

# Random Graphs

hari sundaram

**hs1@illinois.edu**

adapted from Jure Leskovec's slides

# Introduction

What are examples of directed and undirected graphs in real-world datasets?

## Adjacency Matrix



Example: real-world graphs are usually directed, unweighted graphs

Network	Nodes	Edges (Edges)
World-Wide Web (WWW)	118,717	549
Web of Science (WoS)	4,864,048	607
Compendex plus PL	142,2038	111
Cheriton's ENR	12,949	244
Twitter (Twitter)	107,817	18
Protein-Protein	1491	1,8

Adjacency Matrix:  $A_{ij} = 1$  if  $i$  and  $j$  are connected,  $0$  otherwise.  $A$  is symmetric if the graph is undirected.

## Node Degree



## Complete Graph



## Distance



what does a real-world dataset look like?



what is a consequence of this observation?

Are these numbers surprising? Unexpected?



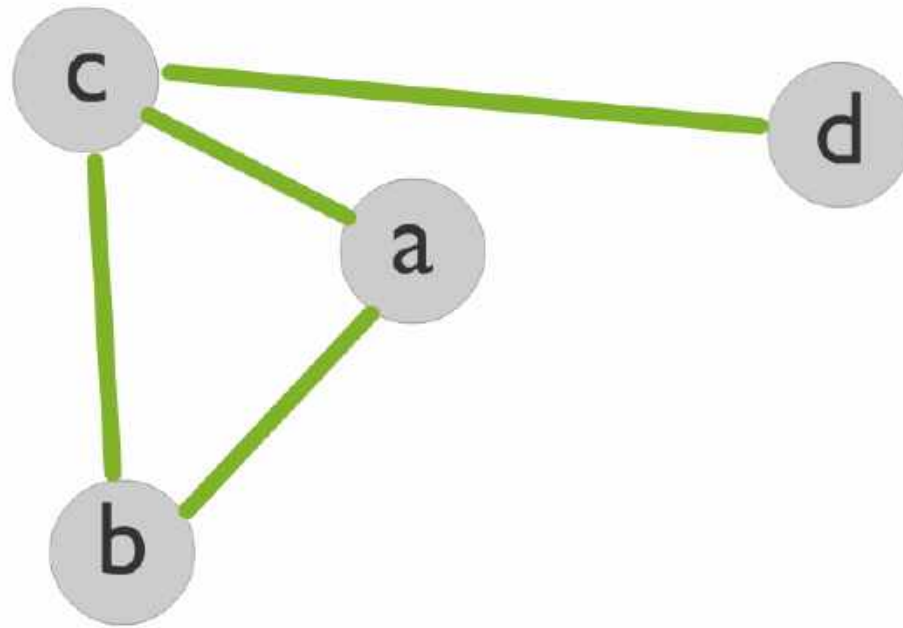
Is the WWW really a "Small World"?

Is the WWW really a "Small World"?

Clustering Coefficient



# Adjacency Matrix



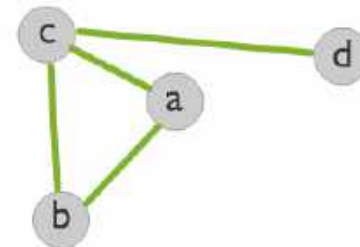
undirected graph

*symmetric*  $A = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$

directed graphs  
do not have  
symmetric  
adjacency  
matrices

What are examples of directed and undirected graphs in real-world datasets?

## Adjacency Matrix



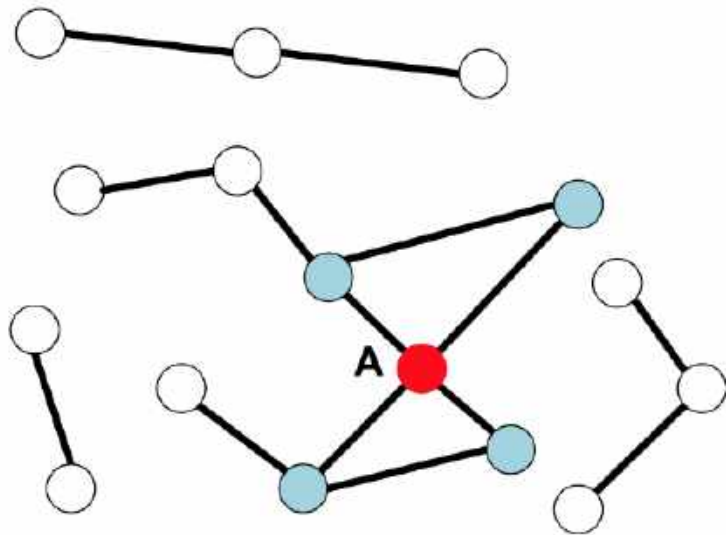
undirected graph

$$A = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

symmetric



# Node Degrees



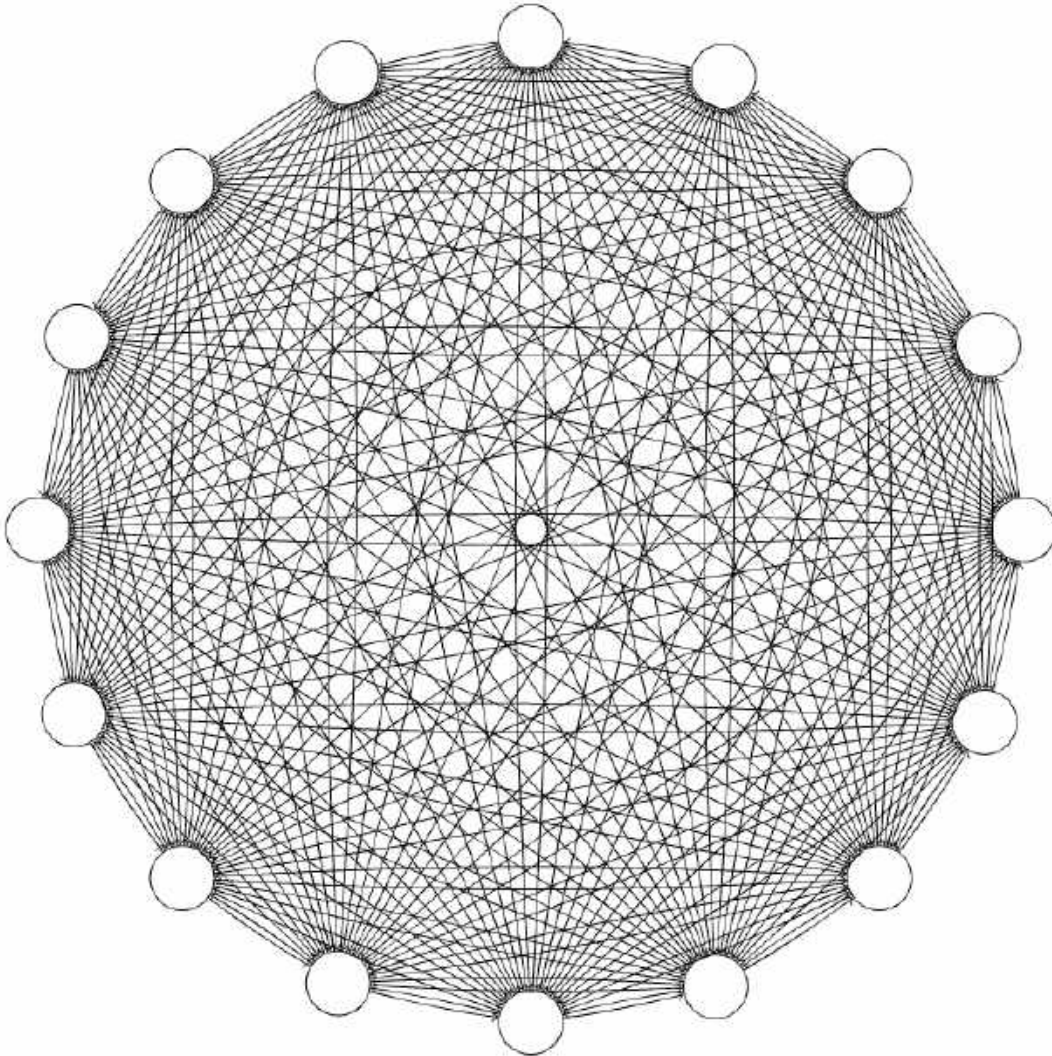
$$k_A = 4$$

$$\bar{k} = \langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i = \frac{2E}{N}$$

↑  
degree of node  $i$

We can compute similar quantities for directed graphs too, except that we will have two numbers.

# Complete Graph



$$E_{\max} = \binom{N}{2} = \frac{n(n-1)}{2}$$

$$\text{degree} = n-1$$

however, most real-world graphs are strikingly different from complete graphs

Network	Nodes	Average degree
WWW (Stanford–Berkeley):	319,717	9.65
Social networks (LinkedIn):	6,946,668	8.87
Communication (MSN IM):	242,720,596	11.1
Coauthorships (DBLP):	317,080	6.62
Roads (California):	1,957,027	2.82
Proteins ( <i>S. Cerevisiae</i> ):	1,870	2.39

matrix density  $\equiv \frac{E}{N^2}$

WWW :  $1.51 \times 10^{-5}$

MSM IM :  $2.27 \times 10^{-8}$

most adjacency matrices are filled with zeros

Most networks are sparse!

$$\bar{k} \ll n - 1$$

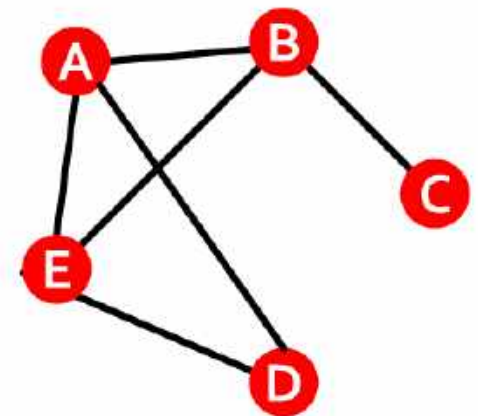
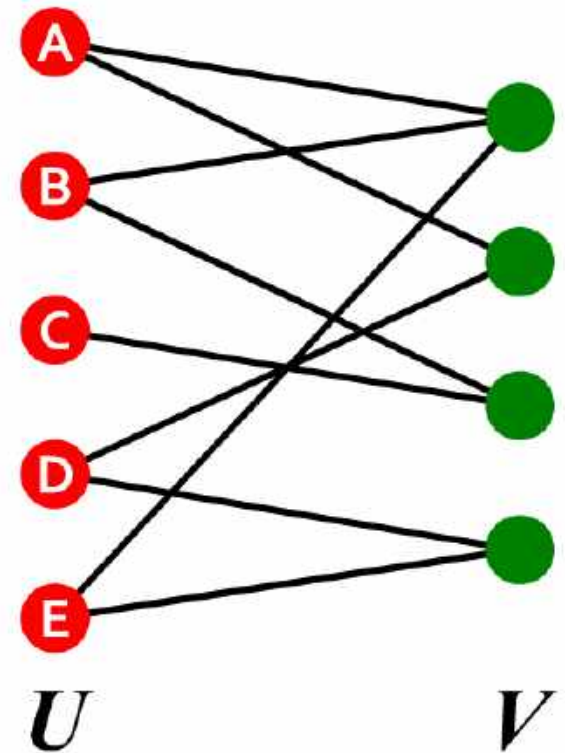


A bipartite graph is a graph whose nodes can be divided into two disjoint sets  $U$  and  $V$  such that every link connects a node in  $U$  to one in  $V$ ; that is,  $U$  and  $V$  are independent sets

That is, it is a set of vertices  $S$  such that for every two vertices in  $S$  there is no edge connecting the two.

Author collaboration networks; Movie co-rating networks

actors-movies; author-papers etc.

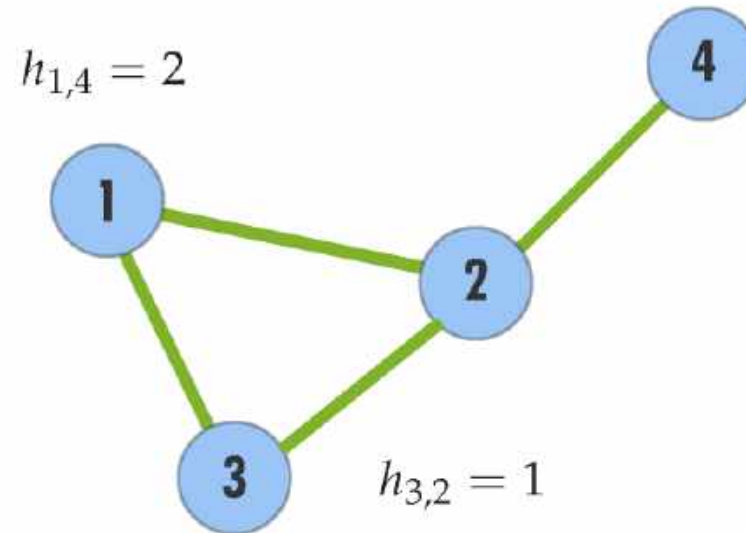


folded graph

# Distance

If the two nodes are disconnected, the distance is usually defined as infinite

**Diameter:** the maximum (shortest path) distance between any pair of nodes in a graph



## BFS

**Distance** (shortest path, geodesic) between a pair of nodes is defined as the number of edges along the shortest path connecting the nodes

$$\bar{h} = \frac{1}{2E_{\max}} \sum_{i,j \neq i} h_{i,j}$$

Average path length for a connected graph (component) or a strongly connected component of a directed graph

Many times we compute the average **only over** the connected pairs of nodes (we ignore “infinite” length paths)


There are **three** useful  
ways to characterize a  
graph





# Degree Distribution

# nodes with degree  $k$


$$P(k) = \frac{N_k}{N}$$

# 2

## Average path length

$$\bar{h} = \frac{1}{2E_{\max}} \sum_{i,j \neq i} h_{i,j}$$

distance between  $i$  and  $j$

$\binom{N}{2}$

# 3

## Clustering coefficient

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

# edges amongst *i*'s neighbors

degree of node *i*

There are **three** useful ways to characterize a graph

1 Degree Distribution

# nodes with degree k

$$P(k) = \frac{N_k}{N}$$

3

Clustering coefficient

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

(degree of node)

2

Average path length

$$\bar{h} = \frac{1}{\binom{N}{2}} \sum_{i,j \neq i} h_{i,j}$$

(distance between i and j)





# MSN Messenger Dataset

June, 2006

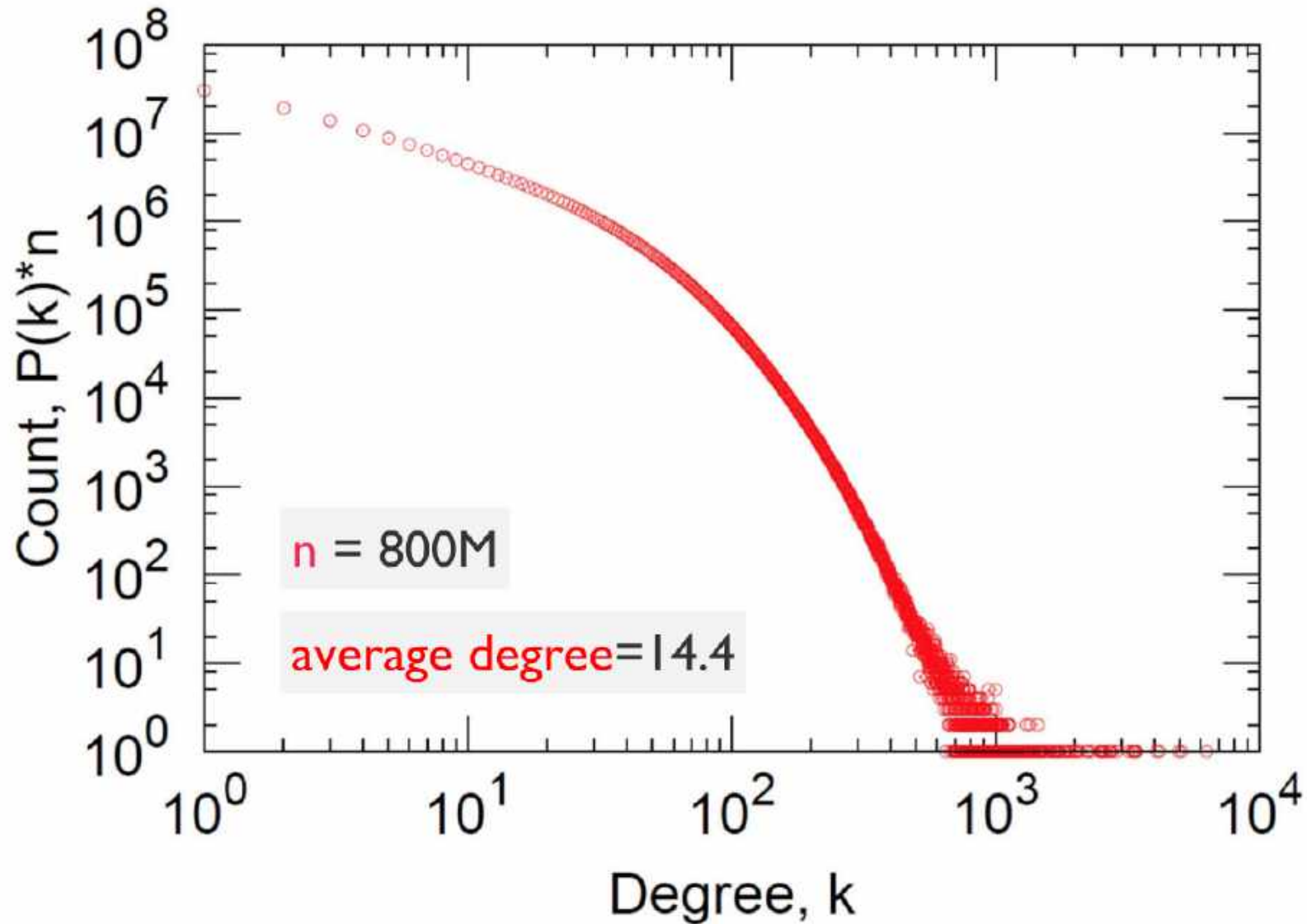
Two individuals are connected if they exchange a message within the last month

- 180M nodes
- 1.3B edges

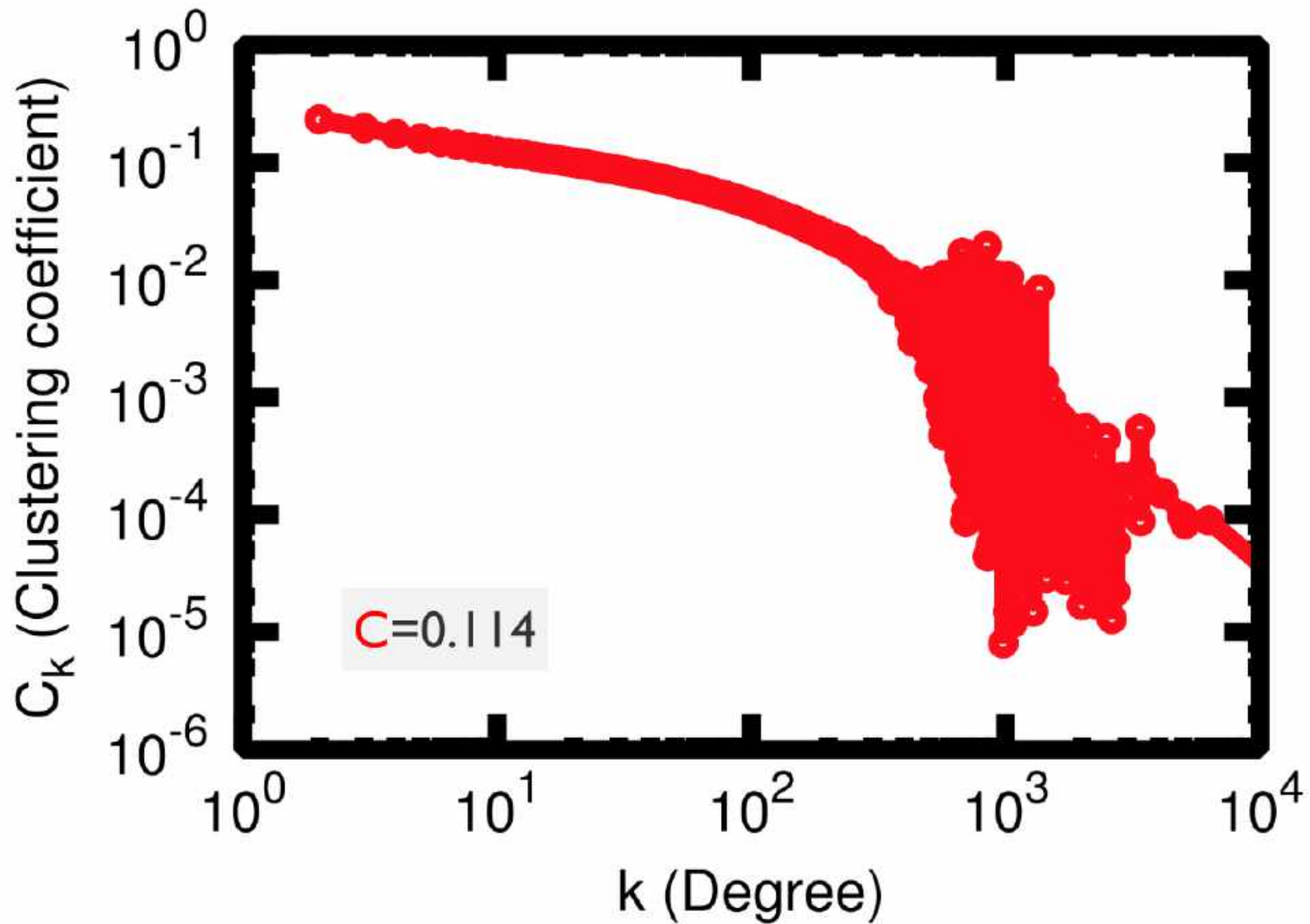
- 4.5Tb / month
- 245 million users logged in
- 180 million users engaged in conversations
- More than 30 billion conversations
- More than 255 billion exchanged messages

J. Leskovec and E. Horvitz. **Planetary-scale views on a large instant-messaging network**. In Proceedings of the 17th international conference on World Wide Web, pages 915–924. ACM, 2008.

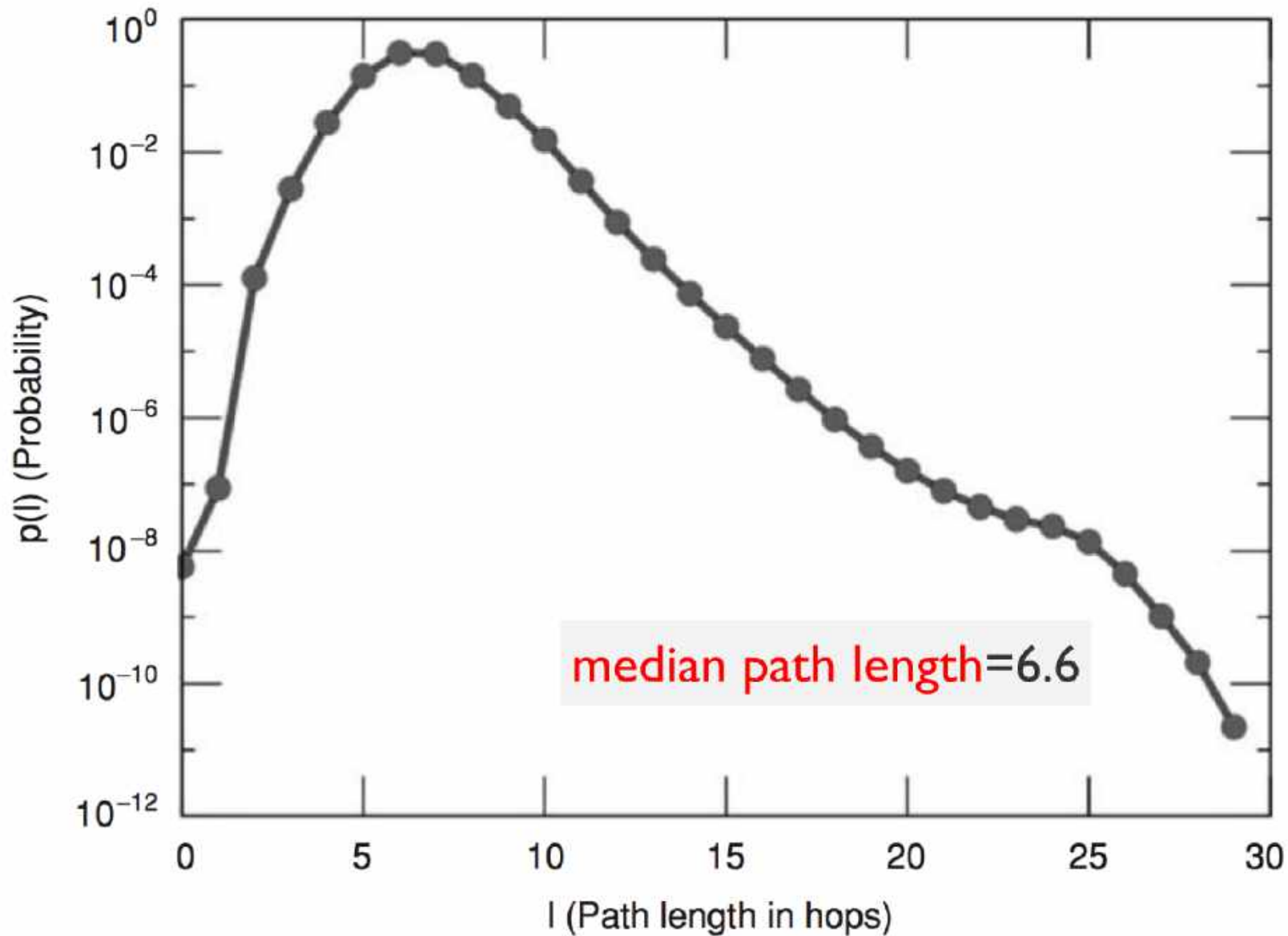
# Degree distribution



# Clustering Coefficient



# Path Length





# MSN Messenger Dataset

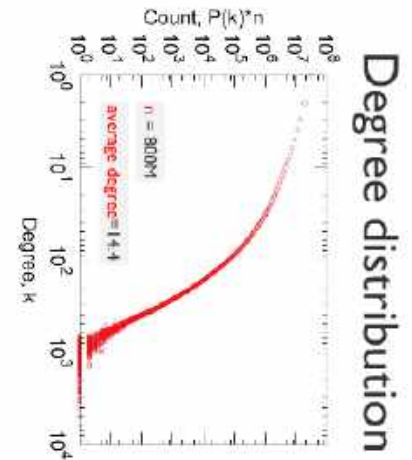
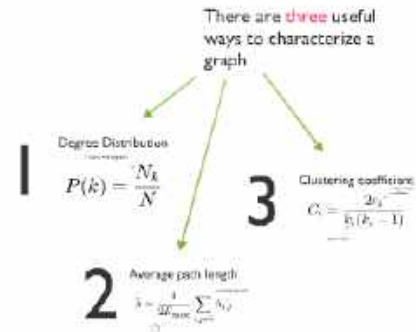
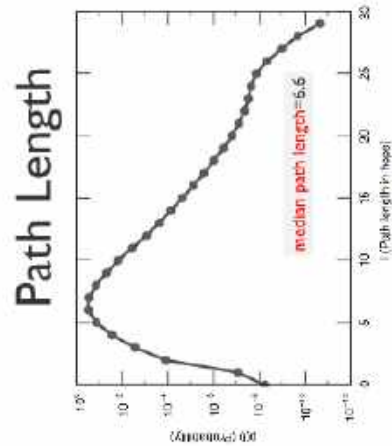
June, 2006

Two individuals are connected if they exchange a message within the last month

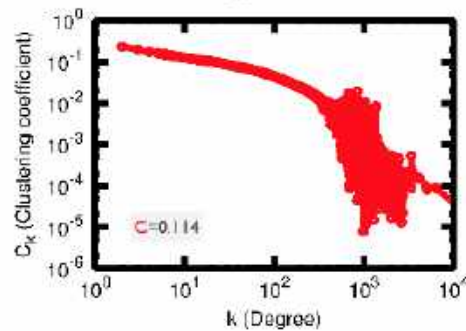
- 180M nodes
- 1.3B edges

- 4.5Tb / month
- 245 million users logged in
- 180 million users engaged in conversations
- More than 30 billion conversations
- More than 255 billion exchanged messages

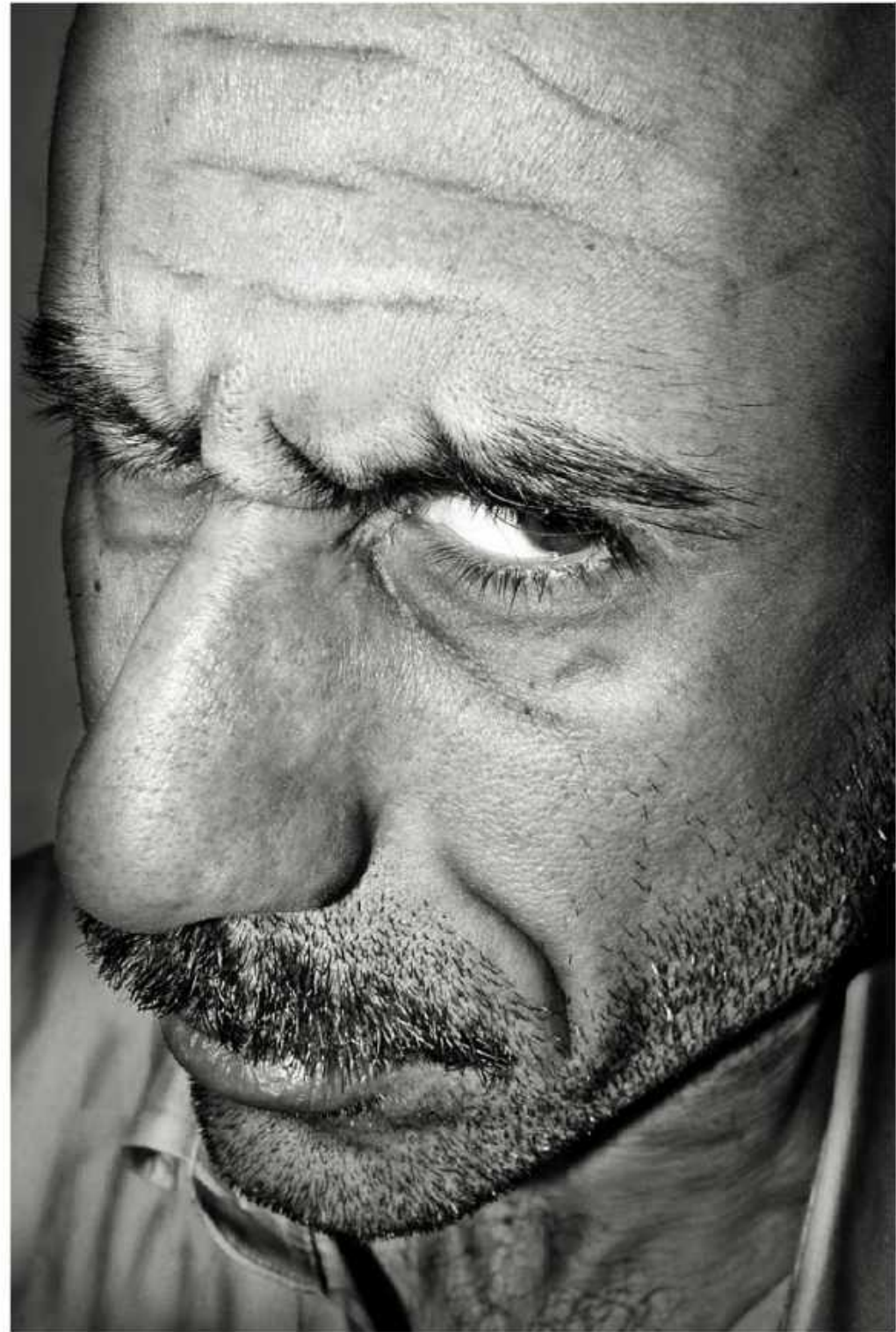
J. Leskovec and E. Horvitz. **Platonic-scale views on a large instant-messaging network.** In Proceedings of the 17th International Conference on World Wide Web, pages 915-924. ACM, 2008.



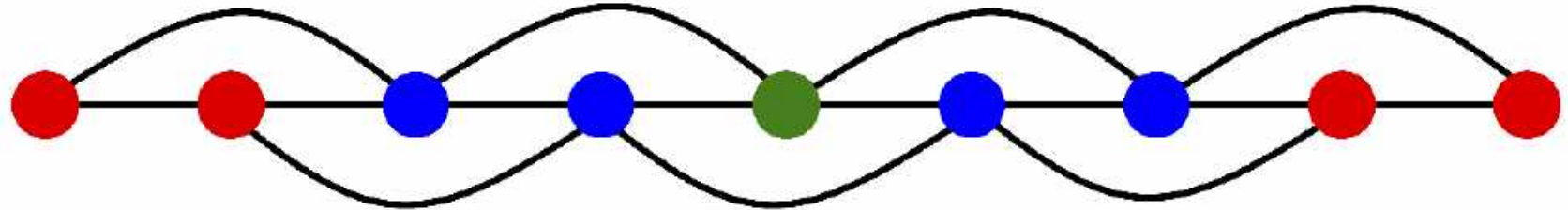
## Clustering Coefficient



**Are these  
numbers  
surprising?  
Unexpected?**



# Is the MSN network like a “Chain”?



$$P(k) = \delta(k - 4) \quad P(k): \text{the degree distribution of a node}$$

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

$$\bar{h} = O(N)$$

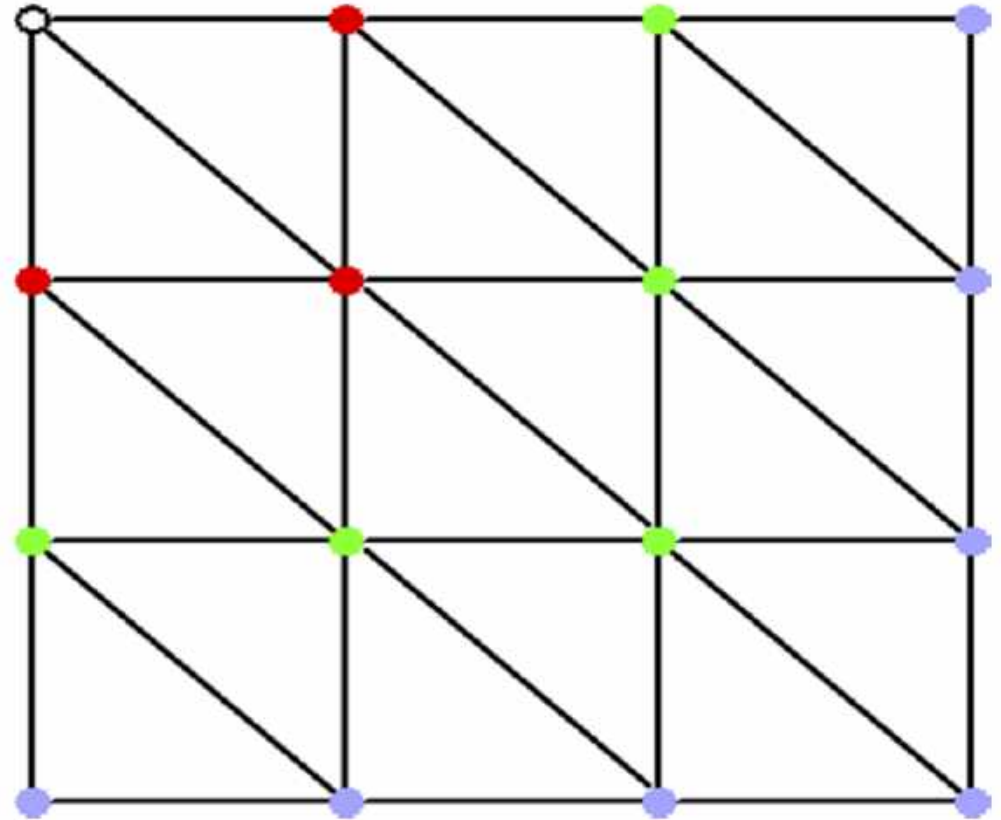
Constant degree; constant average clustering coefficient; average path length linear in  $N$

# Is the MSN network like a "Grid"?

$$P(k) = \delta(k - 6)$$

$$C = \frac{6}{15}$$

$$\bar{h} = O(\sqrt{N})$$



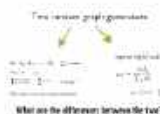
Constant degree, constant clustering coefficient



# Erdős Rényi models



the diameter of the world



What are the differences between the two?

## Poisson models: a reality check

Given components: Clustering coefficient is too low

Average path length is compatible

Degree distributions are unrealistic

Real networks are **not** Poisson!



to what degree is a network random?

# Two variants





$G_{n,p}$

$n$  nodes, edges are i.i.d  
with probability  $p$

2  $G_{n,m}$

$n$  nodes,  $m$  edges are  
picked at random

# Two variants



**1**  $G_{n,p}$

$n$  nodes, edges are i.i.d  
with probability  $p$

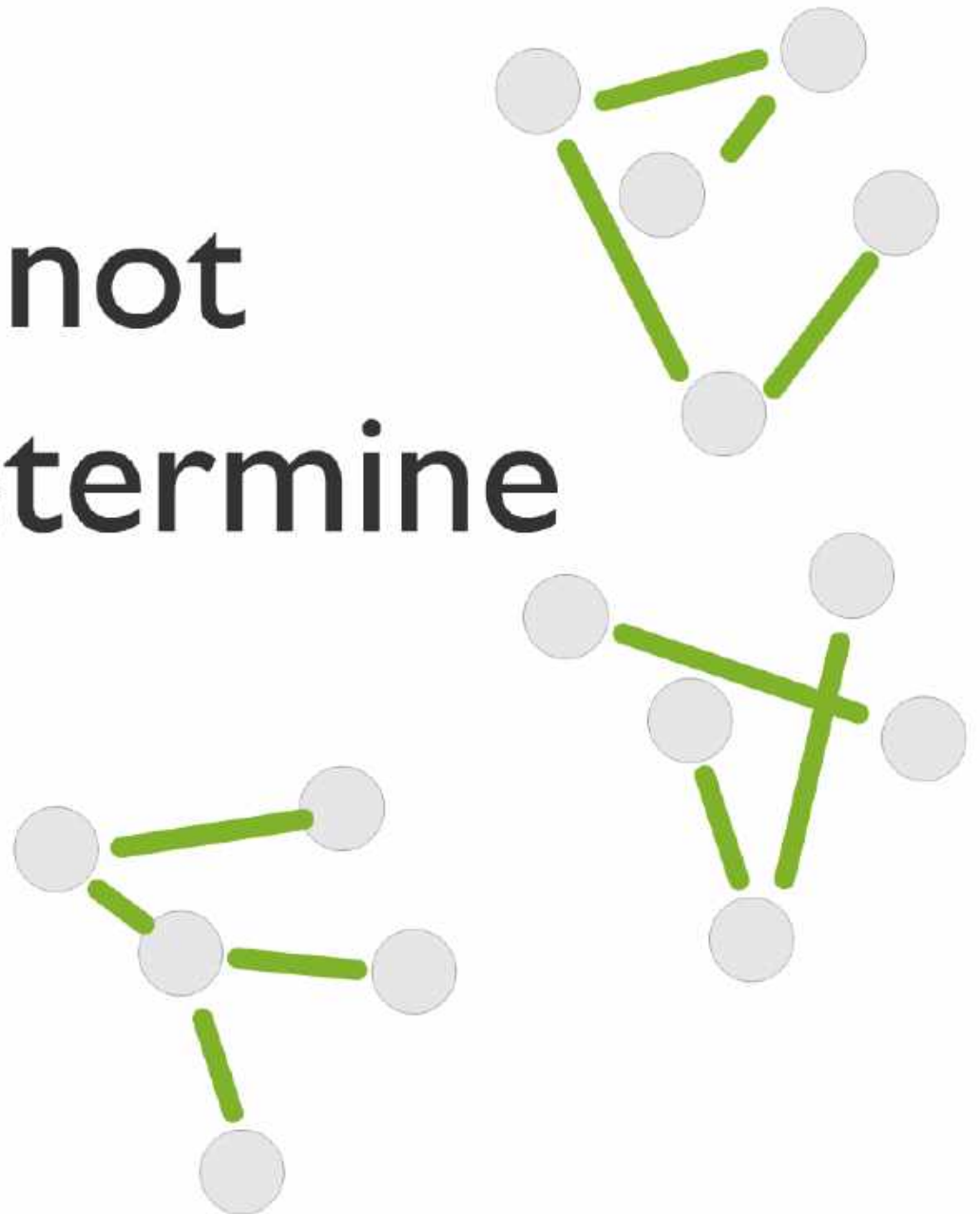
what is the  
difference  
between the  
two?

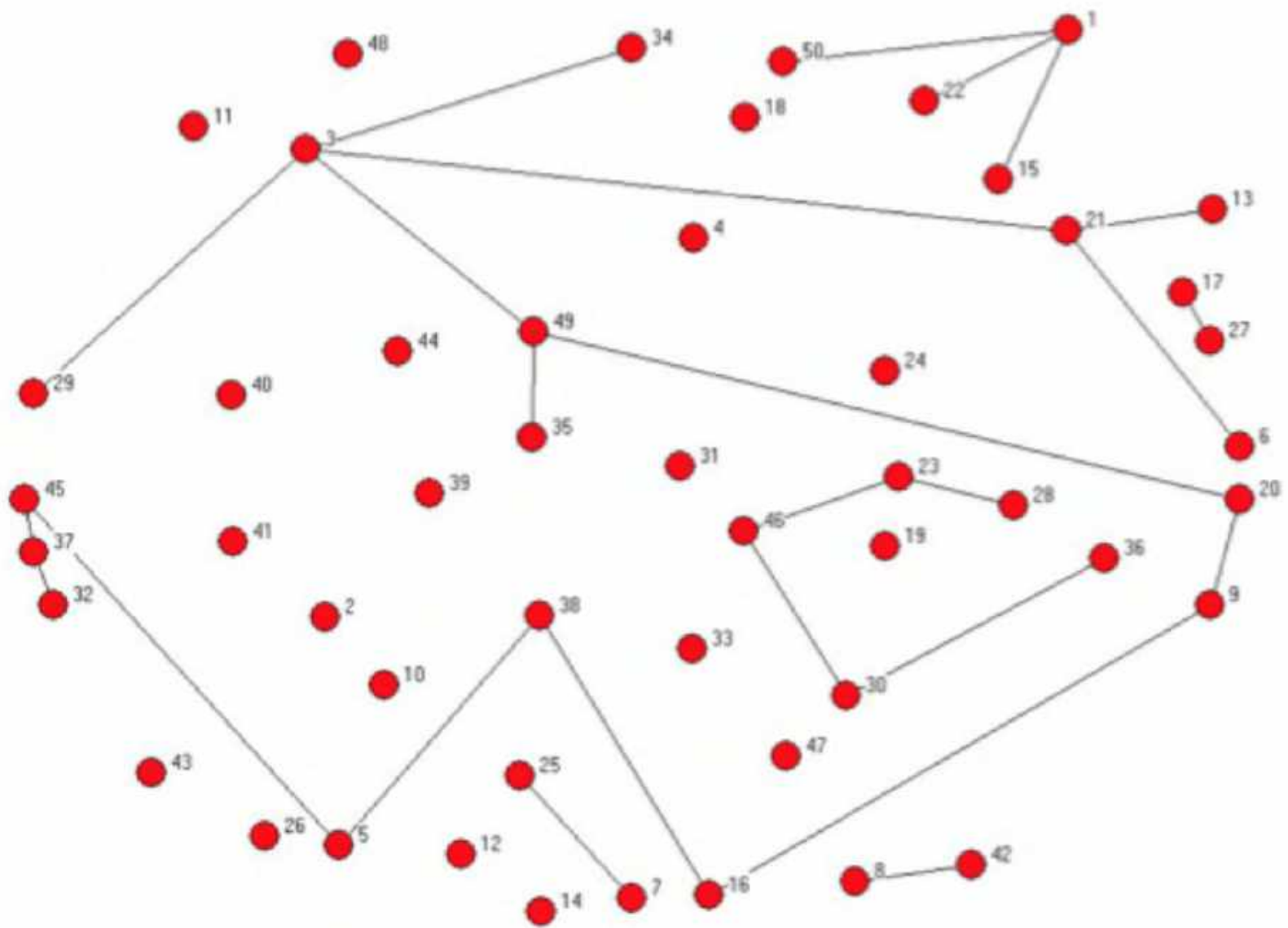


**2**  $G_{n,m}$

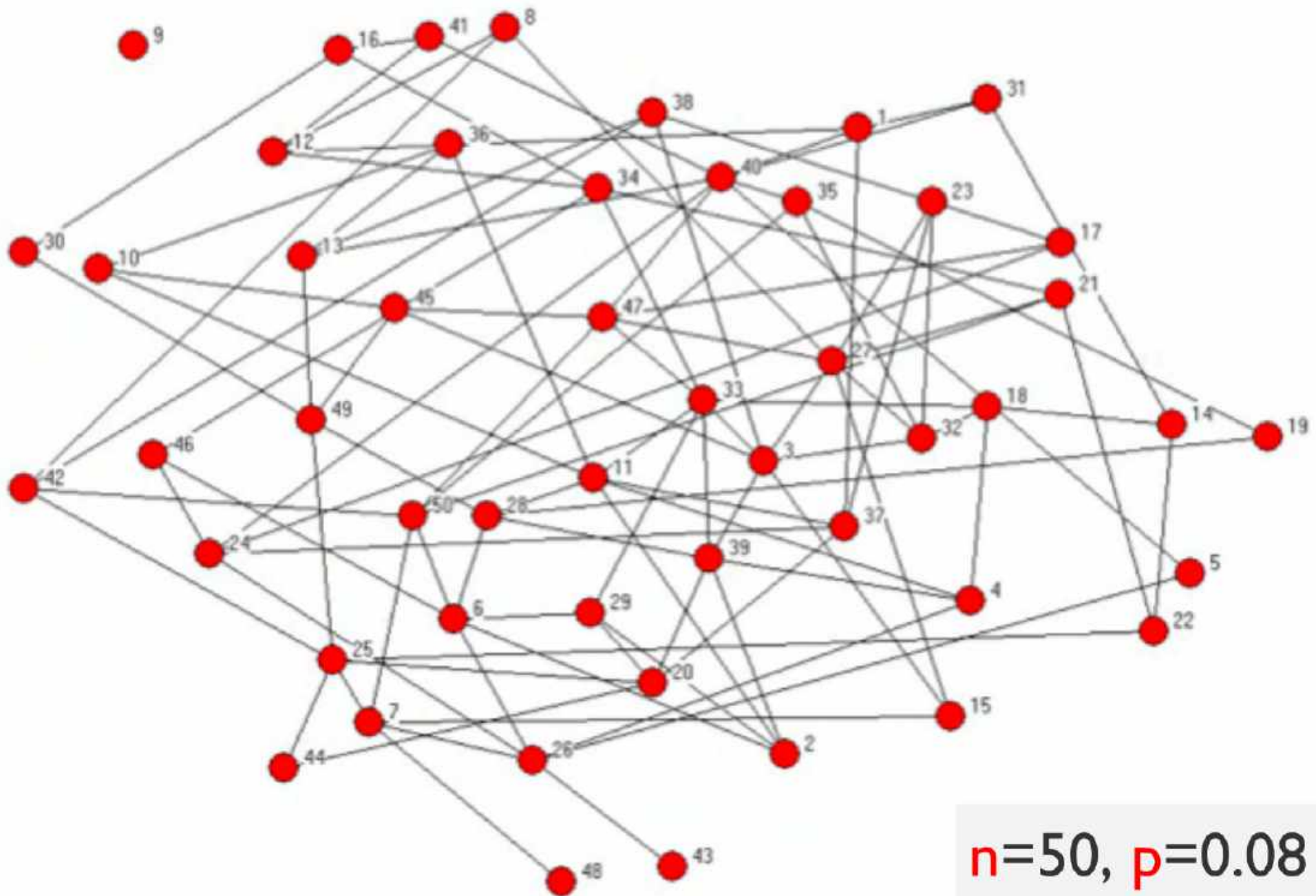
$n$  nodes,  $m$  edges are  
picked at random

**n** and **p** do not  
uniquely determine  
the graph!





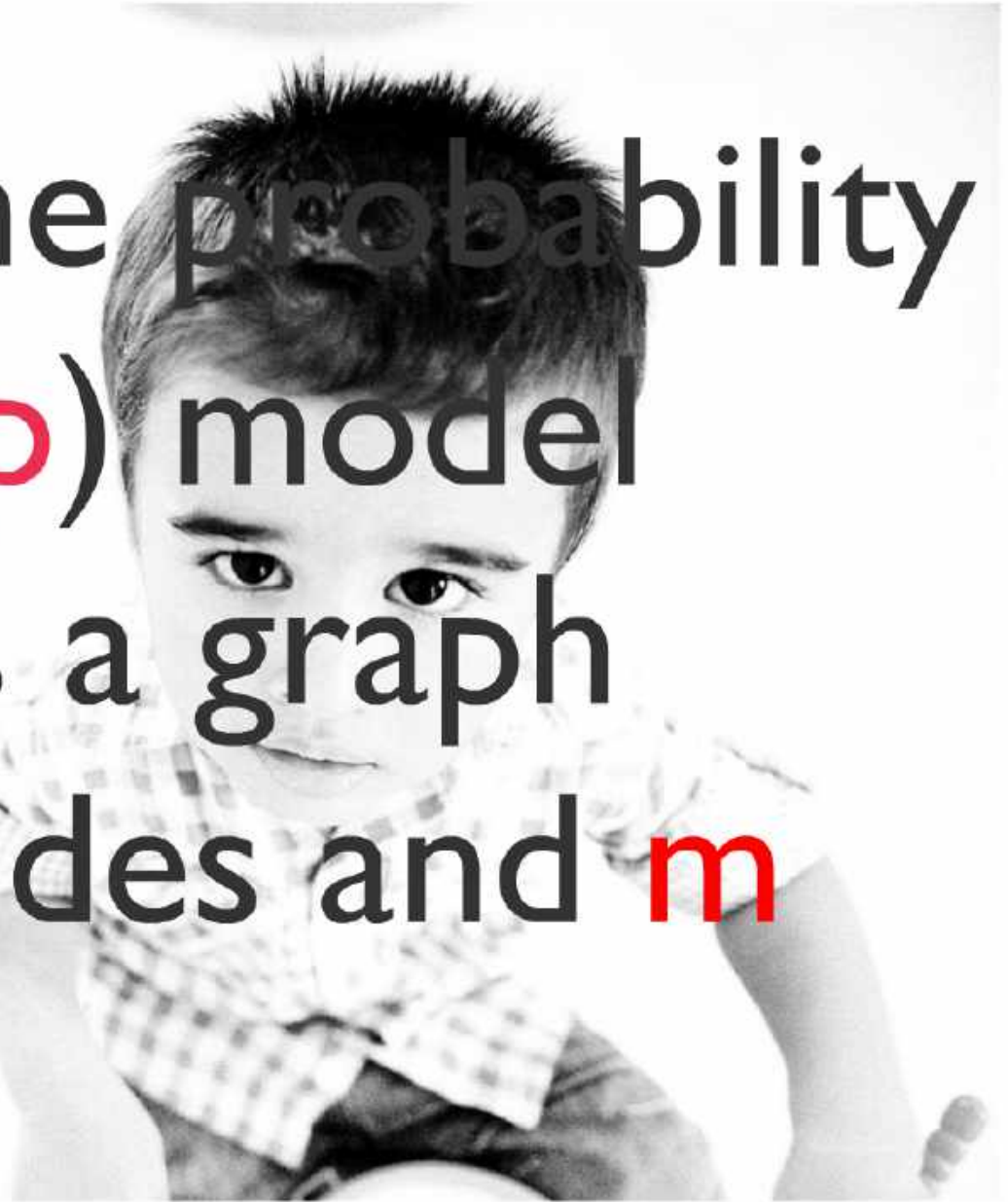
$n=50$ ,  $p=0.02$



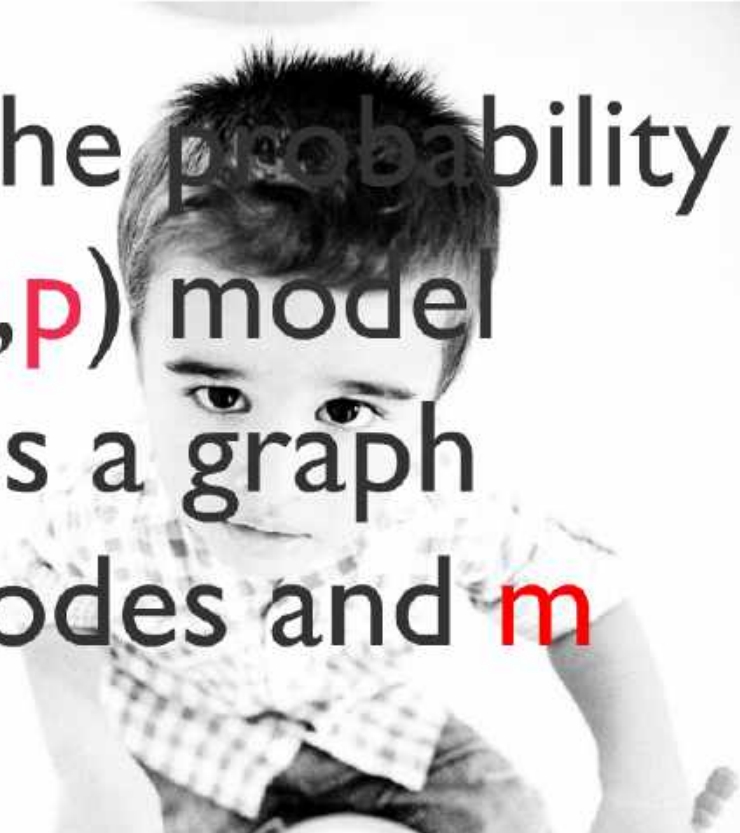
$n=50, p=0.08$



what is the probability  
that  $G(n, p)$  model  
generates a graph  
with  $n$  nodes and  $m$   
edges?

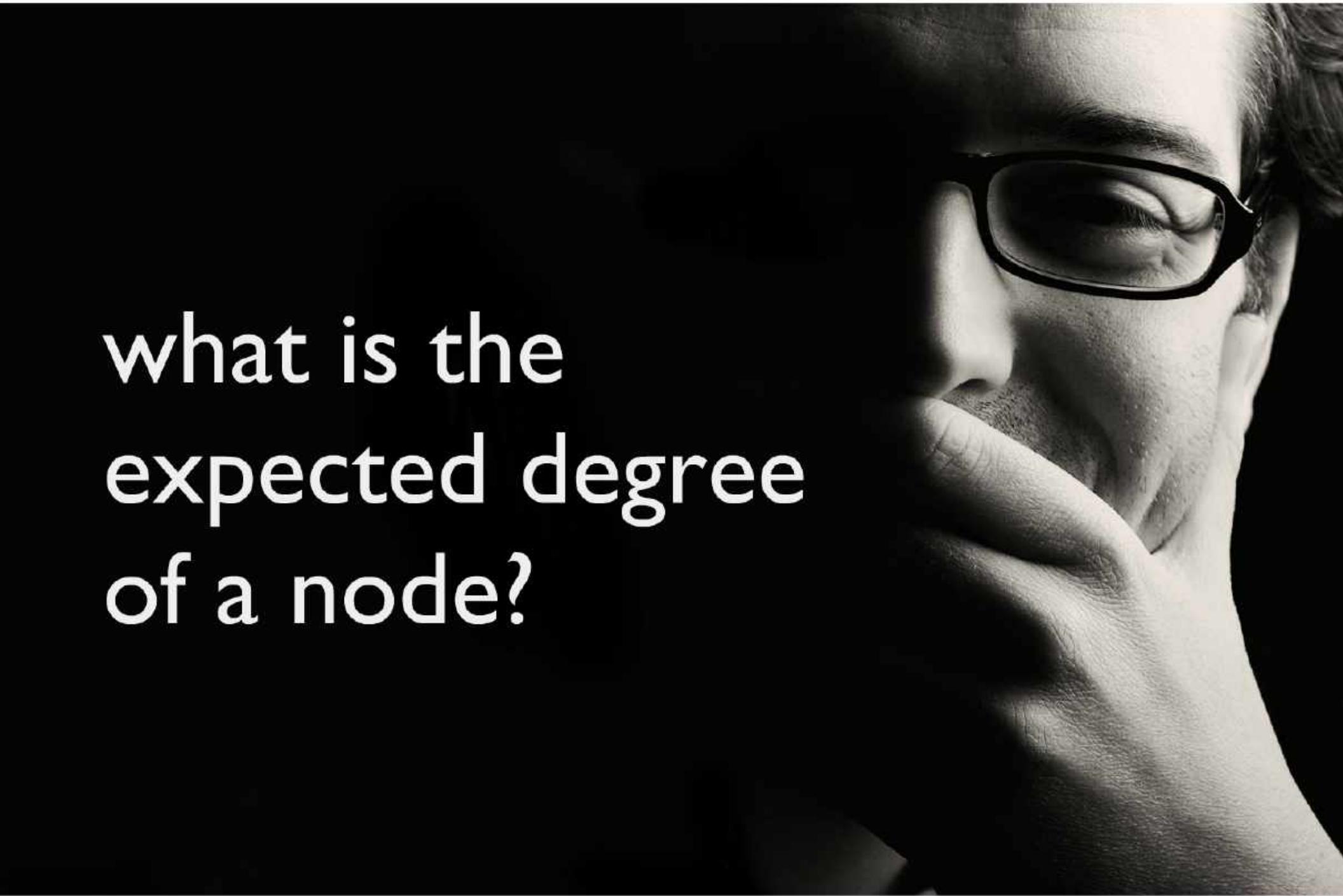


what is the probability  
that  $G(n, p)$  model  
generates a graph  
with  $n$  nodes and  $m$   
edges?



$$\binom{\binom{n}{2}}{m} p^m (1 - p)^{\frac{n(n-1)}{2} - m}$$

what is the  
expected degree  
of a node?



**The probability that a node has degree  $d$**

$$\binom{n-1}{d} p^d (1-p)^{n-1-d}$$

$P(d)$ : the degree distribution of a node

The probability that a node has degree **d**

$$\binom{n-1}{d} p^d (1-p)^{n-1-d}$$

$P(d)$ : the degree distribution of a node

Poisson random networks



$$\approx \frac{e^{-(n-1)p} ((n-1)p)^d}{d!}$$

$$\frac{e^{-\lambda} \lambda^d}{d!} \quad \lambda = np$$

what is the clustering coefficient?



what is the average  
path length?

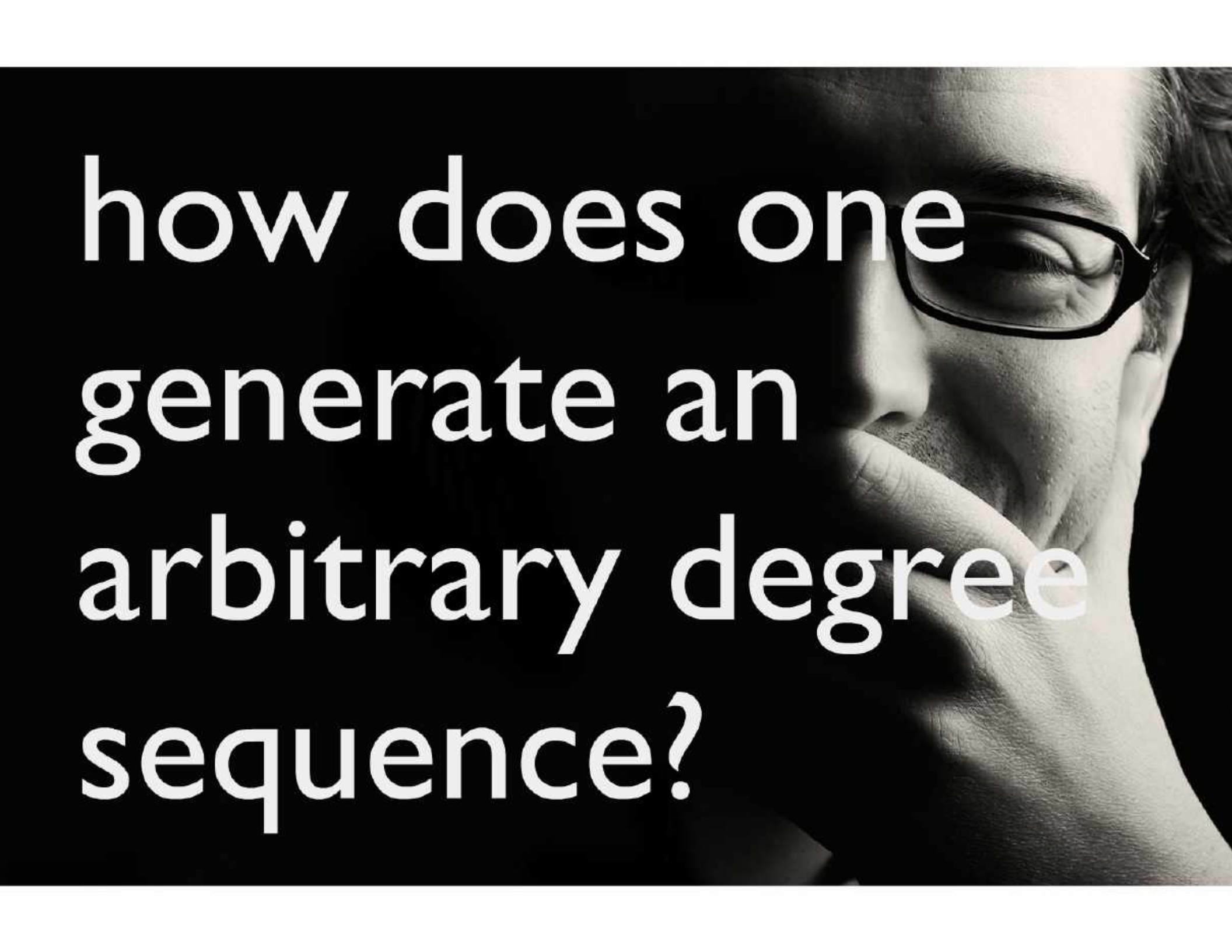


what is the average  
path length?



the diameter grows as  $O(\log n/\alpha)$





how does one  
generate an  
arbitrary degree  
sequence?

$d_1, d_2, d_3, \dots, d_n$

configuration model



$i=3$

assume that each node has the corresponding number of half edges

2

simply pair up edges

$d_1, d_2, d_3, \dots, d_n$

$\sum_i d_i$  is even

configuration model



$i=3$

**1** assume that each node has the corresponding number of half edges

**2** simply pair up edges

delete self loops and multiple edges

What are some assumptions here?

$$d_1, d_2, d_3, \dots, d_n$$

expected degree model

$$p_{ij} \propto \frac{d_i d_j}{\sum_k d_k}$$

where  $\left(\max_i d_i\right)^2 \leq \sum_k d_k$

# Two random graph generators



$$d_1, d_2, d_3, \dots, d_n \quad \sum_i d_i \text{ is even}$$

configuration model

$i=3$

1 assume that each node has the corresponding number of half edges

2 simply pair up edges

delete self loops and multiple edges

$d_1, d_2, d_3, \dots, d_n$   
expected degree model

$$p_{ij} \propto \frac{d_i d_j}{\sum_k d_k}$$

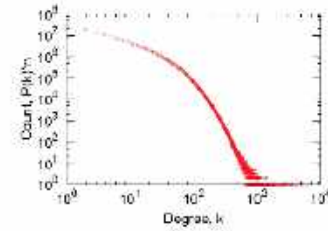
where  $\left(\max_i d_i\right)^2 \leq \sum_k d_k$

What are some assumptions here?

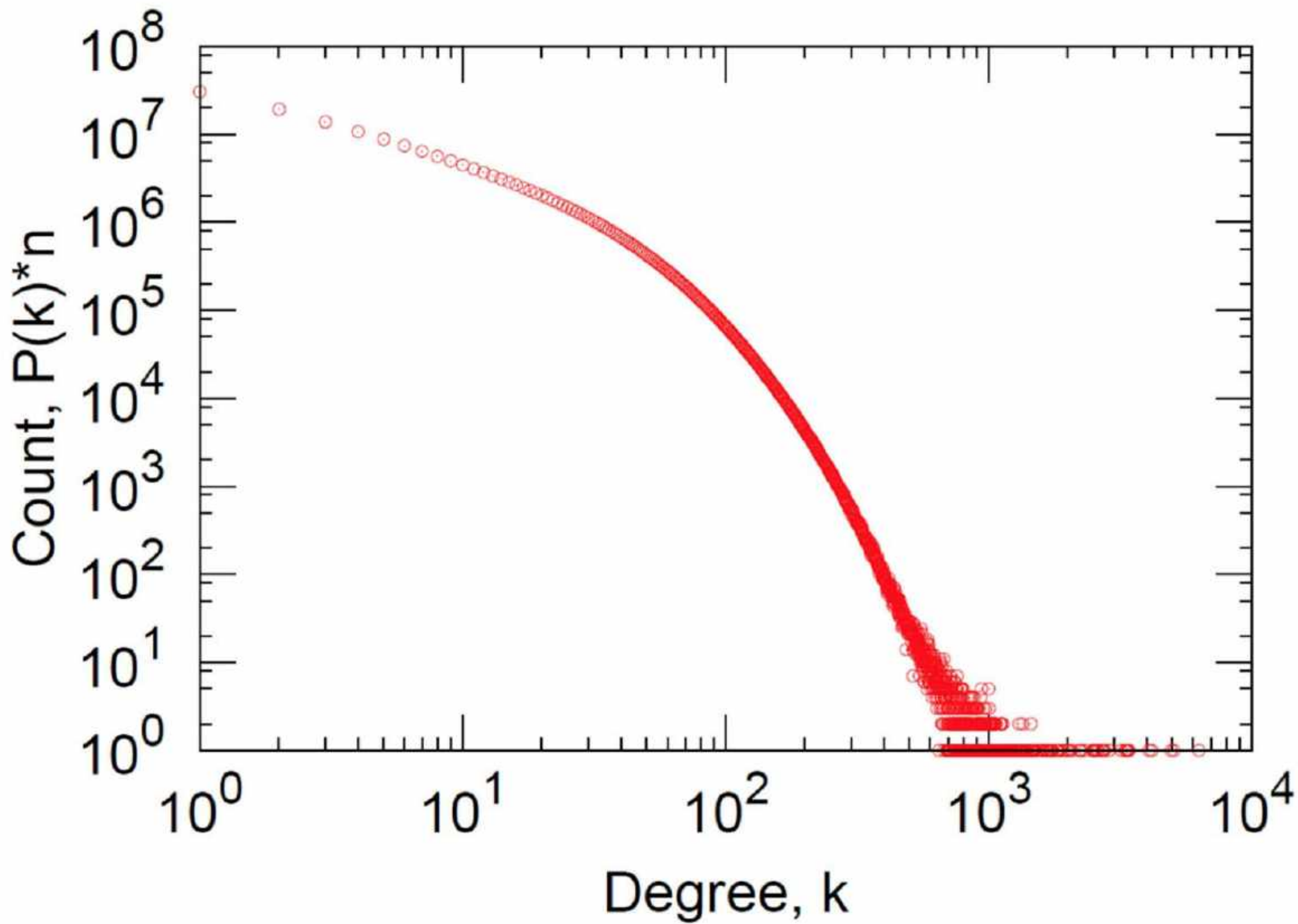
## What are the differences between the two?

how does the  
Poisson model  
stack up to  
reality?





Property	MSN	Poisson
Degree distribution	Power law*	Binomial
Path Length	6.6	$\log(n) \cong 8.8$
Clustering Coefficient	0.11	$p$





# Poisson models: a reality check

Giant component

Clustering coefficient is way off

Average path length is comparable

Degree distributions are unrealistic

**Real networks are *not* Poisson!**

why did we  
study this at  
all!?

reference model

analytically tractable



why did we  
study this at  
all!?

reference model  
analytically tractable



to what degree is a network random?

# Giant Components



The distribution is  $\text{max}(p, q)$



what is happening to the network when we increase  $p$ ?



the network is undergoing a phase transition, we see the emergence of a giant component, for some value of  $p$



A property  $A(N)$  is often specified by simply listing the networks that satisfy the property

$$A(N) = \{g : N_i(g) \neq \emptyset \forall i \in N\}$$

the first giant component



$$(1-p)^{n-1}$$

Particular example with a fixed number of nodes. As graph becomes sparser (smaller  $p$ )

$$P_i(N_i(g) = 1) = p(1-p)^{n-1} \rightarrow 0$$

$$P_i(N_i(g) = 0) = (1-p)^{n-1} \rightarrow 1$$

As a concrete example, consider the case that node  $i$  has at least one link.

$$A(N) = \{g : d_i \geq 1\}$$

probability that  $A(N)$  is true

$$1 - (1-p(n))^{n-1}$$



$$P_i(N_i(g) = 2) = p^2(1-p)^{n-2}$$

then

$$P_i(N_i(g) = 2) = 2p(1-p)^{n-2}$$



$$P_i(N_i(g) = 1) = p(1-p)^{n-1}$$

$$\frac{1}{n^2}$$

the first links emerge



A large component is said to exist when the number of nodes in the component exceeds

$$> n^{2/3}$$

$$\frac{\log(n)}{n}$$

the network becomes connected

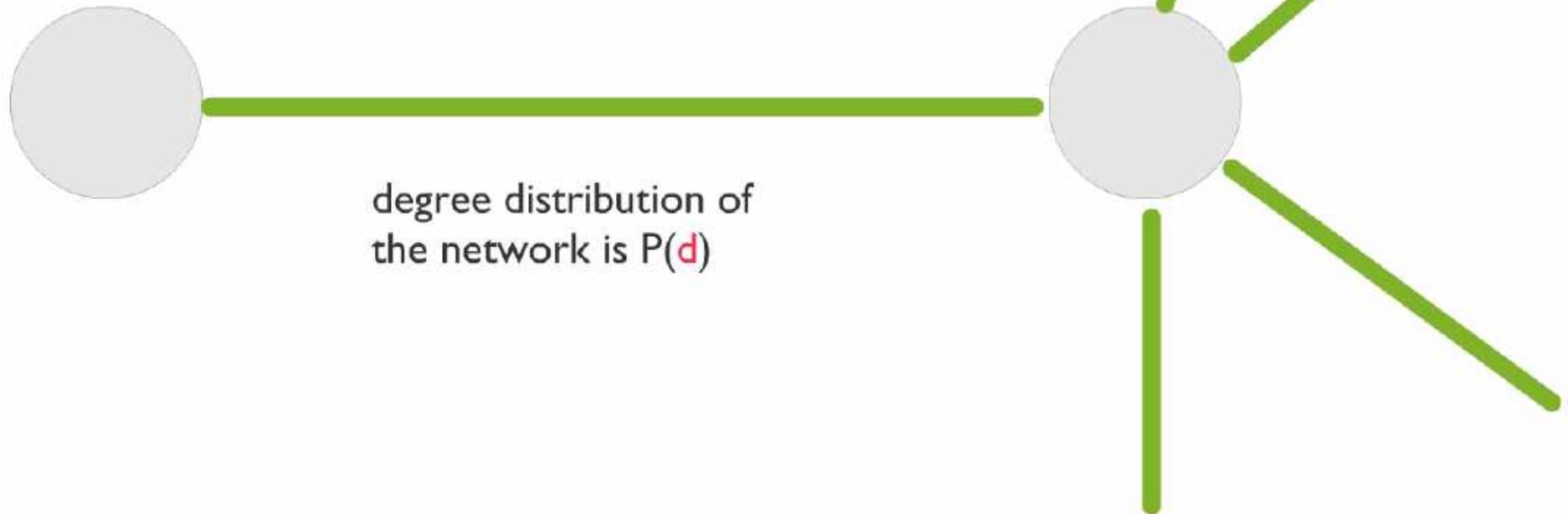


$$\frac{1}{n}$$

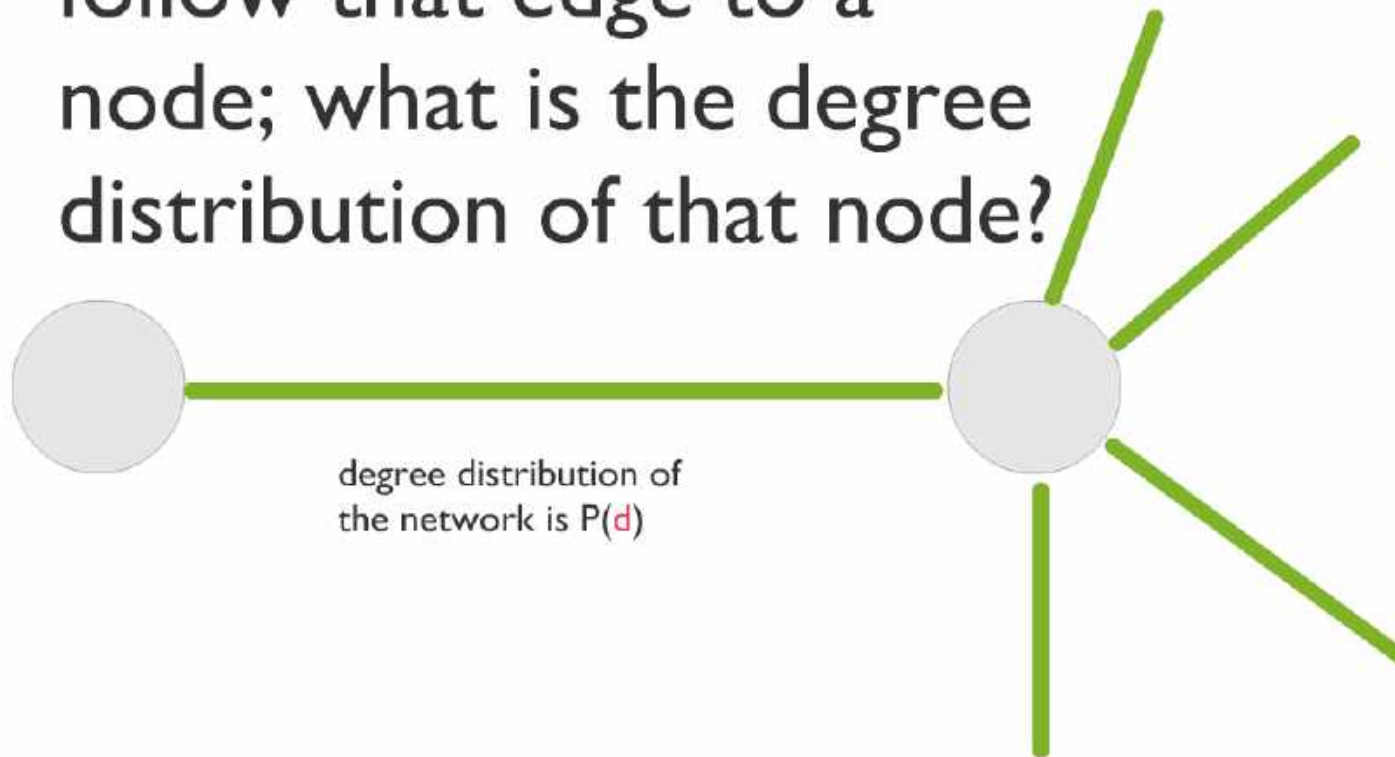
the first cycle emerges giant component begins to form



Let's assume that you pick a random edge and follow that edge to a node; what is the degree distribution of that node?

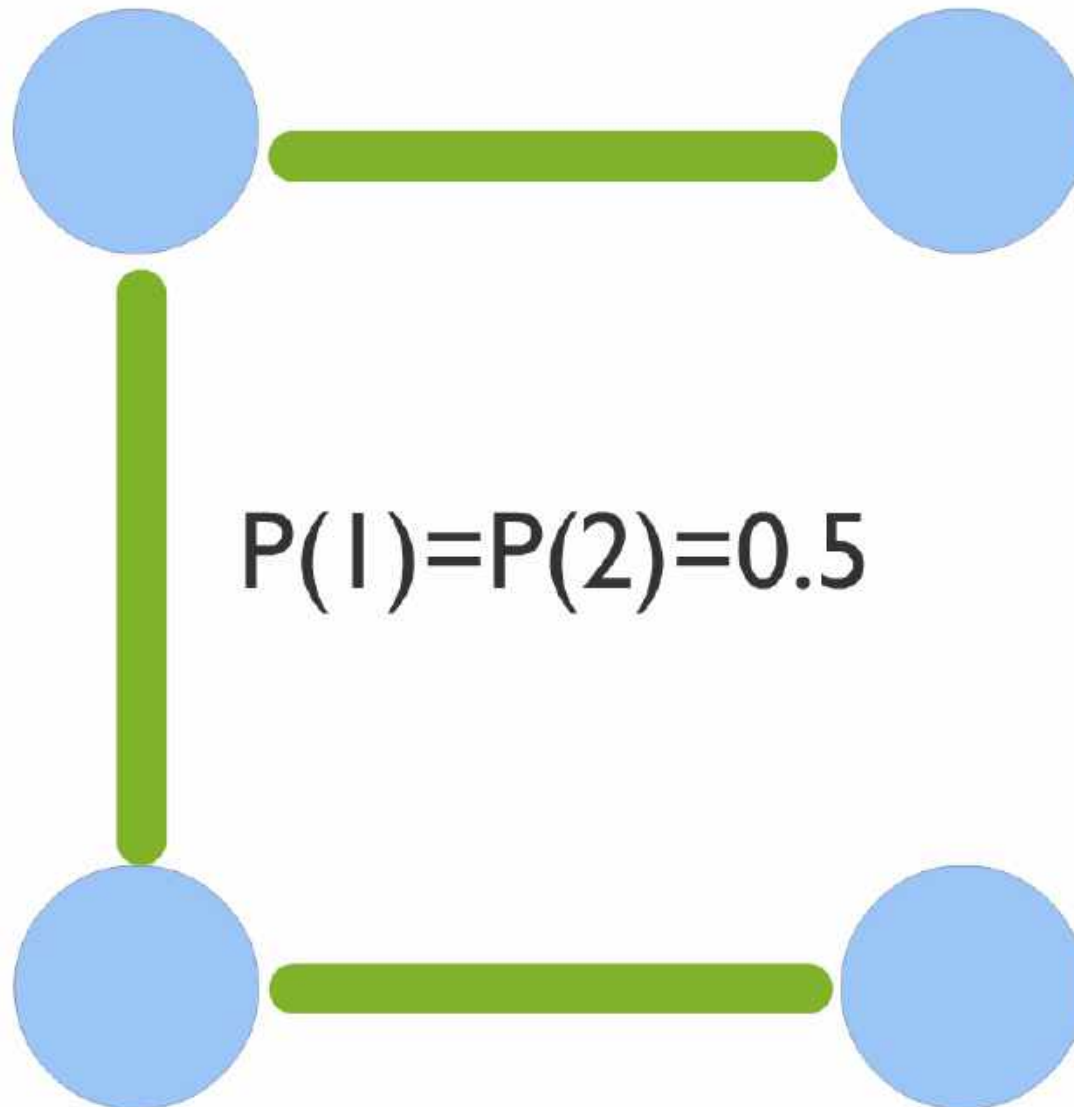


Let's assume that you pick a random edge and follow that edge to a node; what is the degree distribution of that node?



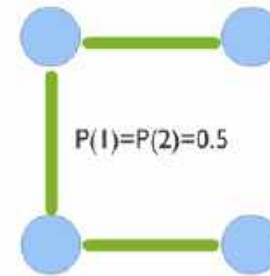
The distribution is **not**  $P(d)$

$P(d)$ : the degree distribution of a node



what is the probability  
of finding a node with  
degree 2?


$P(c)$ : the degree distribution of a node





The distribution is

$$\tilde{P}(d) = \frac{P(d) \cdot d}{\langle d \rangle}$$

  $\langle d \rangle = \sum_d P(d) \cdot d$

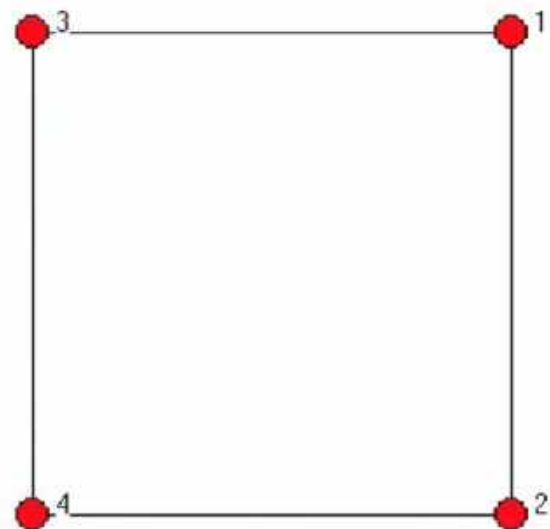
Doesn't hold when the neighboring degrees are highly correlated



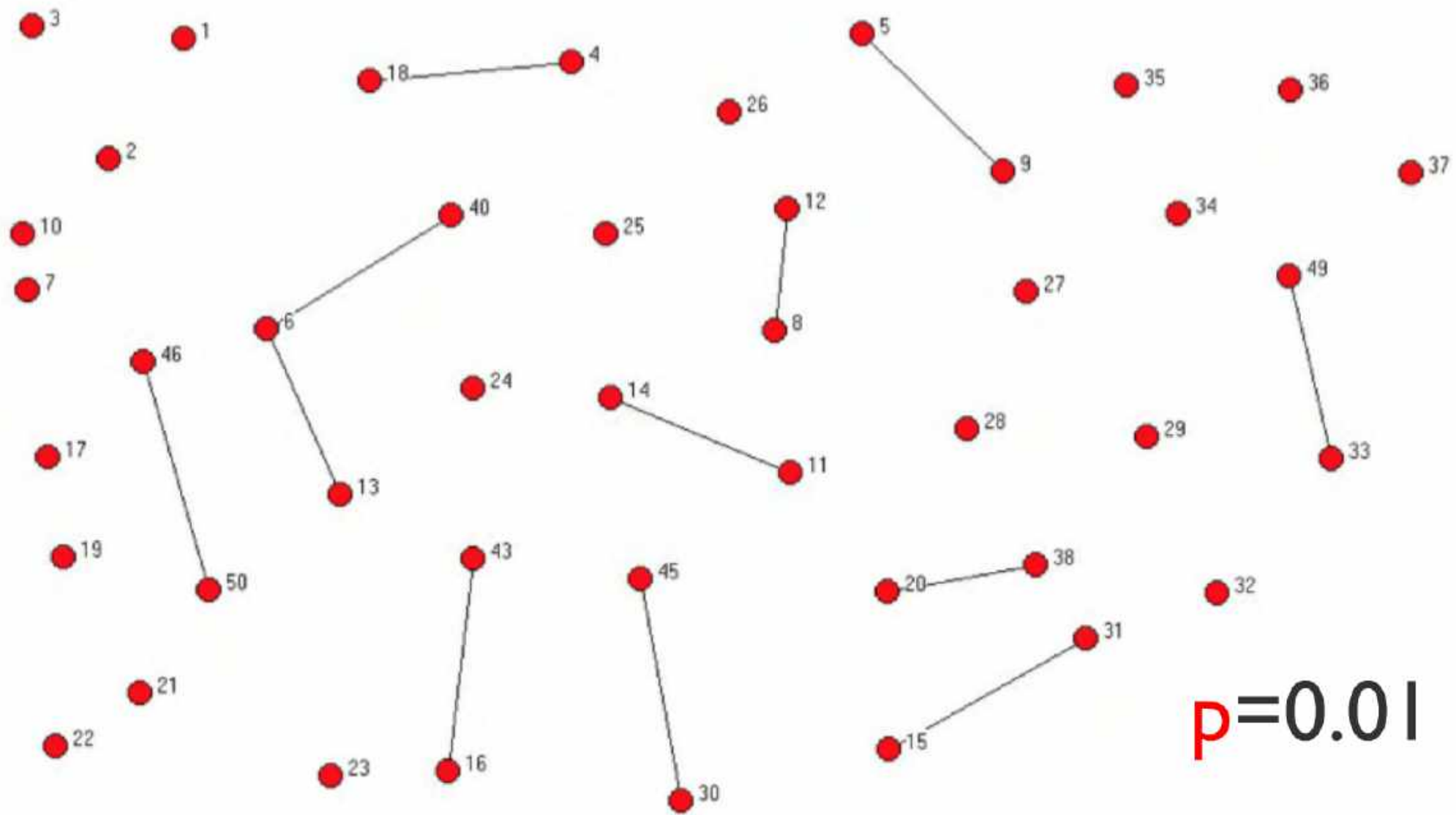
The distribution is

$$\tilde{P}(d) = \frac{P(d) \cdot d}{\langle d \rangle}$$

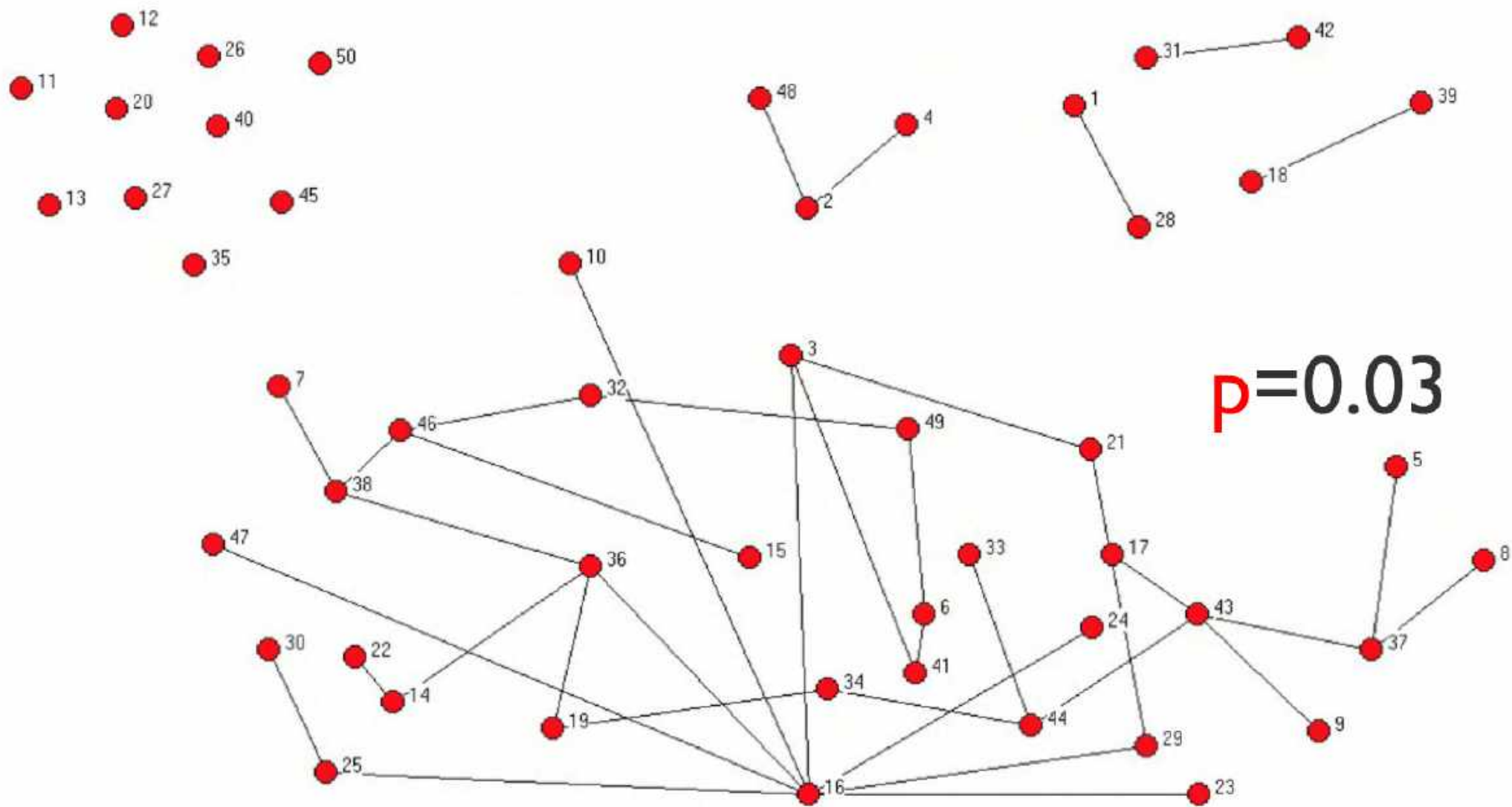
$\langle d \rangle = \sum_i P(d_i) \cdot d_i$

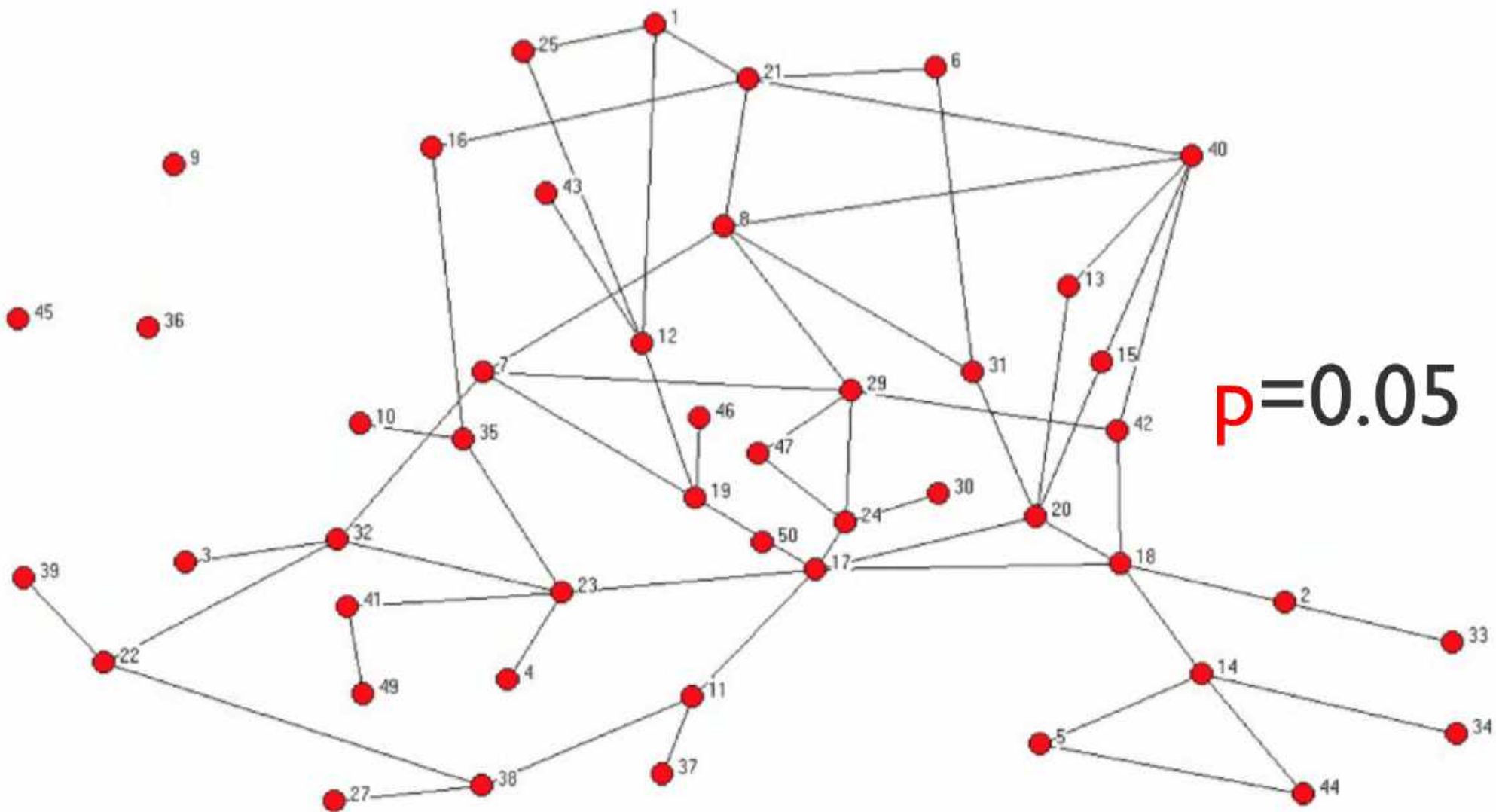


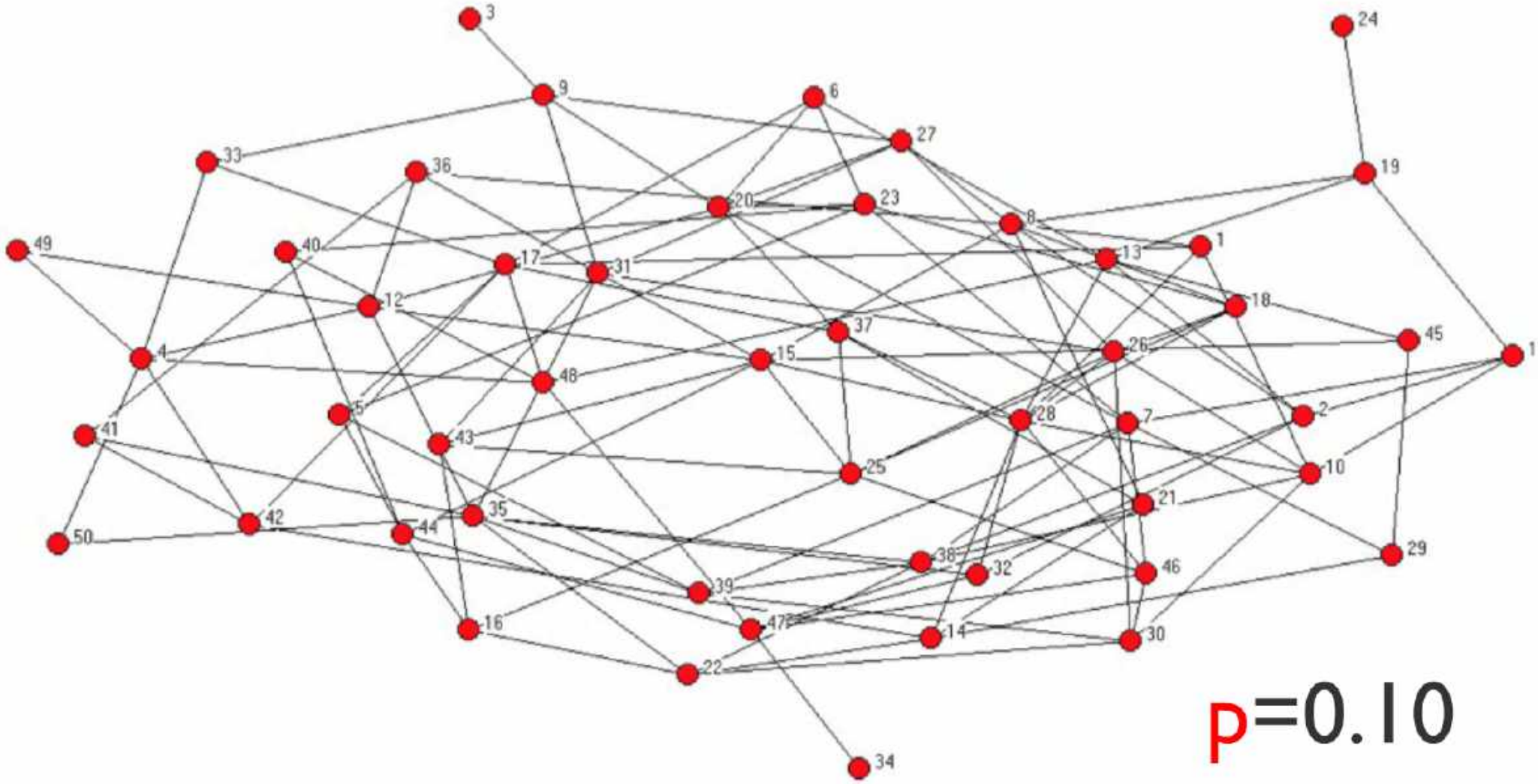
Networks that form due to preferential attachment will also **not** exhibit this behavior



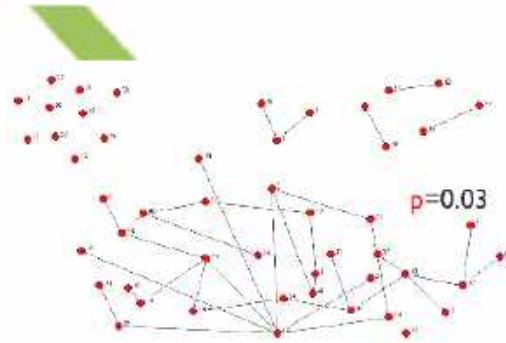
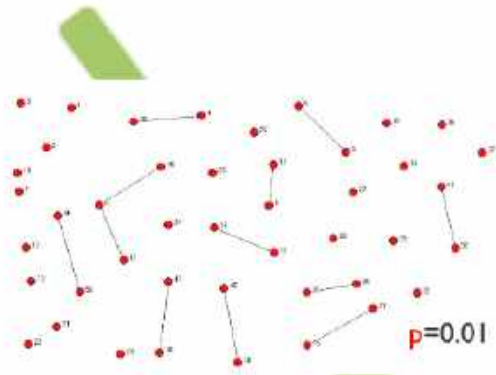
$p=0.01$







$p=0.10$



what is happening to the network when we increase  $p$ ?

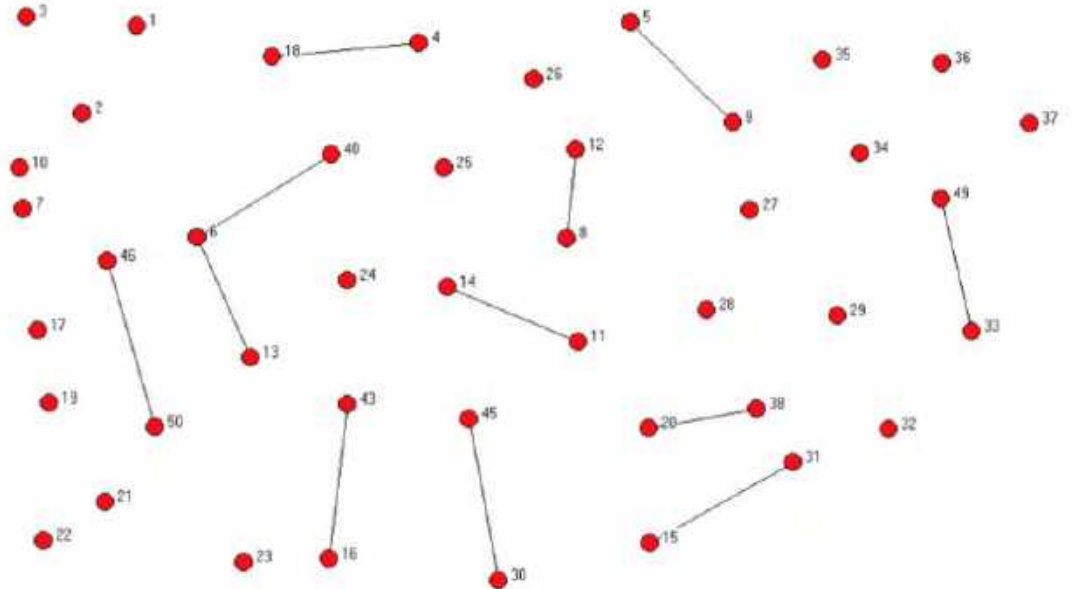
the network is undergoing a phase transition; we see the emergence of a giant component for some value of  $p$

the first links  
emerge

1

---

$n^2$



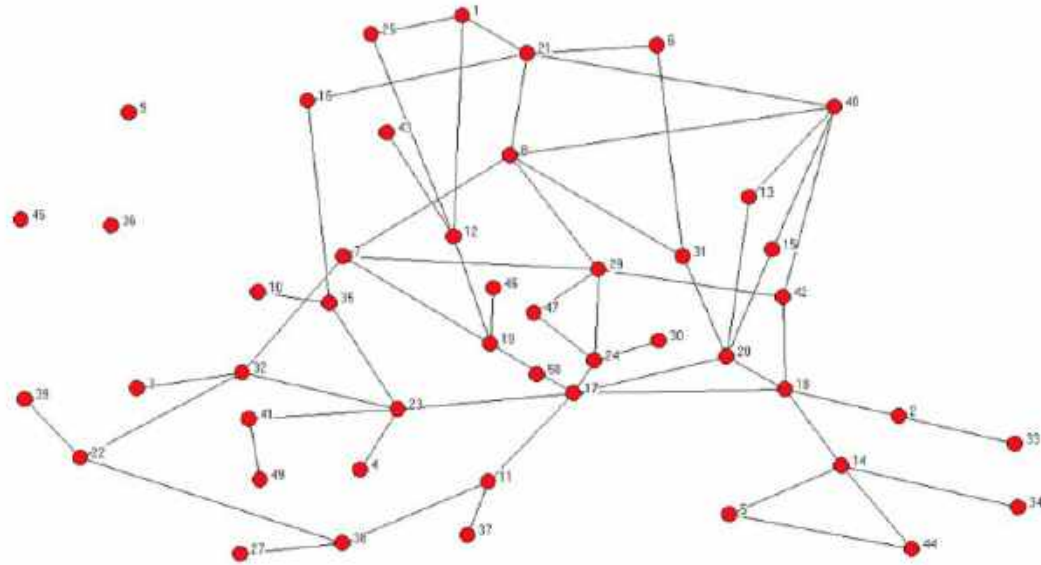


# 1



*n*

the first cycles emerge;  
giant component begins  
to form



# $\log(n)$

---

the network becomes connected

$n$

