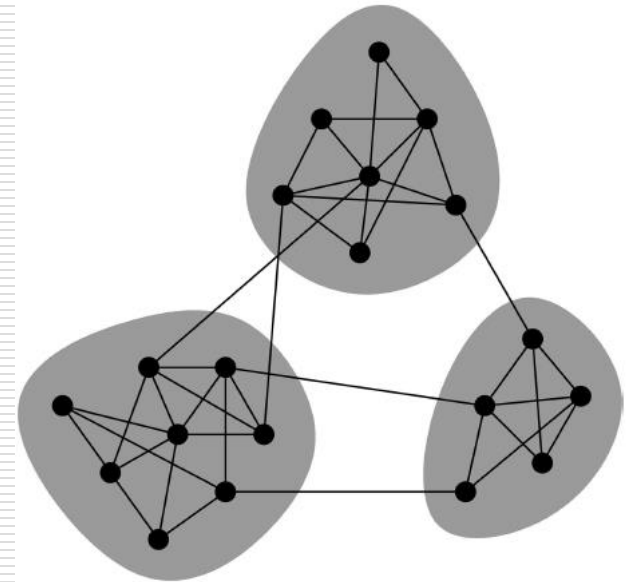
Network Communities

□ **Communities:** sets of tightly connected nodes

□ Define: **Modularity Q**

- A measure of how well a network is partitioned into communities
- Given a partitioning of the network into groups $s \in S$:

$$Q \propto \sum_{s \in S} [(\# \text{ edges within group } s) - (\text{expected } \# \text{ edges within group } s)]$$



Configuration Model

□ Given real G on n nodes and m edges, construct rewired network G'

■ Same degree distribution but random connections

■ Consider G' as a **multigraph**

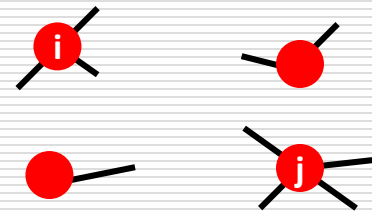
■ The expected number of edges between nodes

i and j of degrees k_i and k_j equals to: $k_i \cdot \frac{k_j}{2m} = \frac{k_i k_j}{2m}$

□ The expected number of edges in (multigraph) G' :

■
$$= \frac{1}{2} \sum_{i \in N} \sum_{j \in N} \frac{k_i k_j}{2m} = \frac{1}{2} \cdot \frac{1}{2m} \sum_{i \in N} k_i \left(\sum_{j \in N} k_j \right) =$$

■
$$= \frac{1}{4m} 2m \cdot 2m = m$$



Note:

$$\sum_{u \in N} k_u = 2m$$

Modularity

□ Modularity of partitioning S of graph G :

- $Q \propto \sum_{s \in S} [(\# \text{ edges within group } s) - (\text{expected } \# \text{ edges within group } s)]$

- $$Q(G, S) = \frac{1}{2m} \sum_{s \in S} \sum_{i \in s} \sum_{j \in s} \left(A_{ij} - \frac{k_i k_j}{2m} \right)$$

Normalizing cost.: $-1 < Q < 1$

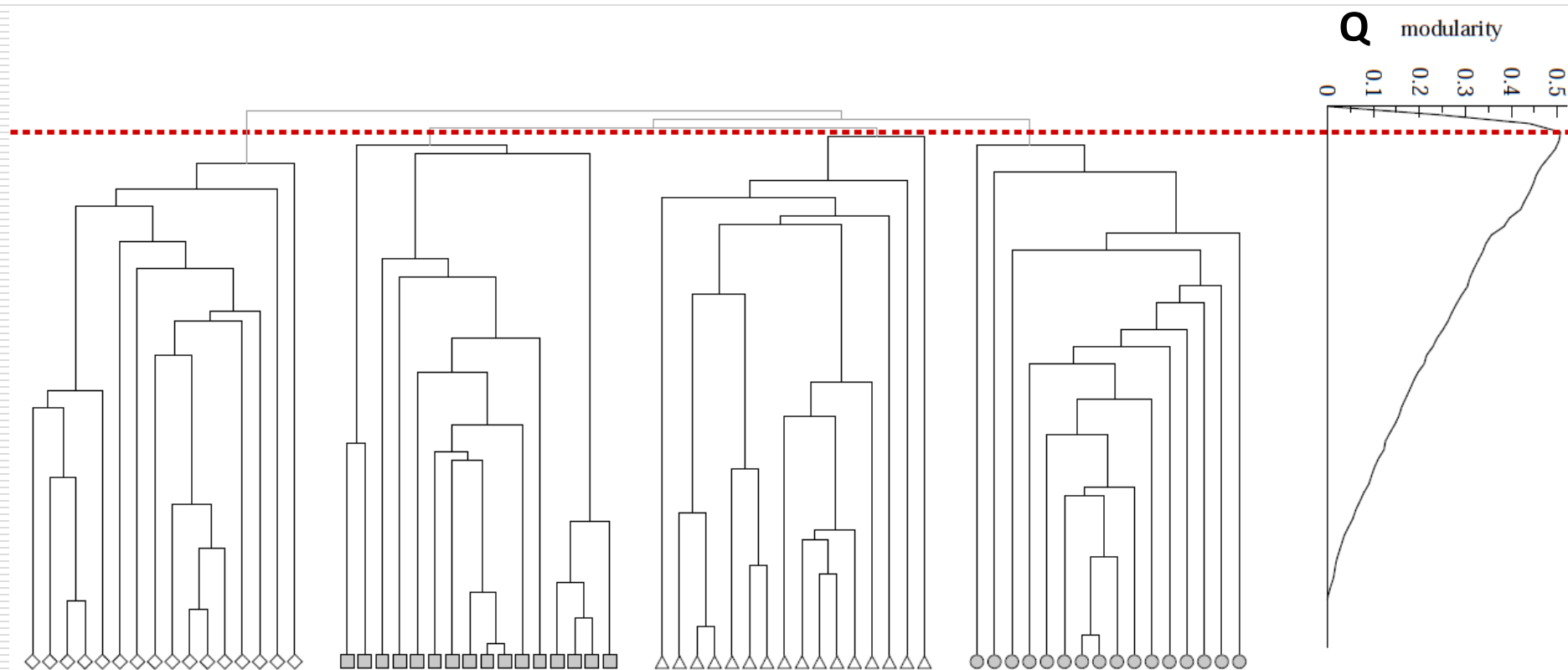
$A_{ij} = 1$ if $i \rightarrow j$,
0 else

□ Modularity values take range $[-1, 1]$

- It is positive if the number of edges within groups exceeds the expected number
- $0.3-0.7 < Q$ means significant community structure

Modularity: Number of clusters

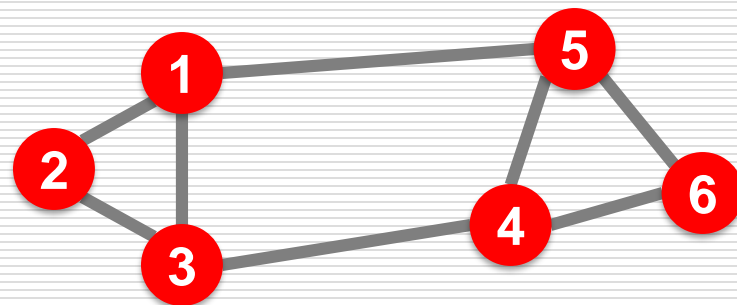
- Modularity is useful for selecting the number of clusters:



Spectral Clustering

Graph Partitioning

□ Undirected graph $G(V, E)$:



□ Bi-partitioning task:

- Divide vertices into two disjoint groups A, B



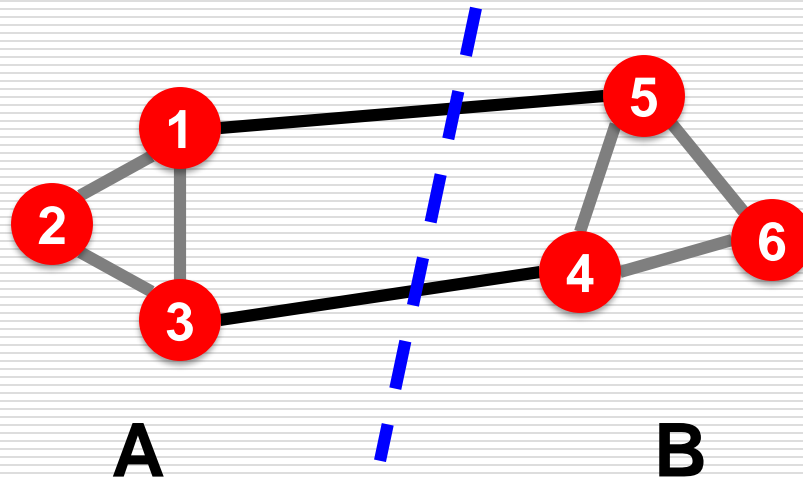
□ Questions:

- How can we define a “good” partition of G ?
- How can we efficiently identify such a partition?

Graph Partitioning

□ What makes a good partition?

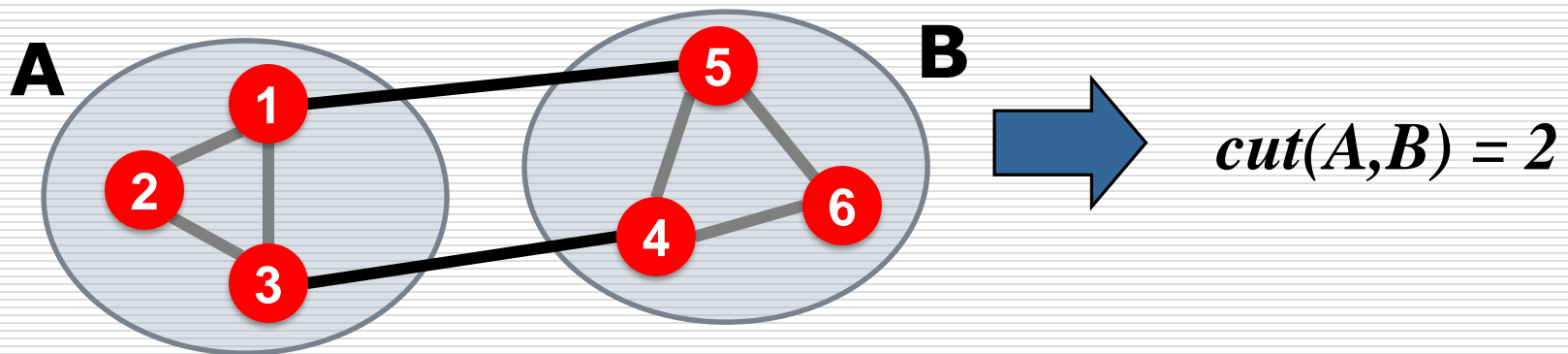
- Maximize the number of within-group connections
- Minimize the number of between-group connections



Graph Cuts

- Express partitioning objectives as a function of the “edge cut” of the partition
- **Cut:** Set of edges with only one vertex in a group:

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



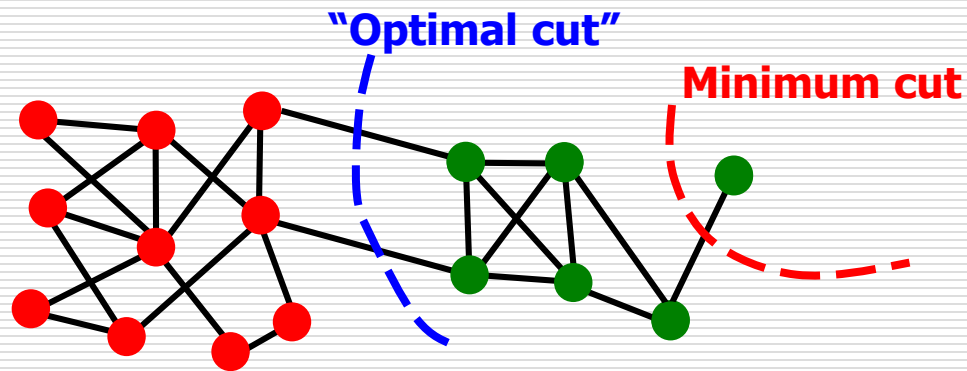
Graph Cut Criterion

□ Criterion: **Minimum-cut**

- Minimize weight of connections between groups

$$\arg \min_{A,B} \text{cut}(A,B)$$

□ Degenerate case:



□ Problem:

- Only considers external cluster connections
- Does not consider internal cluster connectivity

Graph Cut Criteria

□ Criterion: **Normalized-cut** [Shi-Malik, '97]

- Connectivity between groups relative to the density of each group

$$ncut(A, B) = \frac{cut(A, B)}{vol(A)} + \frac{cut(A, B)}{vol(B)}$$

$vol(A)$: total weight of the edges with at least one endpoint in A : $vol(A) = \sum_{i \in A} k_i$

■ **Why use this criterion?**

- Produces more balanced partitions

□ **How do we efficiently find a good partition?**

- **Problem:** Computing optimal cut is NP-hard

Spectral Graph Partitioning

- **A**: adjacency matrix of undirected **G**
 - $A_{ij} = 1$ if (i, j) is an edge, else 0
- **x** is a vector in \Re^n with components (x_1, \dots, x_n)
 - Think of it as a label/value of each node of **G**
- **What is the meaning of $A \cdot x$?**

$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad y_i = \sum_{j=1}^n A_{ij} x_j = \sum_{(i,j) \in E} x_j$$

- **Entry y_i is a sum of labels x_j of neighbors of i**

What is the meaning of Ax ?

- **j^{th} coordinate of $A \cdot x$:**
- Sum of the x -values of neighbors of j
 - Make this a new value at node j
- $$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$
- $$A \cdot x = \lambda \cdot x$$

□ **Spectral Graph Theory:**

- Analyze the “spectrum” of matrix representing G
- **Spectrum:** Eigenvectors x_i of a graph, ordered by the magnitude (strength) of their corresponding eigenvalues λ_i :

$$\Lambda = \{\lambda_1, \lambda_2, \dots, \lambda_n\}$$
$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

Example: d-regular graph

- Suppose all nodes in G have degree d and G is connected
- What are some eigenvalues/vectors of G ?

$A \cdot x = \lambda \cdot x$ What is λ ? What x ?

■ Let's try: $x = (1, 1, \dots, 1)$

■ Then: $A \cdot x = (d, d, \dots, d) = \lambda \cdot x$. So: $\lambda = d$

■ We found eigenpair of G : $x = (1, 1, \dots, 1), \lambda = d$

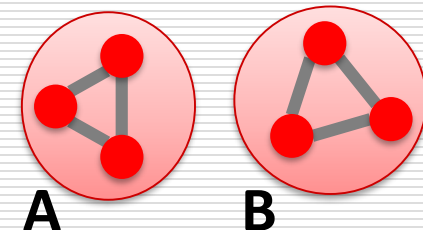
Remember the meaning of $y = A \cdot x$:

$$y_j = \sum_{i=1}^n A_{ij} x_i = \sum_{(j,i) \in E} x_i$$

Example: Graph on 2 components

□ What if G is not connected?

- G has 2 components, each d -regular



□ What are some eigenvectors?

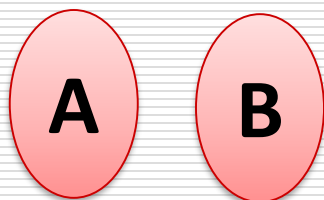
- x = Put all 1s on A and 0s on B or vice versa

- $x' = (\underbrace{1, \dots, 1}_{|A|}, \underbrace{0, \dots, 0}_{|B|})$ then $A \cdot x' = (d, \dots, d, 0, \dots, 0)$

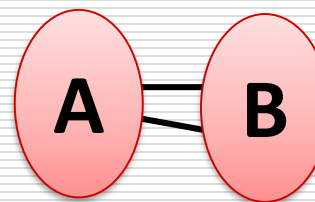
- $x'' = (0, \dots, 0, \underbrace{1, \dots, 1}_{|B|}, \underbrace{0, \dots, 0}_{|A|})$ then $A \cdot x'' = (0, \dots, 0, d, \dots, d)$

- And so in both cases the corresponding $\lambda = d$

□ A bit of intuition:



$$\lambda_n = \lambda_{n-1}$$



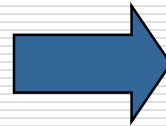
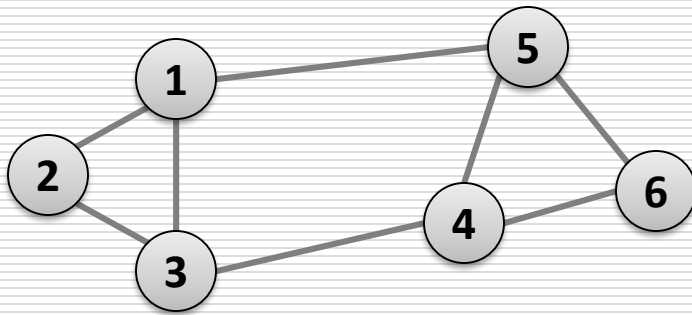
$$\lambda_n - \lambda_{n-1} \approx 0$$

2nd largest eigval.
 λ_{n-1} now has
value very close
to λ_n

Matrix Representations

□ Adjacency matrix (A):

- $n \times n$ matrix
- $A=[a_{ij}]$, $a_{ij}=1$ if edge between node i and j



	1	2	3	4	5	6
1	0	1	1	0	1	0
2	1	0	1	0	0	0
3	1	1	0	1	0	0
4	0	0	1	0	1	1
5	1	0	0	1	0	1
6	0	0	0	1	1	0

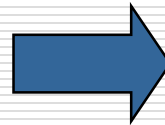
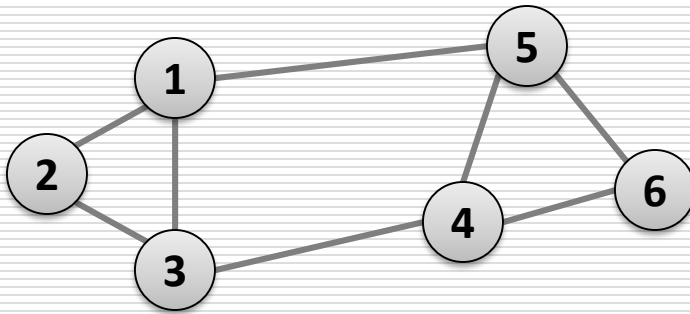
□ Important properties:

- Symmetric matrix
- Eigenvectors are real and orthogonal

Matrix Representations

□ Degree matrix (D):

- $n \times n$ diagonal matrix
- $D=[d_{ij}]$, d_{ij} = degree of node i

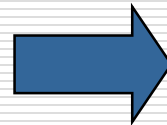
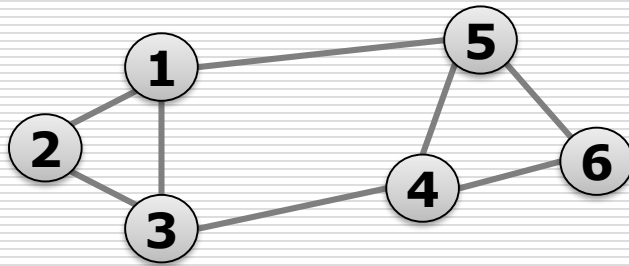


	1	2	3	4	5	6
1	3	0	0	0	0	0
2	0	2	0	0	0	0
3	0	0	3	0	0	0
4	0	0	0	3	0	0
5	0	0	0	0	3	0
6	0	0	0	0	0	2

Matrix Representations

□ Laplacian matrix (L):

- $n \times n$ symmetric matrix



	1	2	3	4	5	6
1	3	-1	-1	0	-1	0
2	-1	2	-1	0	0	0
3	-1	-1	3	-1	0	0
4	0	0	-1	3	-1	-1
5	-1	0	0	-1	3	-1
6	0	0	0	-1	-1	2

$$L = D - A$$

□ What is trivial eigenpair?

- $x = (1, \dots, 1)$ then $L \cdot x = 0$ and so $\lambda = \lambda_1 = 0$

□ Important properties:

- Eigenvalues are non-negative real numbers
- Eigenvectors are real and orthogonal

So far...

☐ How to define a “good” partition of a graph?

- Minimize a given graph cut criterion

☐ How to efficiently identify such a partition?

- Approximate using information provided by the eigenvalues and eigenvectors of a graph

☐ Spectral Clustering

Spectral Clustering Algorithms

□ Three basic stages:

■ 1) Pre-processing

- Construct a matrix representation of the graph

■ 2) Decomposition

- Compute eigenvalues and eigenvectors of the matrix
- Map each point to a lower-dimensional representation based on one or more eigenvectors

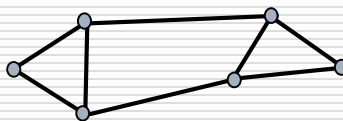
■ 3) Grouping

- Assign points to two or more clusters, based on the new representation

Spectral Partitioning Algorithm

1) Pre-processing:

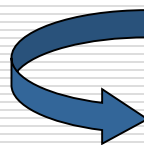
- Build Laplacian matrix L of the graph



	1	2	3	4	5	6
1	3	-1	-1	0	-1	0
2	-1	2	-1	0	0	0
3	-1	-1	3	-1	0	0
4	0	0	-1	3	-1	-1
5	-1	0	0	-1	3	-1
6	0	0	0	-1	-1	2

2) Decomposition:

- Find eigenvalues λ and eigenvectors x of the matrix L
- Map vertices to corresponding components of λ_2

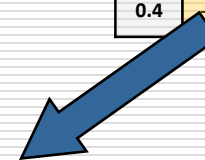


$\lambda =$

0.0
1.0
3.0
3.0
4.0
5.0

$x =$

0.4	0.3	-0.5	-0.2	-0.4	-0.5
0.4	0.6	0.4	-0.4	0.4	0.0
0.4	0.3	0.1	0.6	-0.4	0.5
0.4	-0.3	0.1	0.6	0.4	-0.5
0.4	-0.3	-0.5	-0.2	0.4	0.5
0.4	-0.6	0.4	-0.4	-0.4	0.0



1	0.3
2	0.6
3	0.3
4	-0.3
5	-0.3
6	-0.6

How do we now find the clusters?

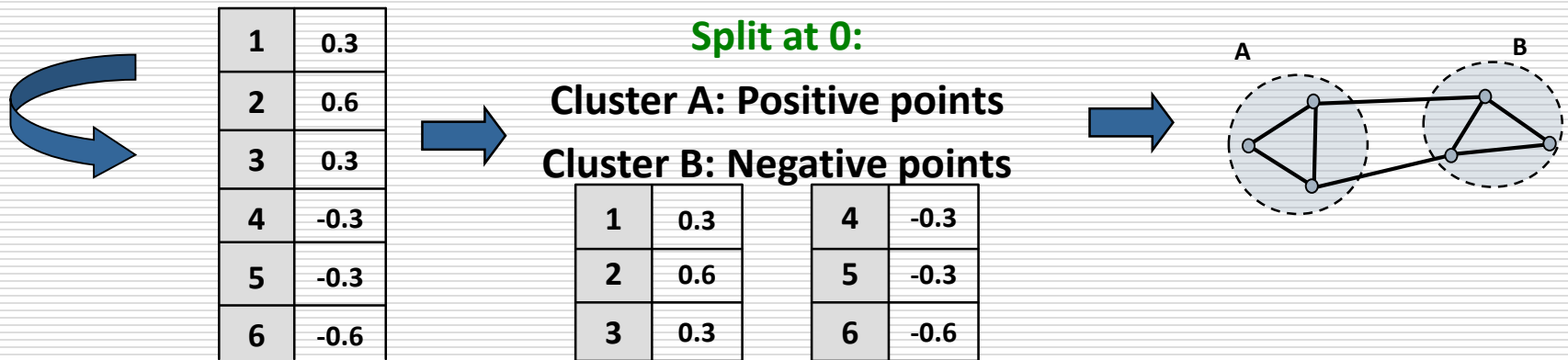
Spectral Partitioning

□ 3) Grouping:

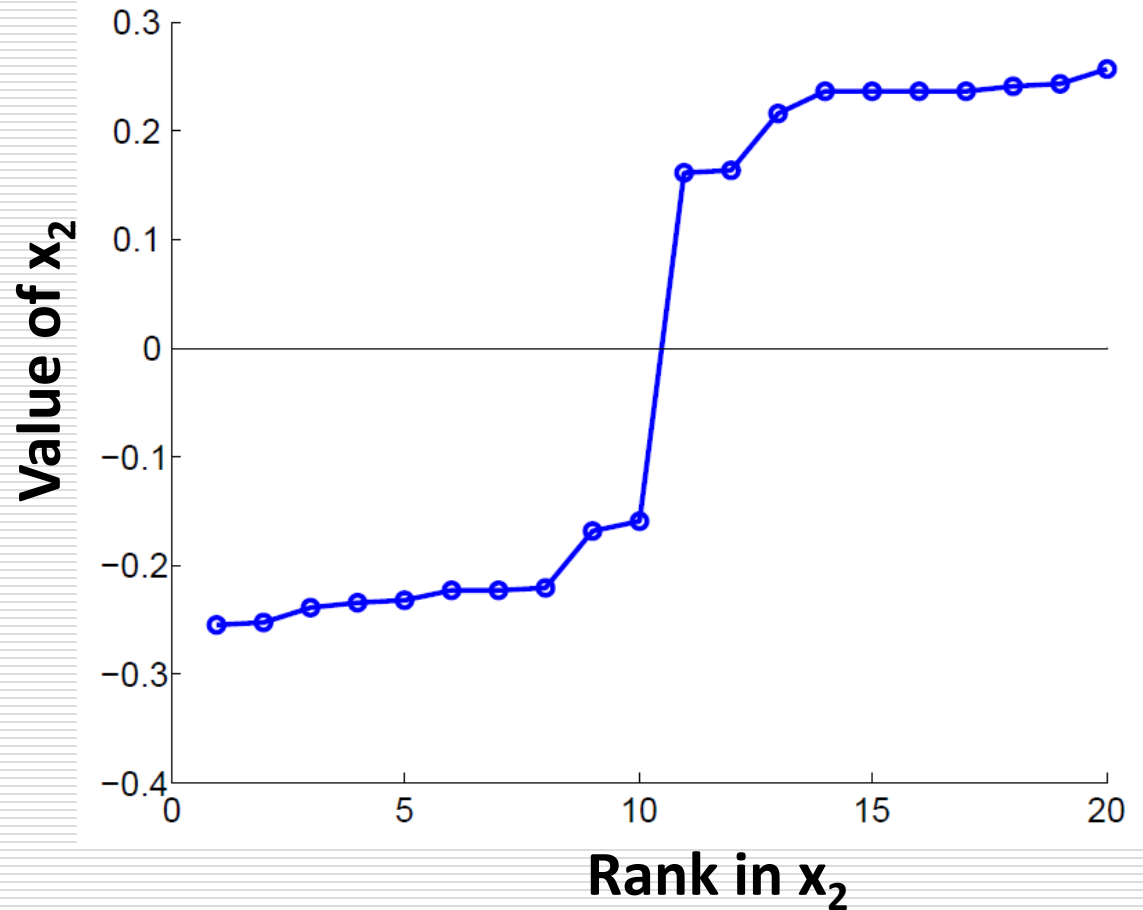
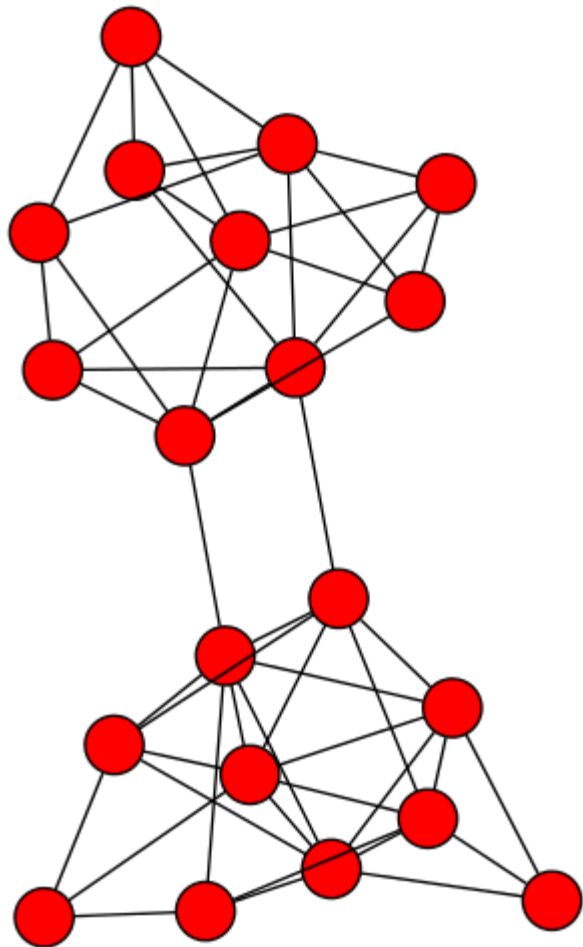
- Sort components of reduced 1-dimensional vector
- Identify clusters by splitting the sorted vector in two

□ How to choose a splitting point?

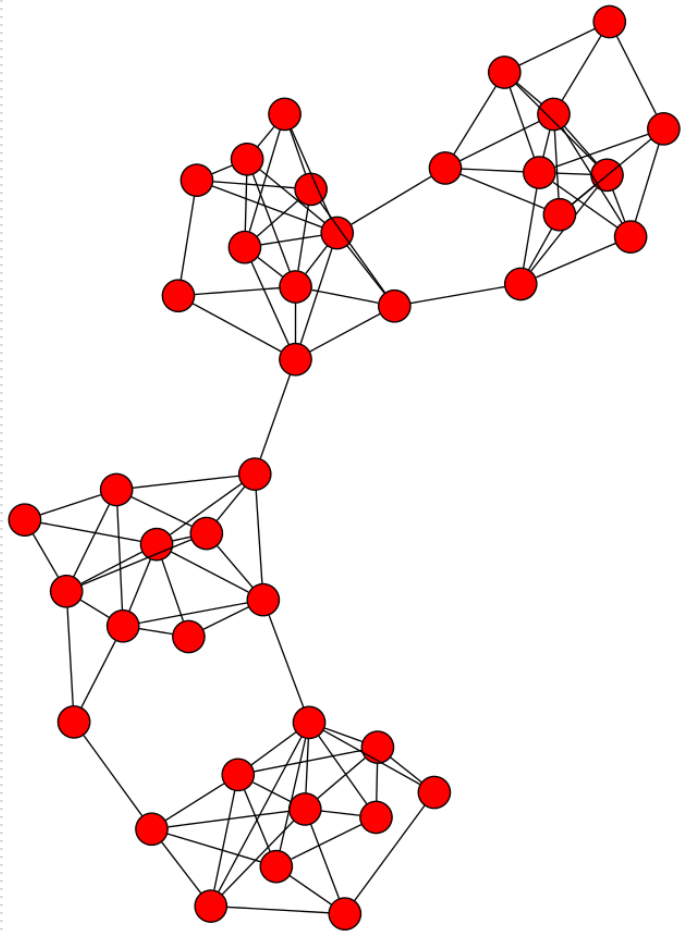
- Naïve approaches:
 - Split at **0** or median value
- More expensive approaches:
 - Attempt to minimize normalized cut in 1-dimension (sweep over ordering of nodes induced by the eigenvector)



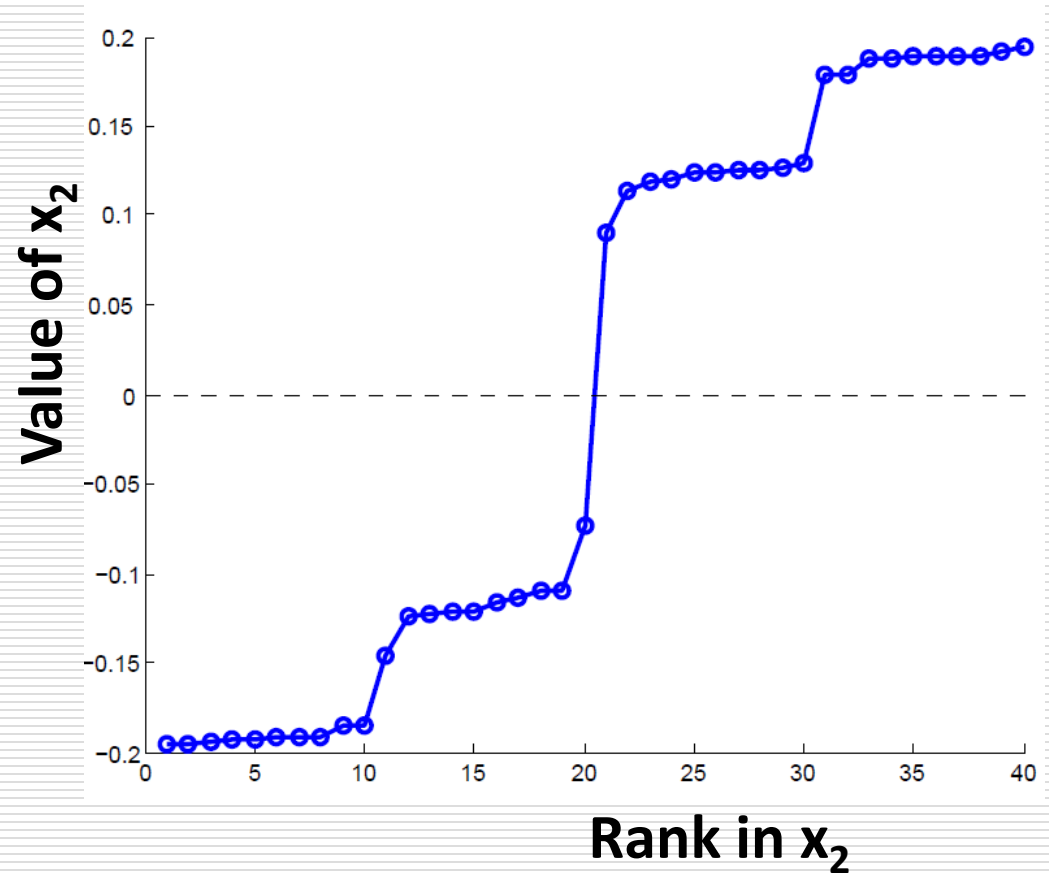
Example: Spectral Partitioning



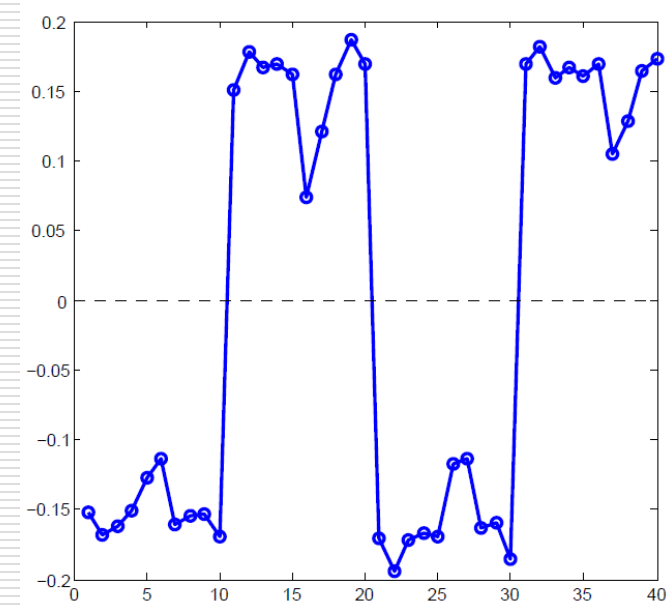
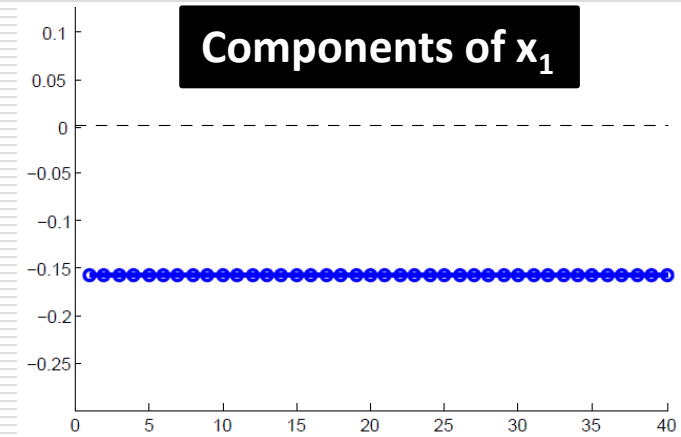
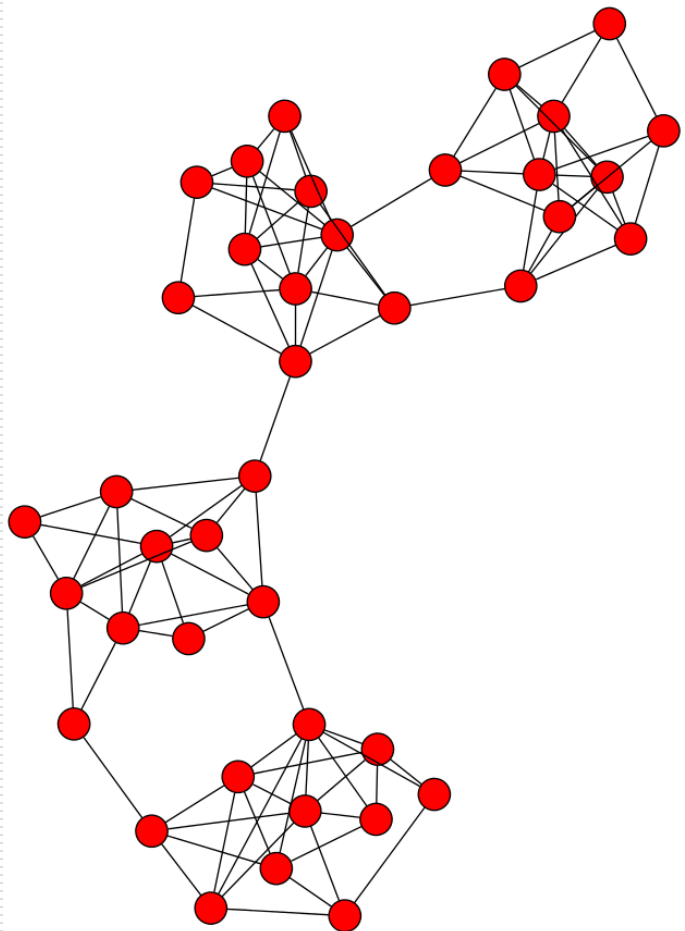
Example: Spectral Partitioning



Components of x_2



Example: Spectral Partitioning



Components of x_3

k-Way Spectral Clustering

□ How do we partition a graph into k clusters?

□ Two basic approaches:

■ **Recursive bi-partitioning** [Hagen et al., '92]

□ Recursively apply bi-partitioning algorithm in a hierarchical divisive manner

□ Disadvantages: Inefficient, unstable

■ **Cluster multiple eigenvectors** [Shi-Malik, '00]

□ Build a reduced space from multiple eigenvectors

□ Commonly used in recent papers

□ A preferable approach...

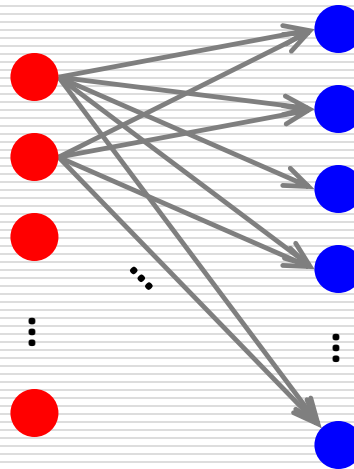
Why use multiple eigenvectors?

- ❑ **Approximates the optimal cut** [Shi-Malik, '00]
 - Can be used to approximate optimal k -way normalized cut
- ❑ **Emphasizes cohesive clusters**
 - Increases the unevenness in the distribution of the data
 - Associations between similar points are amplified, associations between dissimilar points are attenuated
 - The data begins to “approximate a clustering”
- ❑ **Well-separated space**
 - Transforms data to a new “embedded space”, consisting of k orthogonal basis vectors
- ❑ Multiple eigenvectors prevent instability due to information loss

Analysis of Large Graphs: Trawling

Trawling

- Searching for small communities in the Web graph
- What is the signature of a community / discussion in a Web graph?



Dense 2-layer graph

Use this to define “topics”:
What the same people on
the left talk about on the right
Remember HITS!

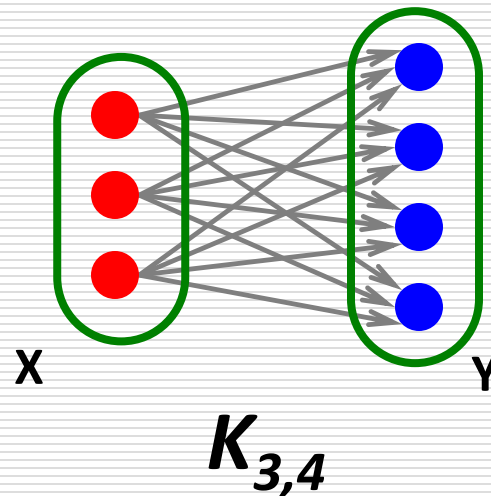
Intuition: Many people all talking about the same things

Searching for Small Communities

□ A more well-defined problem:

Enumerate complete bipartite subgraphs $K_{s,t}$

- Where $K_{s,t}$: s nodes on the “left” where each links to the same t other nodes on the “right”



$$|X| = s = 3$$

$$|Y| = t = 4$$

Fully connected

Frequent Itemset Enumeration

[Agrawal-Srikant '93]

□ **Market basket analysis.** Setting:

- **Market:** Universe U of n items
- **Baskets:** m subsets of U : $S_1, S_2, \dots, S_m \subseteq U$
(S_i is a set of items one person bought)
- **Support:** Frequency threshold f

□ **Goal:**

- Find all subsets T s.t. $T \subseteq S_i$ of at least f sets S_i
(items in T were bought together at least f times)

□ **What's the connection between the itemsets and complete bipartite graphs?**

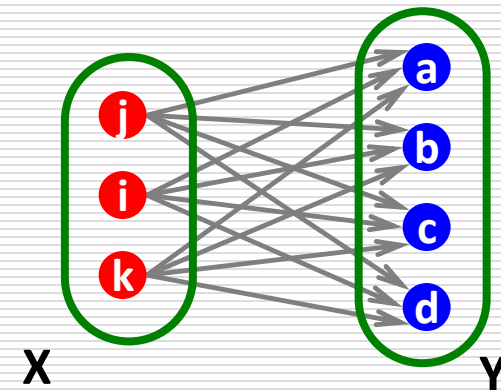
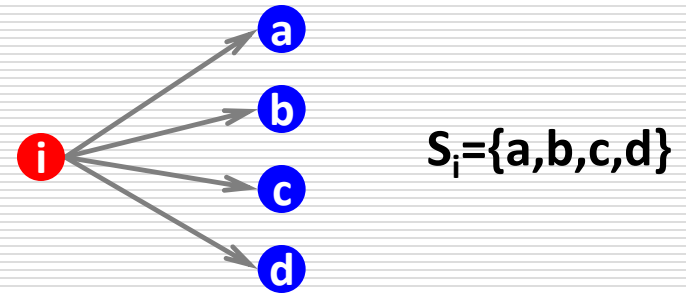
From Itemsets to Bipartite $K_{s,t}$

[Kumar et al. '99]

Frequent itemsets = **complete bipartite graphs!**

□ How?

- View each node i as a set S_i of nodes i points to
- $K_{s,t}$ = a set Y of size t that occurs in s sets S_i
- Looking for $K_{s,t} \rightarrow$ set of frequency threshold to s and look at layer t – all frequent sets of size t

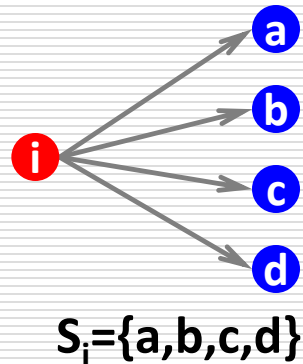


s ... minimum support ($|X|=s$)
 t ... itemset size ($|Y|=t$)

From Itemsets to Bipartite $K_{s,t}$

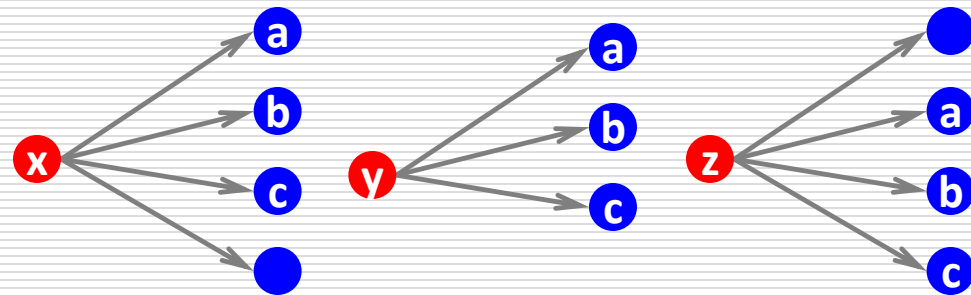
[Kumar et al. '99]

View each node i as a set S_i of nodes i points to



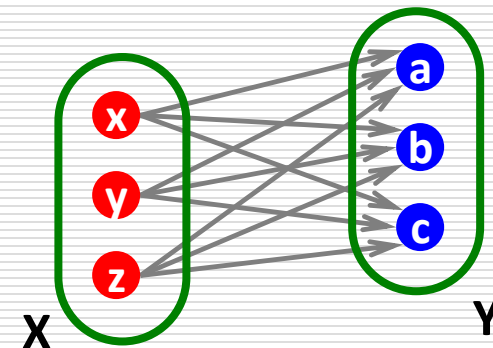
Find frequent itemsets:
 s ... minimum support
 t ... itemset size

Say we find a frequent itemset $Y = \{a, b, c\}$ of supp s
So, there are s nodes that link to all of $\{a, b, c\}$:

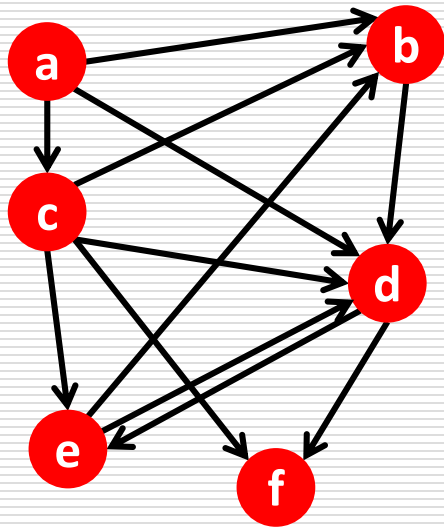


We found $K_{s,t}$!

$K_{s,t}$ = a set Y of size t
that occurs in s sets S_i



Example (1)



□ **Support threshold $s=2$**

■ $\{b,d\}$: support 3

■ $\{e,f\}$: support 2

□ And we just found 2 bipartite subgraphs:

Itemsets:

$a = \{b,c,d\}$

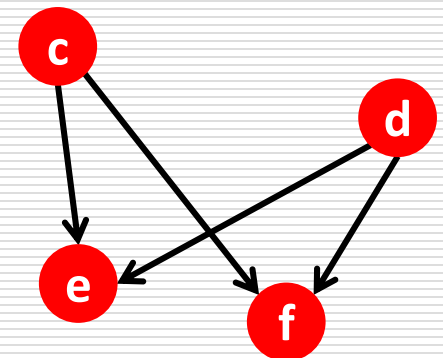
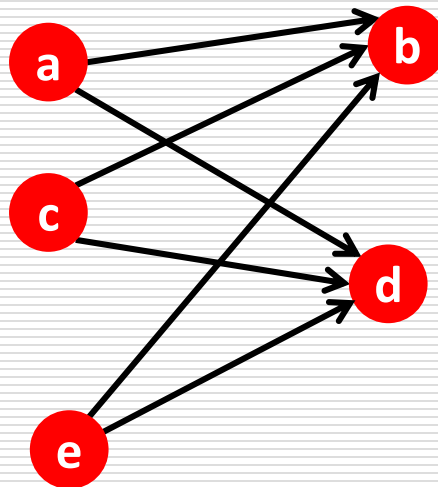
$b = \{d\}$

$c = \{b,d,e,f\}$

$d = \{e,f\}$

$e = \{b,d\}$

$f = \{\}$



Example (2)

□ Example of a community from a web graph

A community of Australian fire brigades

Nodes on the right	Nodes on the left
NSW Rural Fire Service Internet Site	New South Wales Fir...ial Australian Links
NSW Fire Brigades	Feuerwehrlinks Australien
Sutherland Rural Fire Service	FireNet Information Network
CFA: County Fire Authority	The Cherrybrook Rur...re Brigade Home Page
“The National Cente...ted Children’s Ho...	New South Wales Fir...ial Australian Links
CRAFTI Internet Connexions-INFO	Fire Departments, F... Information Network
Welcome to Blackwoo... Fire Safety Serv...	The Australian Firefighter Page
The World Famous Guestbook Server	Kristiansand brannv...dens brannvesener...
Wilberforce County Fire Brigade	Australian Fire Services Links
NEW SOUTH WALES FIR...ES 377 STATION	The 911 F,P,M., Fir...mp; Canada A Section
Woronora Bushfire Brigade	Feuerwehrlinks Australien
Mongarlowe Bush Fire – Home Page	Sanctuary Point Rural Fire Brigade
Golden Square Fire Brigade	Fire Trails “l...ghters around the...
FIREBREAK Home Page	FireSafe – Fire and Safety Directory
Guises Creek Volunt...fficial Home Page...	Kristiansand Firede...departments of th...
[Kumar, Raghavan, Rajagopalan, Tomkins: Trawling the Web for emerging cyber-communities 1999]	

Acknowledgement

- Slides are adapted from:
 - Prof. Jeffrey D. Ullman
 - Dr. Anand Rajaraman
 - Dr. Jure Leskovec