



Network Metrics

CE642: Social and Economic Networks
Maryam Ramezani
Sharif University of Technology
maryam.ramezani@sharif.edu



01

Centrality

Why a centrality measure?

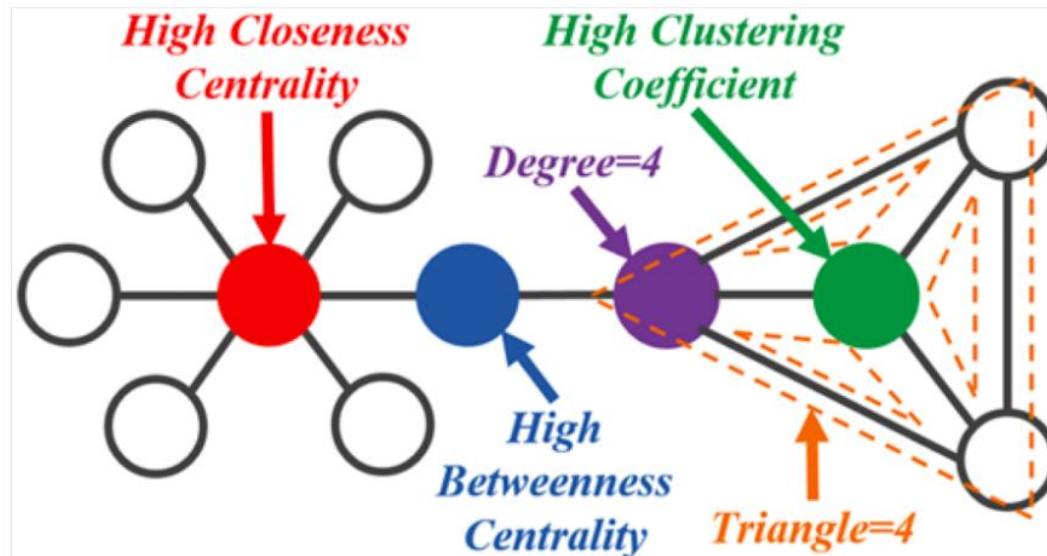
- Often, we are interested in identifying IMPORTANT network components
 - Nodes
 - Edges
- Central components may play critical role in network functions
 - Robustness
 - Collective behavior
 - Synchronization
 - Information spreading
 - Social dynamics
 - ...

Centrality

- Which nodes are most ‘central’?
- Definition of ‘central’ varies by context/purpose.
- Local measure:
 - Degree
- Relative to the rest of the network:
 - Closeness
 - Betweenness
 - Eigenvector (Bonacich power centrality)
- How evenly is centrality distributed among nodes?
 - Centralization

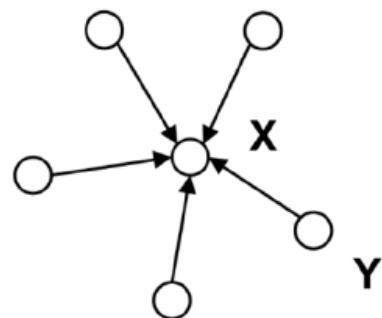
Network Centrality

- Given a social network, which nodes are more **important** or **influential**?
- Centrality measures** were proposed to account for the importance of the nodes of a network

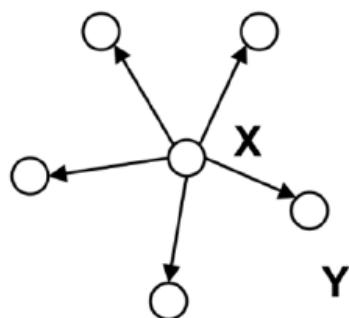


Network Centrality

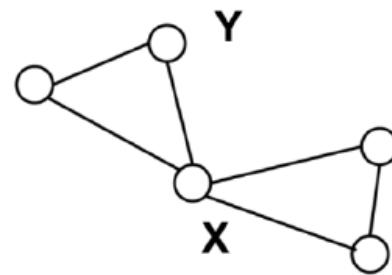
- In each of the following networks, X has higher centrality than Y according to a particular measure



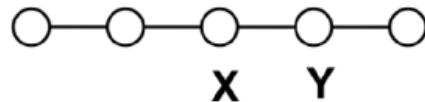
indegree



outdegree



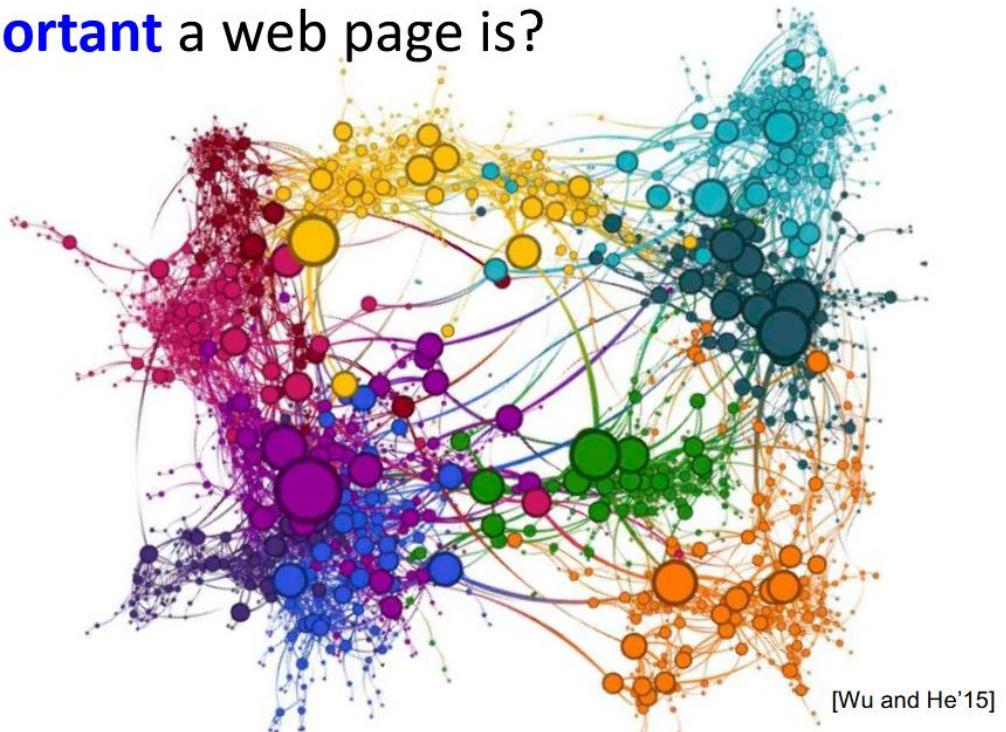
betweenness



closeness

Network Centrality

- **Centrality** is used often for detecting:
 - How **influential** a person is in a social network?
 - How **well used** a road is in a transportation network?
 - How **important** a web page is?



Centrality Measures

■ Geometric Measures:

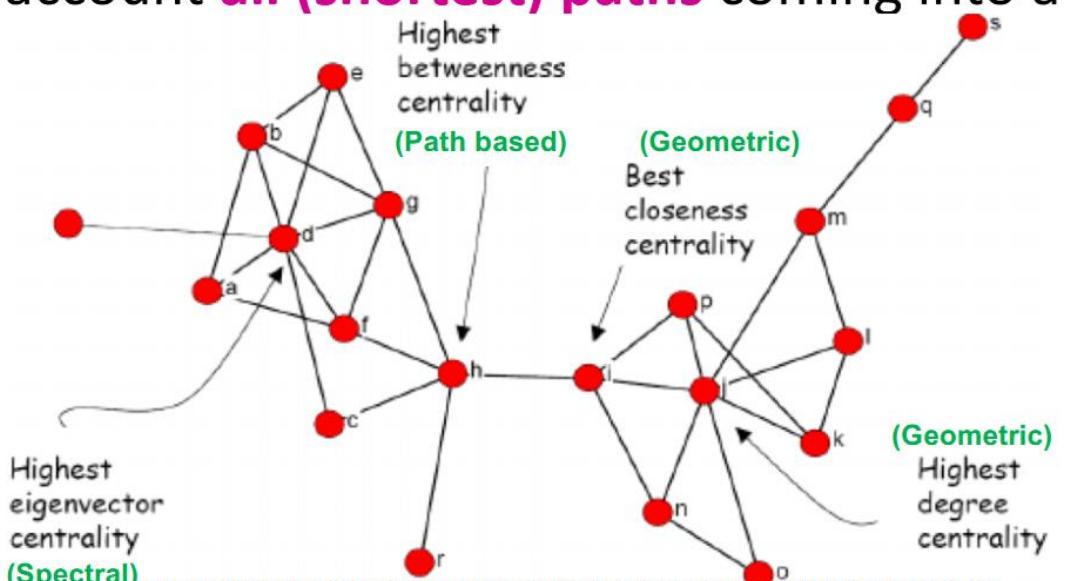
- Importance is a **function of distances** to other nodes.

■ Spectral Measures:

- Based on the **eigen-structure** of some graph-related matrix

■ Path-based Measures:

- Take into account **all (shortest) paths** coming into a node





02

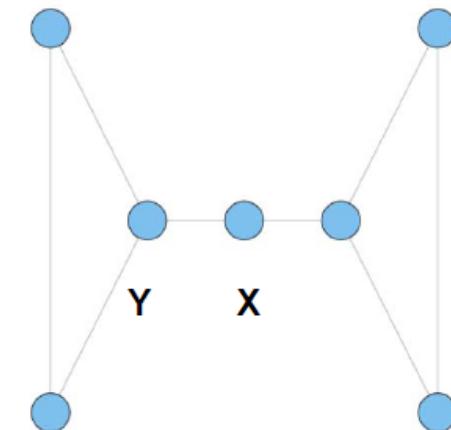
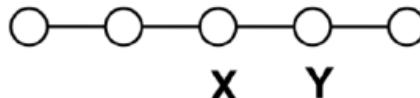
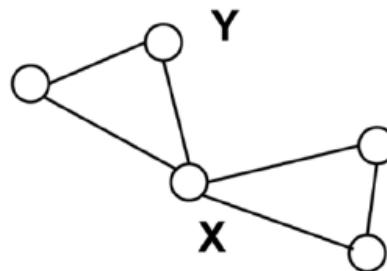
Path-based Measures for Centrality

Betweenness Centrality

- Measure of centrality in a graph based on shortest paths:
 - Edge Betweenness Centrality
 - Node Betweenness Centrality
- In a telecommunications network, a node with higher betweenness centrality would have more control over the network, because more information will pass through that node.

Betweenness Centrality

- Intuition: how many pairs of individuals would have to go through you in order to reach one another in the minimum number of hops?
- Who has higher betweenness, X or Y?



Edge Betweenness Centrality

$$\rho_{ij} = \sum_{p \neq q} \left(\Gamma_{pq}(e_{ij}) / \Gamma_{pq} \right)$$

Γ_{pq} is the number of shortest paths from the p -th to the q -th node

$\Gamma_{pq}(e_{ij})$ is the number of these paths making use of e_{ij} .

- Usually the betweenness is normalized by $[(n-1)(n-2)/2]$

Number of possible edges

Node Betweenness Centrality

$$C_i = \sum_{p \neq i \neq q} \left(\Gamma_{pq}(i) / \Gamma_{pq} \right)$$

Γ_{pq} is the number of shortest paths from the p -th to the q -th node

$\Gamma_{pq}(i)$ is the number of these shortest paths making use of the i -th node (except those that are start or end nodes is i).

$$\lceil (n-1)(n-2)/2 \rceil$$

number of pairs of vertices
excluding the vertex itself

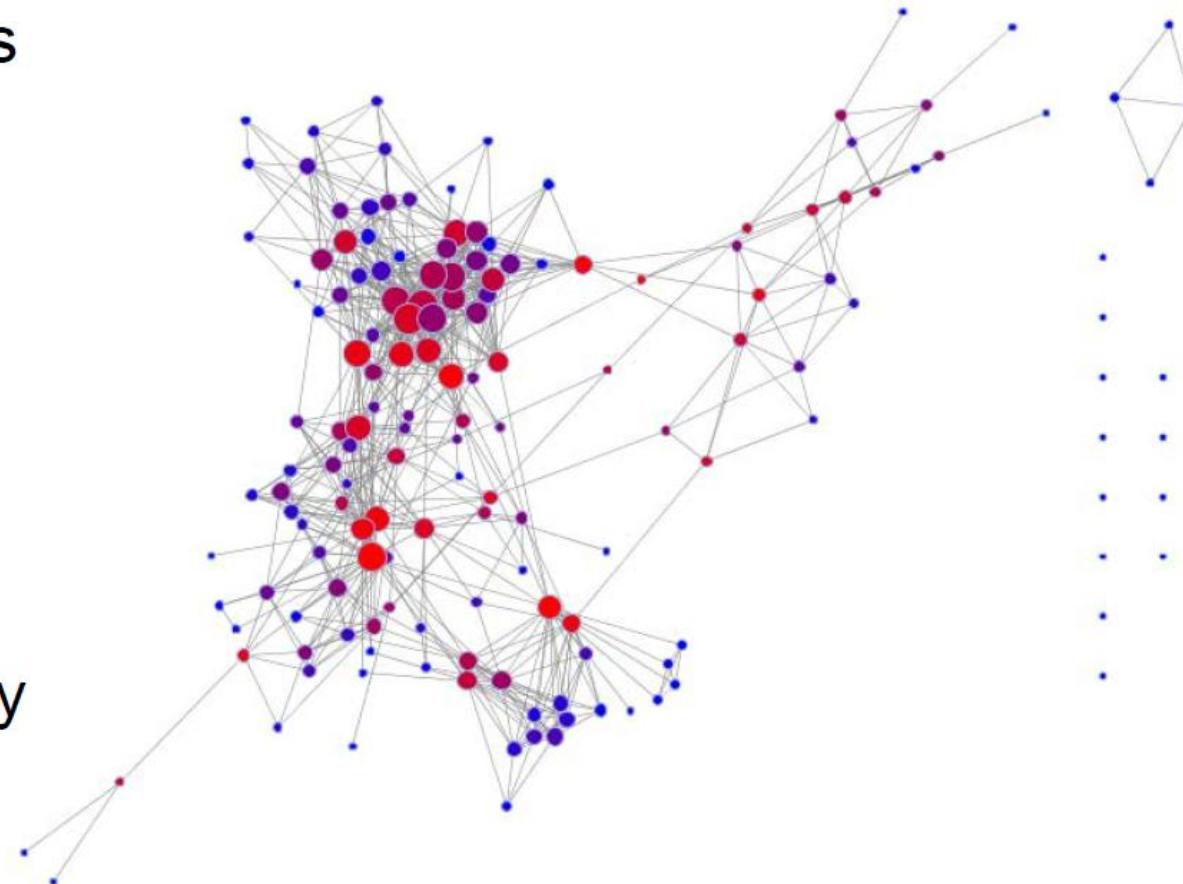
- Usually the betweenness is normalized by

Betweenness centrality: an example

Nodes are sized by degree, and colored by betweenness.

Can you spot nodes with high betweenness but relatively low degree? Explain how this might arise.

What about high degree but relatively low betweenness?



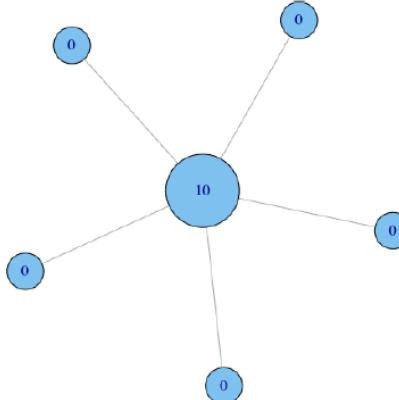
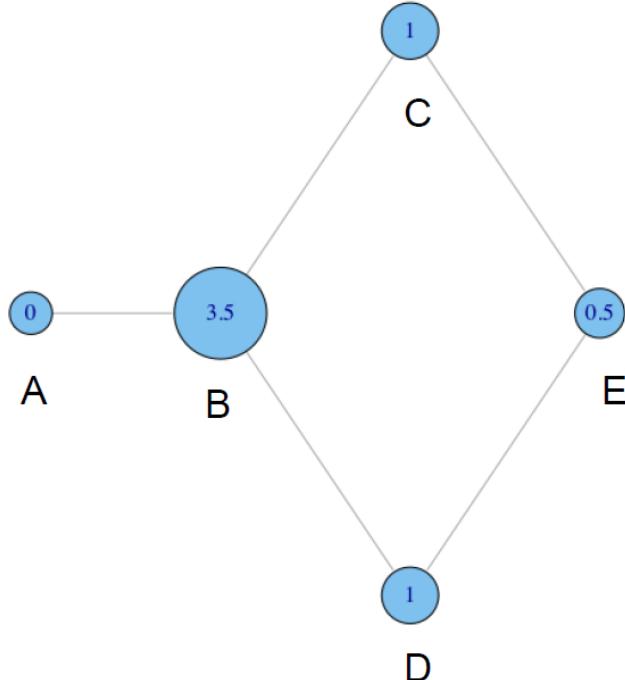
Betweenness Centrality

Why do C and D each have betweenness 1?

They are both on shortest paths for pairs (A,E), and (B,E), and so must share credit:

$$\frac{1}{2} + \frac{1}{2} = 1$$

Can you figure out why B has betweenness 3.5 while E has betweenness 0.5?



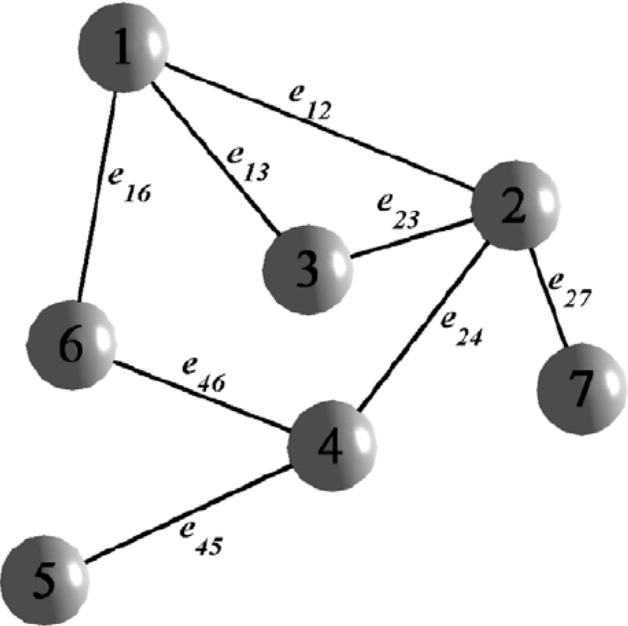
Connection Graph Stability Scores

- In some applications the importance of the shortest paths are also important
- The importance of the shortest paths might be different
- The more important paths making use of an edge the more its importance
- A simple measure of importance would be the length
- The connection graph stability (CGS) method takes into account this issue
- It has application in synchronization analysis
- The CGS-score b_{ij} for the link between the nodes i and j is defined as

$$b_{ij} = \sum_{u=1}^{n-1} \sum_{v>u; e_{ij} \in P_{uv}} |P_{uv}|$$

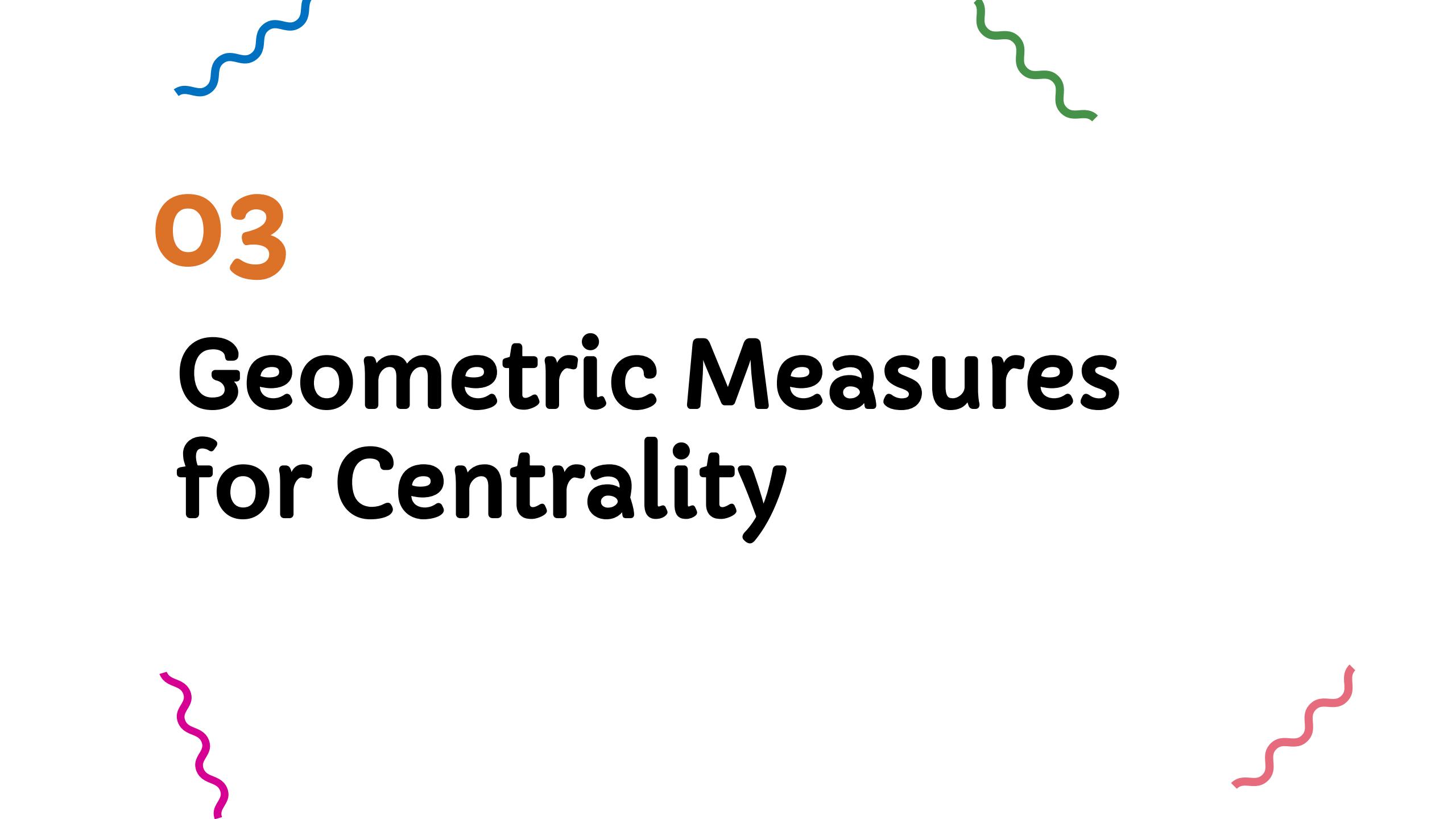
- $|P_{uv}|$: length of path P_{uv} between the nodes u and v

Connection Graph Stability Scores



$$\begin{aligned}P_{12} &= e_{12}, P_{13} = e_{13}, P_{14} = e_{12}e_{24}, \\P_{15} &= e_{16}e_{46}e_{45}, P_{16} = e_{16}, P_{17} = e_{12}e_{77}, \\P_{23} &= e_{23}, P_{24} = e_{24}, P_{25} = e_{24}e_{45}, \\P_{26} &= e_{24}e_{46}, P_{27} = e_{27}, P_{34} = e_{23}e_{24}, \\P_{35} &= e_{23}e_{24}e_{45}, P_{36} = e_{13}e_{16}, P_{37} = e_{23}e_{27}, \\P_{45} &= e_{45}, P_{46} = e_{46}, P_{47} = e_{24}e_{27}, P_{56} = \\&e_{45}e_{46}, P_{57} = e_{45}e_{24}e_{27}, P_{67} = e_{16}e_{12}e_{27}\end{aligned}$$

$$\begin{aligned}b_{12} &= |P_{12}| + |P_{14}| + |P_{17}| + |P_{67}| = 1 + 2 + 2 + 3 = 8, \\b_{13} &= |P_{13}| + |P_{36}| = 1 + 2 = 3, \\b_{16} &= |P_{15}| + |P_{16}| + |P_{36}| + |P_{67}| = 3 + 1 + 2 + 3 = 9, \\b_{23} &= |P_{23}| + |P_{34}| + |P_{35}| + |P_{37}| = 1 + 2 + 3 + 2 = 8, \\b_{24} &= |P_{14}| + |P_{24}| + |P_{25}| + |P_{26}| + |P_{34}| + |P_{35}| + |P_{47}| \\&+ |P_{57}| = 2 + 1 + 2 + 2 + 2 + 3 + 2 + 3 = 17, \\b_{27} &= |P_{17}| + |P_{27}| + |P_{37}| + |P_{47}| + |P_{57}| + |P_{67}| = \\2 + 1 + 2 + 2 + 3 + 3 = 13, \\b_{45} &= |P_{15}| + |P_{25}| + |P_{35}| + |P_{45}| + |P_{56}| + |P_{57}| = \\3 + 2 + 3 + 1 + 2 + 3 = 14, \\b_{46} &= |P_{15}| + |P_{26}| + |P_{46}| + |P_{56}| = 3 + 2 + 1 + 2 = 8.\end{aligned}$$

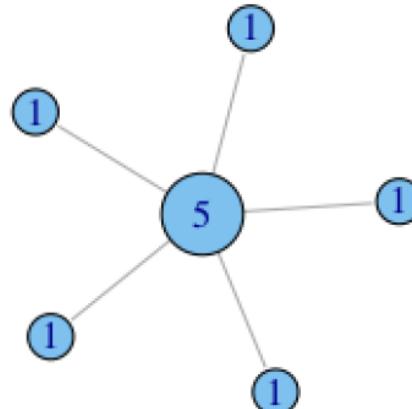


03

Geometric Measures for Centrality

Degree centrality (undirected)

The more the friends the more the importance (the richer the better)

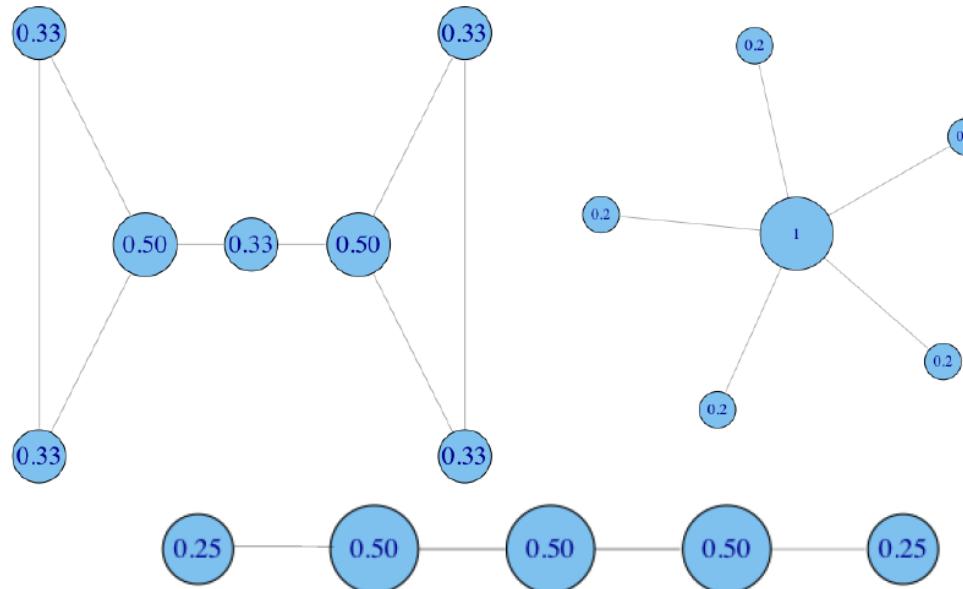


Normalized degree centrality:

Degree is divided by the max. possible, i.e. $(N-1)$

When is the number of connections the best centrality measure?

- people who will do favors for you
- people you can talk to / have a drink with

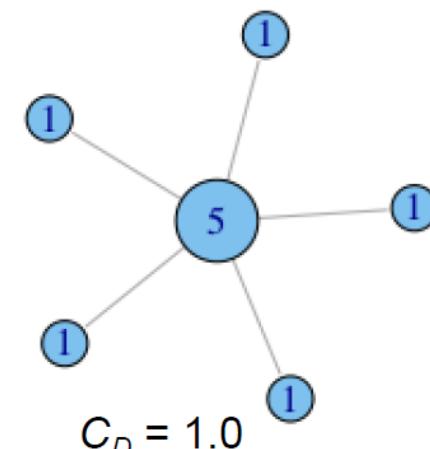
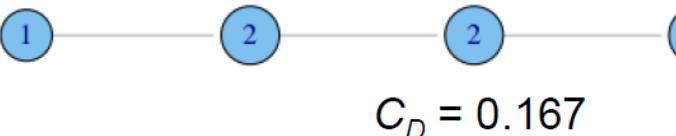
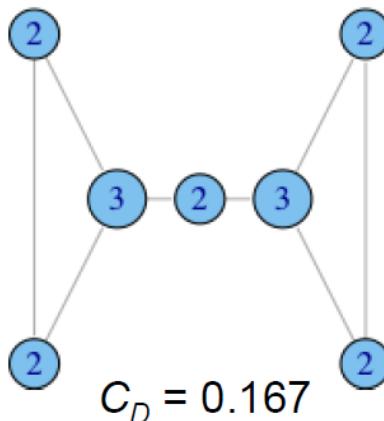


How equal are the nodes?

- How much variation is there in the centrality scores among the nodes?
- Freeman's general formula for centralization (can use other metrics, e.g. Gini coefficient or standard deviation):

$$C_D = \frac{\sum_{i=1}^g [C_D(n^*) - C_D(i)]}{[(N-1)(N-2)]}$$

maximum value in the network



Gini Coefficient (Index)

The bar chart on the left shows a simple distribution of incomes. The total population is split up in 5 parts and ordered from the poorest to the richest 20%. The bar chart shows how much income each 20% part of the income distribution earns.

The chart on the right shows the same information in a different way, both axis show the cumulative shares:

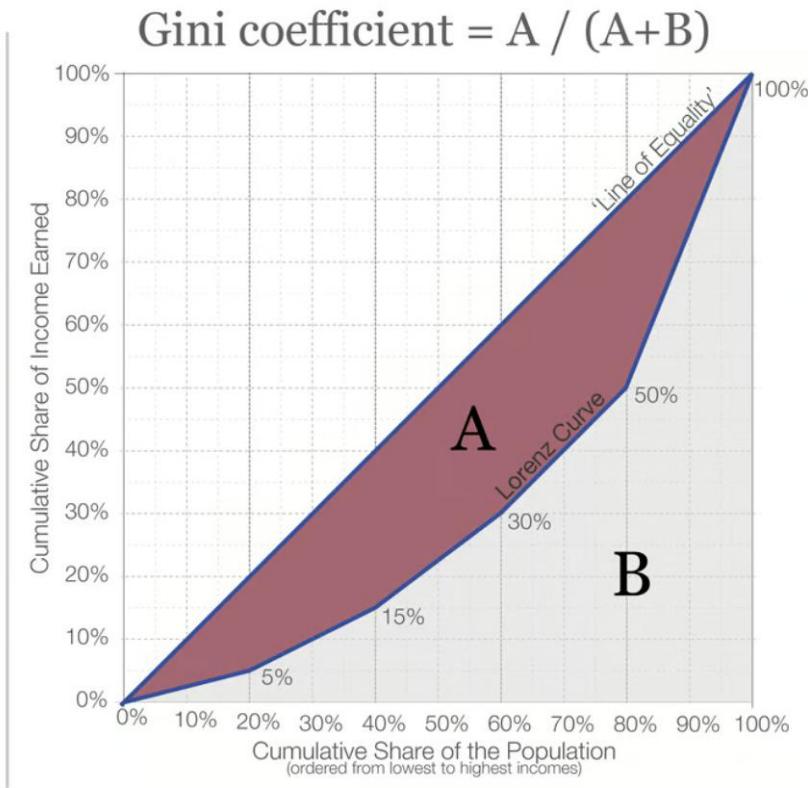
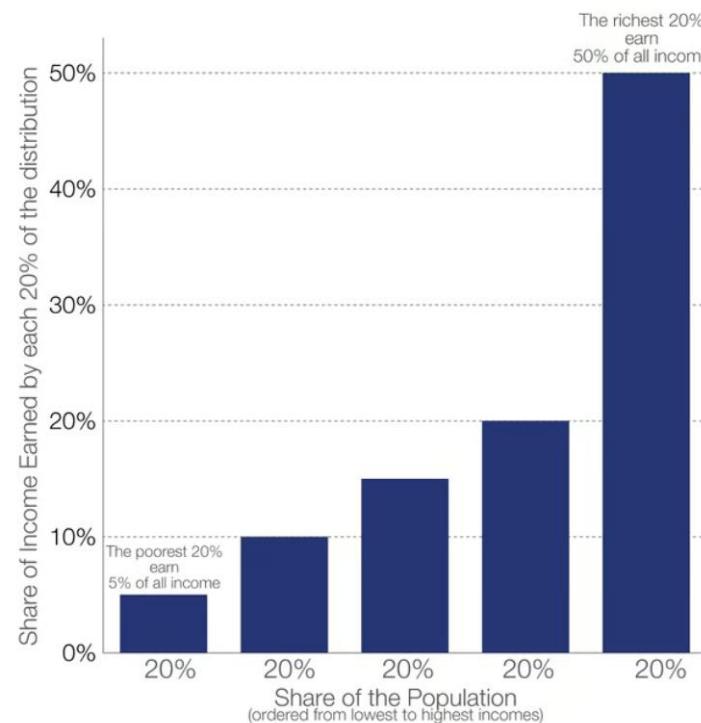
The poorest 20% of the population earn 5% of the total income, the next 20% earn 10% – so that the poorest 40% of the population earn 15% etc.

The curve resulting from this way of displaying the data is called the Lorenz Curve.

If there was no income inequality the resulting Lorenz Curve would be a straight line – the ‘Line of Equality’.

A larger area (A) between the Lorenz Curve and the Line of Equality means a higher level of inequality.

The ratio of $A/(A+B)$ is therefore a measure of inequality and is referred to as the Gini coefficient, Gini index, or simply the Gini.



Gini Coefficient (Index)



Information Sciences

Volume 462, September 2018, Pages 16-39



Sparsity measure of a network graph: Gini index

Swati Goswami ^{a b 1}   , C.A. Murthy ^a, Asit K. Das ^b

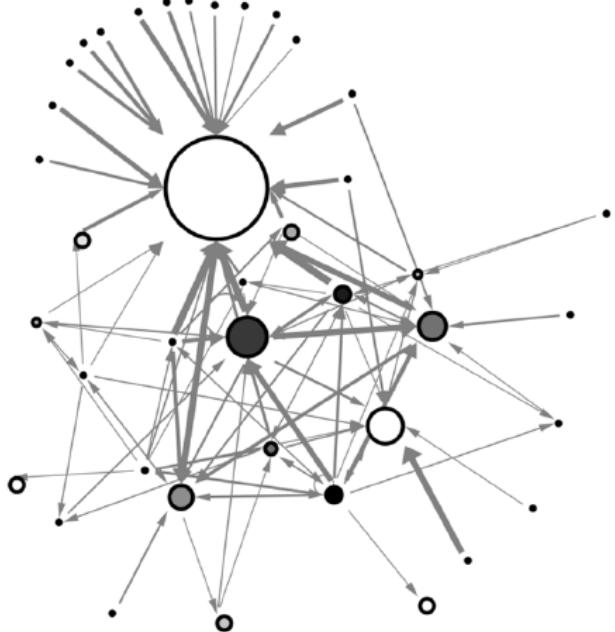
Show more 

 Add to Mendeley  Share  Cite

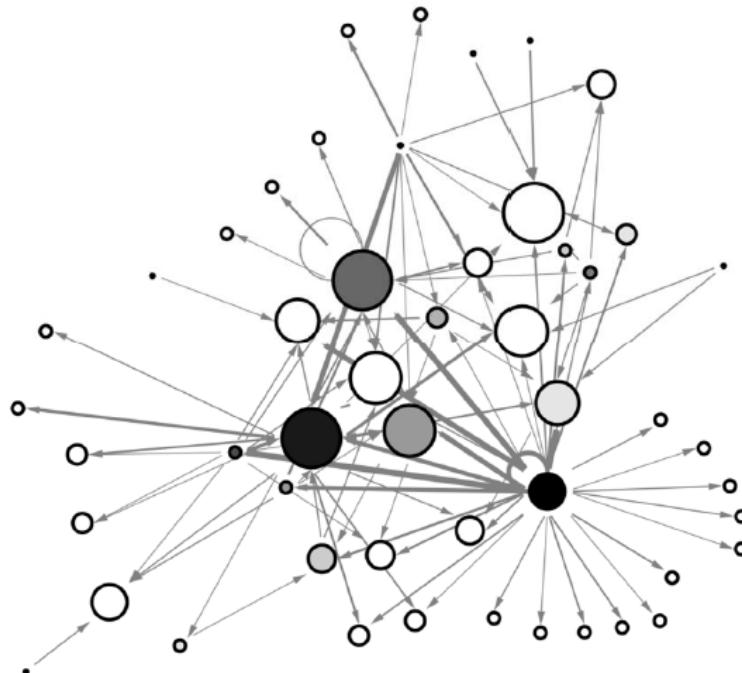
<https://doi.org/10.1016/j.ins.2018.05.044> 

[Get rights and content](#) 

Example: financial trading networks



high centralization: one node
trading with many others



low centralization: trades
are more evenly distributed

Characteristic path length

- A network with N nodes
 - Compute the shortest path (distance) between any two nodes d_{ij}
 - The length of the path is the number of edges (unweighted networks) or the weighted sum of the edges (weighted networks)
 - If the nodes are not connected, the path length between them is set to infinity
 - It is also called average geodesic distance
 - If d_{ij} is infinity, it diverges
 - Many times we compute the average only over the connected pairs of nodes (that is, we ignore “infinite” length paths)
 - $$l = \frac{1}{N(N-1)} \sum_{i,j, i \neq j} d_{ij}$$

Efficiency

- In this way the divergence is avoided
- The inverse of efficiency E is called harmonic mean
- Efficiency is an indicator of the traffic capacity of the network
- The couple of disconnected nodes have a contribution of zero in computing E
- The more the values of E are the more the communication-efficient the network is
- It is also called global efficiency of the network.

$$E = \frac{1}{N(N-1)} \sum_{i,j, i \neq j} \frac{1}{d_{ij}}$$

Efficiency

- Higher Efficiency → Faster communication, better connectivity.
- Why is Efficiency Important?
 -  High Efficiency → Network is Well-Connected
 - Information spreads quickly.
 - Fewer intermediate steps needed.
 - Helps in optimizing transportation, communication, and social interactions.
 -  Low Efficiency → Poor Connectivity
 - Long paths between nodes.
 - Slower communication and bottlenecks.
 - Less effective in handling information flow.

$$E = \frac{1}{N(N-1)} \sum_{i,j, j \neq i} \frac{1}{d_{ij}}$$

The Role of Connectivity in Efficiency

| Network Type | Effect on E |
|--|---|
| Fully Connected (Strongly Connected Component - SCC) | ▲ High E, every node can reach any other node efficiently. |
| Weakly Connected Component (WCC) | ▼ Lower E, some nodes may only be accessible in one direction. |
| Disconnected Components | ▬ Very Low E, as some distances are infinite (ignored in practical calculations). |

The Power of Shortcuts

- A shortcut in a graph is an additional edge that significantly reduces the shortest path distance between two nodes without being essential for the overall connectivity of the network.
- Extra direct connections between distant nodes.
- Reduce the shortest path distance d_{ij} .
- Improve efficiency without requiring full connectivity.
- Real-World Examples:
 - Social Networks → Influencers create bridges between distant groups.
 - Transport Networks → Highways and express routes reduce travel time.
 - Computer Networks → Fast routing through backbone connections (CDNs).
- Fewer hops → Shorter paths → Higher Efficiency.

Vulnerability

- It is important to know which component (nodes or edges) are crucial to the best performance.
- The more the drop in the efficiency by removing a component the more crucial that component.
- Degree (hub node) might be a criterion
 - Only degree is not enough, e.g. all vertices of a binary tree network have equal degree, i.e. no hub, but disconnection of vertices closer to the root and the root itself have a greater impact than of those near the leaves.
- The amount of change in the efficiency (or other network properties) as a component is removed can be an indicator of the vulnerability

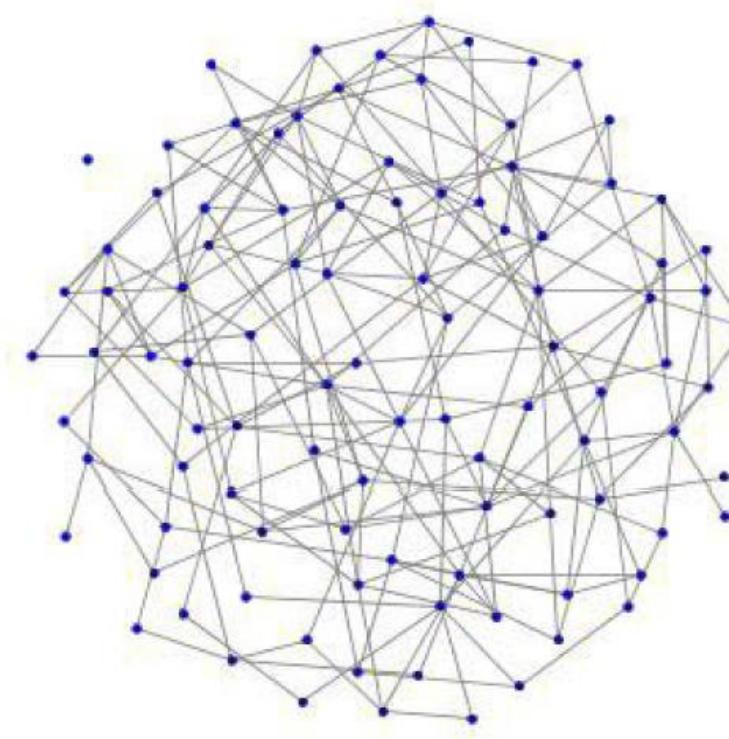
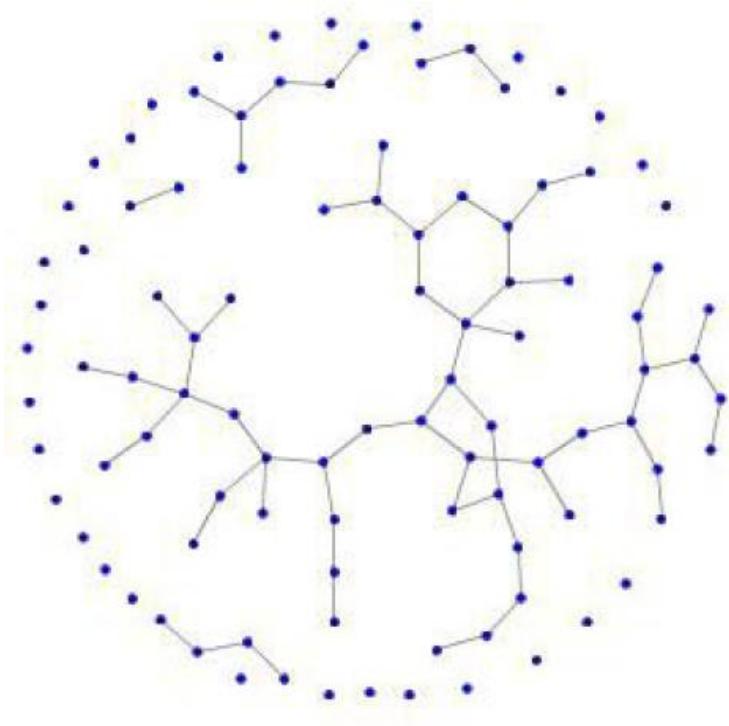
Vulnerability

$$V_i = \frac{E - E_i}{E} \quad V = \max_i V_i$$

- where V_i is the vulnerability of component i and E_i is the efficiency of the networks by removing that component.
- V can be regarded as the vulnerability of the network the ordered distribution of nodes with respect to their vulnerability V_i is related to the network hierarchy.
- The most vulnerable (critical) node occupies the highest position in the network hierarchy.
- The same is also true for the edges.

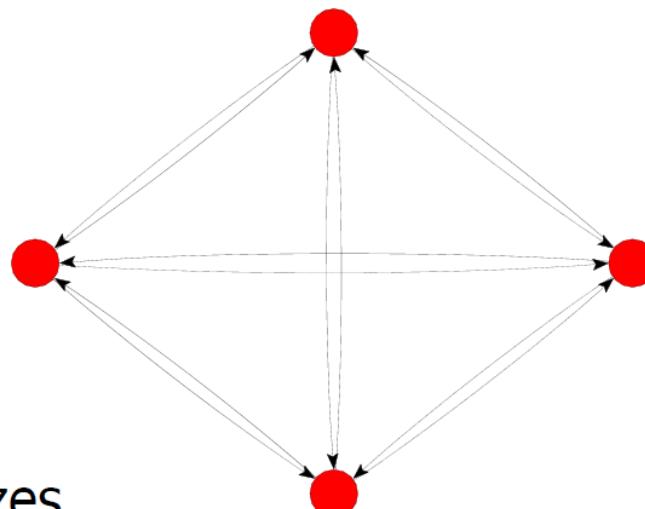
Density

- How dense the networks are?



Density

- Number of the connections that may exist between n nodes
 - directed graph
 $e_{\max} = n*(n-1)$
each of the n nodes can connect to $(n-1)$ other nodes
 - undirected graph
 $e_{\max} = n*(n-1)/2$
since edges are undirected, count each one only once
- What fraction are present?
 - density = e/ e_{\max}
 - For example, out of 12 possible connections, this graph has 7, giving it a density of $7/12 = 0.583$
- Would this measure be useful for comparing networks of different sizes (different numbers of nodes)?



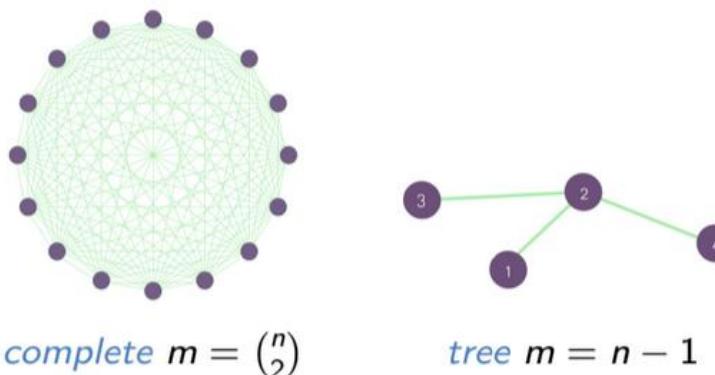
Density

- for *undirected G density* ρ is defined as

$$\rho = \frac{2m}{n(n-1)} = \frac{\langle k \rangle}{n-1}$$

- for *directed G density* ρ^* is defined as

$$\rho^* = \frac{m}{n(n-1)} = \frac{\langle k^* \rangle}{n-1}$$



- G is *dense* if $\rho \rightarrow \text{const.}$ as $n \rightarrow \infty$ thus $\langle k \rangle = \mathcal{O}(n)$
- G is *sparse* if $\rho \rightarrow 0$ as $n \rightarrow \infty$ thus $\langle k \rangle \neq \mathcal{O}(n)$

Closeness

- What if it's not so important to have many direct friends?
 - Degree Centrality is not important
- Or be “between” others
 - Betweenness Centrality is not important
- But one still wants to be in the “middle” of things, not too far from the center.
- Closeness is based on the length of the average shortest path between a node and all other nodes in the network

Example

- **High Degree Centrality**
 - A person who has many direct friends in a social network.
- **High Betweenness Centrality**
 - A person who acts as a bridge between two separate groups.
- **High Closeness Centrality**
 - A person who can reach anyone in the network with the fewest intermediaries, even if they don't have many direct friends.
- Closeness Centrality focuses on being "near" other nodes in terms of short paths across the entire network, rather than just having many direct connections!

Closeness

Formula:

Closeness Centrality:

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

▪ Closeness Centrality:

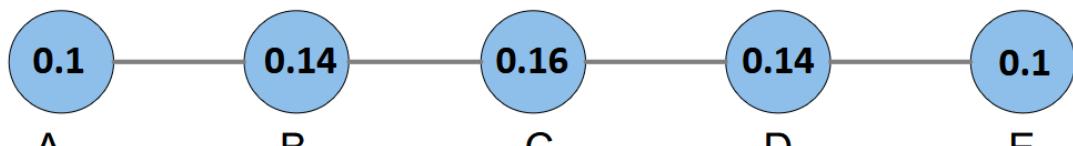
$$c_{\text{clos}}(x) = \frac{1}{\sum_y d(y, x)}$$

length of the shortest path from x to y

Normalized Closeness Centrality

$$C'_c(i) = \left[\sum_{j=1}^N d(i,j) / (N-1) \right]^{-1}$$

- How much a vertex can communicate without relying on third parties for his messages to be delivered

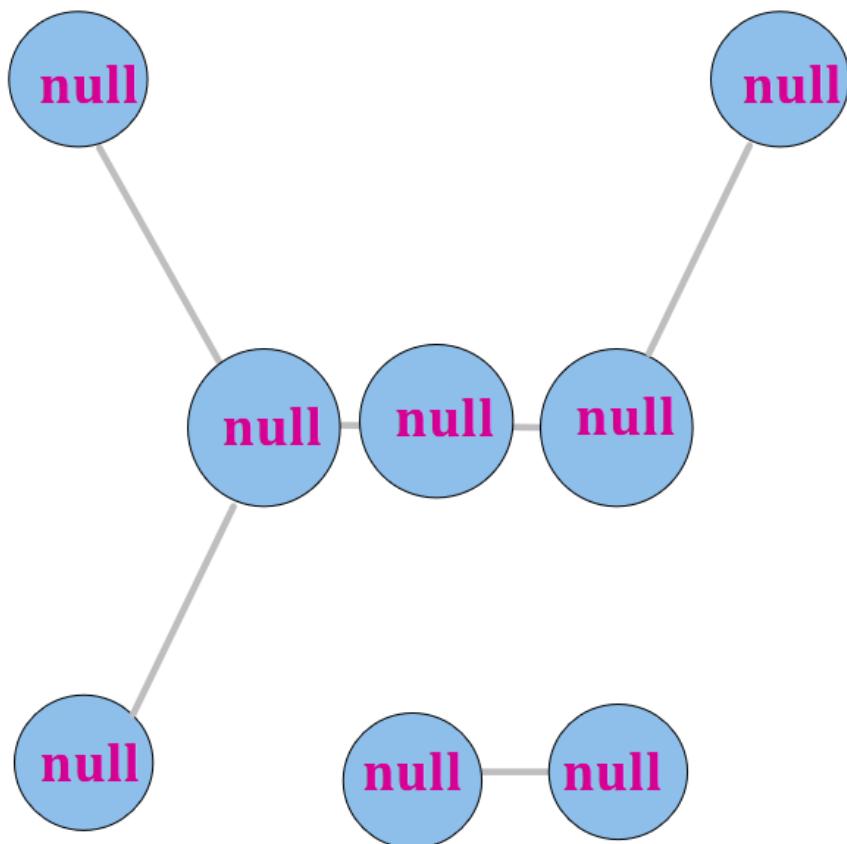


$$c_{\text{clos}}(A) = \frac{1}{1+2+3+4} = 0.1$$

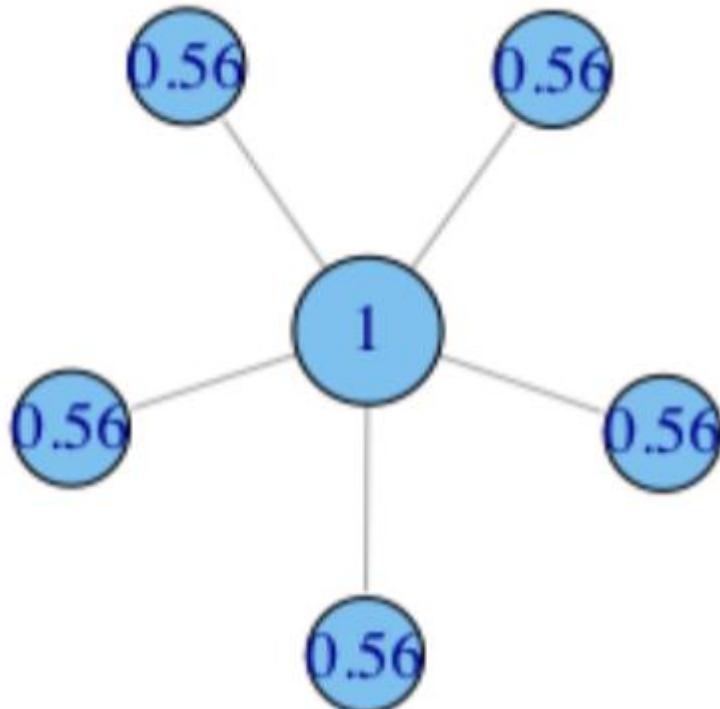
$$C'_c(A) = \left[\frac{\sum_{j=1}^N d(A,j)}{N-1} \right]^{-1} = \left[\frac{1+2+3+4}{4} \right]^{-1} = \left[\frac{10}{4} \right]^{-1} = 0.4$$

- Problem: The graph must be (strongly) connected!

Closeness Example



We get null score for all nodes,
if the graph is not connected!



Harmonic Centrality

■ Geometric measures

■ Harmonic Centrality:

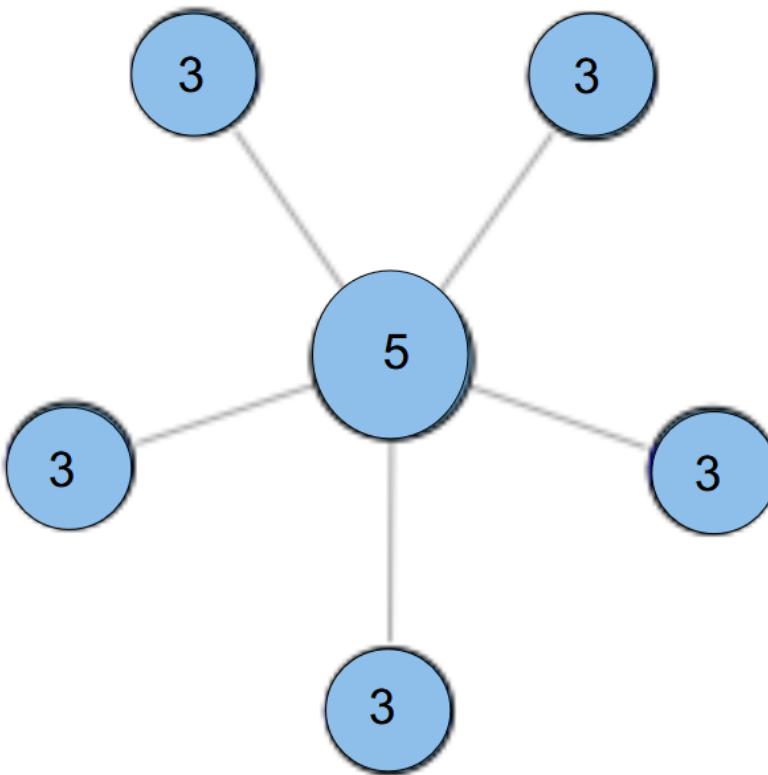
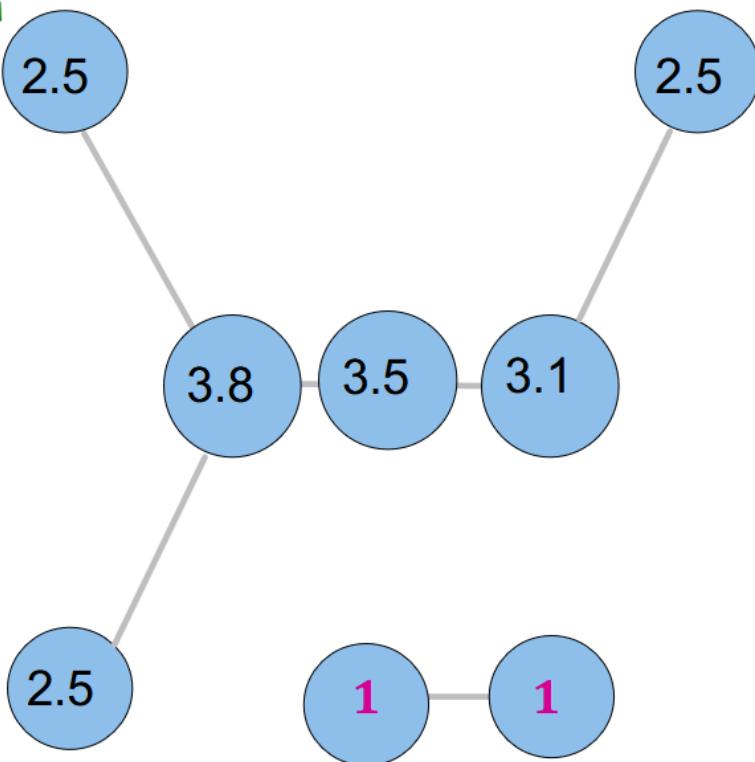
- Replace the average distance with the harmonic mean of all distances.
- The $n(n - 1)$ distances between every pair of distinct nodes:

$$c_{\text{har}}(x) = \frac{\text{Harmonic mean}}{\sum_{y \neq x} \frac{1}{d(y, x)}} = \sum_{d(y,x) < \infty, y \neq x} \frac{1}{d(y, x)}$$

- Strongly correlated to closeness centrality
- Naturally also accounts for nodes y that cannot reach x
- Can be applied to graphs that are **not strongly connected**

Harmonic Centrality Example

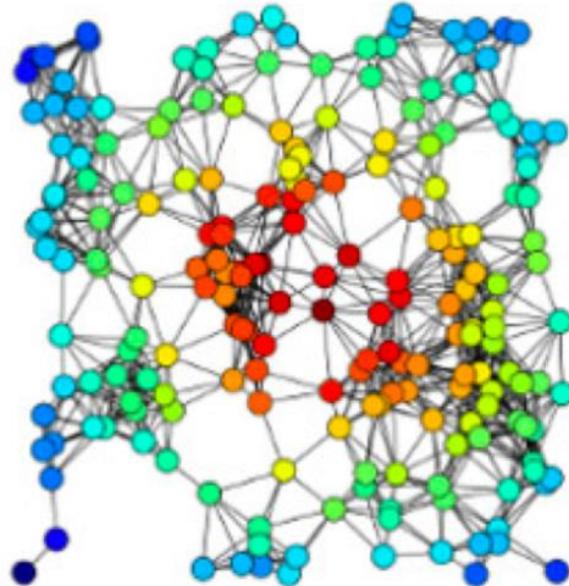
$$c_{harm} = \frac{1}{1} + \frac{1}{2} + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} = 2.5$$



Comparison

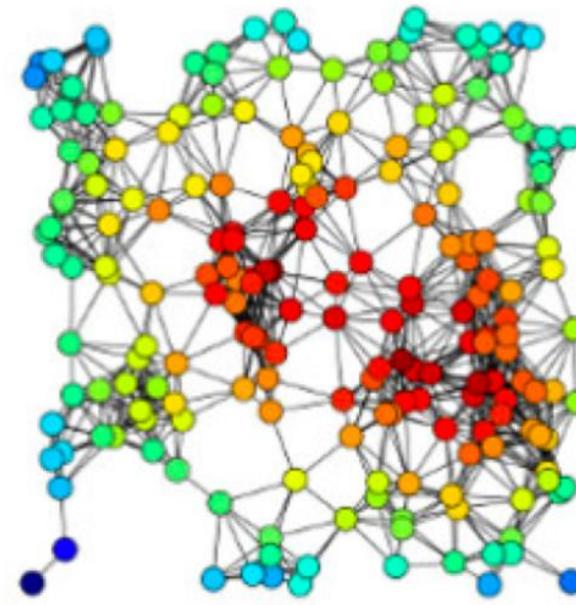
- Closeness Centrality is affected by average distances, while Harmonic Centrality is influenced by nearby nodes.
- A node can be well-positioned (high Closeness) but still have many distant nodes that lower its Harmonic score.
- Harmonic Centrality gives more weight to close neighbors, whereas Closeness considers all distances equally.

Closeness vs Harmonic Centrality



Closeness

**Red nodes are closer to all
the other nodes**



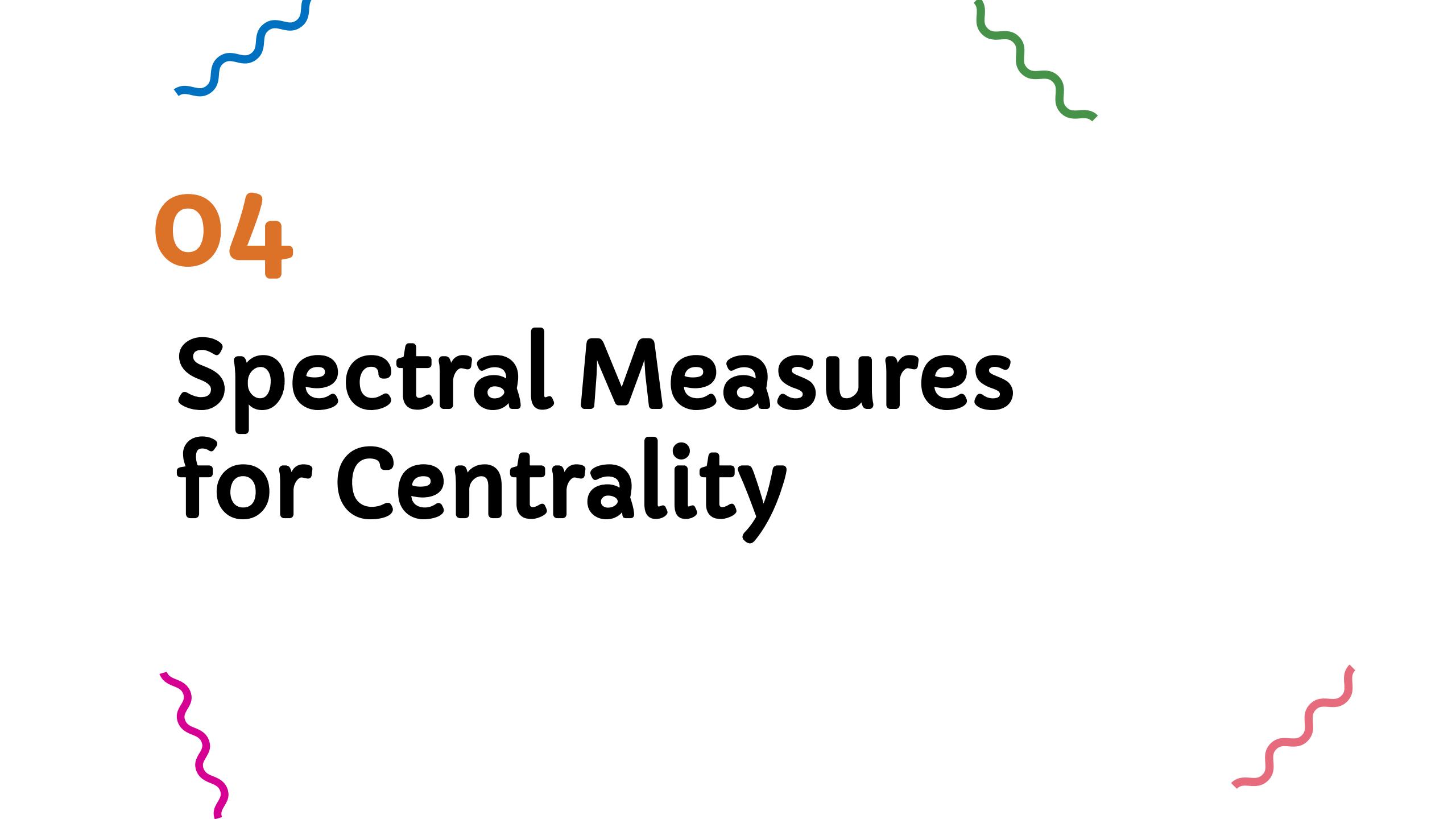
Harmonic

**Red nodes are closer to all the other
nodes, and have larger degrees**

Examples of Closeness centrality, and Harmonic Centrality of the same graph.

Let's Think

- Can a Node Have High Harmonic Centrality but Low Degree?
 - Imagine you have only two friends, but they are well-connected influencers in the network. Your degree is low (only 2 connections). However, because your friends have strong connections, you can quickly reach many people.
- Can a Node Have High Degree but Low Harmonic Centrality?
 - Imagine a node has 10 direct connections, but all these connections are to each other and not to the rest of the network. Degree is high (10 connections), but reaching other parts of the network requires multiple hops.
- Can a Node Have High Closeness but Low Harmonic Centrality?
 - Consider the tree or ring networks.
- Can a Node Have Low Closeness but High Harmonic Centrality?
 - In a tree structure, a node close to a highly connected hub can still have low Closeness (because reaching deep parts of the tree takes many steps). However, its Harmonic Centrality is high because it can reach nearby nodes very efficiently.



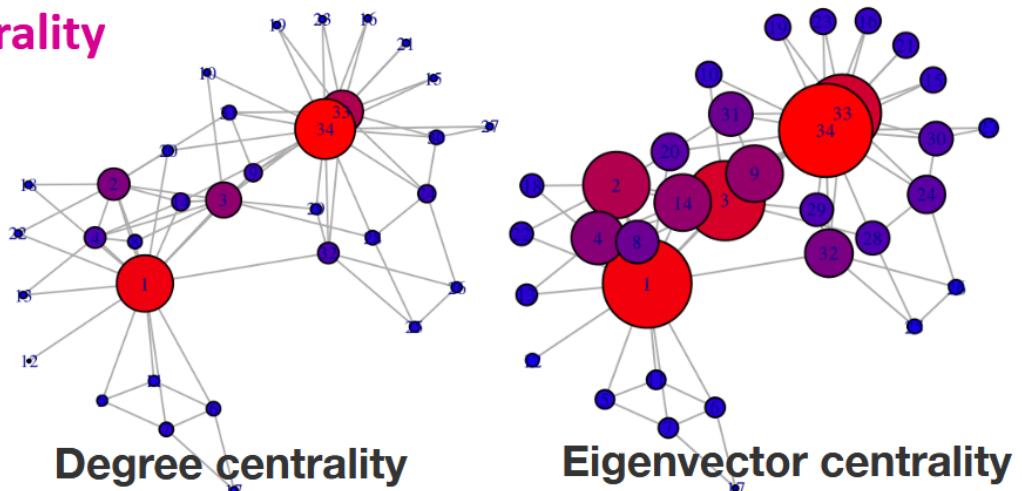
04

Spectral Measures for Centrality

Spectral Measures

■ Spectral measures

- Compute the left dominant eigenvector of some matrix derived from the graph
- **Idea:** A node's centrality is a function of the **centrality of its neighbors**
 - Nodes connected to central nodes has a larger centrality score than those connected to non-central nodes.
 - **Eigenvector Centrality**
 - **Katz's Index**
 - **Page Rank**
 - **Hits**



Eigencentrality

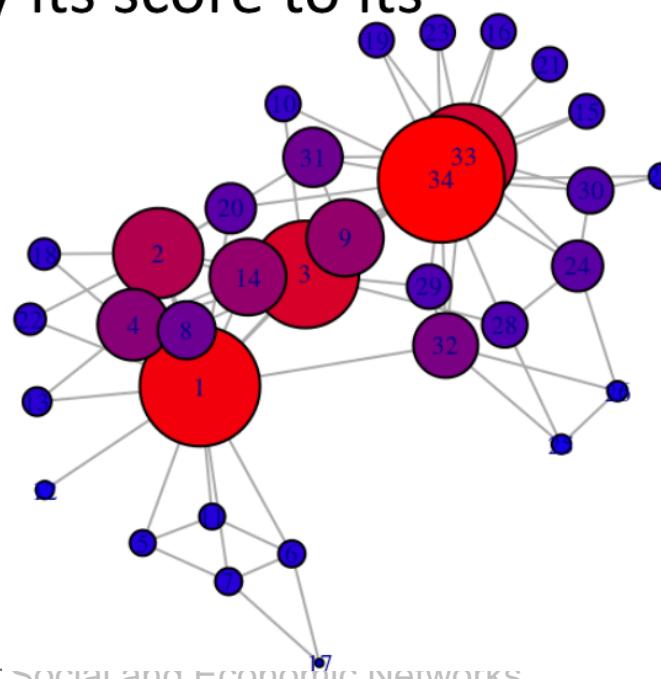
■ Spectral measures

- **Eigenvector Centrality**: Measure of the **influence** of a node in a network
- **Idea**: Every node starts with the same score, and then each node gives away its score to its successors

$$c_{\text{eig}}(x) = \frac{1}{\lambda} \sum_{y \rightarrow x} c_{\text{eig}}(y)$$

Normalization constant = $\|c_{\text{eig}}\|_2$

- **Intuitively**: Degree counts walks of length one, the eigenvalue centrality counts walks of length infinity



Eigencentrality

■ Spectral measures

- **Eigenvector Centrality:** Measure of the **influence** of a node in a network:

$$c_{\text{eig}}(x) = \frac{1}{\lambda} \sum_{y \rightarrow x} c_{\text{eig}}(y)$$

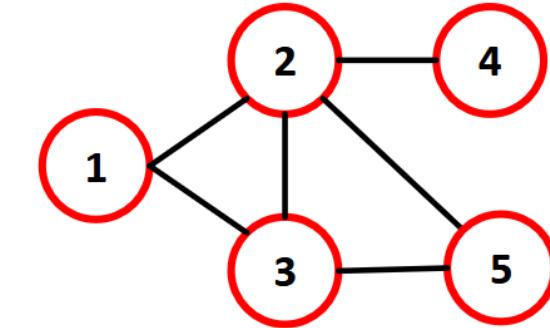
- c_{eig} converges to the dominant eigenvector of adj. matrix A
- λ converges to the dominant eigenvalue of adj. matrix A
- Equivalently, eigencentrality is the eigenvector corresponding to the dominant eigenvalue (λ) of A
$$AX = \lambda X$$
- **Problem:** Graph should be **strongly connected!**

How to compute Eigencentrality?

■ Power Iteration:

- Set $c^{(0)} \leftarrow 1, k \leftarrow 1$
- 1: $c^{(k)} \leftarrow Ac^{(k-1)}$
- 2: $c^{(k)} = c^{(k)}/\|c^{(k)}\|_2$
- 3: If $\|c^{(k)} - c^{(k-1)}\| > \varepsilon$:
- 4: $k \leftarrow k + 1$, goto 1

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix} \quad \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix}$$



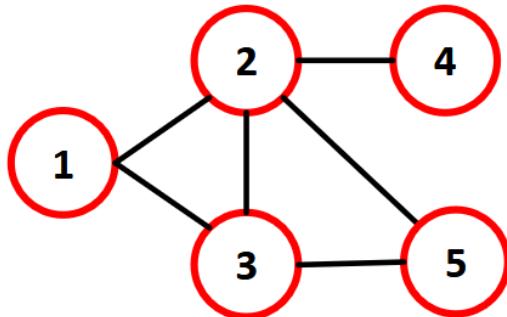
$$c = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

How to compute Eigencentrality?

■ Power Iteration:

Iteration 1

$$\begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 3 \\ 1 \\ 2 \end{bmatrix} \equiv \begin{bmatrix} 0.34 \\ 0.68 \\ 0.51 \\ 0.17 \\ 0.34 \end{bmatrix} \\ A & c^{(0)} & c^{(1)} = Ac^{(0)} & c^{(1)} = c^{(1)}/\|c^{(1)}\|_2 \end{matrix}$$



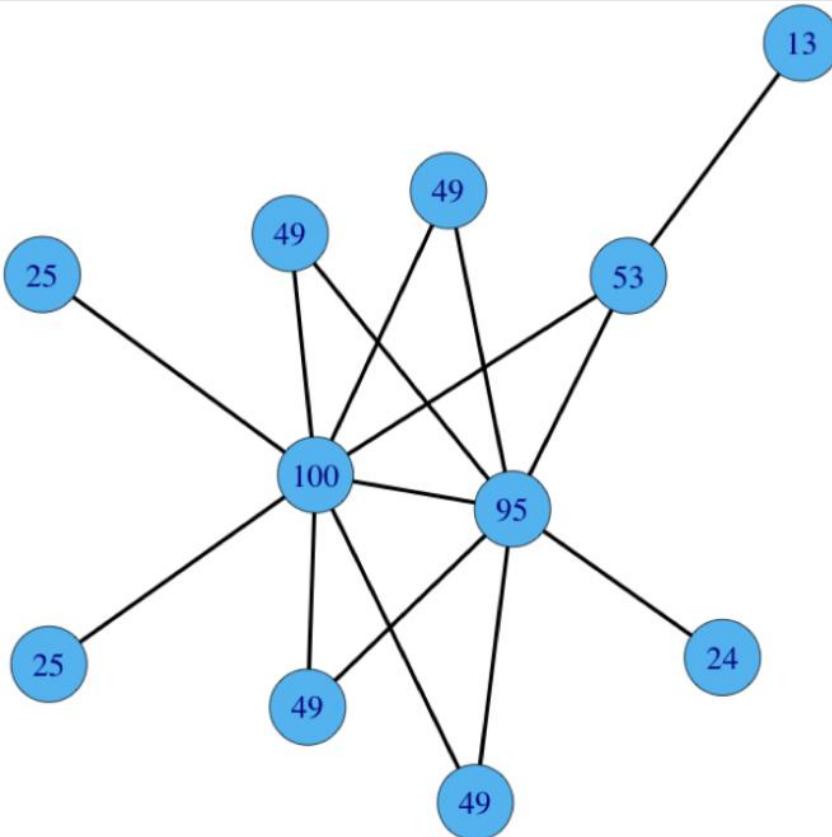
Iteration 2

$$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.34 \\ 0.68 \\ 0.51 \\ 0.17 \\ 0.34 \end{bmatrix} = \begin{bmatrix} 1.19 \\ 1.36 \\ 1.36 \\ 0.68 \\ 1.19 \end{bmatrix} \equiv \begin{bmatrix} 0.45 \\ 0.51 \\ 0.51 \\ 0.25 \\ 0.45 \end{bmatrix}$$

Iteration 3

$$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.45 \\ 0.51 \\ 0.51 \\ 0.25 \\ 0.45 \end{bmatrix} = \begin{bmatrix} 1.02 \\ 1.66 \\ 1.41 \\ 0.51 \\ 1.02 \end{bmatrix} \equiv \begin{bmatrix} 0.38 \\ 0.62 \\ 0.53 \\ 0.19 \\ 0.38 \end{bmatrix} \dots c = \begin{bmatrix} 1 \\ 1.41 \\ 1.27 \\ 0.52 \\ 1 \end{bmatrix}$$

Example



Eigenvalue centrality counts walks of length infinity

Katz's Index

■ Spectral measures

- **Katz's Index:** Measures **influence** by taking into account the **total number of walks** between a pair of nodes

$$c_{\text{katz}}(x) = \beta \sum_{k=0}^{\infty} \sum_{x \rightarrow y} \alpha^k (A^k)_{xy}$$

Total number of walks
of length k between x, y

- α is an attenuation factor in range $(0, \frac{1}{\lambda})$, where λ is the largest eigenvalue of A
- β is to give some nodes more privilege
- **Long paths are weighted less than short ones**

Katz's Index

Spectral measures

- Katz's Index: Measures influence by taking into account the total number of walks between a pair of nodes

$$A^0 = I$$

$$c_{\text{katz}}(x) = \beta \sum_{k=0}^{\infty} \sum_{x \rightarrow y} \alpha^k (A^k)_{xy}$$

$$A^1 = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

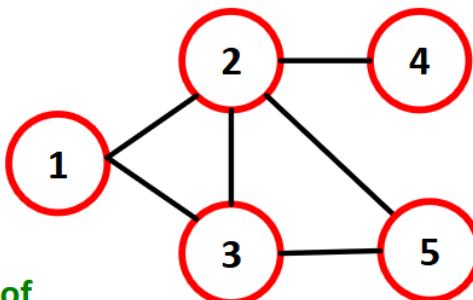
$\alpha < 1$: Long paths are weighted less

Total number of walks of length k between x, y

$$A^2 = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix}^2 = \begin{bmatrix} 2 & 1 & 1 & 1 & 2 \\ 1 & 4 & 2 & 0 & 1 \\ 1 & 2 & 3 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 \\ 2 & 1 & 1 & 1 & 2 \end{bmatrix}$$

$$A^3 = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix}^3 = \begin{bmatrix} 2 & 6 & 5 & 1 & 2 \\ 6 & 4 & 6 & 4 & 6 \\ 5 & 6 & 4 & 2 & 5 \\ 1 & 4 & 2 & 0 & 1 \\ 2 & 6 & 5 & 1 & 2 \end{bmatrix}$$

Number of walks of length 3 between 2, 5
(2,1,3,5), (2,4,2,5), (2,3,2,5),
(2,1,2,5), (2,5,3,5), (2,5,2,5)



How to compute Katz's Index?

■ Spectral measures

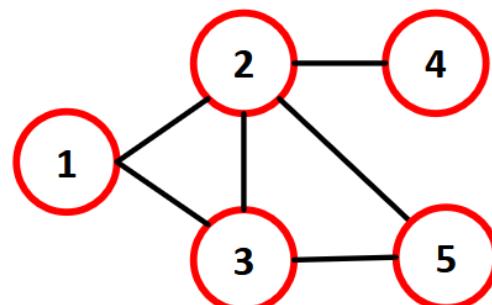
- **Katz's Index:** Give each node a small amount of centrality for free

$$c_{\text{Katz}}(x) = \alpha \sum_{y \rightarrow x} (c_{\text{Katz}}(y) + \beta)$$

Normalization constant

■ Power Iteration:

- Set $\mathbf{c}^{(0)} \leftarrow \mathbf{1}, k \leftarrow 1$
- 1: $\mathbf{c}^{(k)} \leftarrow \alpha A \mathbf{c}^{(k-1)} + \beta \mathbf{1}$
- 2: If $\|\mathbf{c}^{(k)} - \mathbf{c}^{(k-1)}\| > \varepsilon$:
- 3: $k \leftarrow k + 1$, goto 1

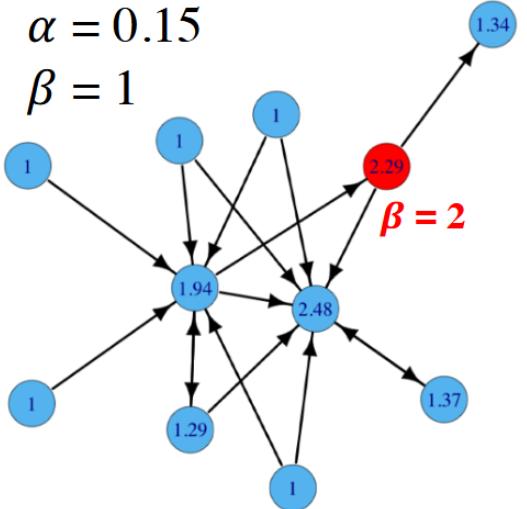
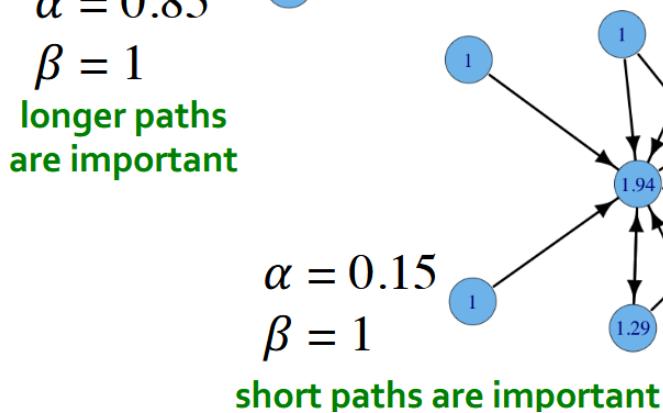
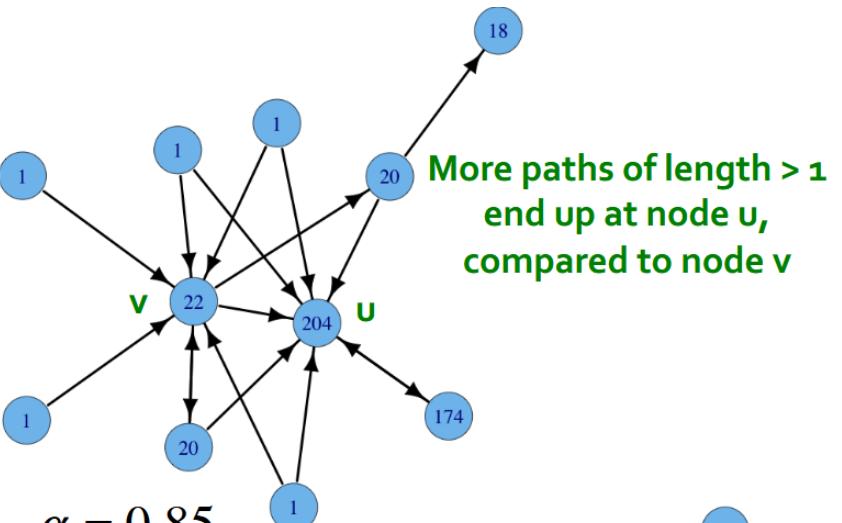


Katz's Index

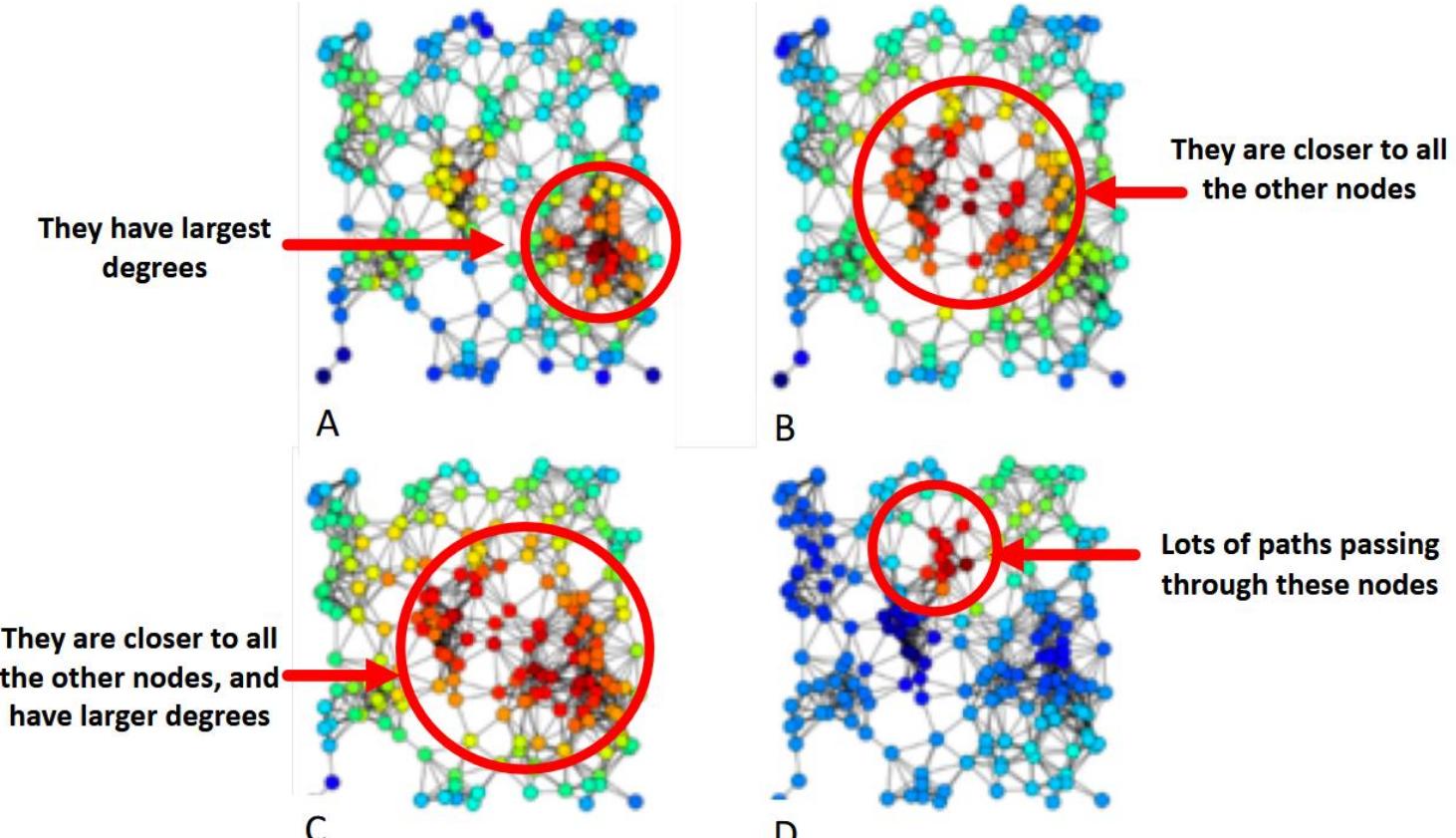
■ Spectral measures

- **Katz's Index**: Suitable for **directed acyclic graphs**
- **How to choose α ?**
 - For α close to 0, the contribution given by paths longer than one rapidly declines, and thus
 - Katz scores are mainly **influenced by short paths** (mostly in-degrees)
 - When the α is large, long paths are devalued smoothly, and
 - Katz scores are more **influenced by topology** of the network
 - The measure diverges at $\alpha > \frac{1}{\lambda}$
 - The dominant eigenvector of A is the limit of Katz centrality as α approaches $1/\lambda$ from below

Example



Example



Examples of A) Degree centrality, B) Closeness centrality, C) Harmonic Centrality and D) Katz centrality of the same graph.

Introduction

- How can we characterize, model, and reason about the structure of social networks?
 - Triadic closure and “the strength of weak ties”
 - Power-laws and scale-free networks, “rich-get-richer” phenomena
 - Small-world phenomena
 - Hubs & Authorities; PageRank
 - Models of network structure

Triangles

- So far we've seen (a little about) how networks can be characterized by their connectivity patterns
- What more can we learn by looking at higher-order properties, such as relationships between triplets of nodes?

Motivation

- Q: Last time you found a job, was it through:
 - A complete stranger?
 - A close friend?
 - An acquaintance?
- A: Surprisingly, people often find jobs through acquaintances rather than through close friends
(Granovetter, 1973)

Motivation

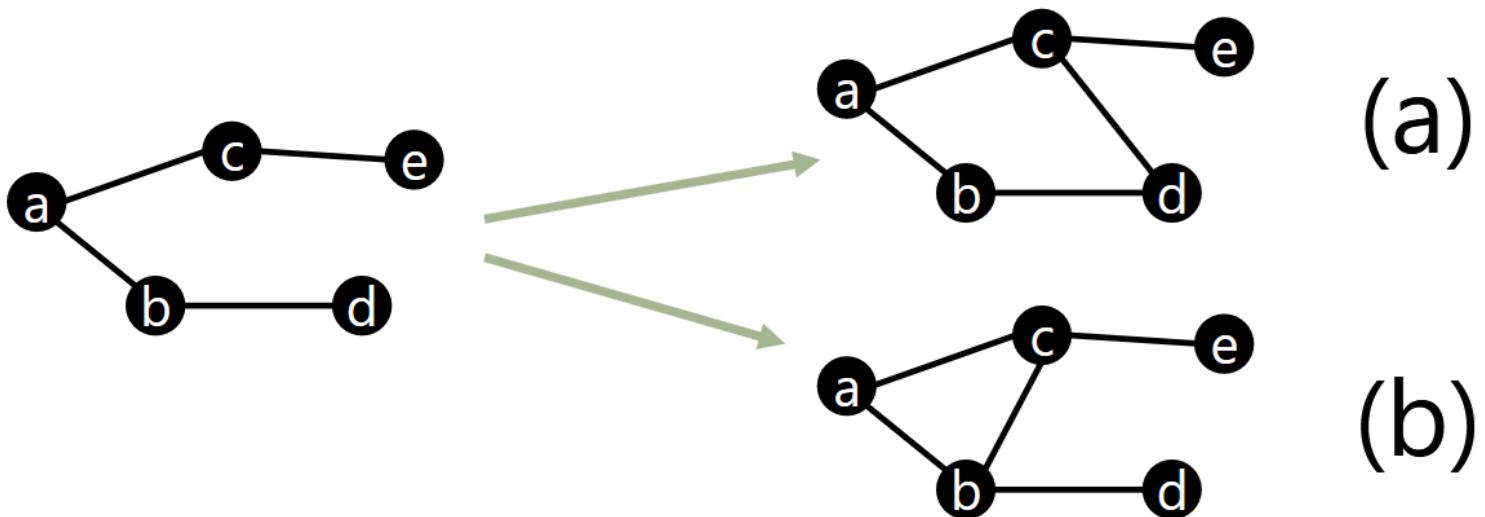
- Your friends (hopefully) would seem to have the greatest motivation to help you.
- But! Your closest friends have limited information that you don't already know about.
- Alternately, acquaintances act as a “bridge” to a different part of the social network, and expose you to new information.
- This phenomenon is known as the **strength of weak ties**.

Motivation

- To make this concrete, we'd like to come up with some notion of “tie strength” in networks
- To do this, we need to go beyond just looking at edges in isolation, and looking at how an edge connects one part of a network to another

Triadic Closure

- Q: Which edge is most likely to form next in this (social) network?



- A: (b), because it creates a triad in the network

Triangles

- “If two people in a social network have a friend in common, then there is an increased likelihood that they will become friends themselves at some point in the future” (Rapoport, 1953)
- Three reasons (see Easley & Kleinberg):
 - Every mutual friend a between bob and chris gives them an opportunity to meet
 - If bob is friends with ashton, then knowing that chris is friends with ashton gives bob a reason to trust chris
 - If chris and bob don’t become friends, this causes stress for ashton (having two friends who don’t like each other), so there is an incentive for them to connect

01

Clustering Coefficient

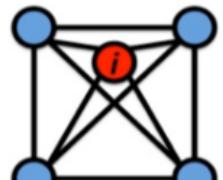
Clustering Coefficient

■ Clustering coefficient:

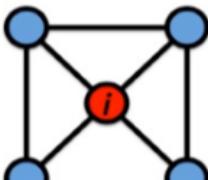
- What portion of i 's neighbors are connected?
- Node i with degree k_i
- $C_i \in [0, 1]$ This ranges between 0 (none of my friends are friends with each other) and 1 (all of my friends are friends with each other)

$$\text{■ } C_i = \frac{2e_i}{k_i(k_i - 1)}$$

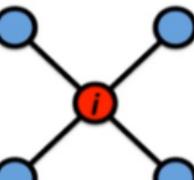
where e_i is the number of edges between the neighbors of node i



$$C_i = 1$$



$$C_i = 1/2$$



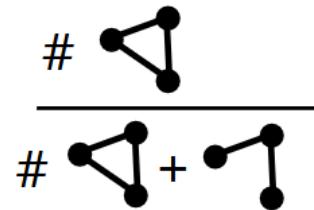
$$C_i = 0$$

- ## ■ Average clustering coefficient: $C = \frac{1}{N} \sum_i C_i$

Clustering Coefficient

- Alternately it can be defined as the fraction of connected triplets in the graph that are closed (these do not evaluate to the same thing!):

$$C = \frac{\# \text{ of closed triplets}}{\# \text{ of connected triplets}}$$



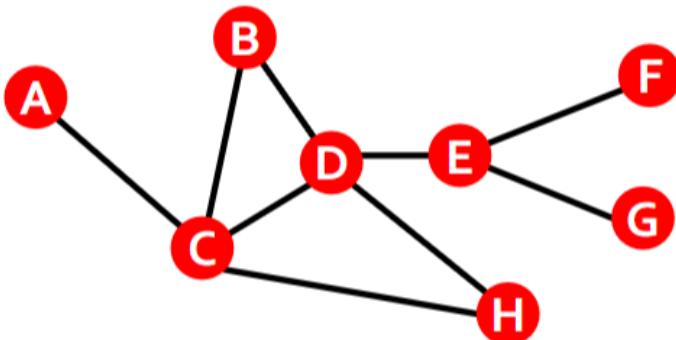
Clustering Coefficient

■ Clustering coefficient:

- What portion of i 's neighbors are connected?
- Node i with degree k_i

- $C_i = \frac{2e_i}{k_i(k_i - 1)}$

where e_i is the number of edges between the neighbors of node i



$$k_B=2, \ e_B=1, \ C_B=2/2 = 1$$

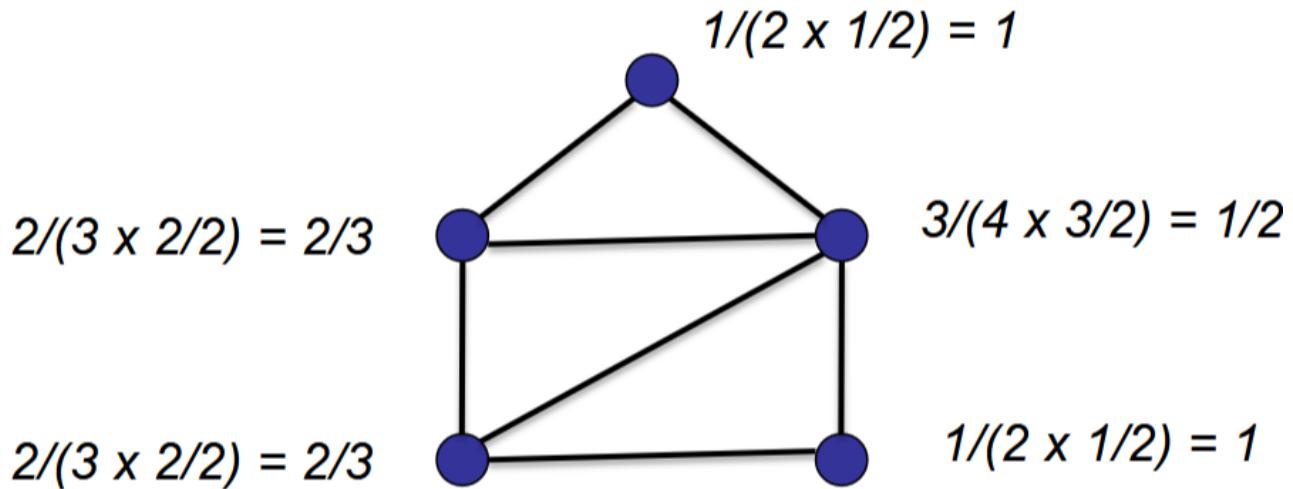
$$k_D=4, \ e_D=2, \ C_D=4/12 = 1/3$$

$$\text{Avg. clustering: } C=0.33$$

Clustering Coefficient

- Clustering coefficient of a graph G :
 $CC(G) = \text{average of } c(u) \text{ over all vertices } u \text{ in } G$
- What Do We Mean By “High” CC?
 - $CC(G)$ measures how likely vertices with a common neighbor are to be neighbors themselves
 - Should be compared to how likely random pairs of vertices are to be neighbors
 - Let p be the edge density of network/graph G : $p = E / (N(N - 1)/2)$
 - Here E = total number of edges in G
 - If we picked a pair of vertices at random in G , probability they are connected is exactly p
 - So, we will say clustering is high if $CC(G) \gg p$

Example



$$C.C. = (1 + \frac{1}{2} + 1 + 2/3 + 2/3)/5 = 0.7666\dots$$

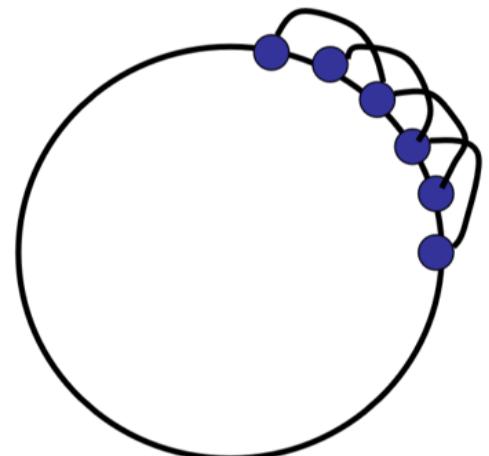
$$p = 7/(5 \times 4/2) = 0.7$$

Not highly clustered

Example

- Network: simple cycle + edges to vertices 2 hops away on cycle

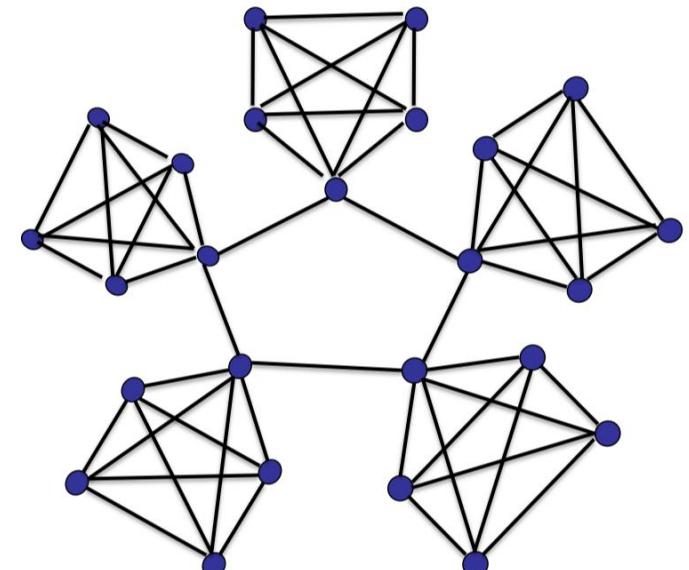
- By symmetry, all vertices have the same clustering coefficient
- Clustering coefficient of a vertex v:
 - Degree of v is 4, so the number of *possible* edges between pairs of neighbors of v is $4 \times 3/2 = 6$
 - How many pairs of v's neighbors actually *are* connected? 3 --- the two clockwise neighbors, the two counterclockwise, and the immediate cycle neighbors
 - So the c.c. of v is $3/6 = 1/2$
- Compare to overall edge density:
 - Total number of edges = $2N$
 - Edge density $p = 2N/(N(N-1)/2) \sim 4/N$
 - As N becomes large, $1/2 \gg 4/N$
 - So this cyclical network is highly clustered



Example

- Divide N vertices into \sqrt{N} groups of size \sqrt{N} (here $N = 25$)
- Add all connections within each group (cliques), connect “leaders” in a cycle
- $N - \sqrt{N}$ non-leaders
 - CC of network as N becomes large?
 - Edge Density?

*Add all connections within each group (cliques), connect “leaders” in a cycle
 $N - \sqrt{N}$ non-leaders have C.C. = 1, so network C.C. $\rightarrow 1$ as N becomes large
Edge density is $p \sim 1/\sqrt{N}$*



Research Study

Higher-order structures such as larger cliques are crucial to the structure and function of complex networks

Higher-order clustering in networks

Hao Yin*

Institute for Computational and Mathematical Engineering, Stanford University, Stanford, CA, 94305, USA

Austin R. Benson[†]

Department of Computer Science, Cornell University, Ithaca, NY, 14850, USA

Jure Leskovec[‡]

Computer Science Department, Stanford University, Stanford, CA, 94305, USA

(Dated: April 30, 2018)

A fundamental property of complex networks is the tendency for edges to cluster. The extent of the clustering is typically quantified by the clustering coefficient, which is the probability that a length-2 path is closed, i.e., induces a triangle in the network. However, higher-order cliques beyond triangles are crucial to understanding complex networks, and the clustering behavior with respect to such higher-order network structures is not well understood. Here we introduce higher-order clustering coefficients that measure the closure probability of higher-order network cliques and provide a more comprehensive view of how the edges of complex networks cluster. Our higher-order clustering coefficients are a natural generalization of the traditional clustering coefficient. We derive several properties about higher-order clustering coefficients and analyze them under common random graph models. Finally, we use higher-order clustering coefficients to gain new insights into the structure of real-world networks from several domains.

<http://snap.stanford.edu/hocc/>

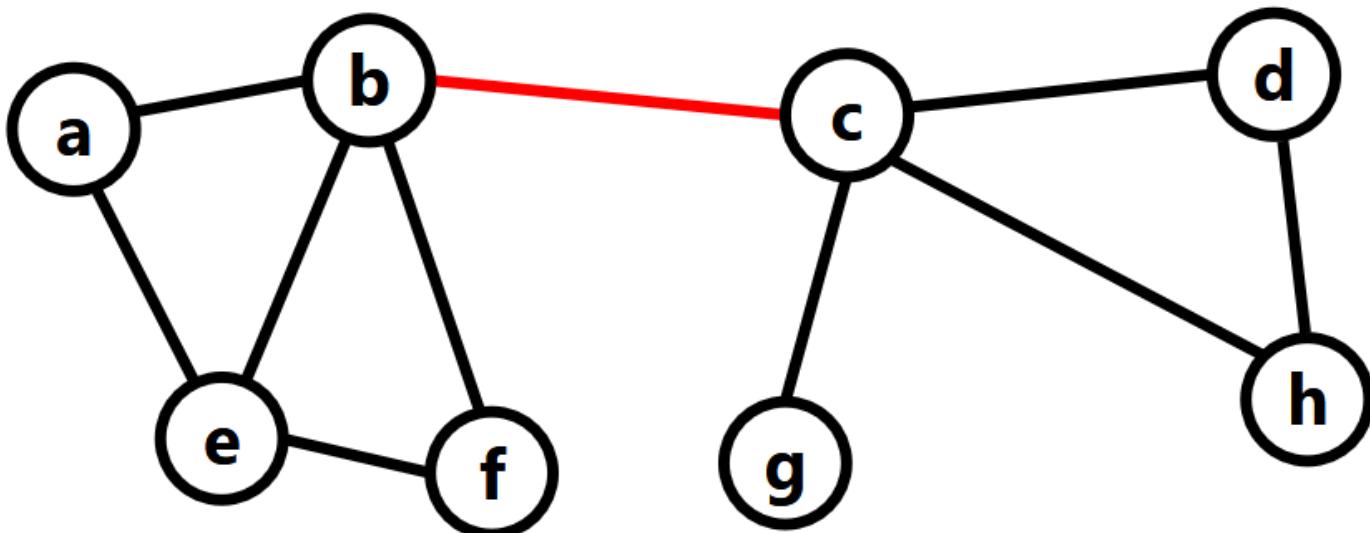


02

Bridges

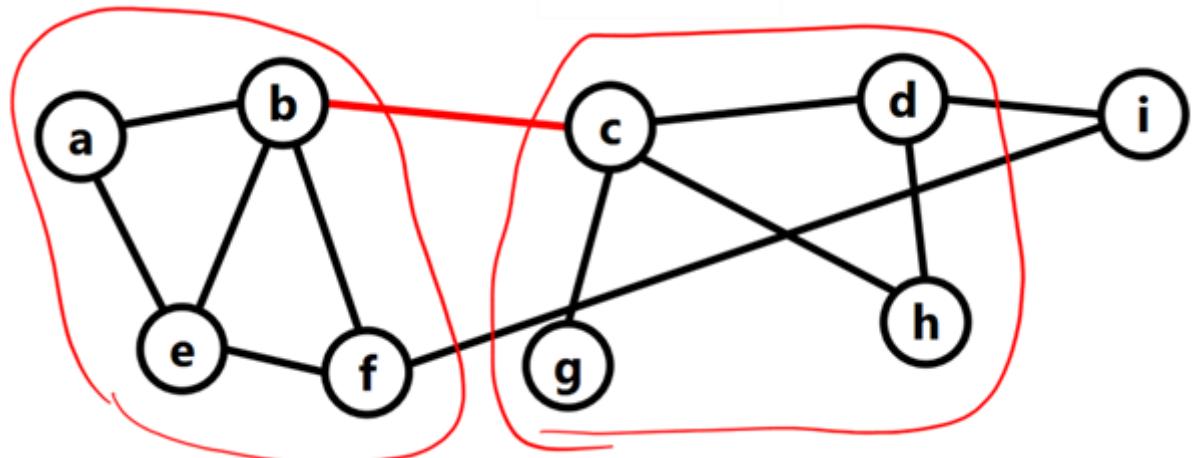
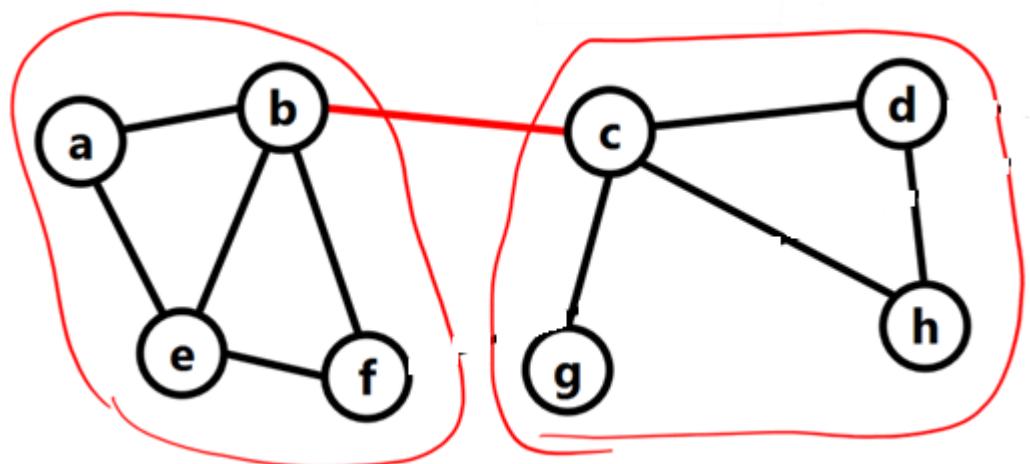
Bridge Edge

- An edge (b,c) is a bridge edge if removing it would leave no path between b and c in the resulting network



Local Bridge Edge

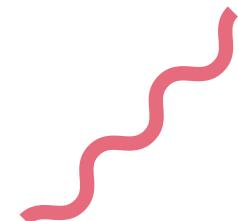
- An edge (b,c) is a local bridge if removing it would leave no edge between b's friends and c's friends (though there could be more distant connections)





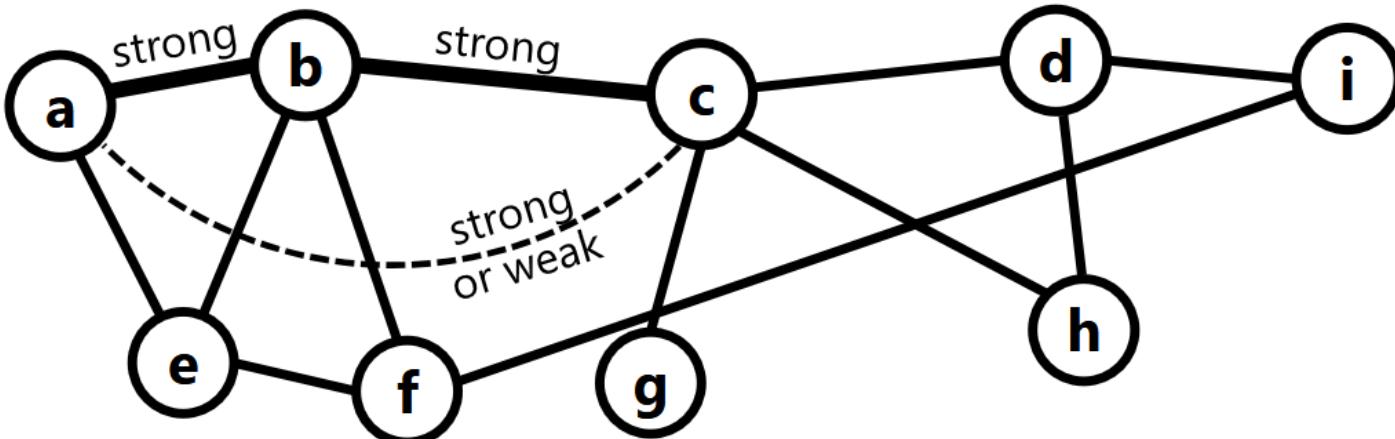
03

Strong & weak ties



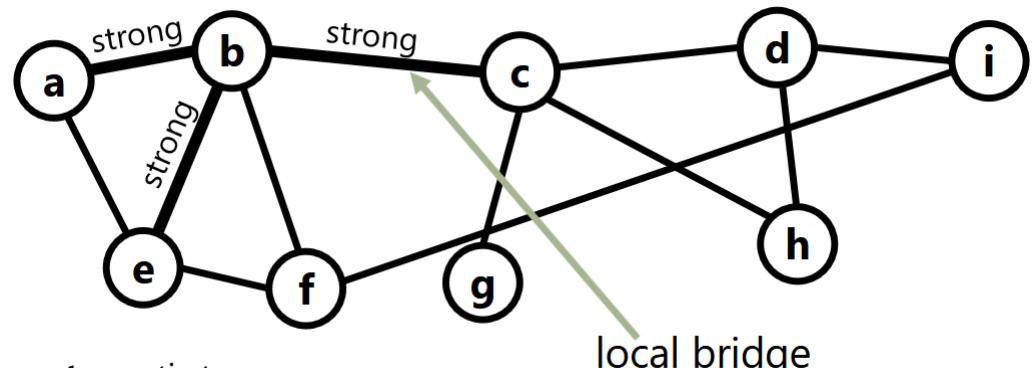
Strong Triadic Closure Property

- If (a,b) and (b,c) are connected by strong ties, there must be at least a weak tie between a and c .
 - Note: (a,c) can be weak or strong!



Granovetter's Theorem

- If the strong triadic closure property is satisfied for a node, and that node is involved in two strong ties, then any incident local bridge must be a weak tie.
 - Proof?



Proof (by contradiction): (1) b has two strong ties (to a and e); (2) suppose it has a **strong** tie to c via a local bridge; (3) but now a tie must exist between c and a (or c and e) due to strong triadic closure; (4) so $b \rightarrow c$ cannot be a bridge

Granovetter's Theorem

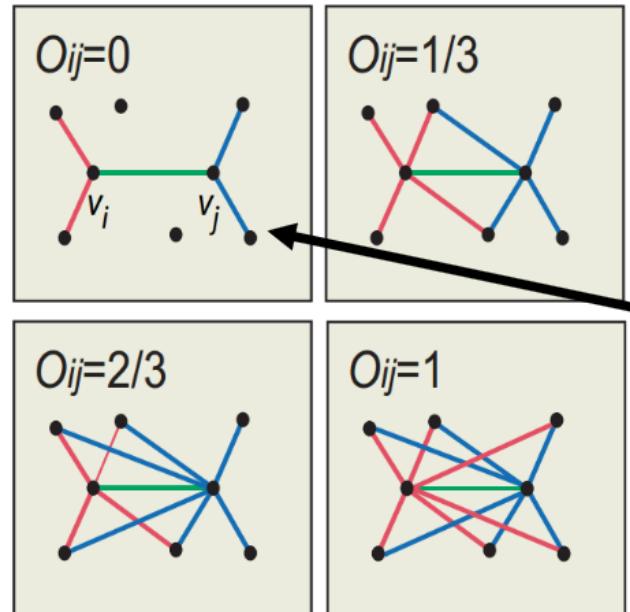
- So, if we're receiving information from distant parts of the network (i.e., via “local bridges”) then we must be receiving it via weak ties.
- Q: How to test this theorem empirically on real data?
 - A: Onnela et al. 2007 studied networks of mobile phone calls.

Study Example #1

Defn. 1: Define the “overlap” between two nodes to be the Jaccard similarity between their connections

$$O_{i,j} = \frac{\Gamma(i) \cap \Gamma(j)}{\Gamma(i) \cup \Gamma(j)}$$

neighbours of i



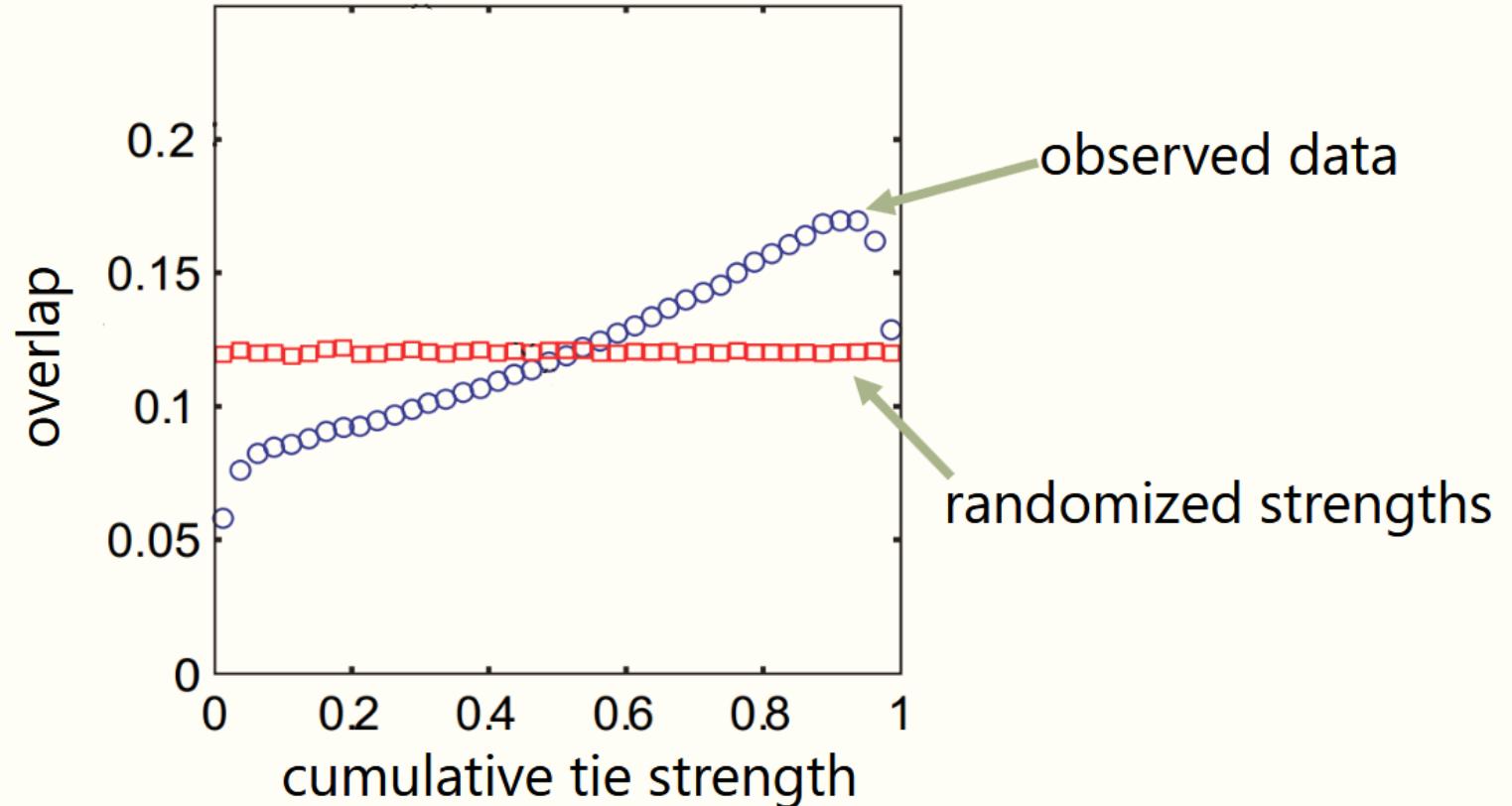
“local bridges” have overlap 0

(picture from Onnela et al., 2007)

Secondly, define the “strength” of a tie in terms of the number of phone calls between i and j

Study Example #1

finding: the “stronger” our tie, the more likely there are to be additional ties between our mutual friends



(picture from Onnela et al., 2007)

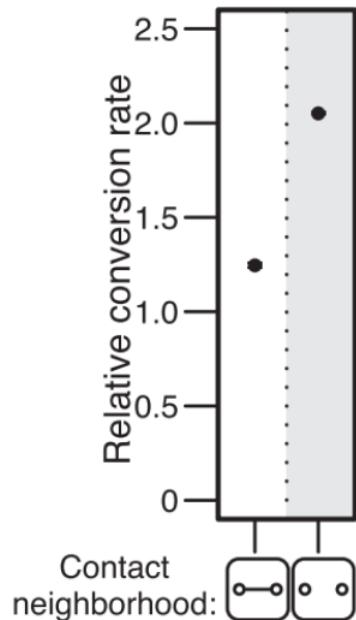
Study Example #2

- Suppose a user receives four e-mail invites to join facebook from users who are already on facebook. Under what conditions are we most likely to accept the invite (and join facebook)?
 - If those four invites are from four close friends?
 - If our invites are from four acquaintances?
 - If the invites are from a combination of friends, acquaintances, work colleagues, and family members?

hypothesis: the invitations are most likely to be adopted if they come from distinct groups of people in the network

Study Example #2

Let's consider the connectivity patterns amongst the people who tried to recruit us

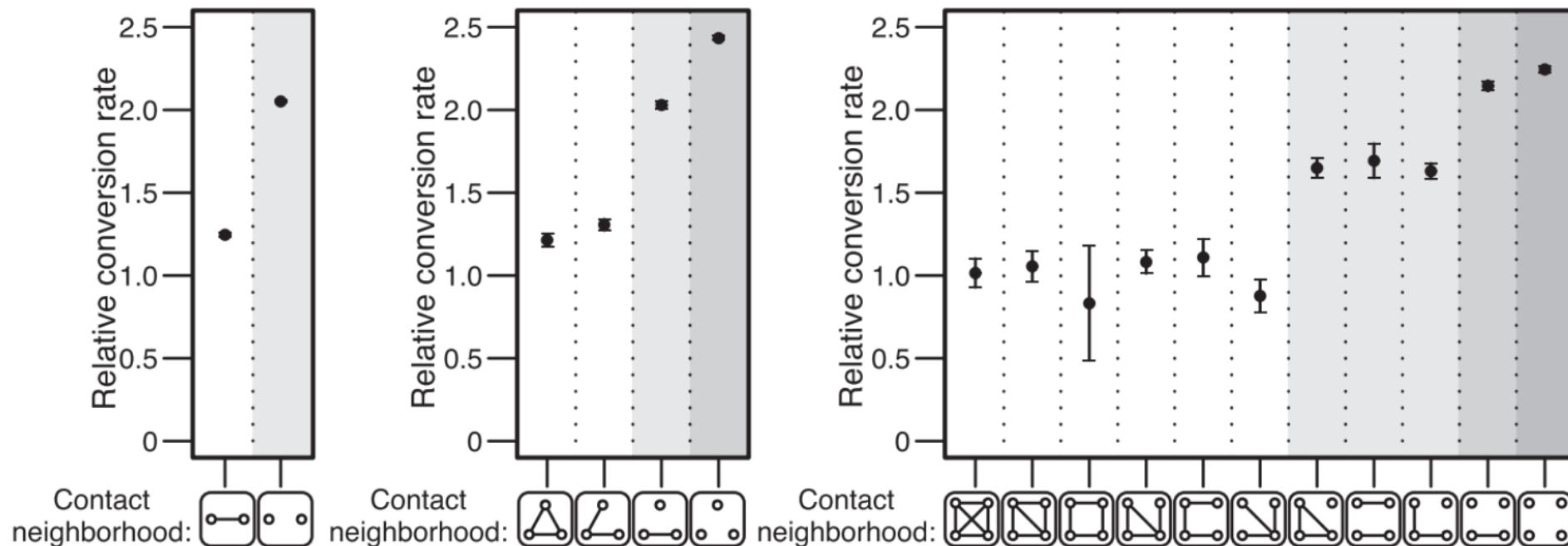


- **Case 1:** two users attempted to recruit
- **y-axis:** relative to recruitment by a single user
- **finding:** recruitments are **more likely to succeed** if they come from friends who are **not connected to each other**

(picture from Ugander et al., 2012)

Study Example #2

Let's consider the connectivity patterns amongst the people who tried to recruit us



error bars are high since this
structure is very very rare

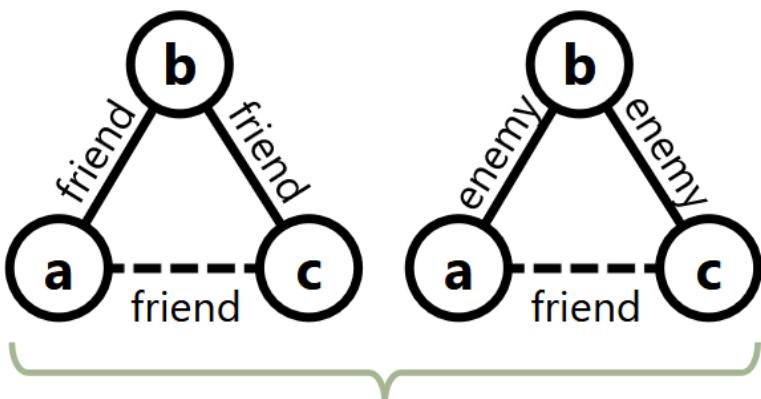
(picture from Ugander et al., 2012)

Conclusion

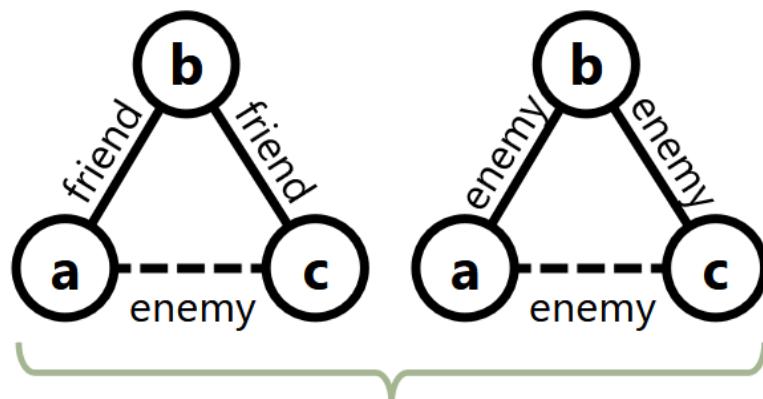
- Important aspects of network structure can be explained by the way an edge connects two parts of the network to each other:
 - Edges tend to close open triads (clustering coefficient etc.)
 - It can be argued that edges that bridge different parts of the network somehow correspond to “weak” connections (Granovetter; Onnela et al.)
 - Disconnected parts of the networks (or parts connected by local bridges) expose us to distinct sources of information (Granovettor; Ugander et al.)

Structural Balance

- Some of the assumptions that we've seen today may not hold if edges have signs associated with them



balanced: the edge
 $a \rightarrow c$ is **likely** to form



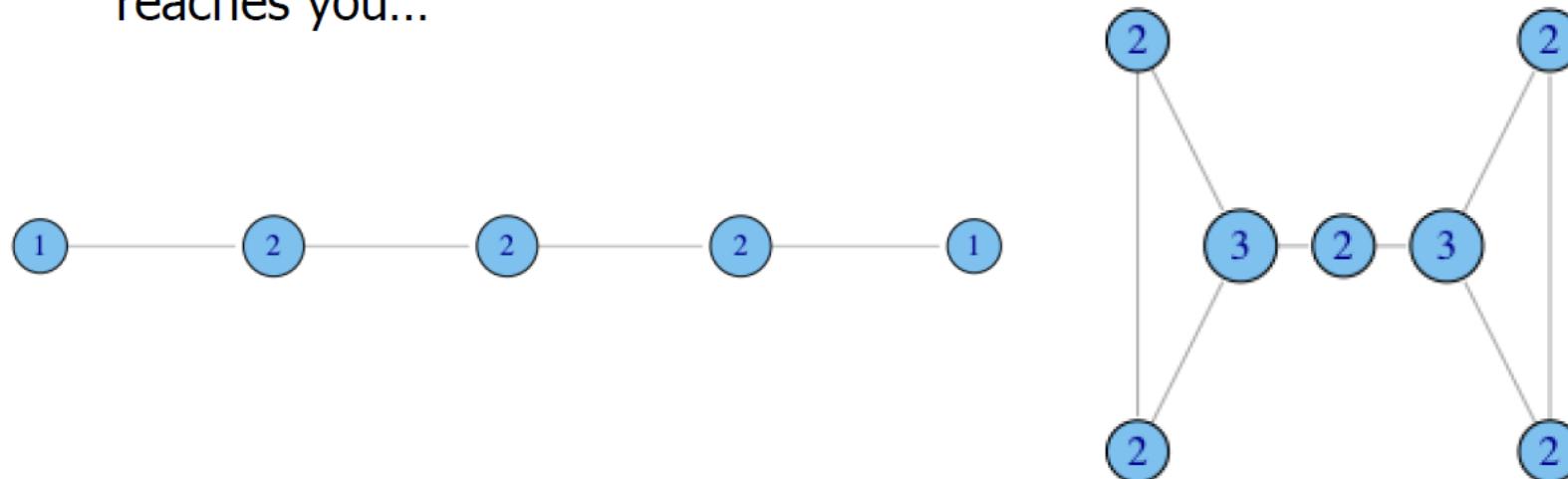
imbalanced: the edge
 $a \rightarrow c$ is **unlikely** to form

06

Degree

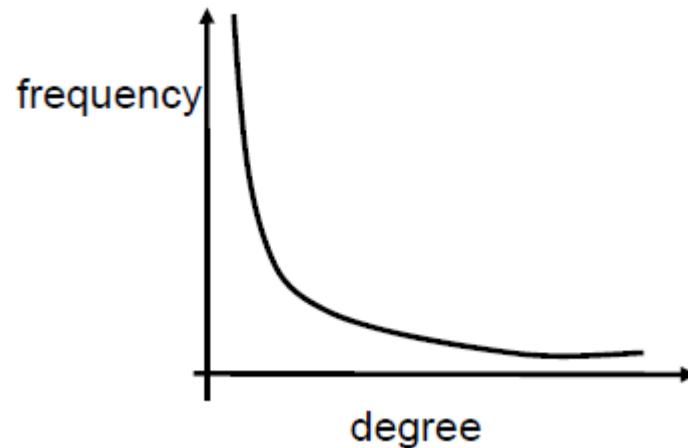
Is degree everything?

- Nodes with the same degree might have different properties
- In what ways does degree fail to capture centrality in the following graphs?
 - ability to broker between groups
 - likelihood that information originating anywhere in the network reaches you...



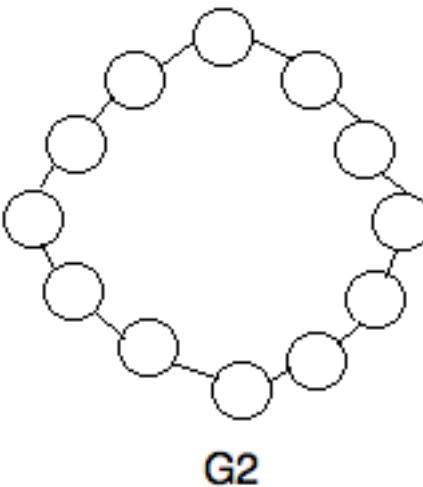
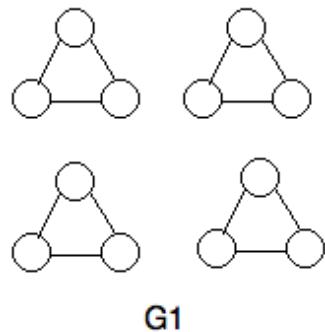
Degree and Degree Distribution

- Degree k_i of node i is a measure of its centrality
- Nodes with high degrees are called hubs
- Maximum degree $k_{max} = \max_i(k_i)$ is also an important measure
- The variance of node-degrees can be an indicator of network heterogeneity, i.e. the more the variance the more the heterogeneity.
- Degree distribution



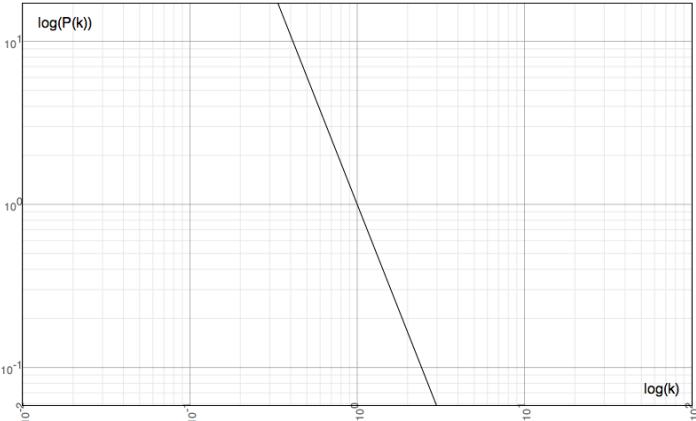
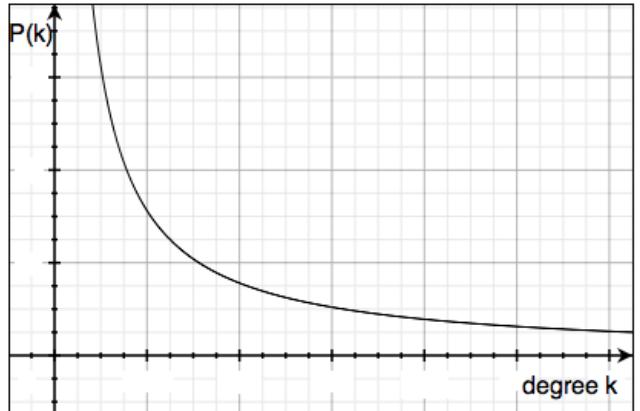
Degree Distribution

- However: degree distribution (and global properties in general) are weak predictors of network structure.
- Illustration:



- G_1 and G_2 are of the same size (i.e., $|G_1|=|G_2|$ -- they have the same number of nodes and edges) and they have same degree distribution, but G_1 and G_2 have very different topologies (i.e., graph structure).

Degree Distribution

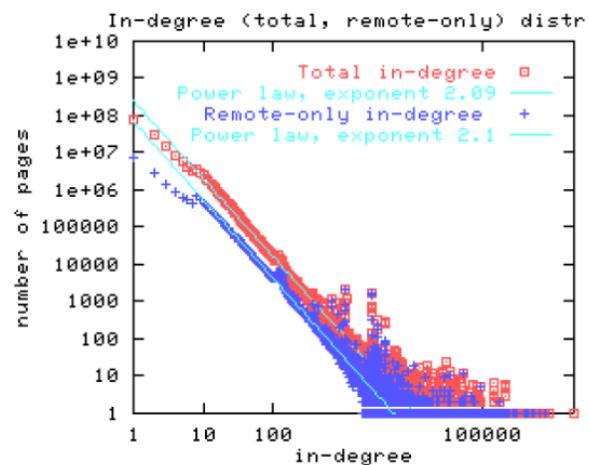


(log-log plot)

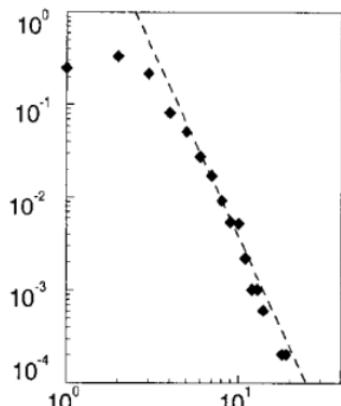
- Here $P(k) \sim k^{-\gamma}$, where often $2 \leq \gamma < 3$. This is a **power-law**, heavy-tailed distribution.
- Networks with power-law degree distributions are called **scale-free networks**. In them, most of the nodes are of low degree, but there is a small number of highly-linked nodes (nodes of high degree) called “hubs.”

Power Laws

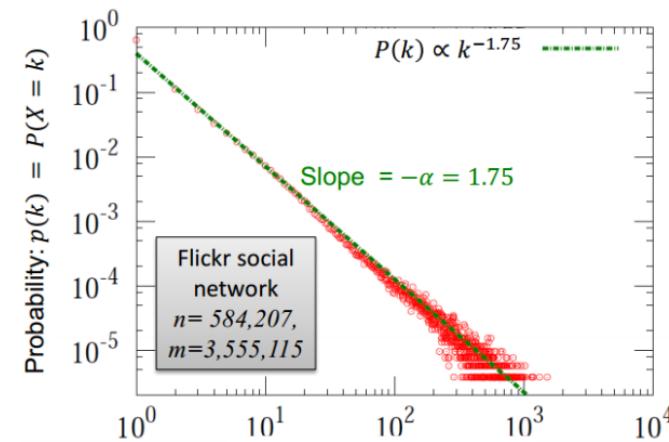
- Social and information networks often follow power laws, meaning that a few nodes have many of the edges, and many nodes have a few edges.



e.g. web graph
(Broder et al.)



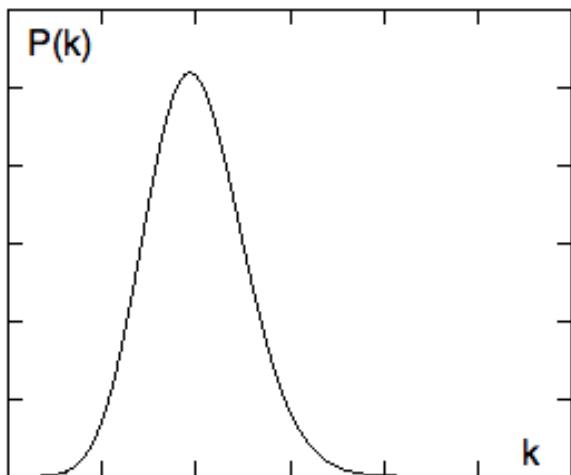
e.g. power grid
(Barabasi-Albert)



e.g. Flickr
(Leskovec)

Degree Distribution

- Here $P(k)$ is a Poisson distribution.



average degree is meaningful

Degree-Degree correlation

- It is important to know if the nodes with degree k are connected to nodes with degree k' . How?
 - 1) Method proposed by Pastor Satoras et al. and plot the mean degree of the neighbors as a function of the degree

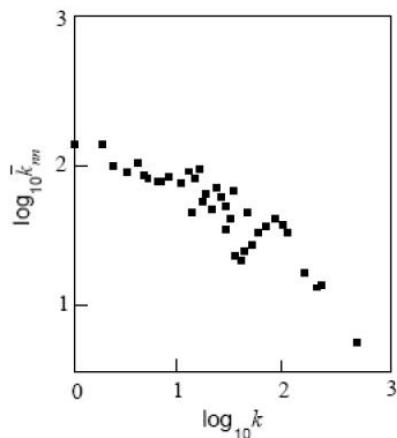


FIG. 3.13. Correlations of the degrees of nearest-neighbour vertices (autonomous systems) in the Internet at the interdomain level (after Pastor-Satorras, Vázquez, and Vespignani 2001). The empirical dependence of the average degree of the nearest neighbours of a vertex on the degree of this vertex is shown in a log-log scale. This empirical dependence was fitted by a power law with exponent approximately 0.5.

Degree-Degree correlation

- It is important to know if the nodes with degree k are connected to nodes with degree k' . How?
 - 2) method proposed by Newman and compute the correlation coefficient

$$r = \frac{\frac{1}{E} \sum_{j>i} k_i k_j a_{ij} - \left[\frac{1}{E} \sum_{j>i} \frac{1}{2} (k_i + k_j) a_{ij} \right]^2}{\frac{1}{E} \sum_{j>i} \frac{1}{2} (k_i^2 + k_j^2) a_{ij} - \left[\frac{1}{E} \sum_{j>i} \frac{1}{2} (k_i + k_j) a_{ij} \right]^2}$$

E is the total number of edges

a_{ij} is the entry (i,j) of the adjacency matrix

k_i is the degree of node i

Degree-Degree correlation

- $r > 0$: the network is called assortative
 - Node with large degree intent to connect to those with large degrees and nodes with low degrees intend to connect to those with low degrees (rich with rich and poor with poor)
- $r < 0$: the network is called disassortative
 - Node with large degree intent to connect to those with low degrees and nodes with low degrees intend to connect to those with high degrees (rich with poor)
- $r = 0$: no correlations
 - There is no specific intention in the connection between the nodes in the sense of their degrees



Any Question?