# Community Detection in Networks

Ahmet Onur Durahim

# How the Class Fits Together

# Network Communities



Communities, clusters, groups, modules

- Granovetter's theory suggest that networks are composed of **tightly connected sets of nodes**

- **Network communities:**
  - Sets of nodes with *lots* of connections *inside* and **few** to **outside** (the rest of the network)

# Finding Network Communities

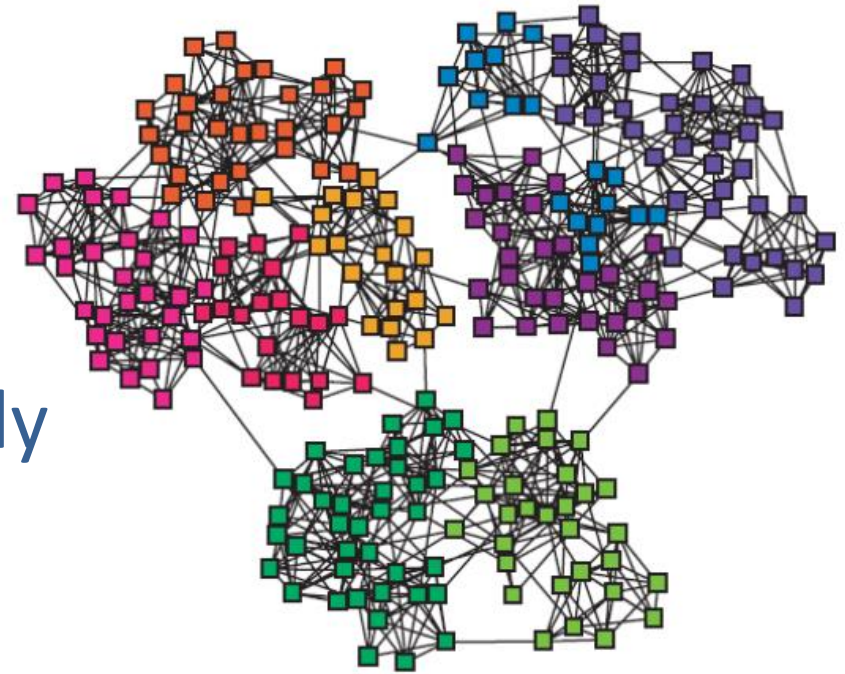- **How to automatically find such densely connected groups of nodes?**

- Ideally such automatically detected clusters would then correspond to real groups

- **For example:**



Communities, clusters, groups, modules

# Social Network Data



- **Zachary's (Ph.D.) Karate club network:**
  - Observe social ties and rivalries in a university karate club
  - During observation, conflicts led the group to split
  - Split could be explained by a minimum cut in the network

# Micro-Markets in Sponsored Search

- **Find micro-markets by partitioning the "Query x Advertiser" graph:**



*Web search:*
keywords and advertisers
bidding with them

# NCAA Football Network

# NCAA Football Network



**NCAA conferences**

- 🟣 Mid American
- 🔴 Big East
- ⚫ Atlantic Coast
- 🌸 SEC
- 🟡 Conference USA
- 🔵 Big 12
- 🟤 Western Athletic
- 🟢 Pacific 10
- 🔵 Mountain West
- 🟢 Big 10
- ⚫ Sun Belt
- ⚪ Independents

**Nodes: Teams**
**Edges: Games played**

# Facebook Ego-network



Can we identify social communities?

**Nodes: Users**
**Edges: Friendships**

# Facebook Ego-network



High school

Company

Stanford (Squash)

Stanford (Basketball)

**Social communities**

**Nodes: Users**
**Edges: Friendships**

# Protein-Protein Interactions



Can we identify functional modules?

Nodes: Proteins
Edges: Interactions

# Protein-Protein Interactions



**Functional modules**

**Nodes: Proteins**
**Edges: Interactions**

# Community Detection

## How to find communities?



We will work with **undirected** (unweighted) networks

# Partitions and Communities

- An algorithm for *community detection*
  - Choose some *score function f*
    - takes *a network A and a partition P of its vertices* as input
    - returns a *scalar value (score)* as output
    $$score = f(G, P)$$
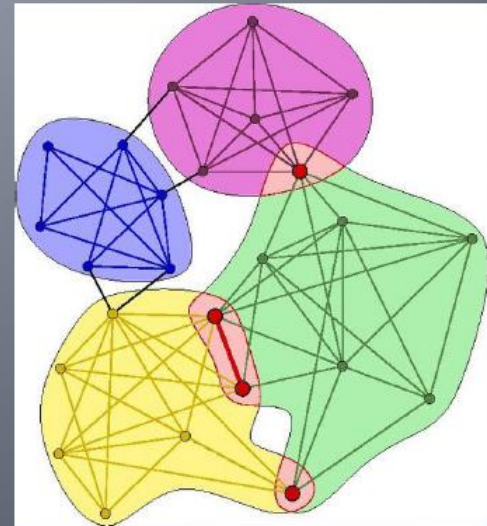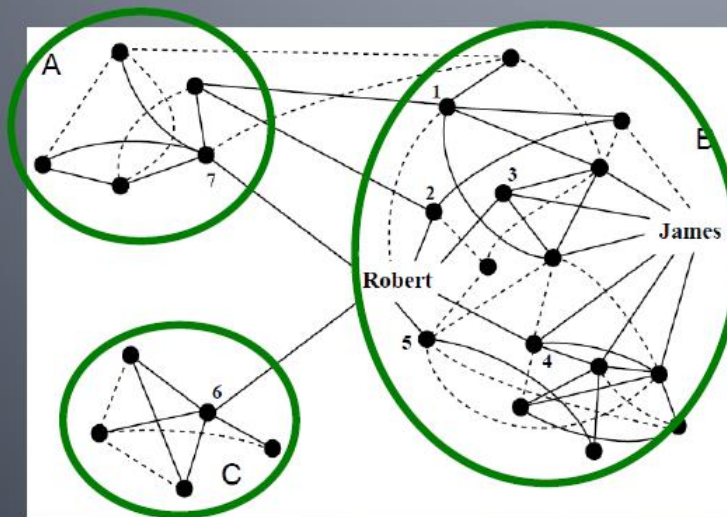  - *Find the partition P that maximizes f*
    - f encodes our beliefs about *what makes a good partition* with respect to representing community structure
  - If *f* prefers assortative structure
    - assigns higher scores to partitions that place more edges within groups than between them
    - maximizing f over all partitions will yield the partition with the most assortative groups possible under f

- *Network partitions (communities):*
  - *a division of a network into k non-empty groups* (communities)
  - *every node v belongs to one and only one group* (community)

# Partitions and Communities

*Combinatorial problem:*

- Number of ways to divide network of n nodes in 2 groups (bi-partition):

$$\frac{n!}{n_1!\,n_2!}, \qquad n = n_1 + n_2$$

- Dividing into k non-empty groups (Stirling numbers of the second kind)

$$S(n, k) = \frac{1}{k!} \sum_{j=0}^{n} (-1)^j C_k^j (k - j)^n$$

there are $S(4, 2) = 7$ possible partitions of a network with $n = 4$ nodes (indexed as 1,2,3,4) partitioned into $k = 2$ groups:

$$\{1\}\{234\}, \ \{2\}\{134\}, \ \{3\}\{124\}, \ \{4\}\{123\}, \ \{12\}\{34\}, \ \{13\}\{24\}, \ \{14\}\{23\}$$

# Partitions and Communities

## *Combinatorial problem:*

- Number of ways to divide network of n nodes in 2 groups (bi-partition):

$$\frac{n!}{n_1!\,n_2!}, \qquad n = n_1 + n_2$$

- Dividing into k non-empty groups (Stirling numbers of the second kind)

$$S(n, k) = \frac{1}{k!} \sum_{j=0}^{n} (-1)^j C_k^j (k - j)^n$$

$$\left\{ {n \atop k} \right\} = \sum_{j=1}^{k} (-1)^{k-j} \frac{j^{n-1}}{(j - 1)!(k - j)!} = \frac{1}{k!} \sum_{j=0}^{k} (-1)^{k-j} \binom{k}{j} j^n$$

there are $S(4, 2) = 7$ possible partitions of a network with $n = 4$ nodes (indexed as 1,2,3,4) partitioned into $k = 2$ groups:

$$\{1\}\{234\}, \ \ \{2\}\{134\}, \ \ \{3\}\{124\}, \ \ \{4\}\{123\}, \ \ \{12\}\{34\}, \ \ \{13\}\{24\}, \ \ \{14\}\{23\}$$

# Partitions and Communities

## *Combinatorial problem:*

- Number of ways to divide network of n nodes in 2 groups (bi-partition):

$$\frac{n!}{n_1! \, n_2!}, \qquad n = n_1 + n_2$$

- Dividing into k non-empty groups (Stirling numbers of the second kind)

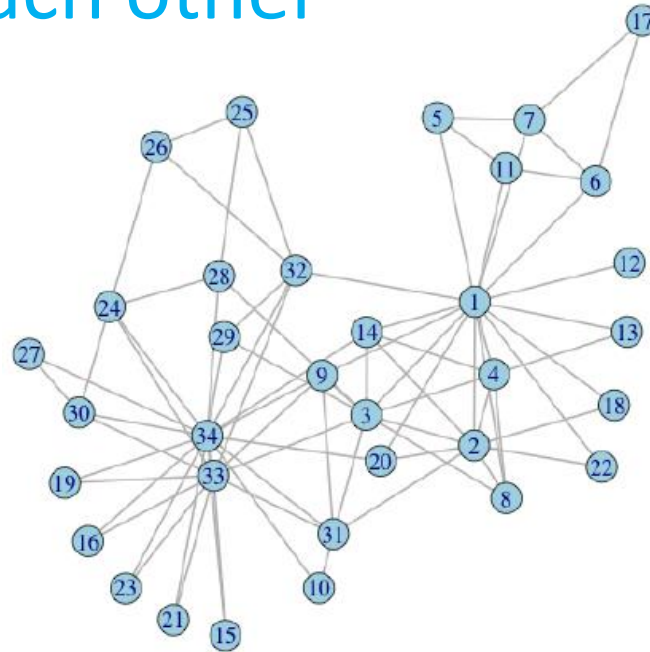$$S(n, k) = \frac{1}{k!} \sum_{j=0}^{n} (-1)^j C_k^j (k - j)^n$$

- Number of all possible partitions (n-th Bell number):

$$B_n = \sum_{k=1}^{n} S(n, k)$$

$B_{20} = 5,832,742,205,057$ ⟹ **Brute force algorithms won't work**
⟹ **Need Heurisric / Greedy Algorithms**

# Network communities

- ***Network communities*** are groups of nodes ***similar*** to each other



- ***Community detection:*** assignment of nodes to communities
  - Non-overlapping communities (every node belongs to a single group)

# Similarity based clustering

## *Similarity based node clustering*

- Define *similarity measure between nodes based on network structure*
  - *Jaccard similarity*
  - *Cosine similarity*
  - *Pearson correlation*
  - *Euclidean distance (dissimilarity)*

- Calculate similarity between all pairs of nodes in the graph (similarity matrix)
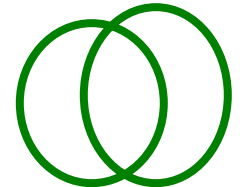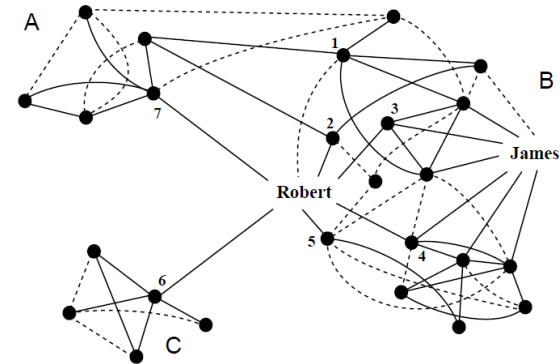- Group together nodes with high similarities

# Clustering and Community Finding

**How to extract groups?**

**Many methods:**

- *Hierarchical methods:*
  – Top-down / bottom-up
- *Graph partitioning methods:*
  – Define "edge counting" metric – conductance, expansion, modularity, etc. – and optimize!
- *Spectral methods:*     **Lecture Notes on Spectral Graph Methods (by Mahoney)**
  – Based on eigenvector decomposition of modified graph adjacency matrix

- *Clique percolation method:*
  – To extract overlapping communities in networks

# Betweenness and Graph Partitioning

- A way of thinking **networks** in terms of
  - **tightly-knit regions** and **weaker ties** that link them together
    - Clustering coefficient and local bridges
  - What we mean by tightly-knit region?
- **Graph Partitioning:**
  - a method that *takes a network and break it down* into a set of tightly-knit regions, with sparser interconnections between the regions
- What we might **hope** from this method?
  - Need a general way to pull (visually observed) groups out of the data, beyond using just visual intuition
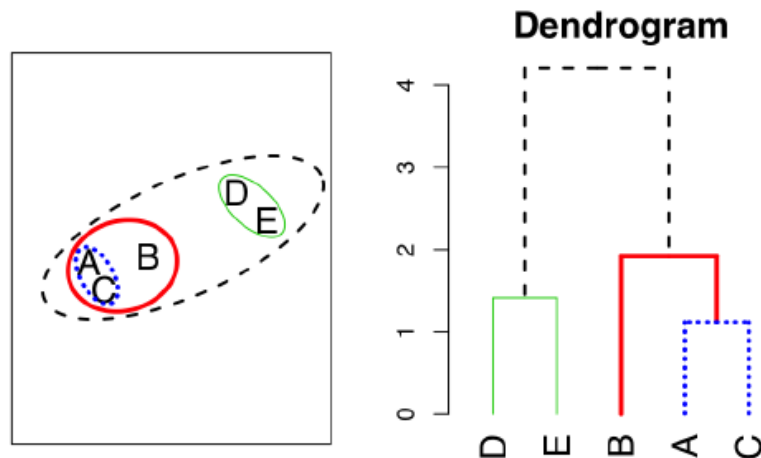
# Methods of Graph Partitioning

- **Agglomerative methods** of graph partitioning
  - find nodes that are likely to belong to the same region and merge them together
    - network consists of a large number of *merged chunks*, each containing the seeds of a *densely-connected region*
  - process looks for chunks that should be further merged together
  - regions are assembled "bottom-up"
- **Divisive methods** of graph partitioning
  - Identify and remove the "*spanning links*" between *densely-connected regions*
    - the network begins to fall apart into large pieces
  - spanning links are further identified and removed within these connected components
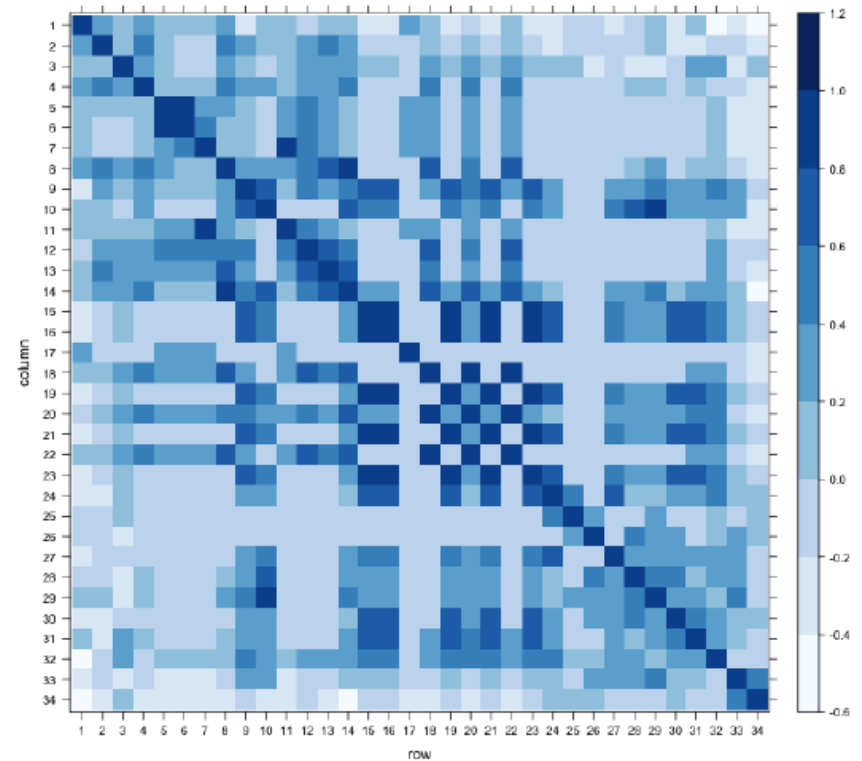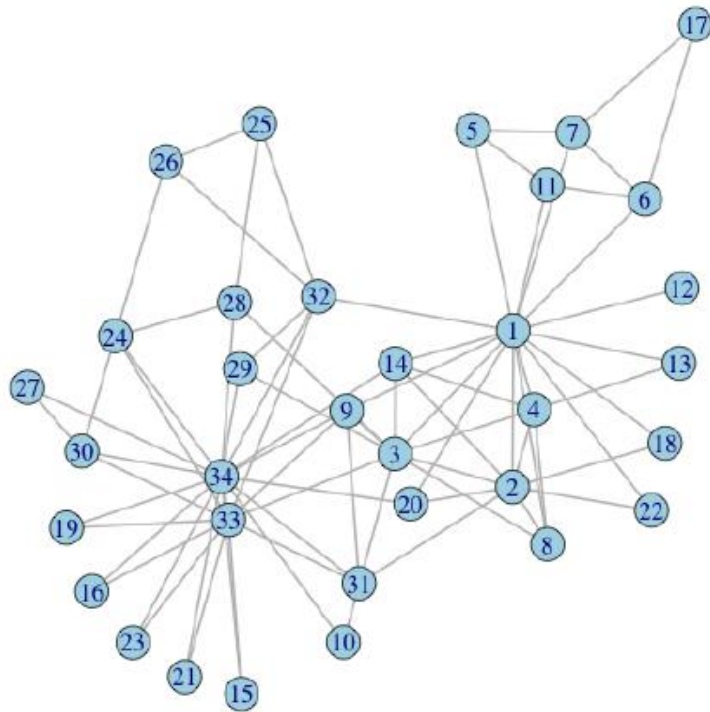  - process continues "top-to-bottom"

# Hierarchical clustering

## Agglomerative clustering:

- Assign each node to a group of its own
- Find two groups with the ***highest similarity*** and ***join*** them in a single group
- Calculate ***similarity between groups***:
  - ***single-linkage clustering:*** most similar in the group
  - ***complete-linkage clustering:*** least similar in the group
  - ***average-linkage clustering:*** mean similarity between groups
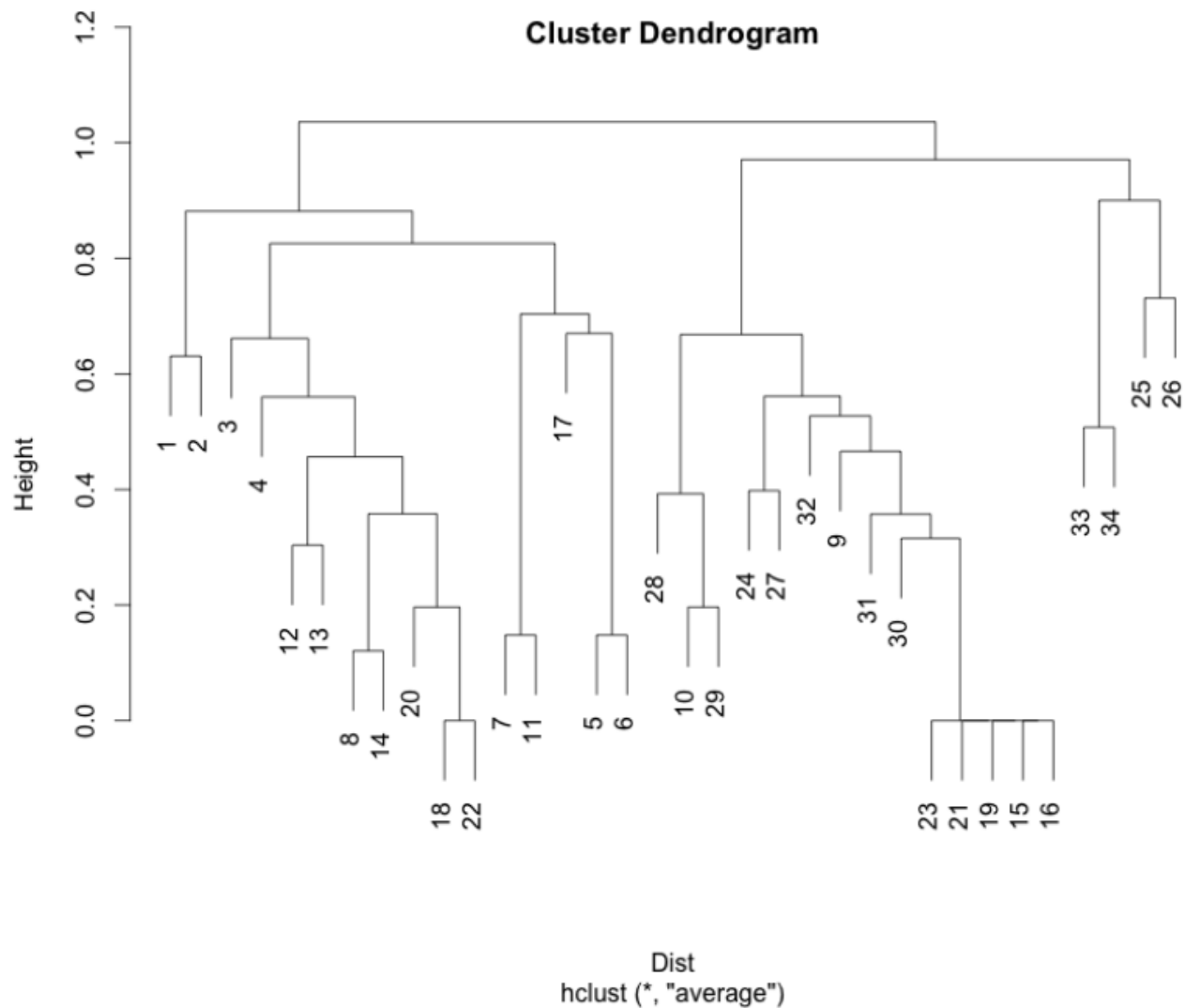- Repeat until all joined into single group



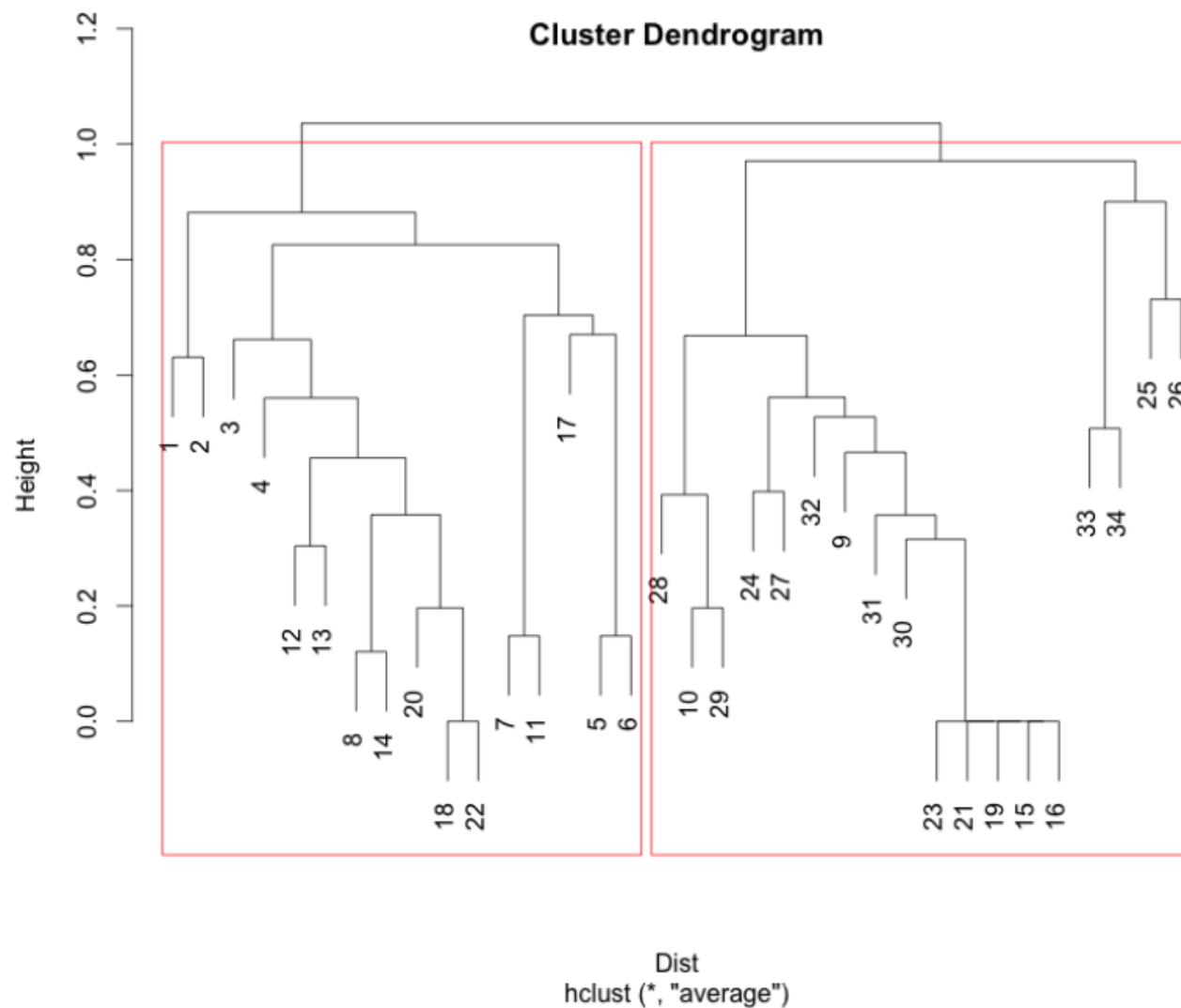Dendrogram

# Similarity matrix
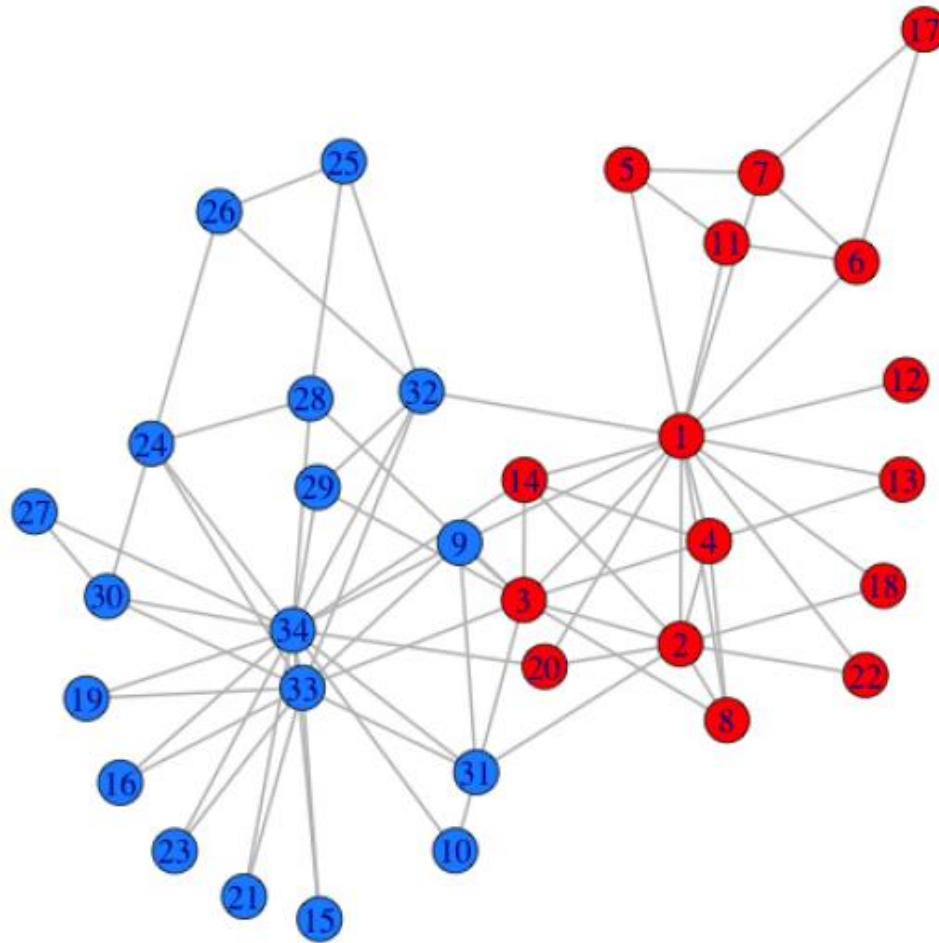
## Zachary karate club



*Remember the notions of node equivalences*
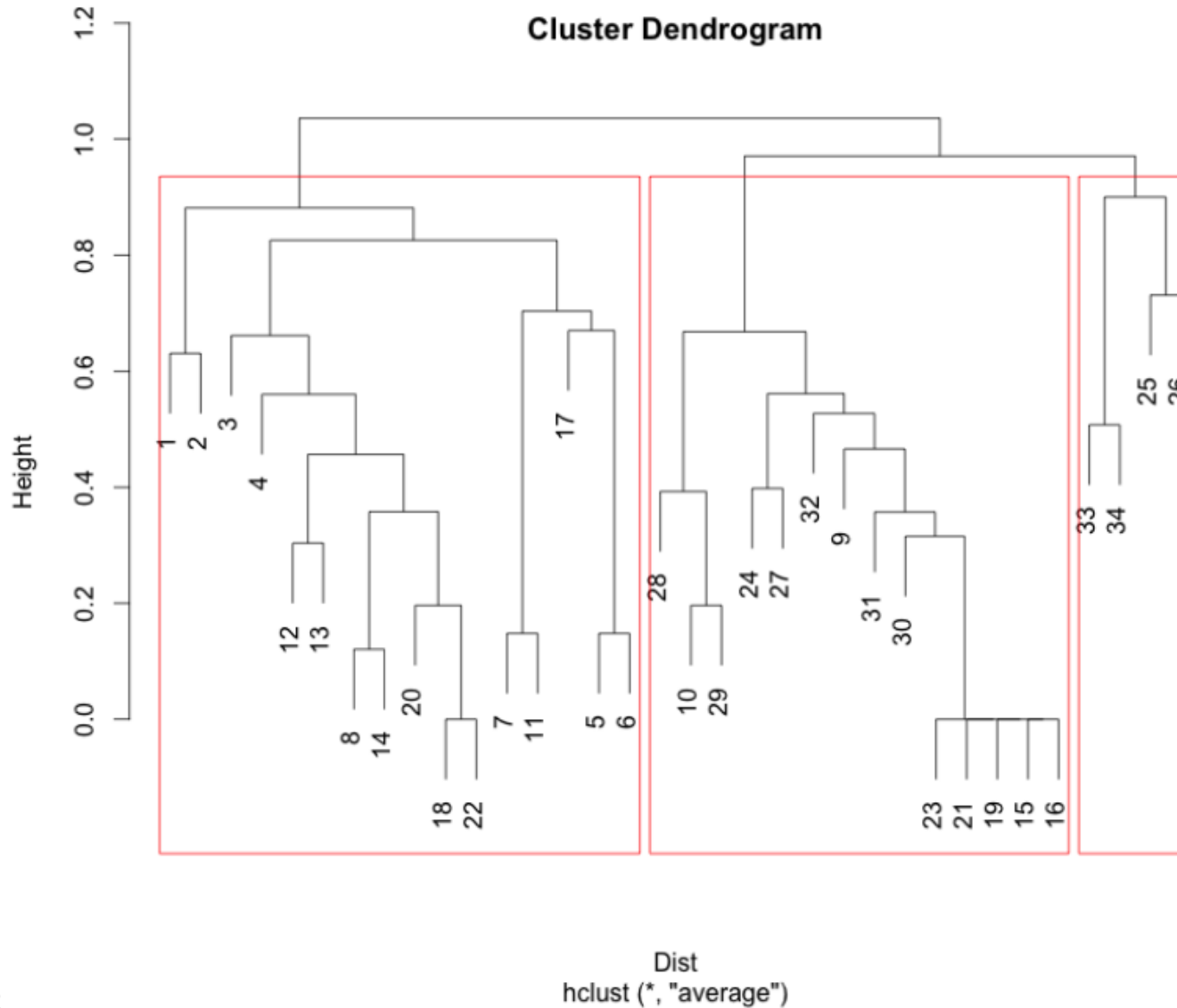
# Hierarchical clustering



Cluster Dendrogram

Height

Dist
hclust (*, "average")

# Hierarchical clustering



**Cluster Dendrogram**

Dist
hclust (*, "average")

# Hierarchical clustering

# Hierarchical clustering



**Cluster Dendrogram**

Height

Dist
hclust (*, "average")
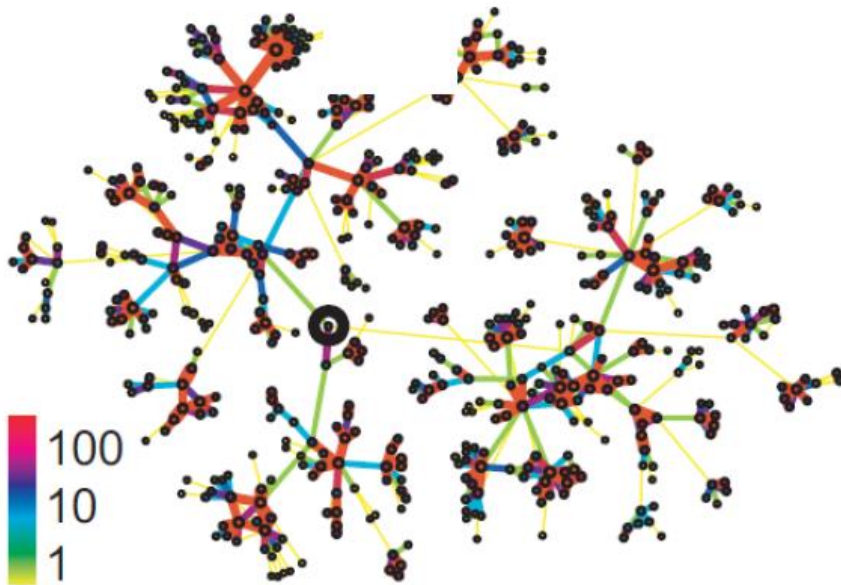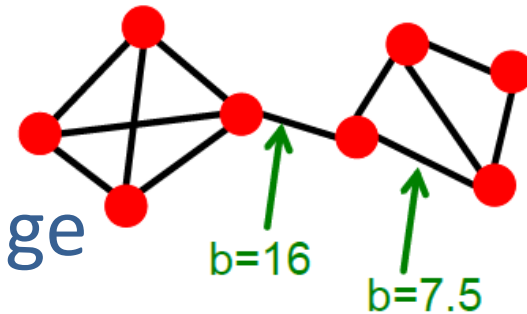
# Hierarchical clustering

# Method 2: *Girvan-Newman*

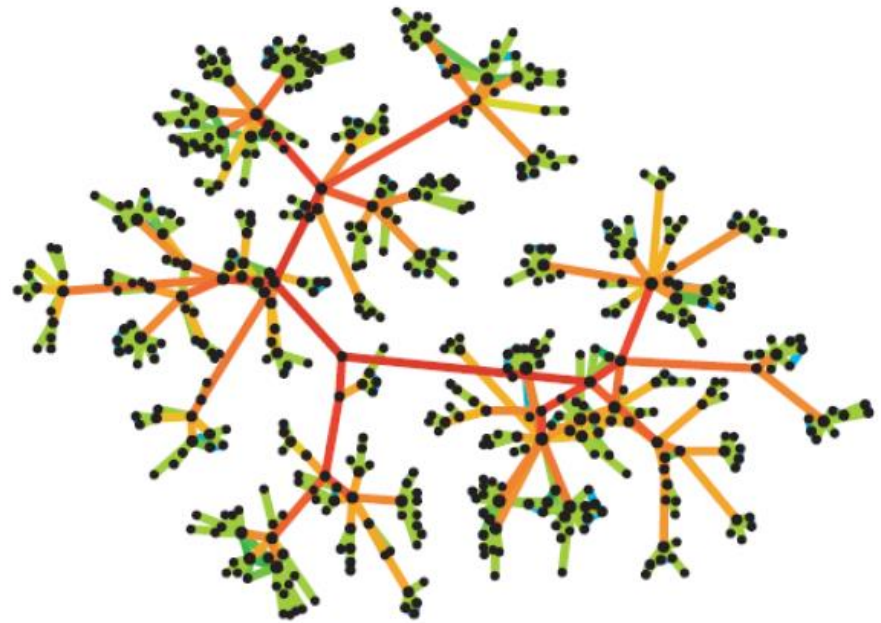- **Divisive hierarchical clustering** based on the notion of *edge betweenness*:
  - **Number of shortest paths passing through the edge**
- *Girvan-Newman Algorithm:*
    - **Undirected unweighted networks**
  - **Repeat until no edges are left:**
    - Calculate betweenness of edges
    - Remove edges with highest betweenness
  - *Connected components* are *communities*
  - Gives a hierarchical decomposition of the network

# Strength of Weak Ties



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

- **Intuition:**



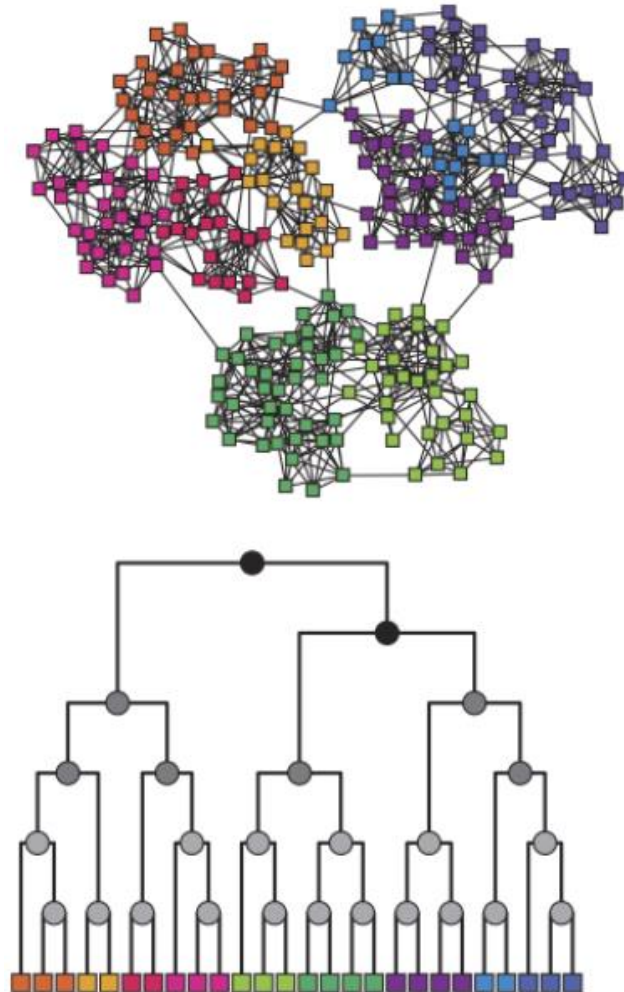Edge strengths (call volume) in a real network

Edge betweenness in a real network

# Edge betweenness

- ***Hierarchical algorithm,*** dendrogram

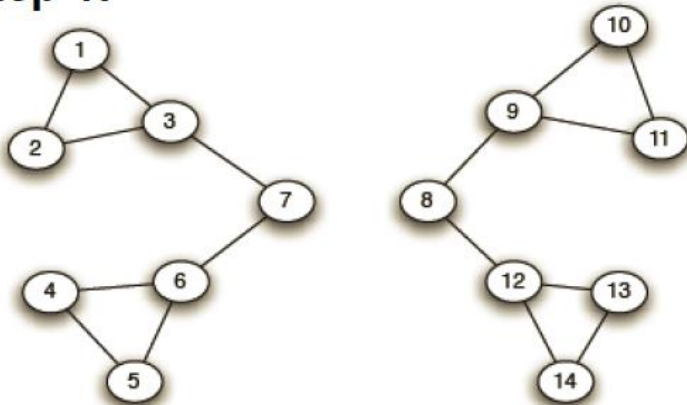# Girvan-Newman: Example



Need to re-compute
betweenness at
every step

# Girvan-Newman: Example

# 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?**

# Determine the flow from one node to all other nodes in the graph

- Perform a breadth-first search of the graph, starting at a node (node A)

- Determine the number of shortest paths from node A to all the other nodes

- Based on these numbers, determine the amount of flow from node A to all other nodes that uses each edge

# How to Compute Betweenness?

- **Want to compute betweenness of paths starting at node $A$ efficiently**

- **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+∑child edges
  - -- split the flow up based on the parent value
- Repeat the BFS procedure for each starting node $U$

1+1 paths to H
Split evenly

1+0.5 paths to J
Split 1:2

1 path to K.
Split evenly

# 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+∑child edges
  -- split the flow up based on the parent value
- Repeat the BFS procedure for each starting node $U$

10 nodes to be reached below this level

6 nodes to be reached below this level

1+1 paths to H
Split evenly

3 nodes to be reached below this level

1+0.5 paths to J
Split 1:2

1 path to K.
Split evenly

# 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 $s \in S$:

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

**Need a null model!**

# Null Model: 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: $k_i \cdot \dfrac{k_j}{2m} = \dfrac{k_i k_j}{2m}$

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

- $= \dfrac{1}{2} \sum_{i \in N} \sum_{j \in N} \dfrac{k_i k_j}{2m} = \dfrac{1}{2} \cdot \dfrac{1}{2m} \sum_{i \in N} k_i \left( \sum_{j \in N} k_j \right) =$
- $= \dfrac{1}{4m} 2m \cdot 2m = m$

Note:

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

# Modularity

- **Modularity of partitioning S of graph G:**

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

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

    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

# Dendrogram and modularity score

5/7/2022

80

# Network communities

- Zachary karate club

# Modularity: Number of Clusters

- *Modularity is useful for selecting the number of clusters:*



**Next time: Why not optimize Modularity directly?**

# Modularity

- **Modularity of partitioning $S$ of graph $G$:**

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

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

    $= \underbrace{\frac{1}{2m}}_{\text{Normalizing const.: } -1 < Q < 1} \sum_{s \in S} \sum_{i \in s} \sum_{j \in s} \left( A_{ij} - \frac{k_i k_j}{2m} \right)$    $A_{ij} = 1$ if $i \rightarrow j$, 0 else

- **Modularity $Q$ tells us whether $S$ represents any significant community structure**
  - So, let's find $S$ that maximizes modularity itself!

# Method 3: *Modularity Optimization*

- *Let's split the graph into 2 communities*

- **Want to directly optimize modularity:**

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

- **Community membership vector s:**
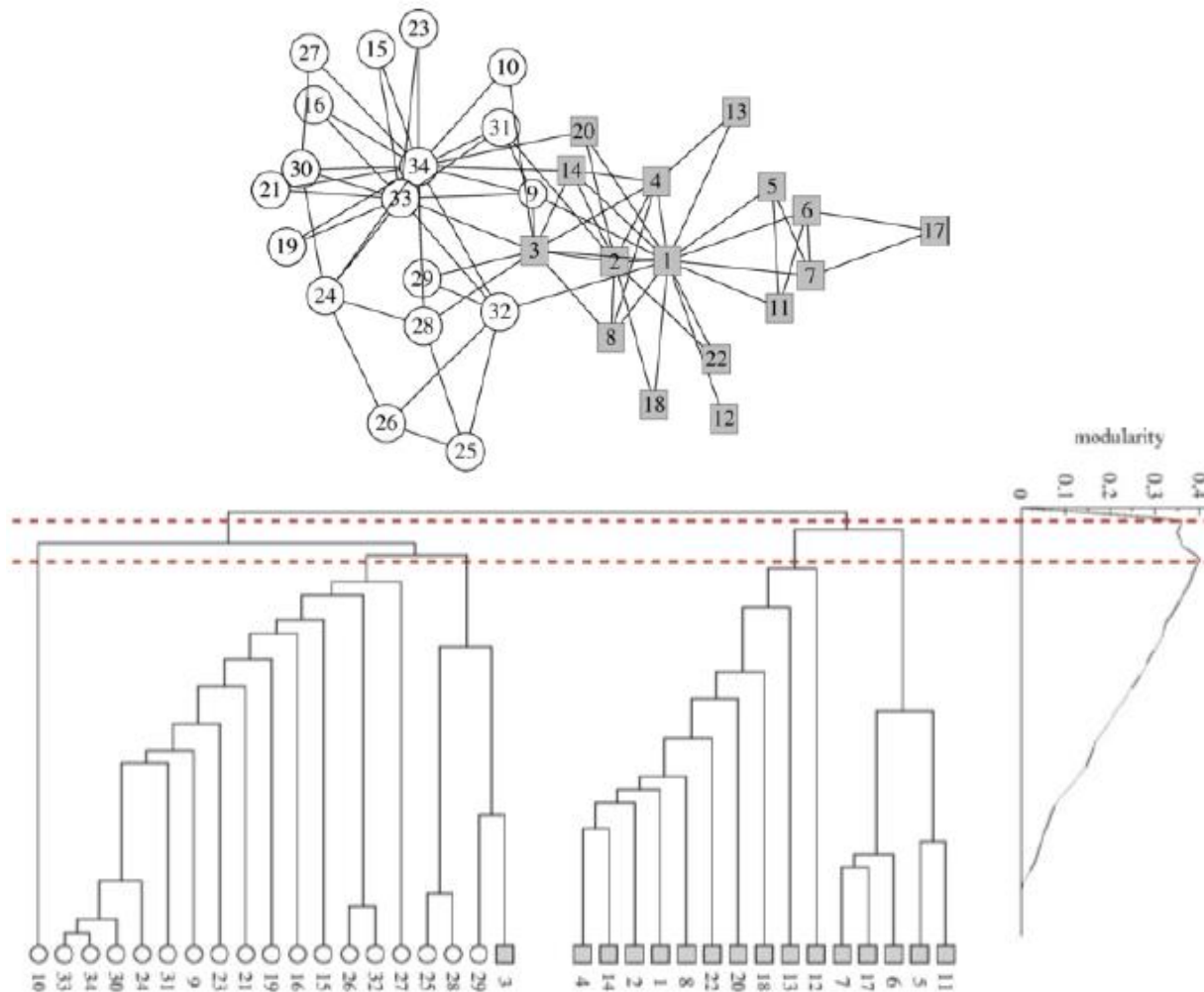
  - $s_i = 1$ if node $i$ is in community **1**
    **-1** if node $i$ is in community **-1**

  $$\frac{s_i s_j + 1}{2} = \begin{array}{l} 1.. \text{ if } s_i = s_j \\ 0.. \text{ else} \end{array}$$

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

  $$= \frac{1}{4m} \sum_{i,j \in N} \left( A_{ij} - \frac{k_i k_j}{2m} \right) s_i s_j$$

# Modularity Matrix

- **Define:**
  - **Modularity matrix:** $B_{ij} = A_{ij} - \frac{k_i k_j}{2m}$
  - **Membership:** $s = \{-1, +1\}$

- **Then:** $Q(G, s) = \frac{1}{4m} \sum_{i \in N} \sum_{j \in N} \left( A_{ij} - \frac{k_i k_j}{2m} \right) s_i s_j$

  $= \frac{1}{4m} \sum_{i,j \in N} B_{ij} s_i s_j$

  $= \frac{1}{4m} \sum_i s_i \underbrace{\sum_j B_{ij} s_j}_{= B_{i\cdot} \cdot s} = \frac{1}{4m} s^T B s$

- **Task:** Find $s \in \{-1,+1\}^n$ that maximizes **Q(G,s)**

# Quick Review of Linear Algebra

$$\begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

- **Symmetric** matrix **A**
  - is Positive Semidefinite: $A = U \cdot U^T$
- Then solutions $\lambda, x$ to equation $A \cdot x = \lambda \cdot x$:
  - **Eigenvectors** $x_i$ ordered by the magnitude of their corresponding **eigenvalues** $\lambda_i$ ($\lambda_1 \leq \lambda_2 \ldots \leq \lambda_n$)
  - $x_i$ are **orthonormal** (orthogonal and unit length)
  - $x_i$ form a coordinate system (basis)
  - If $A$ is Positive Semidefinite: $\lambda_i \geq 0$ (and they always exist)
- **Eigendecomposition theorem:** Can rewrite matrix $A$ in terms of its eigenvectors and eigenvalues:

$$A = \sum_i x_i \cdot \lambda_i \cdot x_i^T$$

# Modularity Optimization

- **Rewrite:** $Q(G, s) = \frac{1}{4m} s^{\mathrm{T}} B s$ **in terms of its eigenvectors $x$ and eigenvalues $\lambda$:**

$$= s^{\mathrm{T}} \left[ \sum_{i=1}^{n} x_i \lambda_i x_i^{T} \right] s = \sum_{i=1}^{n} s^{T} x_i \lambda_i x_i^{T} s = \sum_{i=1}^{n} (s^{T} x_i)^2 \lambda_i$$

- **So, if there would be no other constraints on $s$ then to maximize $Q$, we make $s = x_n$**

  - **Why?** Because $\lambda_n \geq \lambda_{n-1} \geq \cdots$
    - Remember $s$ has fixed length ($\|s\| = 1$)!
  - Assigns all weight in the sum to $\lambda_n$ (largest eigenvalue)
    - All other $s^T x_i$ terms are **zero** because of orthonormality

# Finding the vector *s*

- **Let's consider only the first term in the summation (because $\lambda_n$ is the largest):**
  $$\max_s Q(G,s) = \sum_{i=1}^{n}(s^T x_i)^2 \lambda_i \approx (s^T x_n)^2 \lambda_n$$

- **Let's maximize:** $\sum_{j=1}^{n} s_j \cdot x_{n,j}$ where $s_j \in \{-1,+1\}$

- To do this, we set:

  - $s_j = \begin{cases} +1 & \text{if } x_{n,j} \geq 0 \text{ (j–th coordinate of } x_n \geq 0) \\ -1 & \text{if } x_{n,j} < 0 \text{ (j–th coordinate of } x_n < 0) \end{cases}$

- **Continue the bisection hierarchically**

# Summary: Modularity Optimization

- **Fast Modularity Optimization Algorithm:**
  - Find leading eigenvector $x_n$ of modularity matrix **B**
  - Divide the nodes by the signs of the elements of $x_n$
  - Repeat hierarchically until:
    - If a proposed split does not cause modularity to increase, declare community indivisible and do not split it
    - If all communities are indivisible, stop
- **How to find $x_n$? Power method!**
  - Start with random $v^{(0)}$, repeat :
  - When converged ($v^{(t)} \approx v^{(t+1)}$), set $x_n = v^{(t)}$

$$v^{(t+1)} = \frac{Bv^{(t)}}{\left\| Bv^{(t)} \right\|}$$

# Summary: Modularity based Algorithms

- **Girvan-Newman**
  - Based on the "strength of weak ties"
  - Remove edge of highest betweenness

- **Modularity:**
  - Overall quality of the partitioning of a graph
  - Use to determine the number of communities

- **Fast Modularity Optimization:**
  - Transform the modularity optimization into an eigenvalue problem

# 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***

# Many other partitioning methods

- **METIS:**
  - Heuristic but works really well in practice
  - http://glaros.dtc.umn.edu/gkhome/views/metis

- **Graclus:**
  - Based on kernel k-means
  - http://www.cs.utexas.edu/users/dml/Software/graclus.html

- **Louvain:**
  - Based on Modularity optimization
  - http://perso.uclouvain.be/vincent.blondel/research/louvain.html

- **Clique percolation method:**
  - For finding overlapping clusters
  - http://angel.elte.hu/cfinder/