

Random Graph Theory

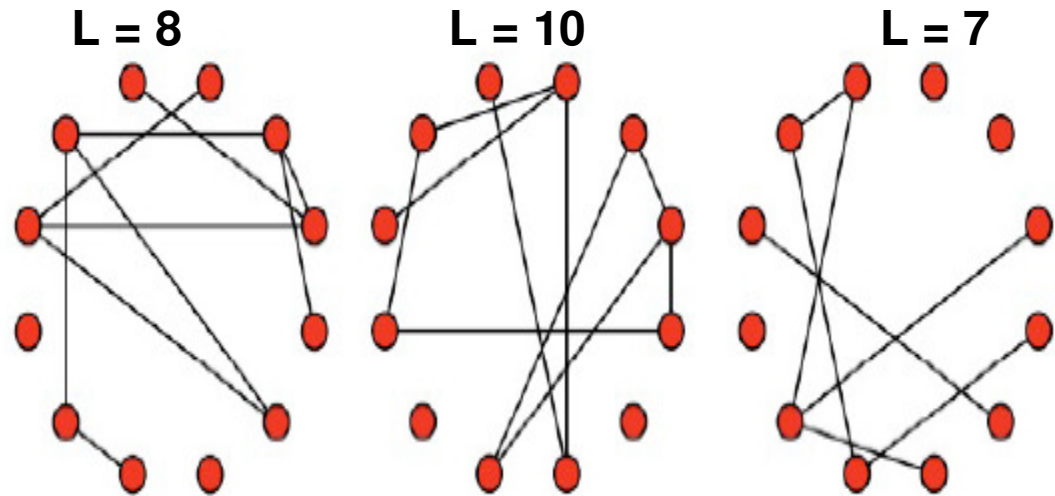
Dr. Natarajan Meghanathan
Associate Professor
Department of Computer Science
Jackson State University, Jackson, MS
E-mail: natarajan.meghanathan@jsums.edu

Introduction

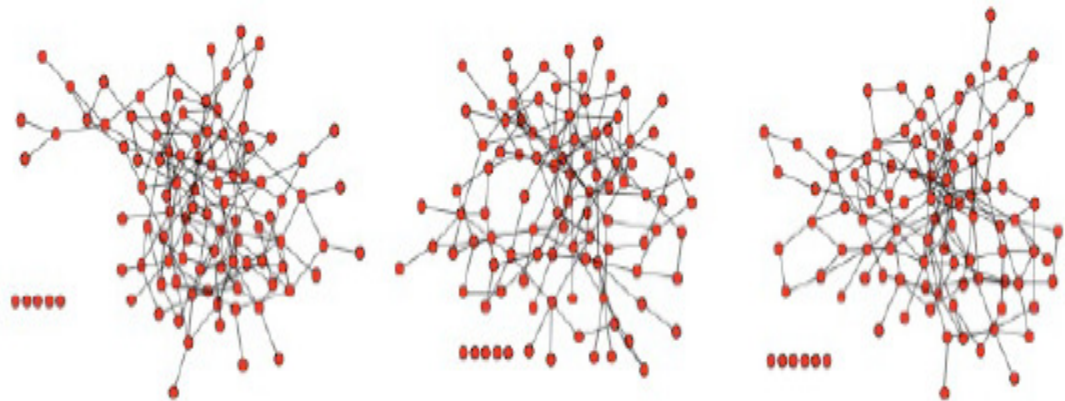
- At first inspection, most real-world networks look as if they are spun randomly.
- To model such networks that are *truly random*, the principle behind “Random Graph Theory” is:
 - Place the links randomly between nodes to reproduce the complexity and apparent randomness of real-world systems.
- Two definitions of random networks
 - $G(N, L)$ model: N labeled nodes are connected with L randomly placed links
 - $G(N, p)$ model: Each pair of N labeled nodes are connected with a probability p .
- Though the average degree for a node is simply $2L/N$ in a $G(N, L)$ model, the other key network characteristics are easier to calculate in the $G(N, p)$ model.
 - The construction of the $G(N, p)$ model is closer to the way real systems develop. The total number of links in a network is rarely fixed.

Constructing a $G(N, p)$ Network

- Step 1: Start with N isolated nodes
- Step 2: For a particular node pair (u, v) , generate a random number r . If $r \leq p$, then, add the link (u, v) to the network.
- Repeat Step 2 for each of the $N(N-1)/2$ node pairs.
- Each random network we generate with the same parameters (N, p) will look slightly different.
 - The number of links L is likely to be different.



N = 12 nodes, p = 1/6



N = 100 nodes, p = 1/6

Source: Figure 3.3a: Barabasi

Review of Binomial Distribution

- Let there be N independent experiments with two possible outcomes (in each experiment: success or failure): with the probability of one outcome (say success) is p and of the other is $1-p$.
- The binomial distribution provides the probability p_x that we obtain exactly x successes in a sequence of N experiments.

The binomial distribution has the form

$$p_x = \binom{N}{x} p^x (1-p)^{N-x}.$$

The mean of the distribution (first moment) is

$$\langle x \rangle = \sum_{x=0}^N x p_x = Np.$$

Its second moment is

$$\langle x^2 \rangle = \sum_{x=0}^N x^2 p_x = p(1-p)N + p^2 N^2,$$

providing its standard deviation as

$$\sigma_x = \left(\langle x^2 \rangle - \langle x \rangle^2 \right)^{\frac{1}{2}} = [p(1-p)N]^{\frac{1}{2}}.$$

$C(N, x) = \binom{N}{x}$ is the different combinations of the results of the N experiments in which there will be X successes and $N-X$ failures.

Binomial Distribution: Tossing a Coin

- Prob[Head] = Prob[Tail] = $\frac{1}{2}$.
- Probability of getting exactly X Heads in a sequence of N tossing of a coin is:

$$p_X = \binom{N}{X} p^X (1-p)^{N-X} \quad p_X = \binom{N}{X} (1/2)^X (1-1/2)^{N-X}$$
$$p_X = \binom{N}{X} (1/2)^N$$

- $C(5,0) = 1$; $C(5,1) = 5$; $C(5,2) = 10$; $C(5,3) = 10$; $C(5,4) = 5$; $C(5,5) = 1$
- $P_0 = 1 \cdot [1/2]^5$; $p_1 = 5 \cdot [1/2]^5$; $p_2 = 10 \cdot [1/2]^5$; $p_3 = 10 \cdot [1/2]^5$; $p_4 = 5 \cdot [1/2]^5$; $p_5 = 1 \cdot [1/2]^5$.

- Avg. # Heads:

$$\langle X \rangle = \sum_{X=0}^N X \cdot p_X = \sum_{X=0}^5 X \cdot p_X$$

$$\langle X \rangle = 0 \cdot p_0 + 1 \cdot p_1 + 2 \cdot p_2 + 3 \cdot p_3 + 4 \cdot p_4 + 5 \cdot p_5$$

$$\langle X \rangle = [1/2]^5 * [1 * 5 + 2 * 10 + 3 * 10 + 4 * 5 + 5 * 1]$$

$$\langle X \rangle = [1/2]^5 * 80 = 2.5 = (N)(p) = (5)(1/2)$$

Links in a G(N, p) Network

- Let L be the number of links arising out of a random network generated according to the $G(N, p)$ model.
- To determine the **Average Number of Links $\langle L \rangle$** , we need to model the probability that there will be exactly L links among the total number of node pairs $N(N-1)/2$ considered to have a link; each node pair has a probability of p to form a link. Let $L_{\max} = N(N-1)/2$.

$$p_L = \binom{L_{\max}}{L} p^L (1-p)^{L_{\max}-L}$$

$$\langle L \rangle = \sum_{L=0}^{L_{\max}} L * p_L = (L_{\max}) * p$$

$$\langle L \rangle = p * \frac{N(N-1)}{2}$$

Average Degree of a Node $\langle k \rangle$

$$\begin{aligned} \langle K \rangle &= \frac{2 * \langle L \rangle}{N} \\ &= \frac{2 * p * N(N-1)}{2 * N} \\ &= p * (N-1) \end{aligned}$$

Degree Distribution

- For a **random network** of N nodes, each node can have potentially $N-1$ links.
- The probability p_k that a node has exactly k links is given by the **binomial distribution**:

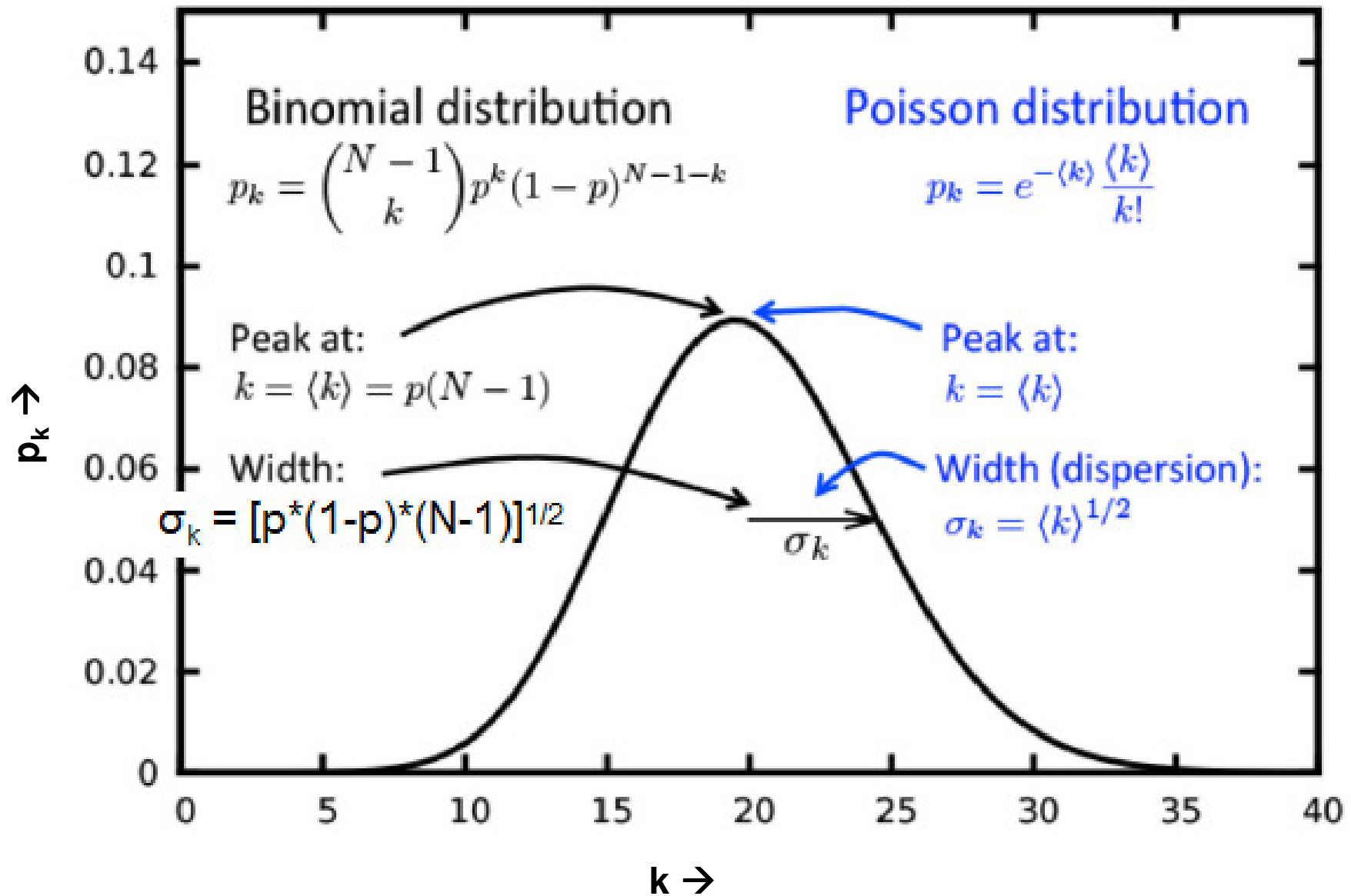
$$p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k}$$

- Using the above binomial distribution to find the average node degree for a random network, we obtain $\langle k \rangle = p^*(N-1)$ and the standard deviation for the node degree is $\sigma_k = [p^*(1-p)^*(N-1)]^{1/2}$.
- **For sparse networks** (for which $\langle k \rangle \ll N$), the probability of finding a node with k neighbors is given by the **Poisson distribution**:

$$p_k = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$

- Using the above Poisson distribution to find the average node degree for a random network, we obtain $\langle k \rangle = k$ and the standard deviation for the node degree is $\sigma_k = (k)^{1/2}$.

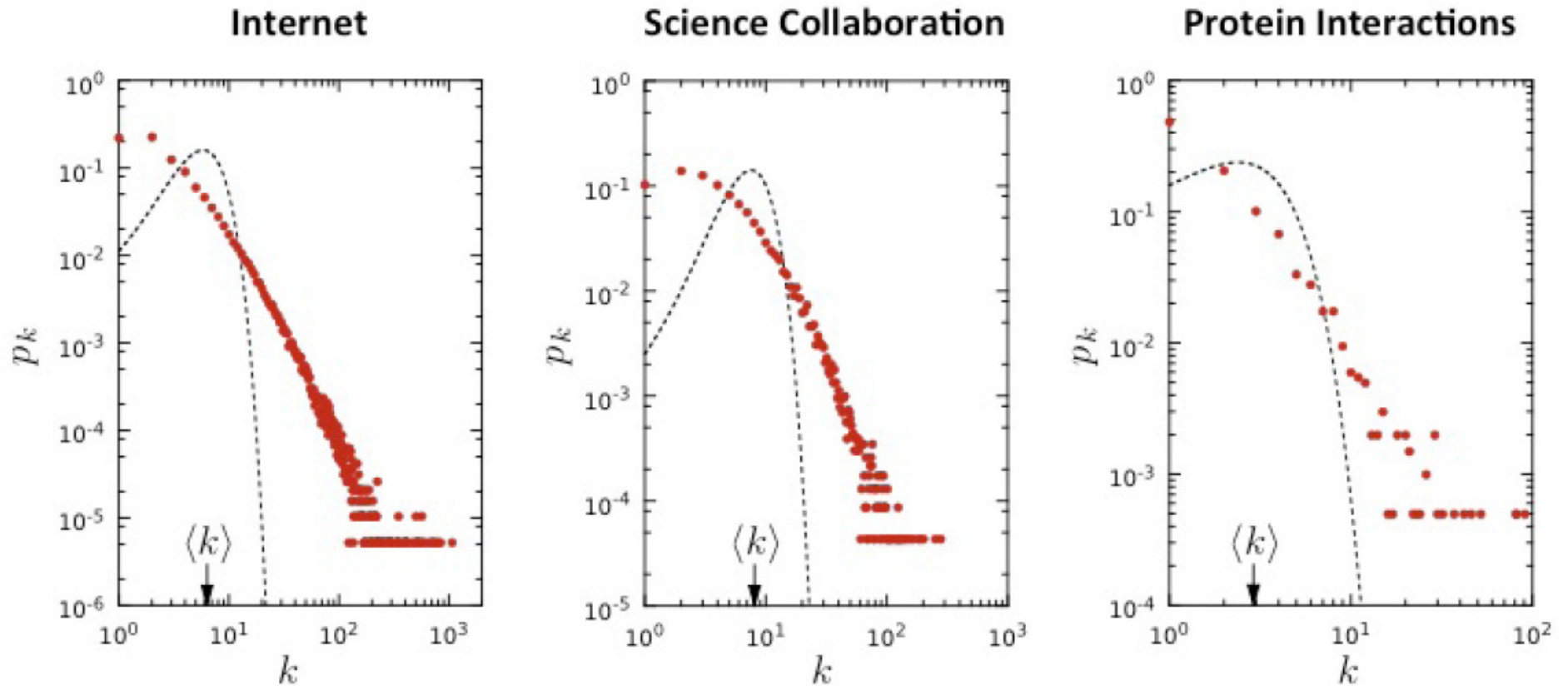
Degree Distribution



Real Networks do not have a Poisson degree distribution

- Let us assume that the world's social network (typically, $N = 10^9$ nodes and average node degree $\langle k \rangle = 1000$) follows a random network model.
- Using the results obtained for random networks, the above values for the global social network corresponds to:
 - Dispersion (std. dev.) = $\langle k \rangle^{1/2} = 31.62$.
- The above results indicate that in the global social network, the degree of most nodes is in the vicinity of $\langle k \rangle$.
 - However, we have people with number of contacts significantly larger than 1000 and significantly lower than 1000 too.
- The random network cannot be used to model a network with few extremely popular individuals (hubs) and networks with large differences in node degrees.

Degree Distribution of Real Networks



Source: Figure 3.5: Barabasi

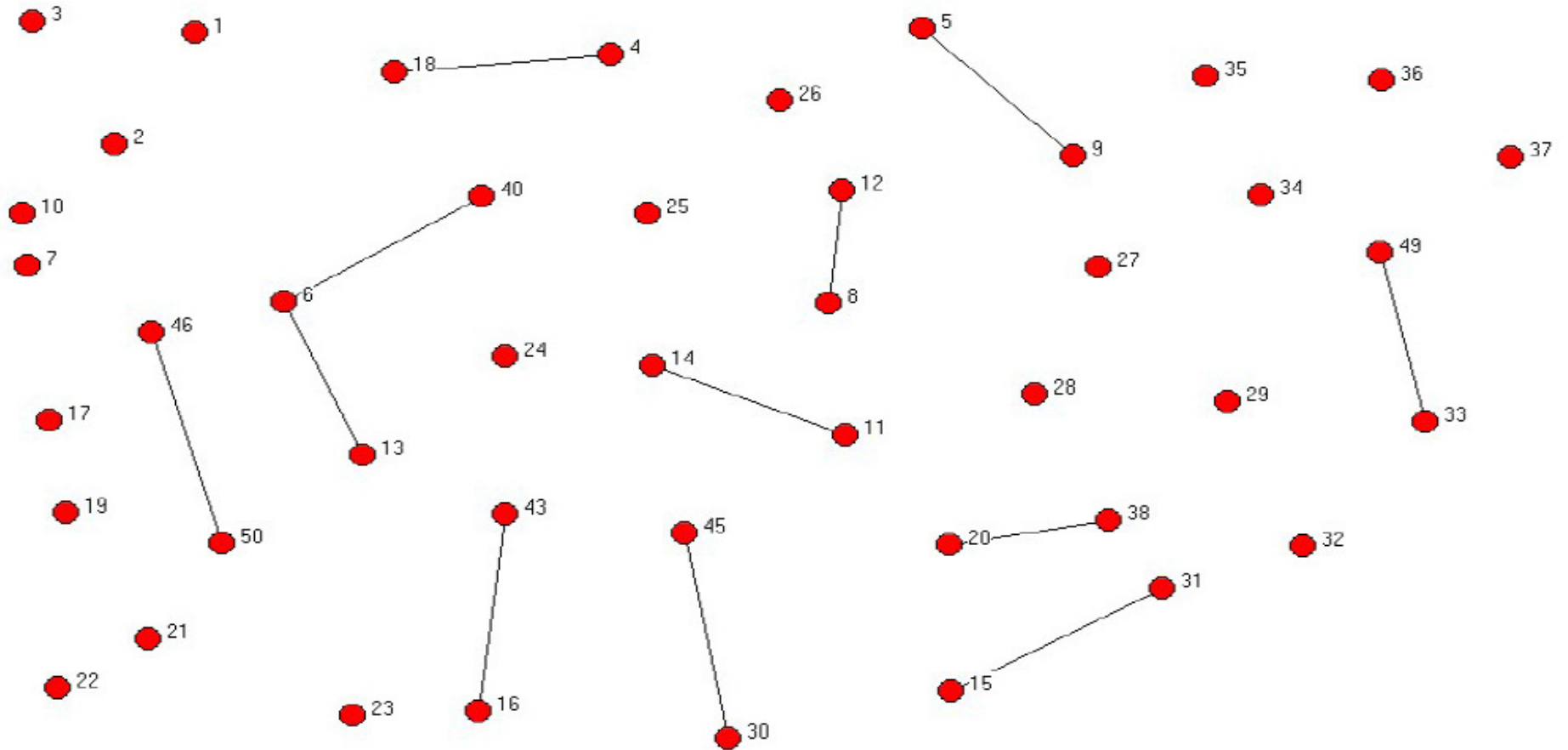
The Poisson distribution underestimates the presence of nodes with larger degrees. For example, the maximum degree for a node in the Internet (according to the random model) is expected to be 20; there are nodes with degrees close to 1000. Likewise, the dispersion predicted under the random model is 2.52 (much smaller than the measured value of 14.44).

Phase Transitions in Random Networks

- If $p \geq 1/n^2$, the network has some links (avg. deg. $1/n$)
- If $p \geq 1/n^{3/2}$, the network has a component with at least three links (avg. deg. $1/n^{1/2}$)
- If $p \geq 1/n$, the network has a cycle; the network has a unique giant component: a component with at least n^a nodes (for some fixed $a < 1$); (avg. deg. 1)
- If $p \geq \log(n)/n$, then the network is connected; (avg. deg. $\log(n)$)

Phase Transitions in Random Networks

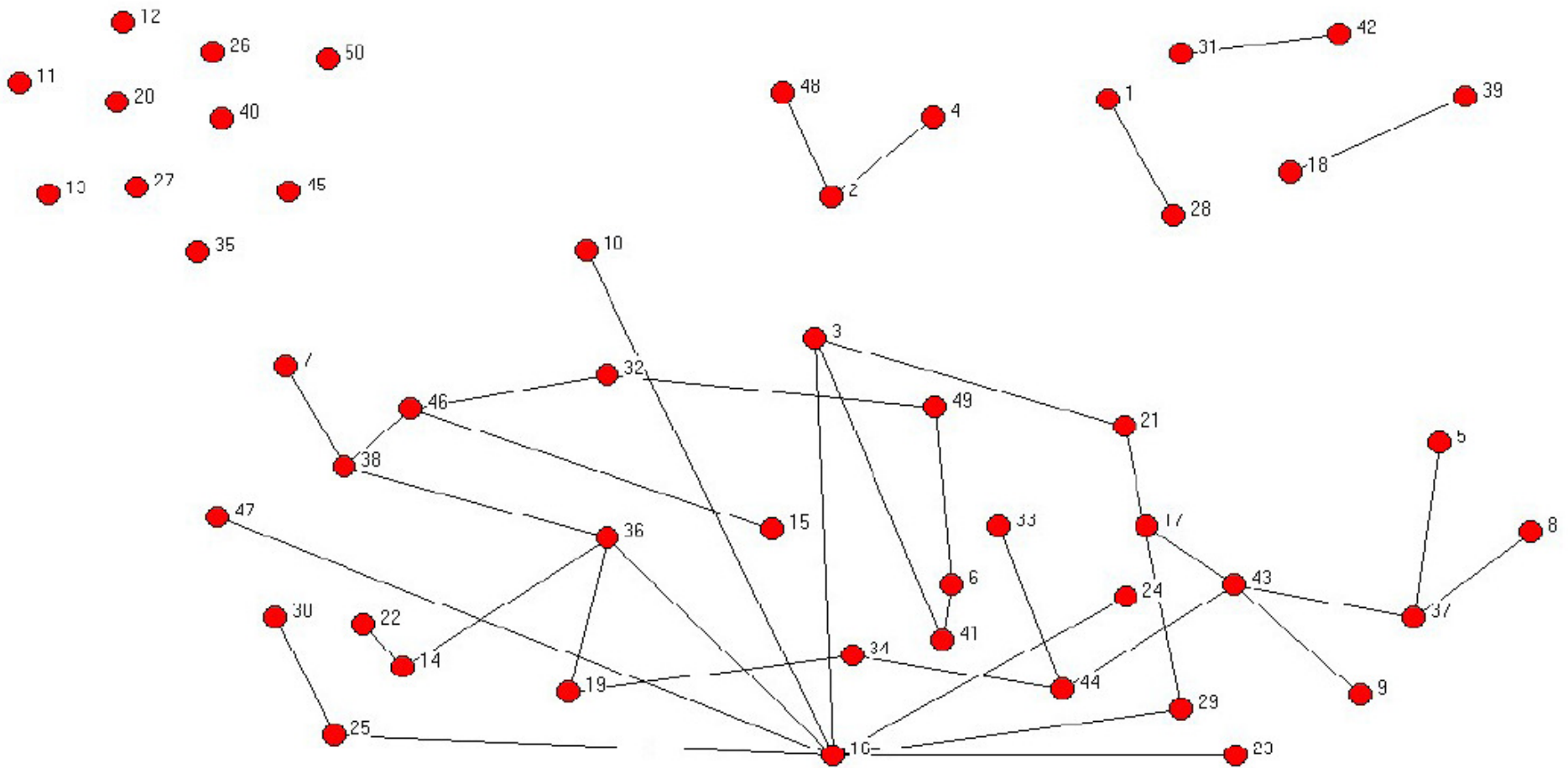
$p = 0.01$; 50 nodes (1)



Phase Transitions in Random Networks

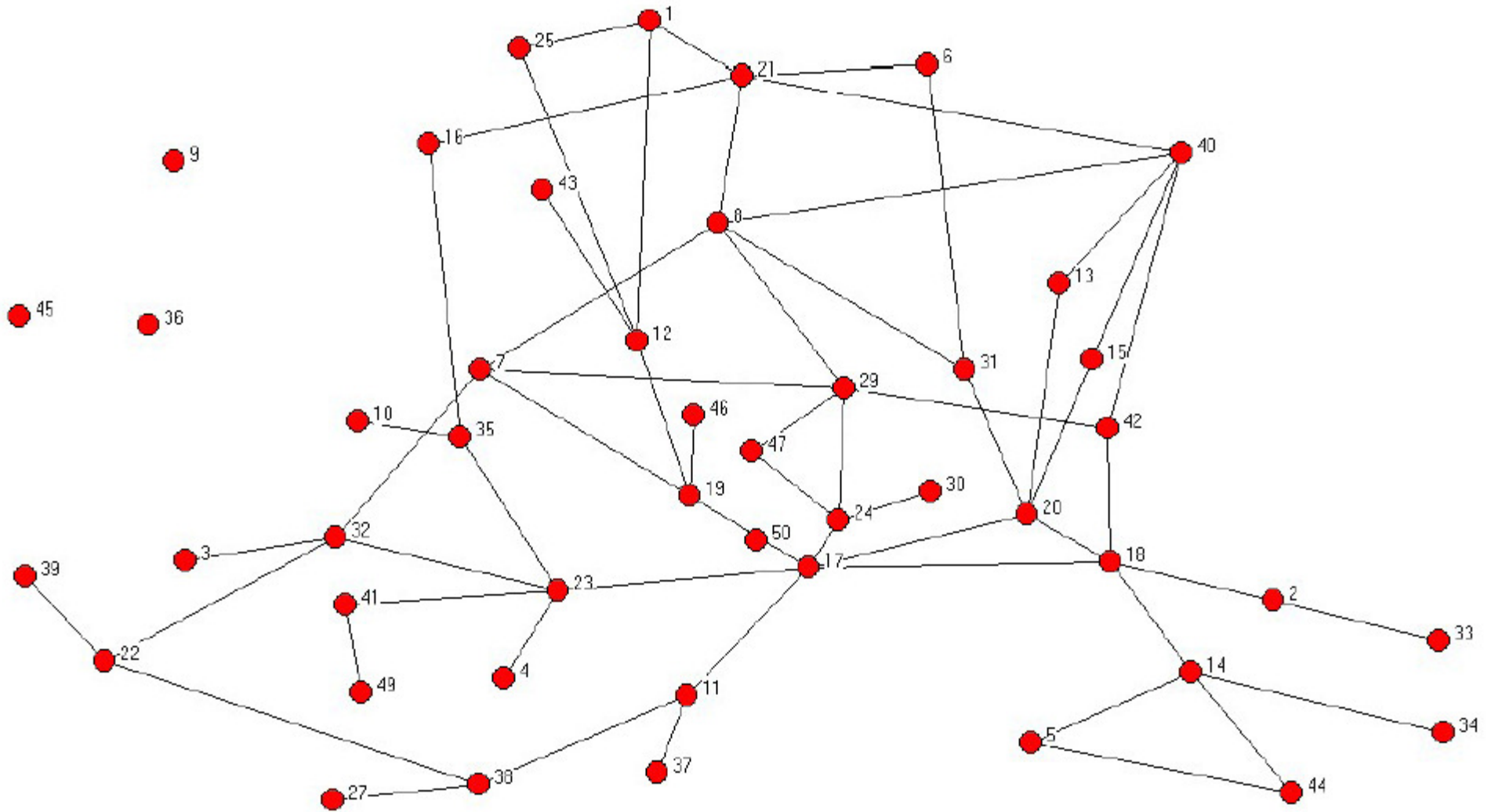
$p = 0.03$; 50 nodes (2)

$p = 0.02$ for the emergence of a cycle and a giant component



Phase Transitions in Random Networks

$p = 0.1$; 50 nodes (3)



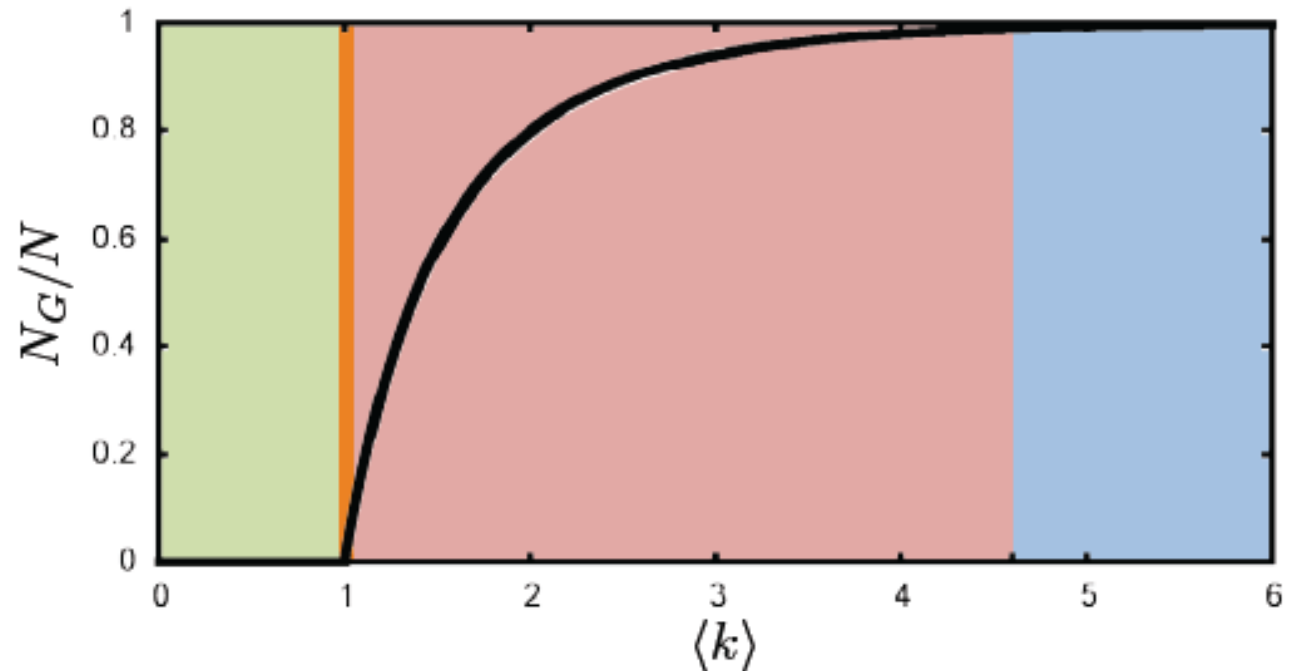
Evolution of a Random Network

- Giant component is the largest cluster within the network.
- The size of the giant component (N_G) varies with the average degree $\langle k \rangle$.
 - For $p = 0$, we have $\langle k \rangle = 0$. Hence, we observe only isolated nodes. Hence, $N_G = 1$ and $N_G/N \rightarrow 0$ for large N .
 - For $p = 1$, we have $\langle k \rangle = N-1$. Hence, the network is a complete graph and all nodes belong to a single cluster. Hence, $N_G = N$ and $N_G/N = 1$.
- One would expect that the giant component will grow gradually from $N_G = 1$ to $N_G = N$ if we increase $\langle k \rangle$ from 0 to $N-1$.
 - However, as observed from theoretical analysis studies, N_G/N remains 0 for small $\langle k \rangle$. Once $\langle k \rangle$ exceeds a critical value (1), N_G/N increases rapidly signaling the emergence of a giant component.
 - We have a giant component if and only if when each node has on average more than one link.

Evolution of a Random Network

- We know that $\langle k \rangle = p(N-1)$.
- For the critical value of $\langle k \rangle = 1$ when the giant component emerges, $p_c(N-1) = 1$.
 - $p_c = 1/(N-1) \approx 1/N$.
 - This indicates: Larger the network, the smaller the value of p for the emergence of a giant component.

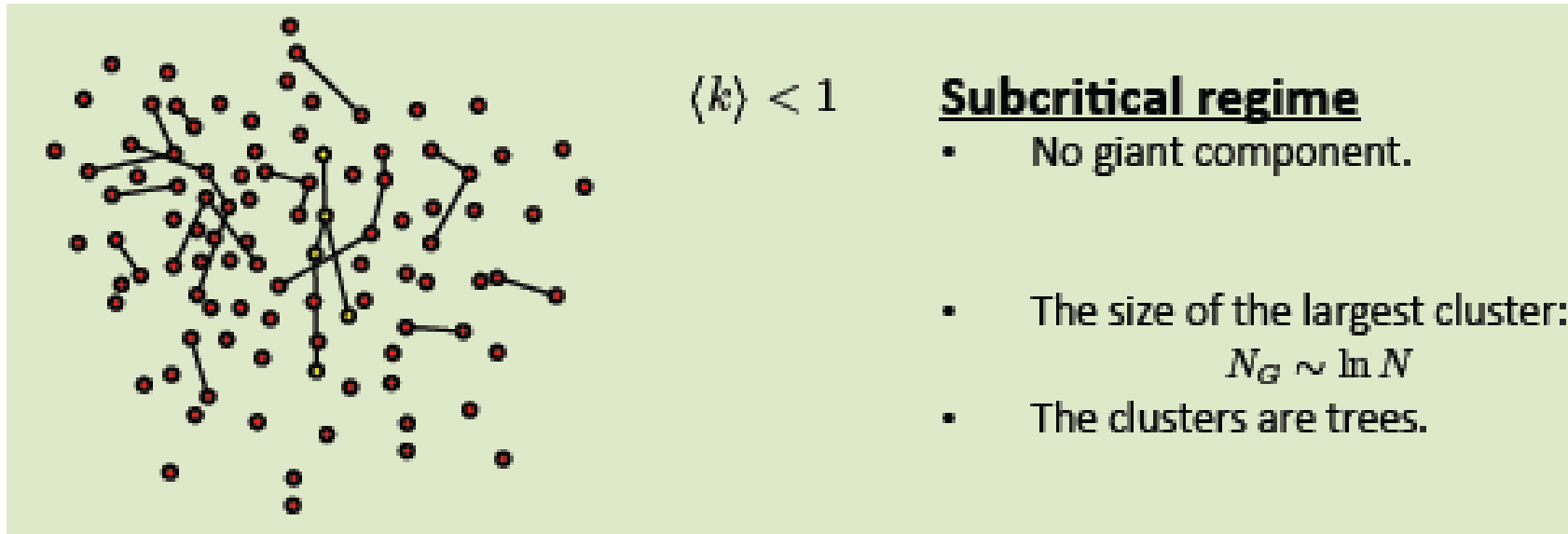
where $S = N_G/N$



Source: Figure 3.6a
Barabasi

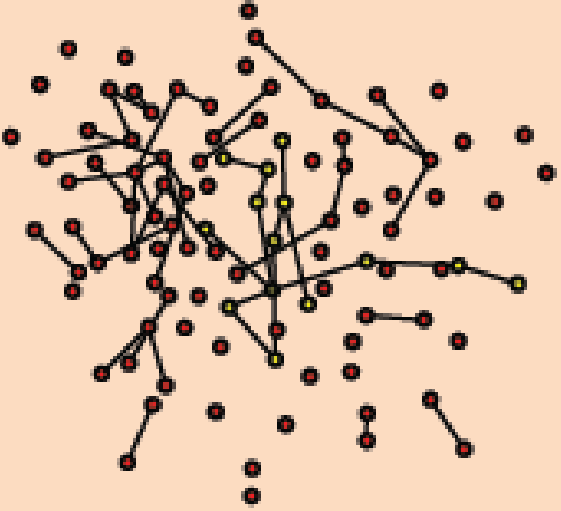
Evolution: Topological Transitions

Consider creating a random network according to the $G(N, L)$ model



- Subcritical regime:
 - $0 < \langle k \rangle < 1$ and $p < 1/N$
 - The largest cluster is expected to be a tree with $\ln N$ nodes. Hence, $N_G/N = \ln N/N \rightarrow 0$ in the $N \rightarrow \infty$ limit, indicating that the largest component is tiny compared to the size of the network.
 - Components have comparable sizes, lacking a clear winner to be designated as a giant component.

Evolution: Topological Transitions



$\langle k \rangle = 1$

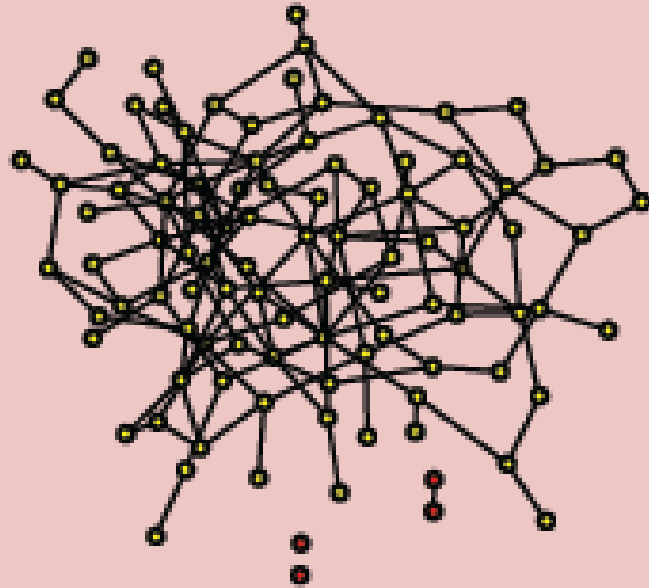
Critical point

- No giant component.
- Size of the largest cluster:
 $N_G \sim N^{2/3}$
- The clusters may contain loops.

- **Critical Point:**

- $\langle k \rangle = 1$ and $p = 1/N$
- The largest cluster is expected to be of size $N^{2/3}$ and contain loops, while the smaller clusters are typically trees.
- The largest cluster is still tiny compared to the network size. $N_G/N = N^{(-1/3)} \rightarrow 0$ as $N \rightarrow \infty$.

Evolution: Topological Transitions

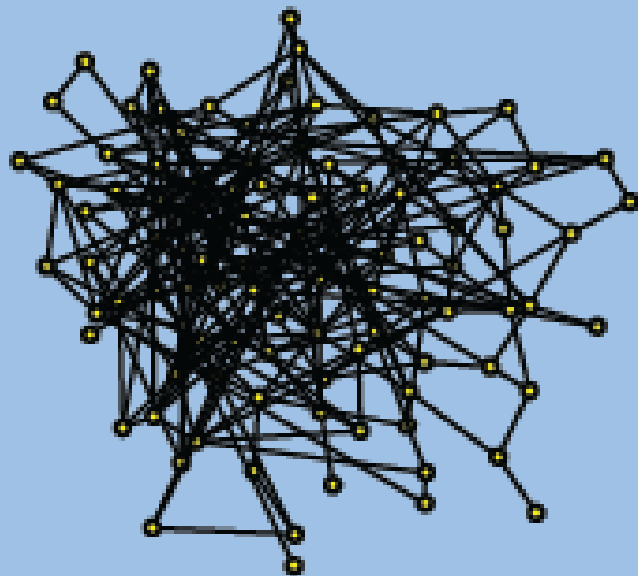


$$\langle k \rangle > 1$$

$$p > 1/N$$

Supercritical regime

- Single giant component.
- Size of the giant component:
 $(N_G / N) \sim (p - p_c)$
- The small clusters are trees.
- GC has loops.



$$\langle k \rangle \geq \ln N$$

$$p \geq (\ln N)/N$$

Fully connected regime

- Single giant component.
- No isolated nodes or clusters.
- Size of the giant component:
 $N_G = N$
- GC has many loops.

Prediction of Random Network Theory: Real Networks are Supercritical

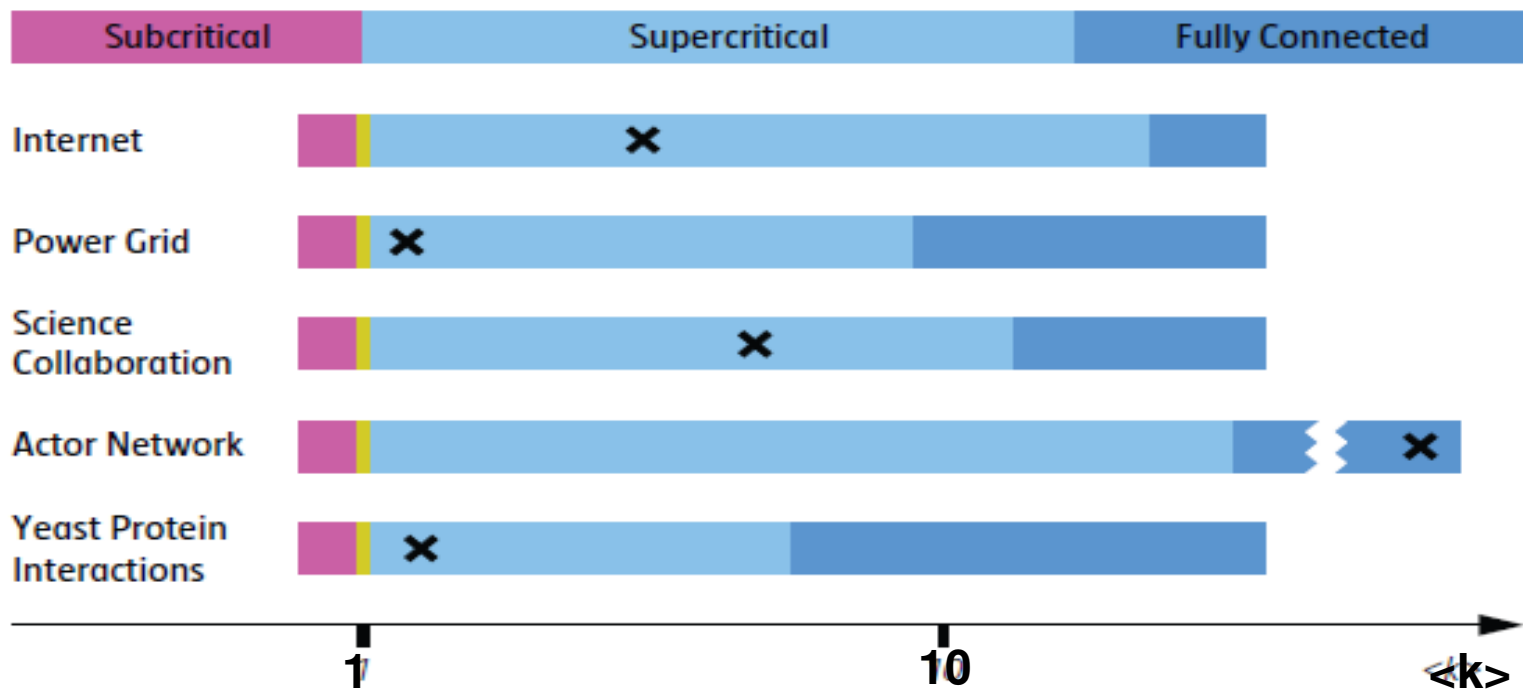
- The theoretical thresholds uncovered for random networks are:
 - For $\langle k \rangle > 1$, a giant component emerges that contains a finite fraction of all nodes.
 - For $\langle k \rangle > \ln N$, all components are absorbed by the giant component, resulting in a single connected network.

Source:
Table 3.1
Barabasi

Network	N	L	$\langle k \rangle$	$\ln N$	$\frac{\langle k \rangle}{\ln N}$	$\frac{\ln N}{\langle k \rangle}$
Internet	192,244	609,066	6.34	12.17	0.52	6.59
Power Grid	4,941	6,594	2.67	8.51	0.31	8.67
Science Collaboration	23,133	186,936	8.08	10.04	0.80	4.81
Actor Network	212,250	3,054,278	28.78	12.27	2.35	3.65
Yeast Protein Interactions	2,018	2,930	2.90	7.61	0.38	7.15

Prediction of Random Network Theory: Real Networks are Supercritical

- Just based on the N and L values for the real networks, we could predict (according to the principles of Random Network Theory) that:
 - All real networks should have a giant component (since their $\langle k \rangle$ exceeds 1)
 - For most real networks (except the actor network), the giant component does not absorb all the nodes (components) as their $\langle k \rangle$ value is less than $\ln N$. Hence, most real networks according to Random Network theory are in the supercritical topology regime.



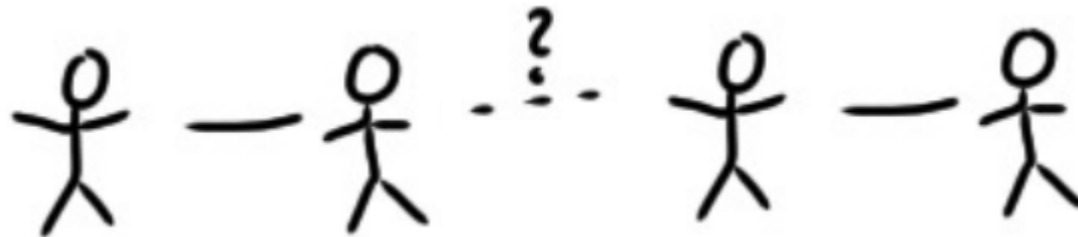
Source:
Figure 3.8
Barabasi

Giant Components: Intuitive Idea



If your friend starts getting connected to someone other than yourself, then you are more likely to belong to a larger component.

The emergence of the giant component sets in when each node has degree of at least 1. Any new edge added to the network is more likely to merge two disconnected groups. Hence, the giant component is very likely to emerge if the average degree of a node exceeds 1.



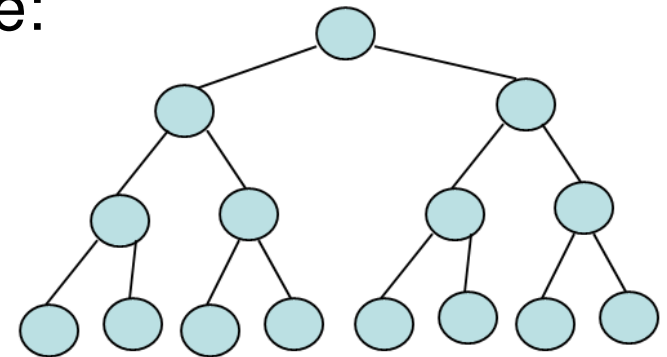
As the network evolves, there cannot be two giant components. The addition of new edges is likely to merge two giant components and evolve them as one single giant component.

Try this applet: <http://ccl.northwestern.edu/netlogo/models/run.cgi?GiantComponent.884.534>

Small World Property

- Distance between two randomly chosen nodes in a network is surprisingly short.
- Consider a random network with average degree $\langle k \rangle$. A node in this network has on average:

- $\langle k \rangle$ nodes at distance one ($d = 1$).
- $\langle k \rangle^2$ nodes at distance two ($d = 2$).
- $\langle k \rangle^3$ nodes at distance three ($d = 3$).
-
- $\langle k \rangle^d$ nodes at distance d .



- The expected number of nodes up to distance d from the starting node is:

$$N(d) = 1 + \langle k \rangle + \langle k \rangle^2 + \dots + \langle k \rangle^d = \frac{\langle k \rangle^{d+1} - 1}{\langle k \rangle - 1}$$

Small World Property

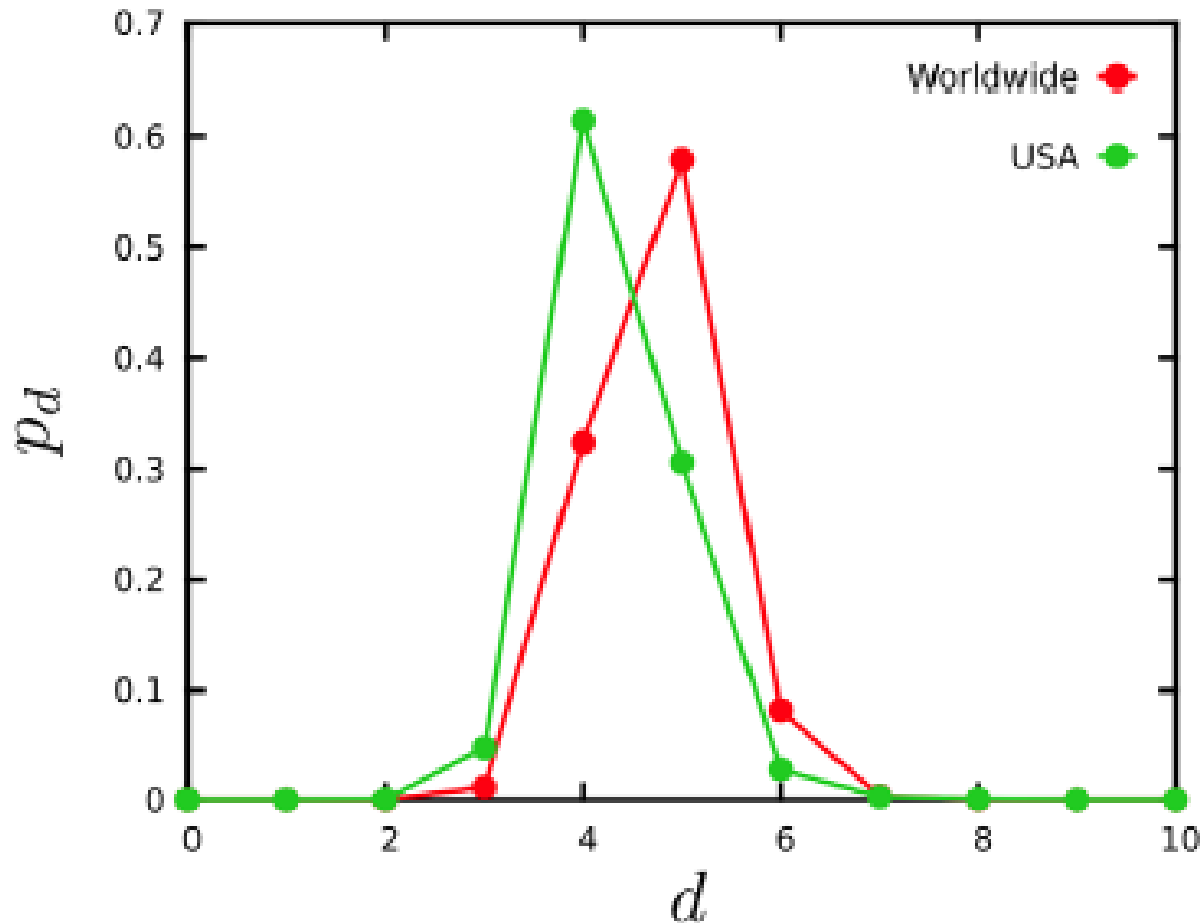
- Let d_{max} denote the maximum distance (the network diameter) at which $N(d)$ reaches N . That is, $N(d_{max}) = N$.
- Assuming that $\langle k \rangle \gg 1$,
 - $\langle k \rangle^{d_{max}} \approx N$.
 - $d_{max} = \ln N / \ln \langle k \rangle$
- As seen from the results for real networks, $\ln N / \ln \langle k \rangle$ approximates more better for the average distance between two randomly chosen nodes.
 - This is because d_{max} is often dominated by a few extreme paths, while $\langle d \rangle$ is averaged over all node pairs, a process that diminishes the fluctuations.
- Thus, the average distances $\langle d \rangle$ in a random network are proportional to $\ln N$, rather than N .
- The $1/(\ln \langle k \rangle)$ term implies that denser the network, the smaller is the distance between the nodes.

Small World Property

<i>Network Name</i>	N	L	$\langle k \rangle$	$\langle d \rangle$	d_{max}	$\frac{\log N}{\log \langle k \rangle}$
Internet	192,244	609,066	6.34	6.98	26	6.59
WWW	325,729	1,497,134	4.60	11.27	93	8.32
Power Grid	4,941	6,594	2.67	18.99	46	8.66
Mobile Phone Calls	36,595	91,826	2.51	11.72	39	11.42
Email	57,194	103,731	1.81	5.88	18	18.4
Science Collaboration	23,133	186,936	8.08	5.35	15	4.81
Actor Network	212,250	3,054,278	28.78	-	-	-
Citation Network	449,673	4,707,958	10.47	11.21	42	5.55
E Coli Metabolism	1,039	5,802	5.84	2.98	8	4.04
Yeast Protein Interactions	2,018	2,930	2.90	5.61	14	7.14

Source: Table 3.2: Barabasi

Small World Property: Facebook



For Facebook,
 $N = 7 \times 10^9$ users
 $\langle k \rangle = 1000$

$$\langle d \rangle = \frac{\ln 7 \times 10^9}{\ln(10^3)} = 3.28$$

Based on the actual
Facebook connections,
 $\langle d \rangle = 4.74$

Clustering Coefficient

- The local clustering coefficient C_i captures the density of links in node i 's immediate neighborhood.
 - $C_i = 0$ implies there are no links between i 's neighbors
 - $C_i = 1$ implies that each of node i 's neighbors link to each other.
- Let k_i be the degree of node i .
- Max. number of possible links between the k_i neighbors of node i are $k_i(k_i - 1)/2$.
- If p is the probability that any two nodes in a network are connected, then the number of links between the k_i neighbors of node i is:

$$\langle L_i \rangle = p \frac{k_i(k_i - 1)}{2}$$

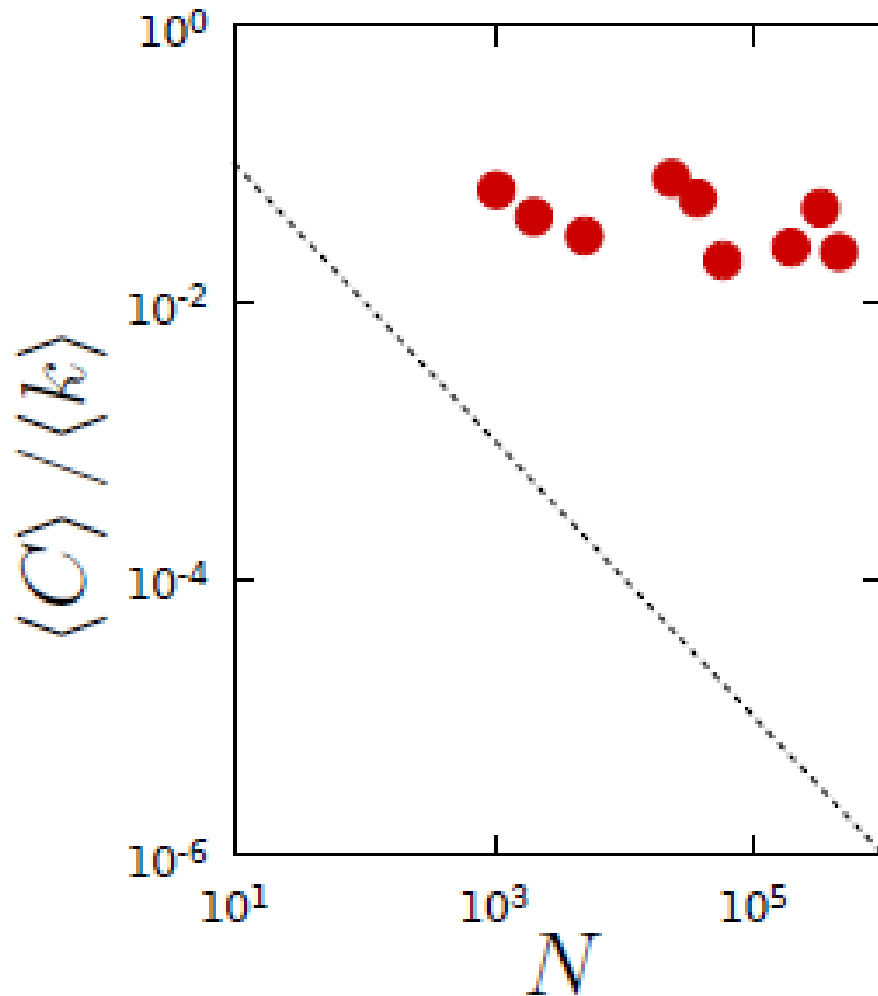
- Local clustering coefficient of node i :

$$C_i = \frac{2\langle L_i \rangle}{k_i(k_i - 1)} = p = \frac{\langle k \rangle}{N}$$

Clustering Coefficient

- Observations based on Random Network Theory
- For fixed $\langle k \rangle$, the larger the network, the smaller is a node's clustering coefficient.
 - Thus, the network's average clustering coefficient $\langle C \rangle$ is expected to decrease as $1/N$.
- The local clustering coefficient of a node is independent of the node's degree

Clustering Coefficients for Real Networks



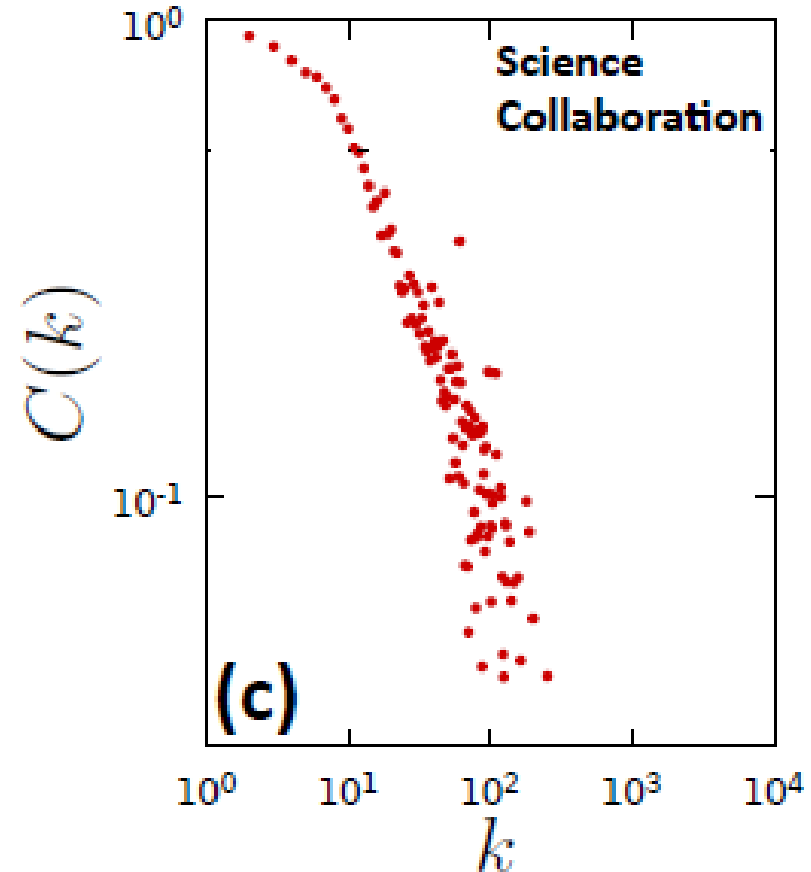
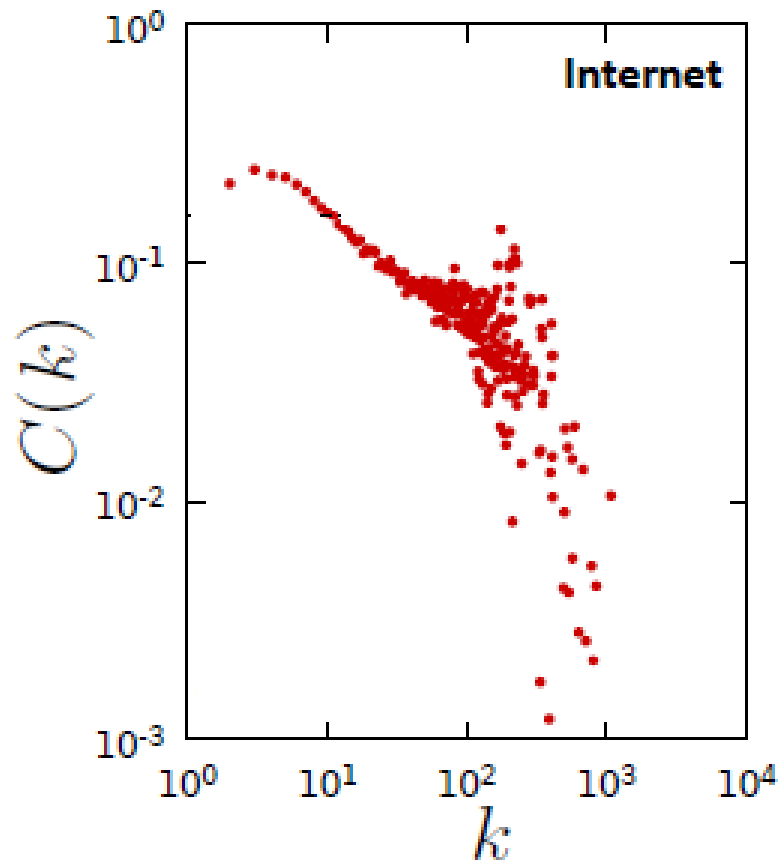
Each circle corresponds to a real network.

Directed networks were made undirected to calculate C .

For random networks, the average clustering coefficient decreases as $1/N$. In contrast, for real networks, $\langle C \rangle$ has only a weak dependence on N .

Real networks have a much higher Clustering coefficient than expected for a random network of similar N and L .

Clustering for Real Networks



$C(k)$ is measured by averaging the local clustering coefficient of all nodes with the same degree k .

According to the Random Network theory model, $C(k)$ is independent of the individual node degrees. However, we find that $C(k)$ decreases as k increases.

Nodes with fewer neighbors have larger local clustering coefficients and vice-versa

Clustering Coeff. Real Networks

• Networks	Actual	Random, $G(n, p)$
– Prison		
Friendships	0.31	0.0134
Co-authorships		
Math	0.15	0.00002
Biology	0.09	0.00001
Economy	0.19	0.00002
WWW		
Web links	0.11	0.002

Real Networks are not Random

- Degree distribution:
 - Random networks – binomial distribution, in general, and Poisson distribution for $k \ll N$.
 - Highly connected nodes (hubs) are effectively forbidden.
 - Real networks: More highly connected nodes, compared to that predicted with random model.
- Connectedness:
 - Random networks: One single giant component exists only if $\langle k \rangle > \ln N$.
 - Real networks: One single giant component exists for several networks with $\langle k \rangle < \ln N$.
- Average Path Length (small world property):
 - For both random and real networks, the average path length scales as $\log N / \log \langle k \rangle$.
- Clustering coefficient:
 - Random model: Local clustering coefficient is independent of the node's degree and $\langle C \rangle$ depends on the system size as $1/N$.
 - Real networks: C decreases with node degrees and is largely independent of the system size.

Real Networks are not Random

- Except for the small world property, the properties observed for real-world networks are not matching with that observed for random networks.
- Then why study random graph theory?
- If a certain property is observed for real-world networks, we can refer to the random graph theory and analyze whether the property is observed by chance (like the small world property).
- If the property observed does not coincide with that of the random networks (like the local clustering coefficient), we need to further analyze the real-world network for the existence of the property because it did not just happen by chance.
- Establish useful benchmarks (e.g., for component structure, diameter, degree distribution, clustering, etc)

Simulating a Random Network

ER Model

- Let S be the set of all node pairs
- Until S gets empty
 - Pick a node u randomly in the network.
 - If this node has at least one node in the set S that it is not yet considered for a possible edge, then randomly select a node v among these candidate nodes.
 - Generate a random number r
 - If the value of $r \leq p$, the probability for an edge, then connect the two nodes $u-v$.
 - Else do not connect them
 - Either way, remove the node pair $u-v$ from set S

Realistic Variations of the Random Network Model

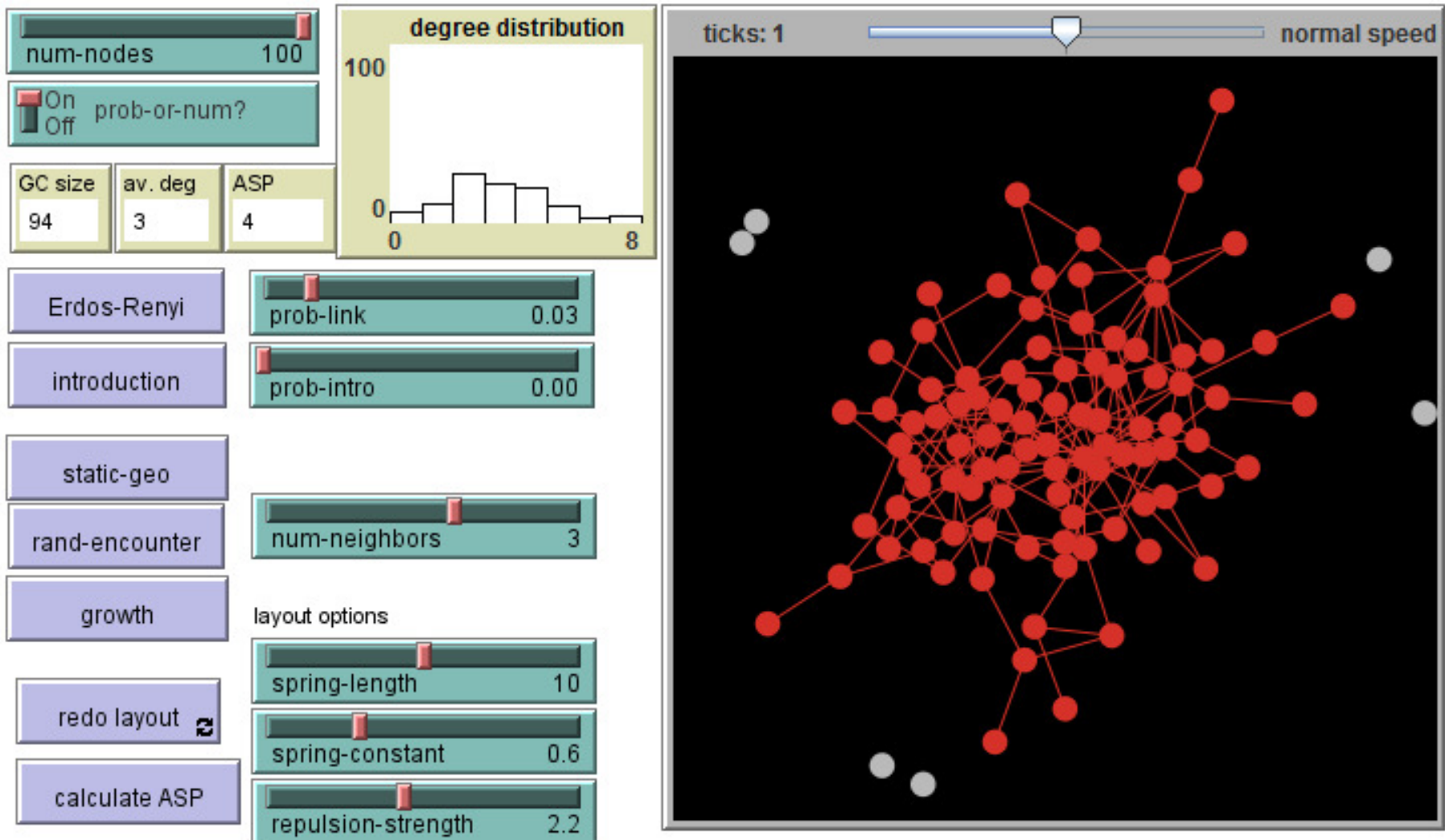
- **Introduction Model:** A node has higher chances of establishing a link with a neighbor of its neighbor (e.g., with the friend of a friend) rather than with an arbitrarily selected node.
 - Operate with a probability, p -intro, the probability that a node prefers to connect to the neighbor of a neighbor node.
- Visit: <http://www.ladamic.com/netlearn/nw/RandomGraphs.html>
- **Key Observations:**
 - Smaller Giant Component Size for smaller p ;
 - Larger average shortest path length;
 - Uneven node degree distribution;

Simulating a Random Network

Introduction Model

- Let S be the set of all node pairs
- Until S gets empty
 - Pick a node u randomly in the network.
 - If this node has at least one node in the set S that it is not yet considered for a possible edge
 - Generate a random number r -intro.
 - If r -intro \leq p -intro, the set of candidate nodes that are chosen for connection are the unconnected neighbors of neighbor nodes.
 - Else, the set of candidate nodes are all the unconnected nodes in the network.
 - Among the chosen candidate nodes, the node connects to a randomly chosen node v with a probability p .
 - » Generate a random number r
 - » If the value of $r \leq p$, the probability for an edge, then connect the two nodes u - v .
 - » Else do not connect them
 - Either way, remove the node pair u - v from set S

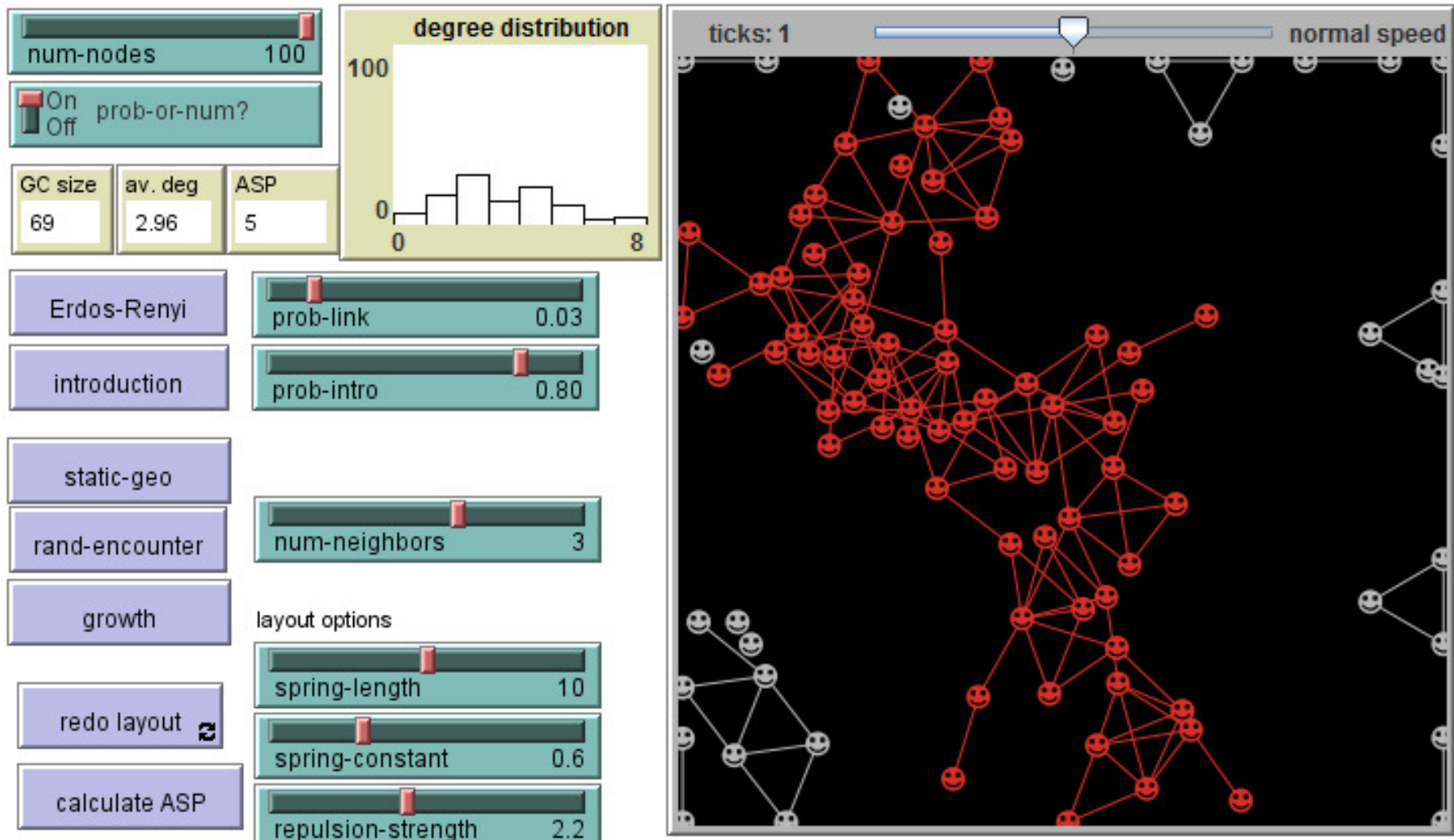
ER Model



Num Nodes = 100; $p = 0.03$; $p\text{-intro} = 0$

GC Size – 94; Avg. Degree = 3; Avg. Shortest Path Length = 4

Introduction Model



Num Nodes = 100; $p = 0.03$; $p\text{-intro} = 0.80$

GC Size – 69; Avg. Degree = 2.96; Avg. Shortest Path Length = 5

Problem Example 1

- Consider a random network generated according to the $G(N, p)$ model where the total number of nodes is 12 and the probability that there are links between any two nodes is 0.20. Determine the following:
 - The probability that there are exactly 60 links in the network
 - The average number of links in the network
 - The average node degree
 - The standard deviation of the node degree
 - The average path length (distance between any two nodes in the network)
 - The average local clustering coefficient for any node in the network.
 - The local clustering coefficient for a node that has exactly 5 neighbors.

Problem Example 1: Solution (1)

- There are $N = 12$ nodes
- Prob[link between any two nodes] = $p = 0.2$

Max. possible number of links between any two nodes is
 $(N)(N-1)/2 = (12*11/2) = 66$

(1) Prob[there are exactly 60 links in the network]

$$= C(66, 60) * p^{60} * (1-p)^{(66-60)}$$

$$C(66, 60) = 66! / (60! * 6!)$$

$$= 60! * 61 * 62 * 63 * 64 * 65 * 66 / (60! * 1 * 2 * 3 * 4 * 5 * 6)$$

$$= 90858768$$

Prob[there are exactly 60 links in the network]

$$= 90858768 * (0.2)^{60} * (0.8)^6$$

$$= 2.75 * 10^{-35}$$

Problem Example 1: Solution (2)

- There are $N = 12$ nodes
- $\text{Prob}[\text{link between any two nodes}] = p = 0.2$

Max. possible number of links between any two nodes is $(N)(N-1)/2 = (12*11/2) = 66$

(2) The average number of links in the network = $p * N(N-1)/2$
 $= 0.2 * 66 = 13.2$

(3) Average node degree = $p*(N-1) = 0.2 * 11 = 2.2$

(4) Standard deviation of node degree = $\sqrt{p(1-p) * (N-1)}$
 $= \text{sqrt}(0.2*0.8*11) = 1.33$

(5) Average path length = $\ln N / \ln \langle k \rangle = \ln(12) / \ln(2.2) = 3.15$

(6) Avg. Local clustering coefficient for any node in the network = $p = 0.2$.

(7) The local clustering coefficient for a node in a random network is independent of its number of neighbors. Hence, the answer is 0.2

Problem Example 2

- Consider the evolution of a random network according to the $G(N, p)$ model, where the total number of nodes is 100 and $p = 0.03$. Consider adding (randomly) one link at a time to the network. The total number of links added is sufficiently large enough to create one single connected component of the entire network. Determine the following:
 - The critical value of the probability (of the number of links) that a giant component emerges for the above network and the average size of the giant component at that value?
 - The minimum value of the average degree per node in the giant component of the fully connected regime.
 - The maximum value for the average path length between any two nodes in the giant component that encompasses all the nodes in the network..

Problem Example 2: Solution (1)

- There are $N = 100$ nodes
 - The critical value of the probability (of the number of links) that a larger cluster for the above network?
$$p_c = 1/N = 1/100 = 0.01$$

Supercritical regime $N_G = (p - p_c) * N = (0.05 - 0.01) * 100 = 4$
 - In the fully connected regime, the average node degree has to be at least $\ln N$. That is, $\langle k \rangle \geq \ln N$.
 - $\text{Min } \langle k \rangle = \ln N = \ln(100) = 4.61$
 - The average path length is given by: $\ln N / \ln \langle k \rangle$
 - Using the minimum value of $\ln \langle k \rangle$ in the above expression, we obtain the max. average path length to be: $\ln(100)/\ln(4.61) = 3.01$.