

# Basics of CNA: Complex Network Analysis



---

# Overview of CNA

---

- ❖ different roles: hubs, weak ties, bridges, betweenness
- ❖ network heterogeneity
- ❖ robustness and immunization
- ❖ weighted and directed networks
- ❖ communities
- ❖ homophily
- ❖ the emergence of social clusters and segregation
- ❖ information and misinformation

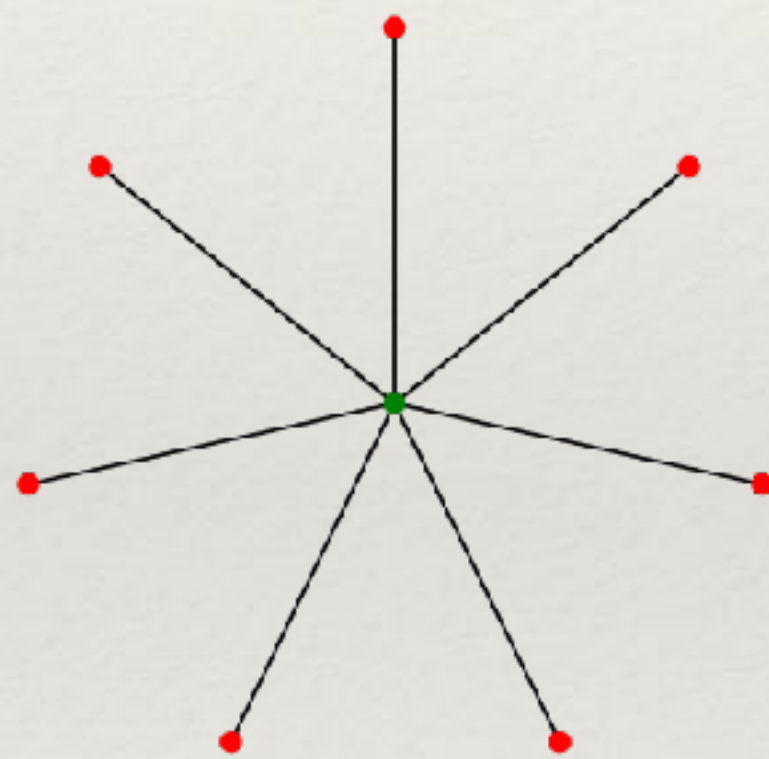


---

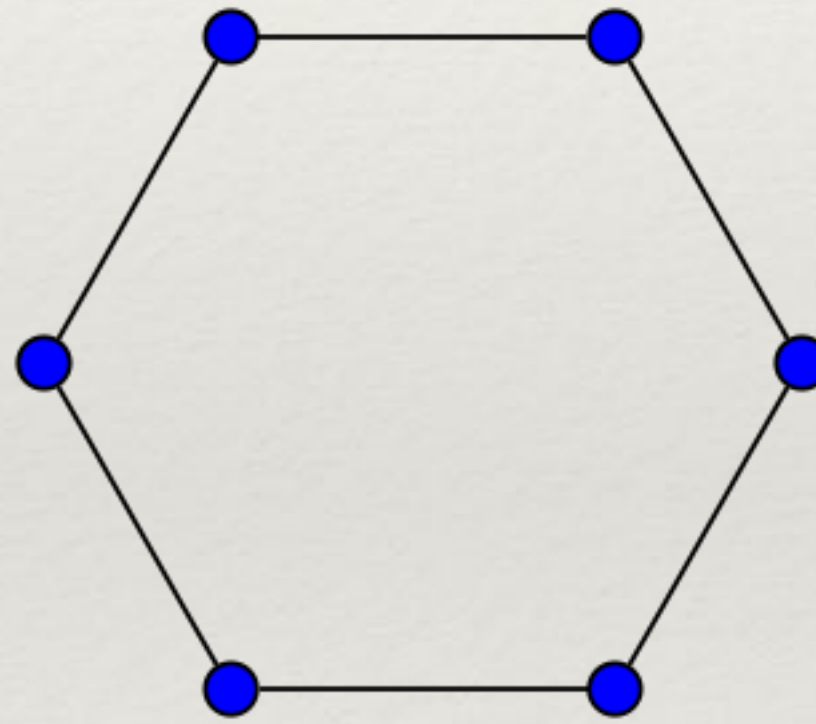
# Networks structural aspects

---

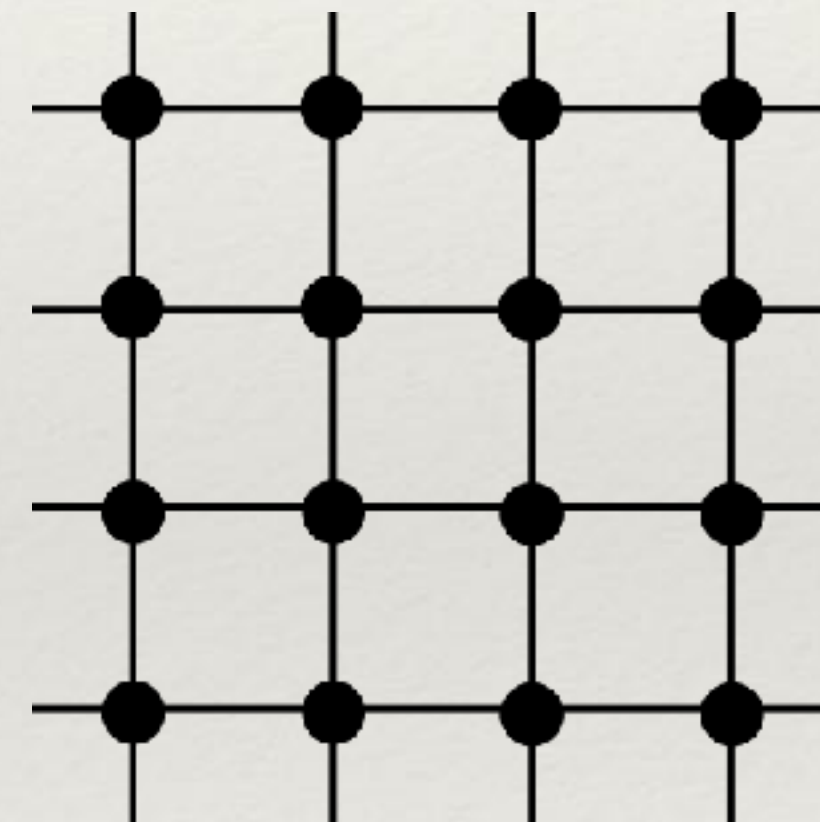
- ❖ "trivial" representation of a complex system
- ❖ **Simple** networks: few characteristics describe the network



Star



Ring



Grid

- ❖ We need a **language** and a **framework** to describe complex networks



---

# Basic Definitions

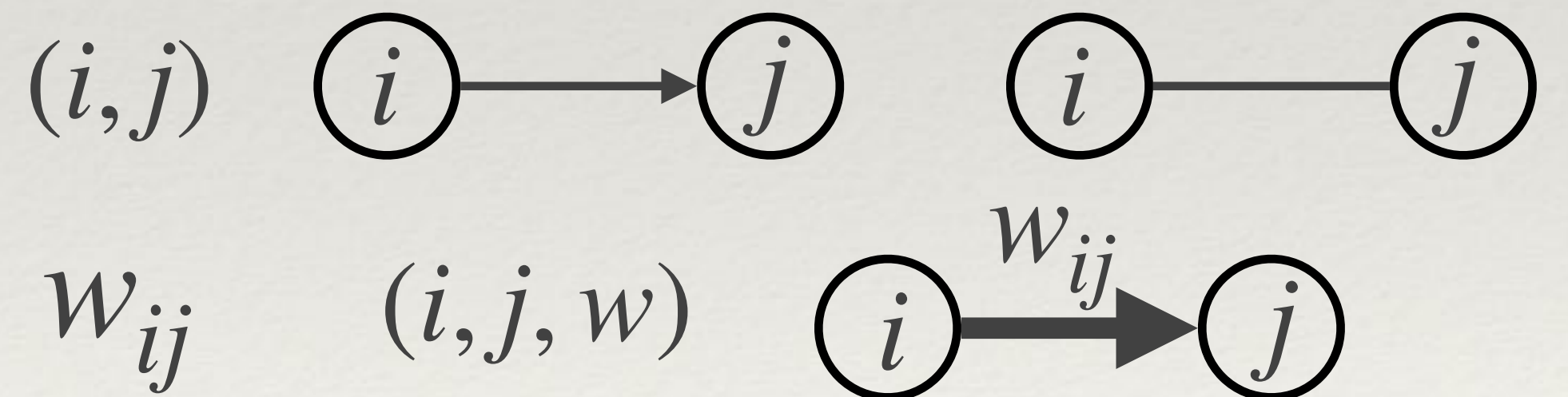
---

- ❖ A **graph** (or a network) is made of nodes and links
- ❖ **nodes** (or vertices)
- ❖ **links** (or edges, or arcs)
- ❖ Graphs can be **directed** or **undirected**
- ❖ Graphs can be **weighted** or **unweighted**

$$G = (N, L)$$

$$N = \{n_1, n_2, \dots, n_l\} = \{1, 2, \dots, l\}$$

$$L = \{(i, j) : i, j \in N\}$$

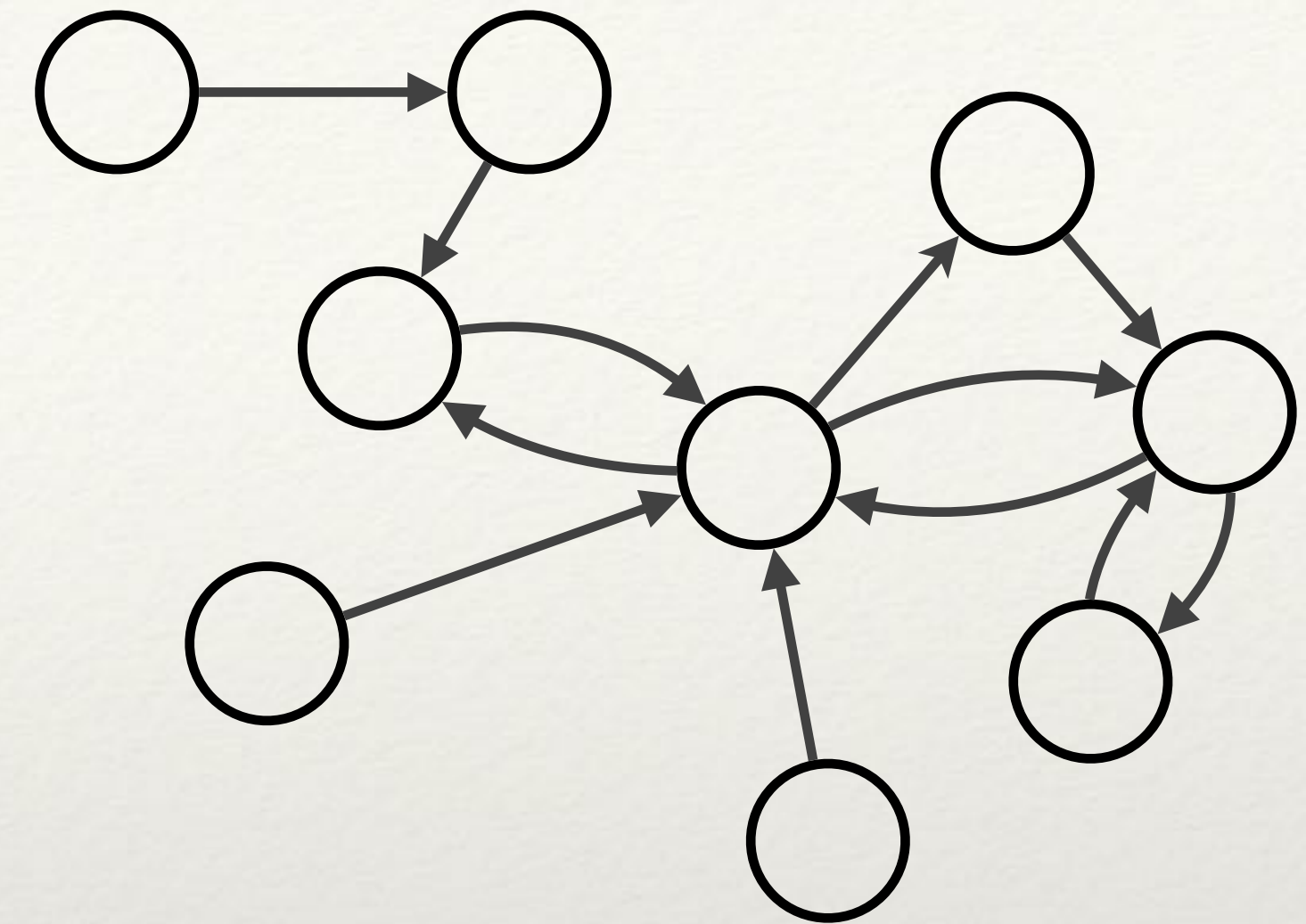
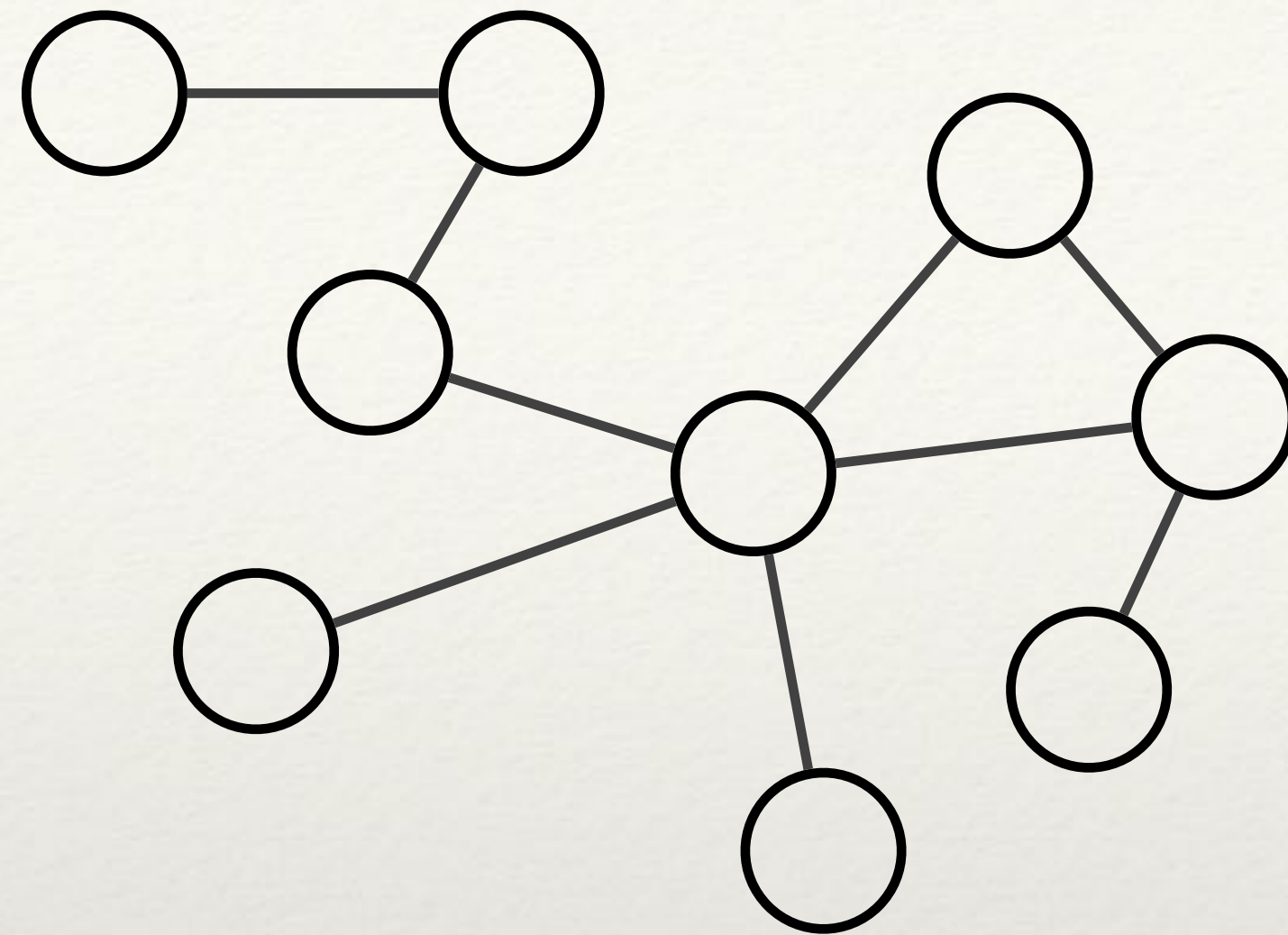




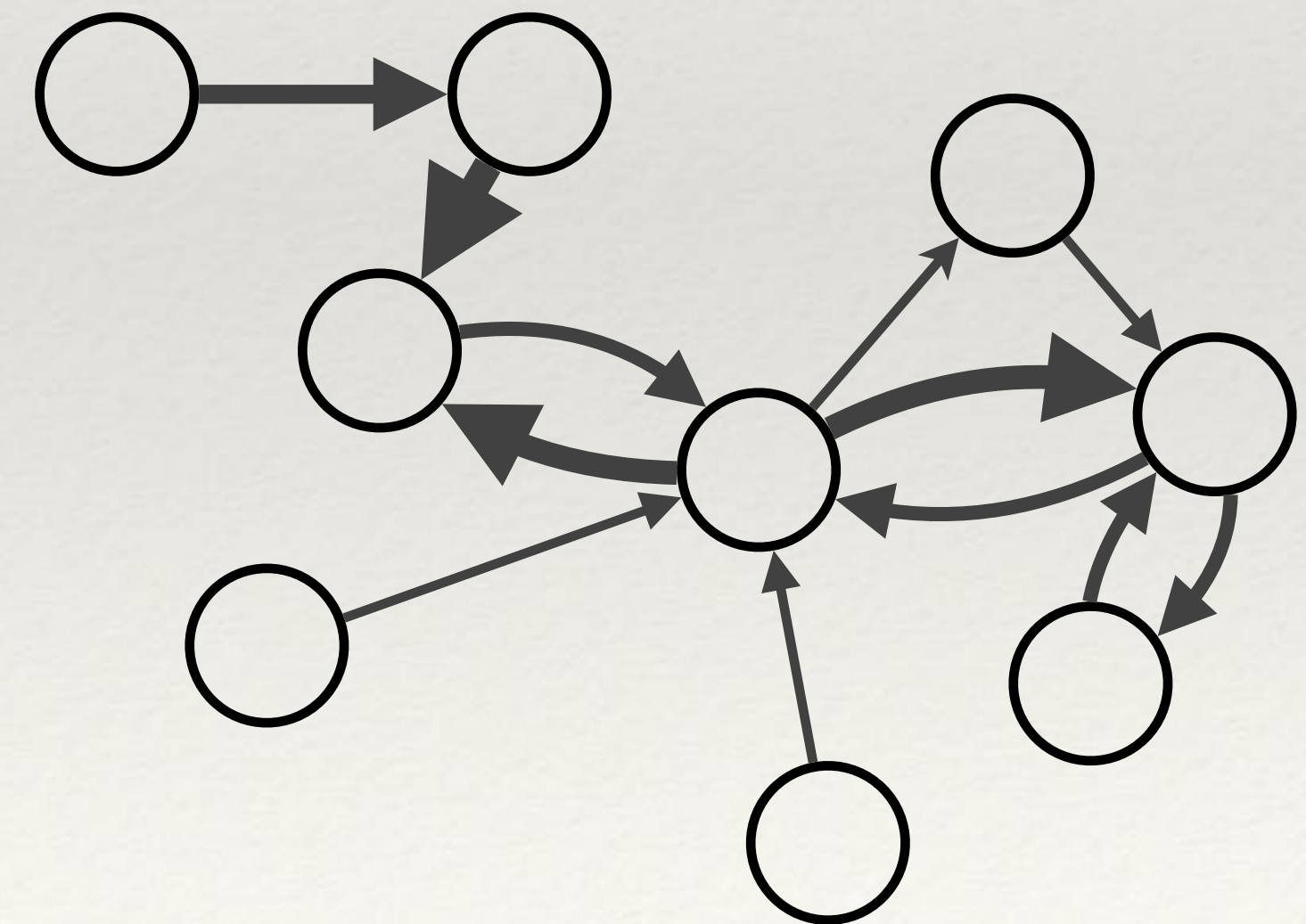
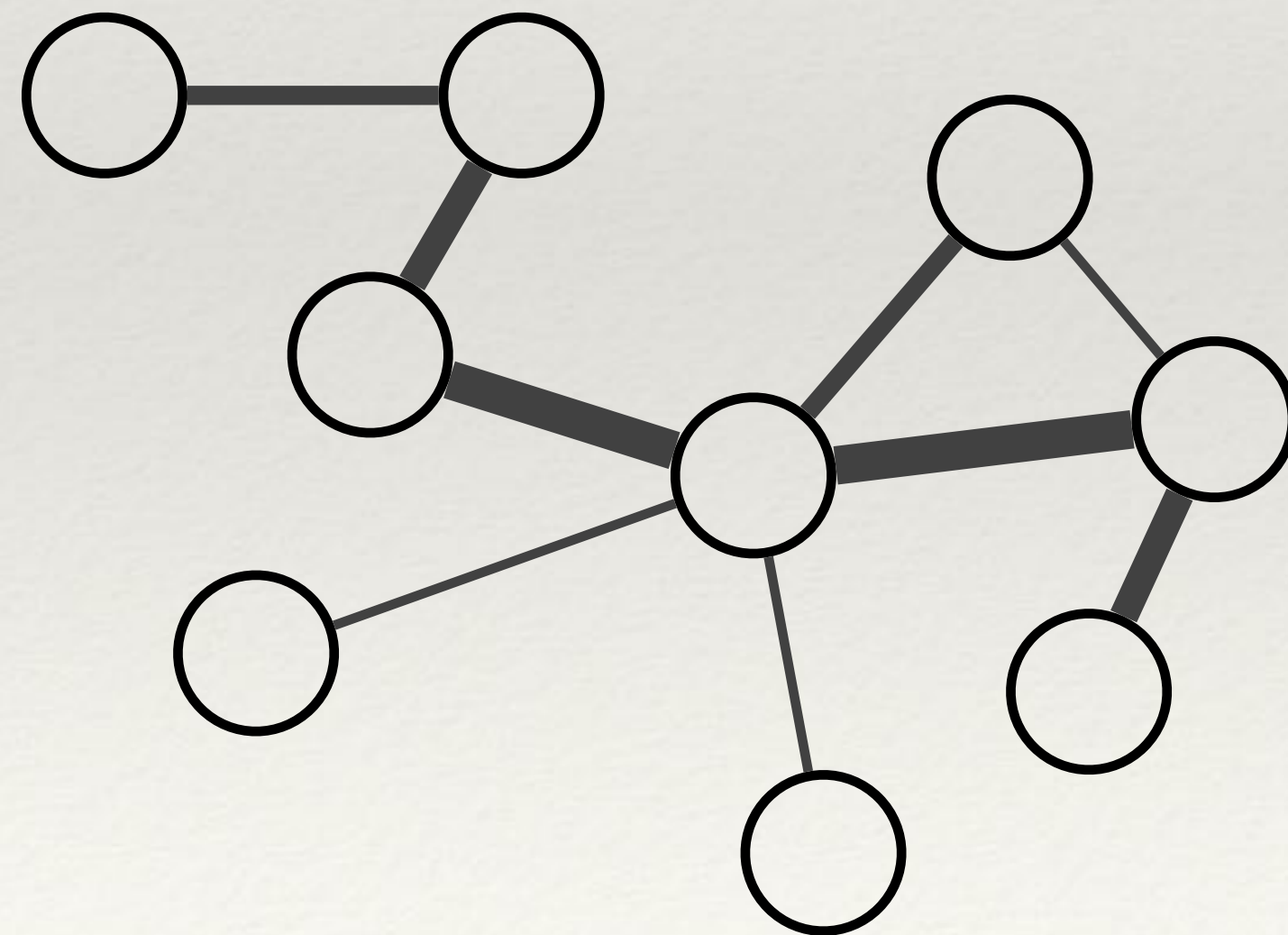
Undirected

Directed

Unweighted



Weighted





---

# Degree

---

Number of links (or neighbors)

$$i \rightarrow N_i \quad k_i = |N_i| \text{ degree}$$

Singleton: a node whose degree is zero

$$N_i = \{\}, k_i = 0$$

In directed networks

$$k_i^{in} = |P_i| \text{ in-degree}$$

$$k_i^{out} = |S_i| \text{ out-degree}$$

$$k_i = k_i^{in} + k_i^{out}$$



---

# Strength

---

Strength: Weighted degree

$$s_i = \sum_{j \in N_i} w_{ij}$$

in-strength

$$s_i^{in} = \sum_{j \in P_i} w_{ji}$$

out-strength

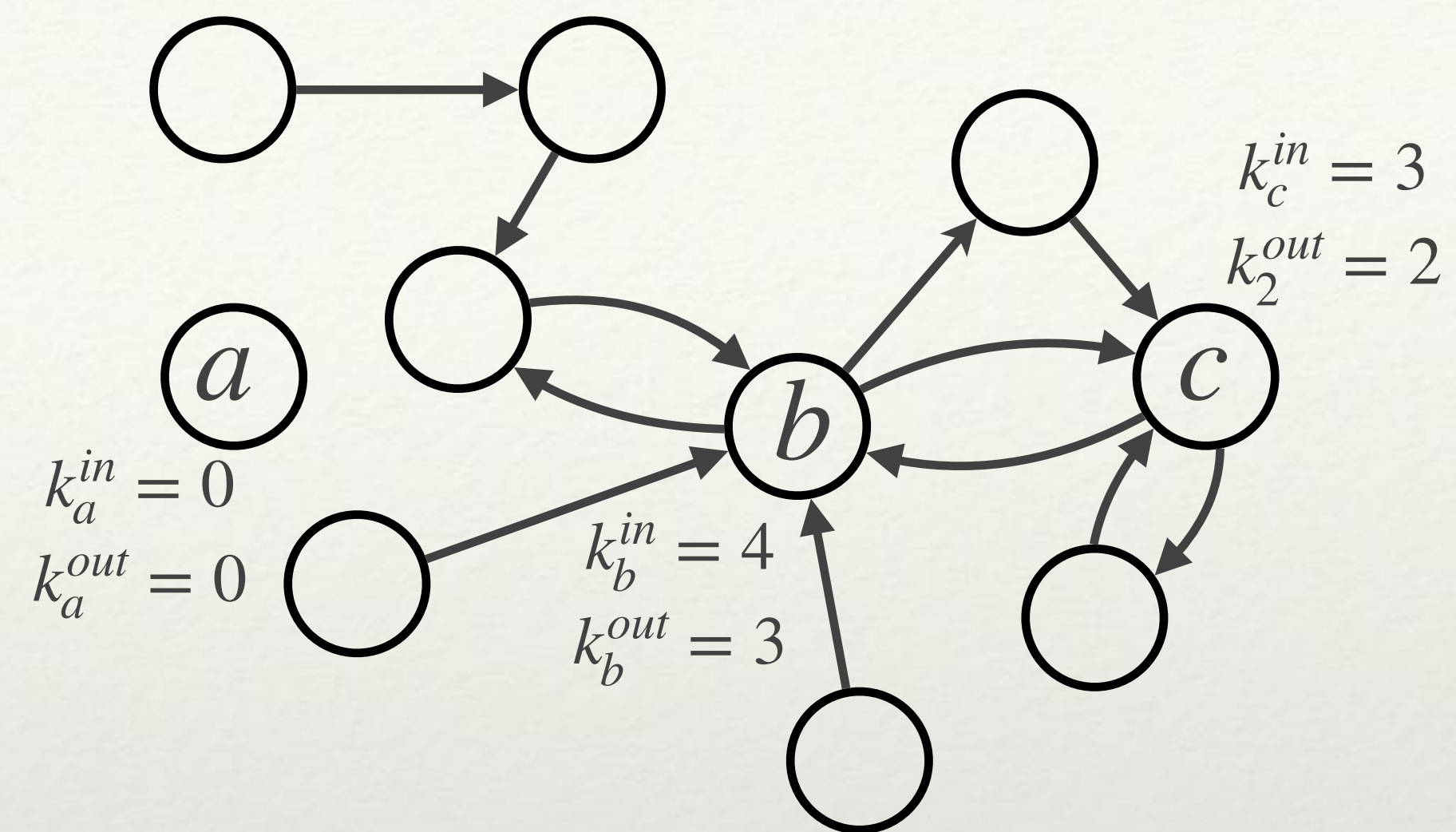
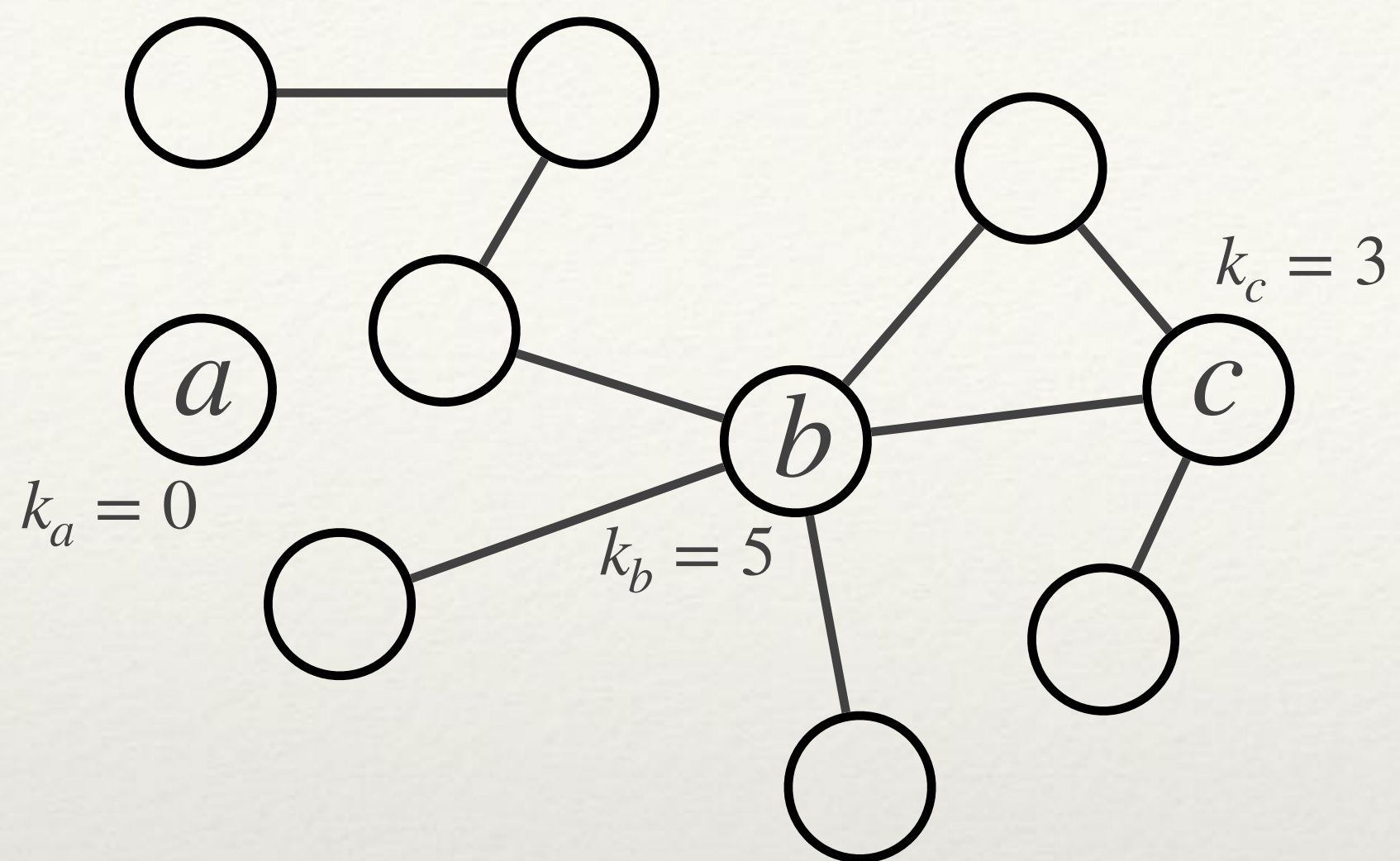
$$s_i^{out} = \sum_{j \in S_i} w_{ij}$$



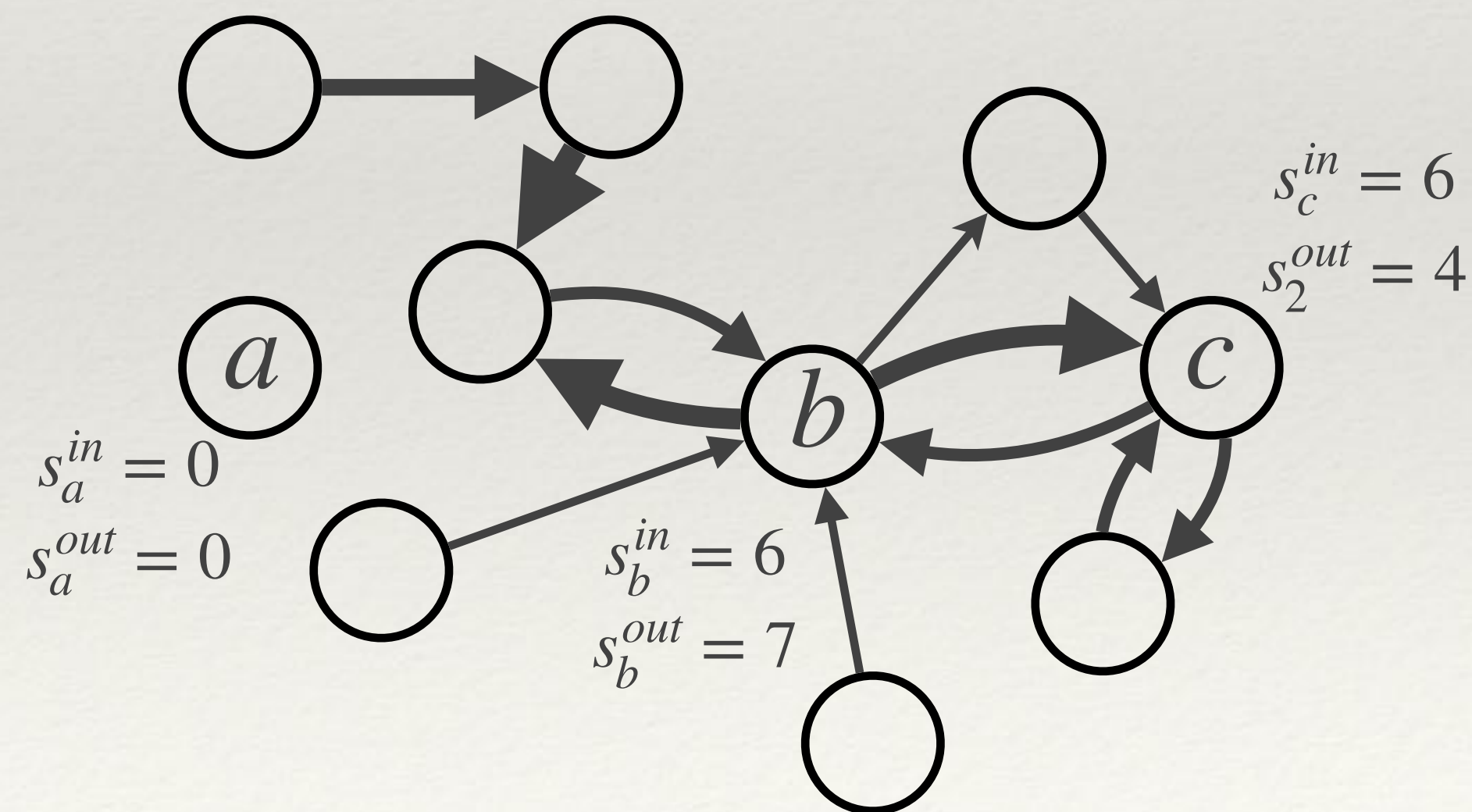
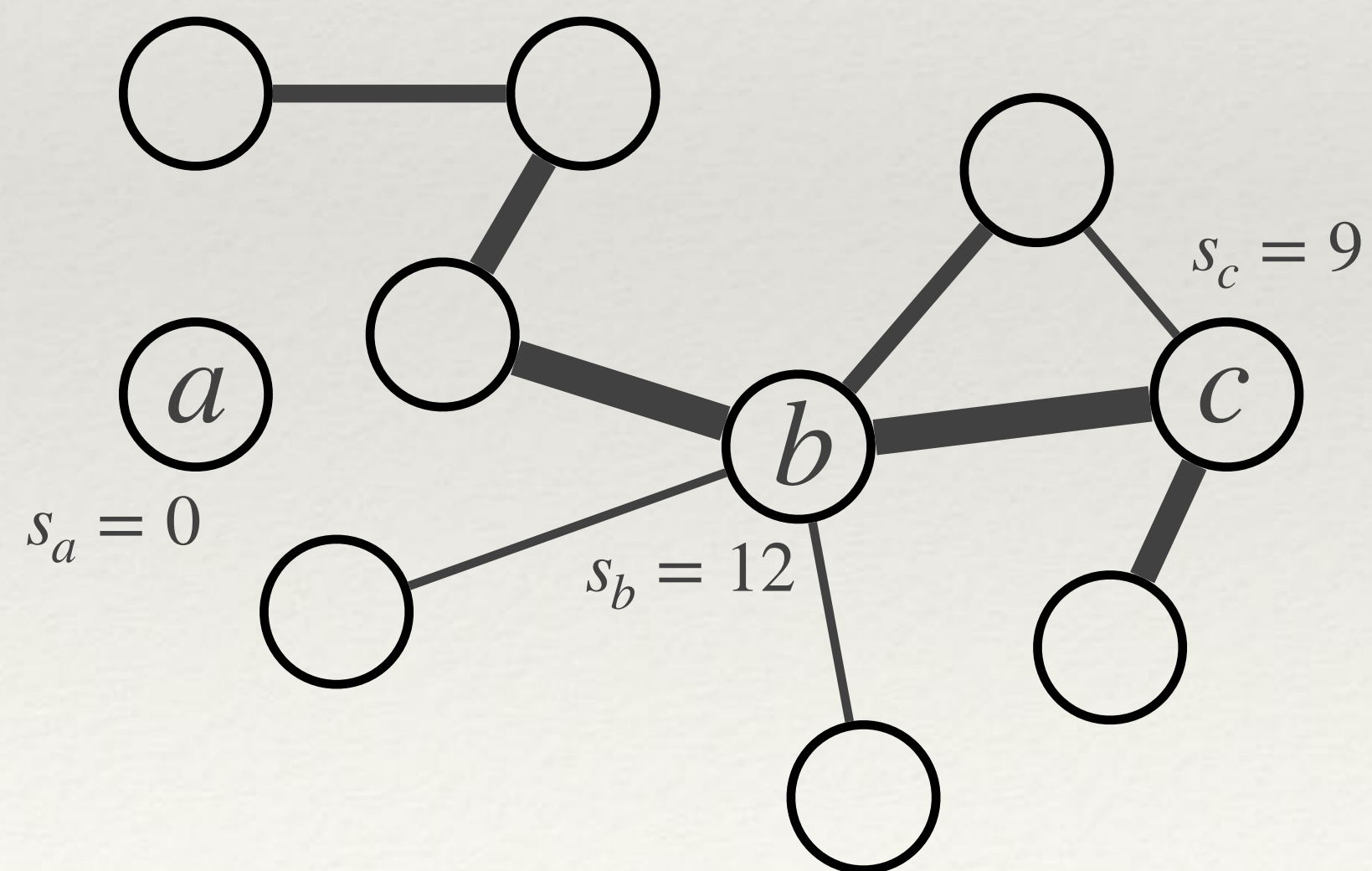
### Undirected

### Directed

Unweighted



Weighted





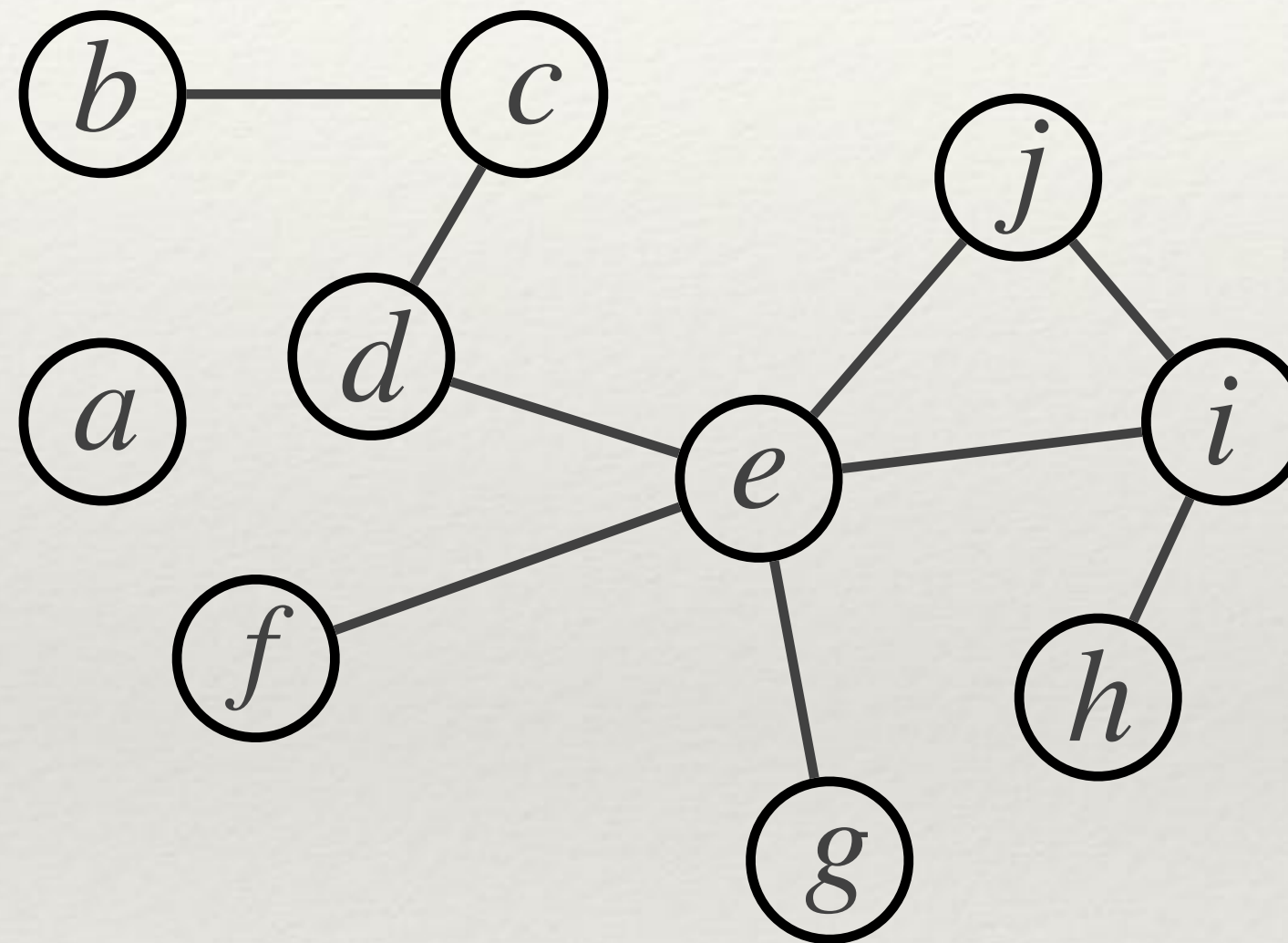
# Network representations

## Adjacency Matrix

$N \times N$  matrix

$$a_{ij} = \begin{cases} 0 & \text{no edge} \\ 1 & (i, j) \in L \end{cases}$$

Undirected network:  $a_{ij} = a_{ji}$



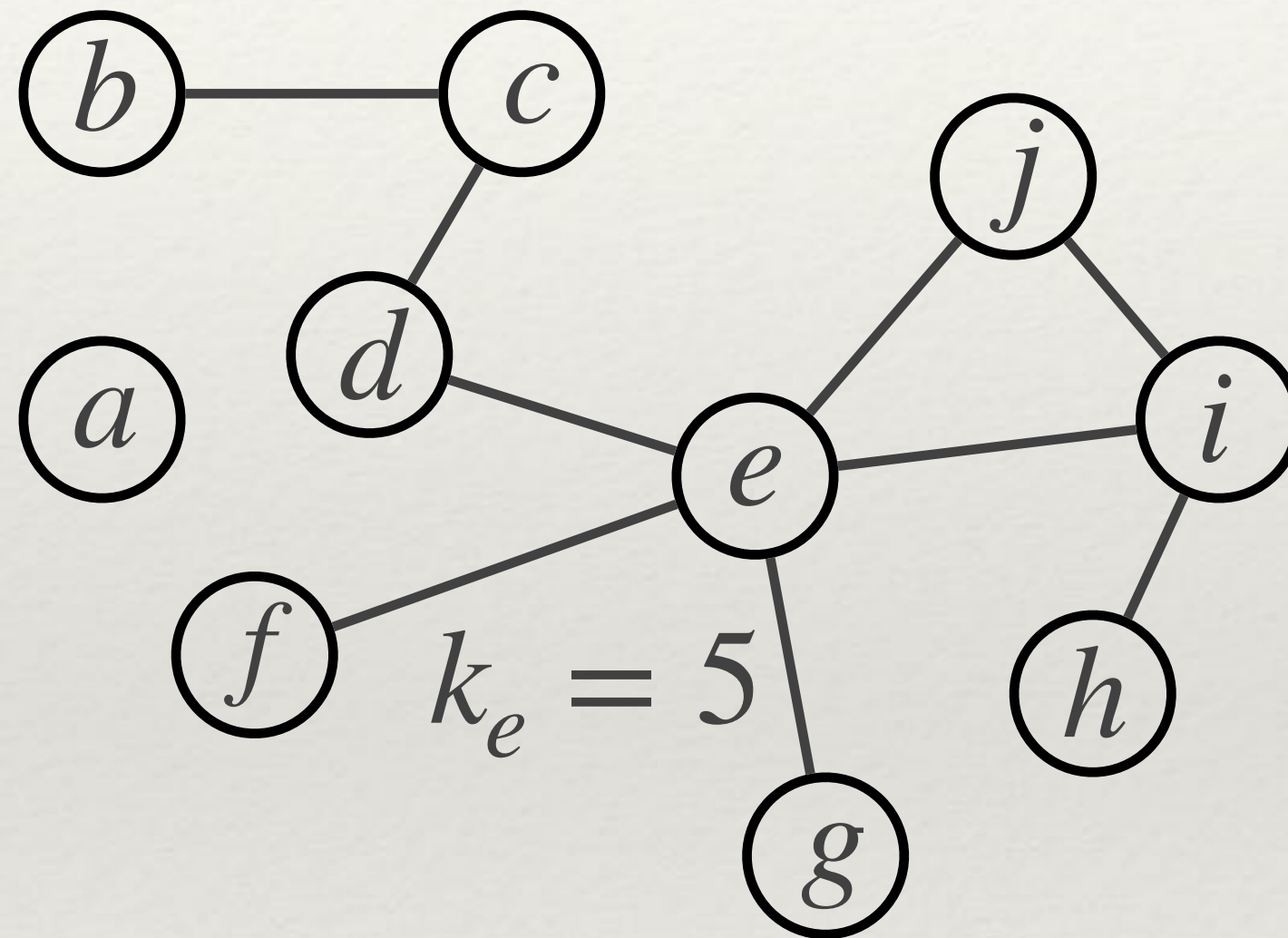
	a	b	c	d	e	f	g	h	i	j
a	0	0	0	0	0	0	0	0	0	0
b	0	0	1	0	0	0	0	0	0	0
c	0	1	0	1	0	0	0	0	0	0
d	0	0	1	0	1	0	0	0	0	0
e	0	0	0	1	0	1	1	0	1	1
f	0	0	0	0	1	0	0	0	0	0
g	0	0	0	0	1	0	0	0	0	0
h	0	0	0	0	0	0	0	0	1	0
i	0	0	0	0	1	0	0	1	0	1
j	0	0	0	0	1	0	0	0	1	0



# Network representations

Adjacency Matrix  
degree

$$k_i = \sum_j a_{ij} = \sum_j a_{ji}$$



	a	b	c	d	e	f	g	h	i	j
a	0	0	0	0	0	0	0	0	0	0
b	0	0	1	0	0	0	0	0	0	0
c	0	1	0	1	0	0	0	0	0	0
d	0	0	1	0	1	0	0	0	0	0
e	0	0	0	1	0	1	1	0	1	1
f	0	0	0	0	1	0	0	0	0	0
g	0	0	0	0	1	0	0	0	0	0
h	0	0	0	0	0	0	0	0	1	0
i	0	0	0	0	1	0	0	1	0	1
j	0	0	0	0	1	0	0	0	1	0











---

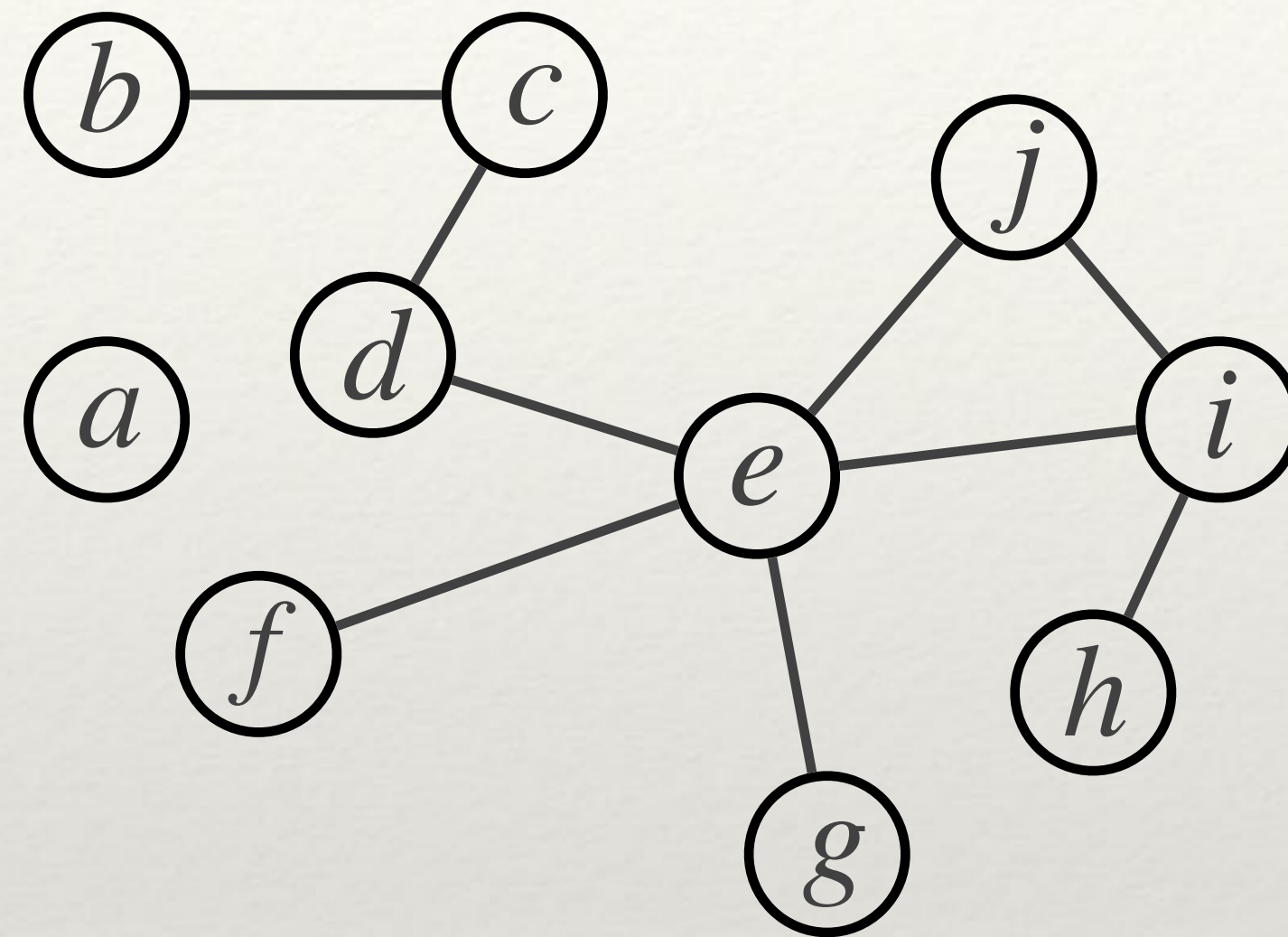
# Sparse network representations

---

- ❖ The memory / disk storage needed by an adjacency matrix is proportional to  $N^2$
- ❖ In sparse networks (most real-world networks), this is terribly inefficient: most of the space is wasted storing zeros (non-links); for very large networks, adjacency matrices are unfeasible
- ❖ It is much more efficient, often necessary, to store only the actual links, and assume that if a link is not listed it means it is not present
- ❖ There are two commonly used sparse networks representations:
  - ❖ **Adjacency list**
  - ❖ **Edge list**



# Adjacency list



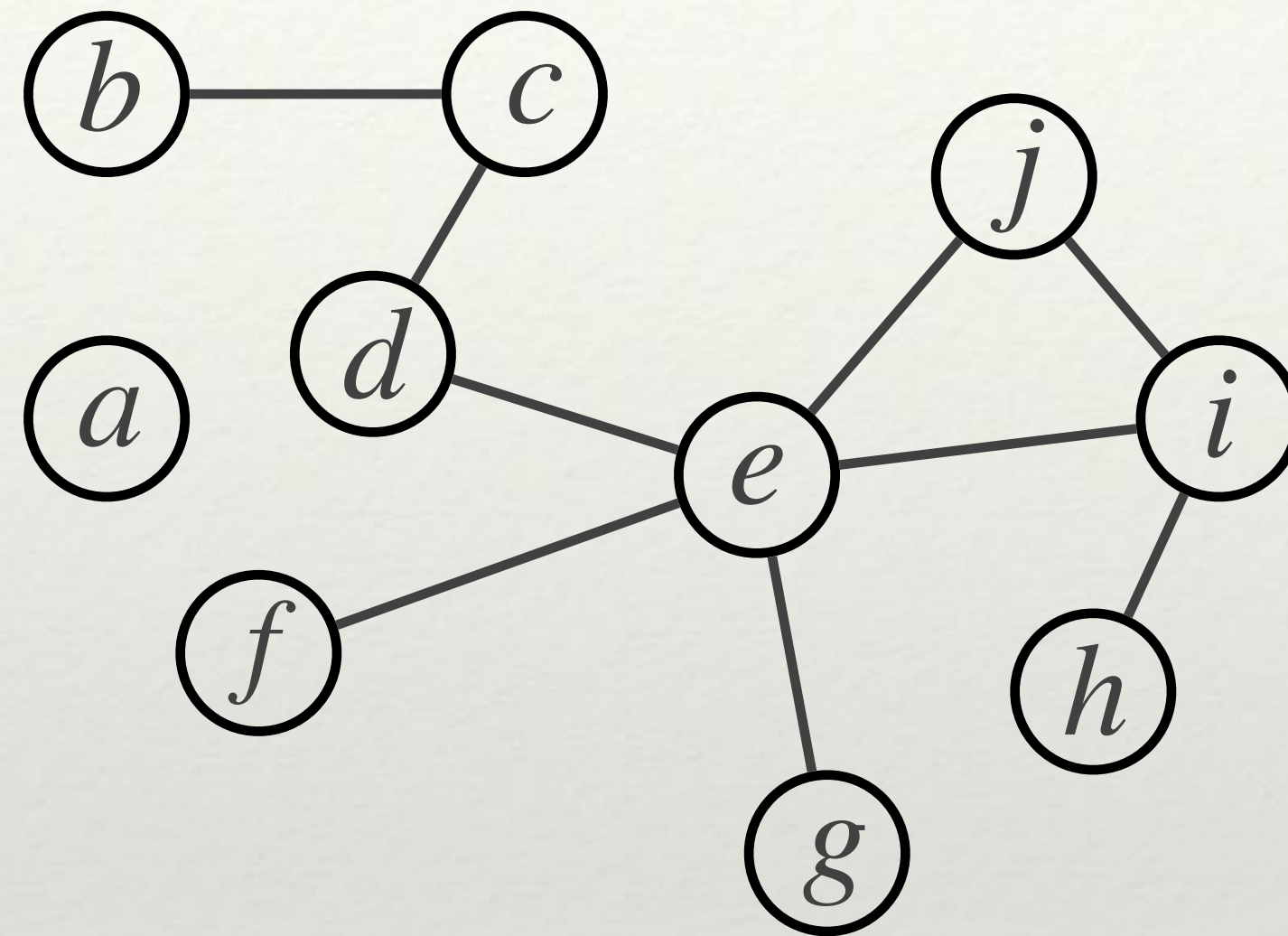
a									
b	c								
c	b	d							
d	c	e							
e	d	f	i	j	g				
f	e								
g	e								
h	i								
i	h	j							
j	e	i							

Undirected network: list each link twice

Directed network: list only existing links



# Edge list

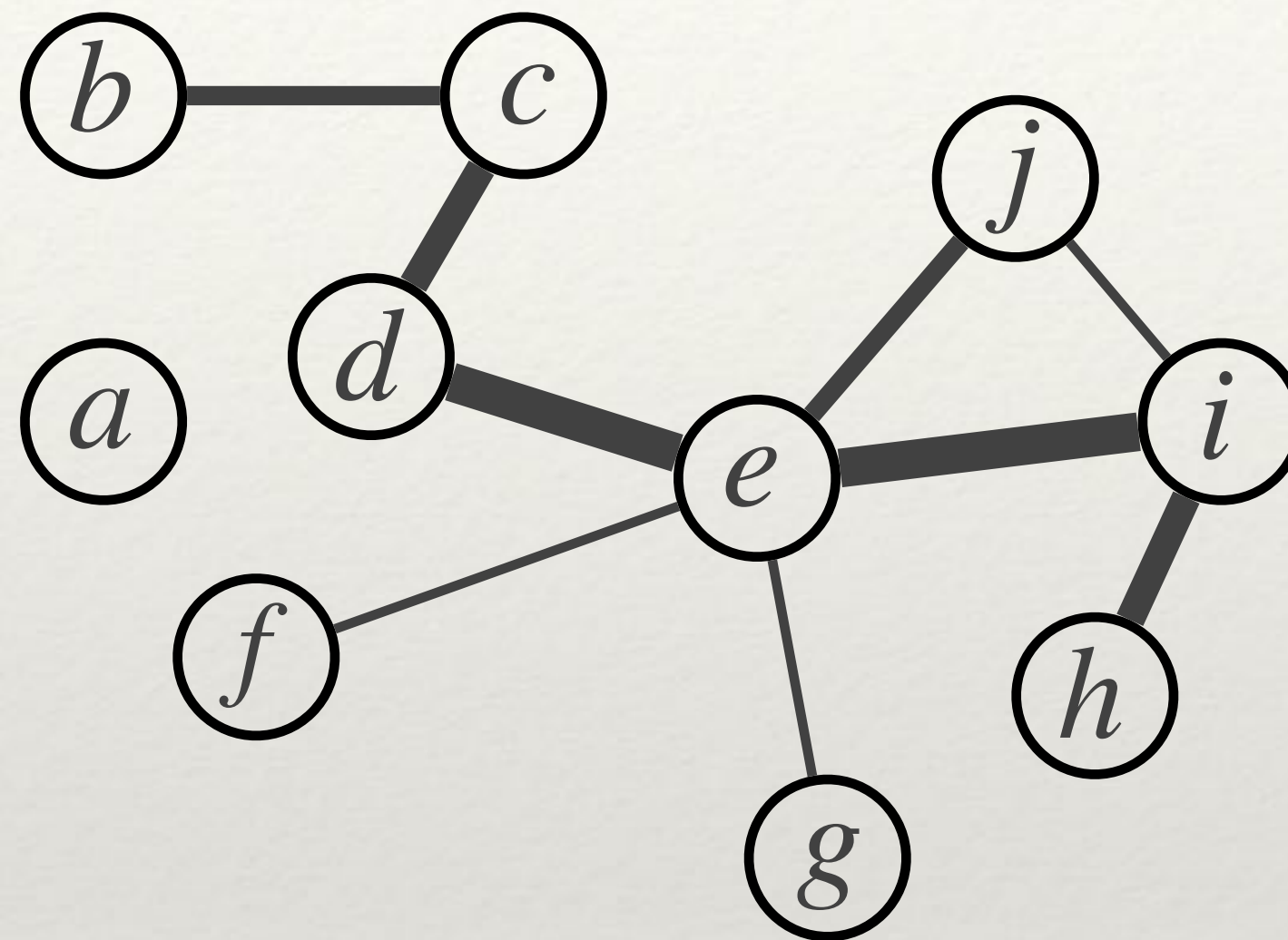


<i>b</i>	<i>c</i>
<i>c</i>	<i>d</i>
<i>d</i>	<i>e</i>
<i>e</i>	<i>f</i>
<i>e</i>	<i>g</i>
<i>e</i>	<i>i</i>
<i>e</i>	<i>j</i>
<i>h</i>	<i>i</i>
<i>i</i>	<i>j</i>

L



# Edge (weighted) list

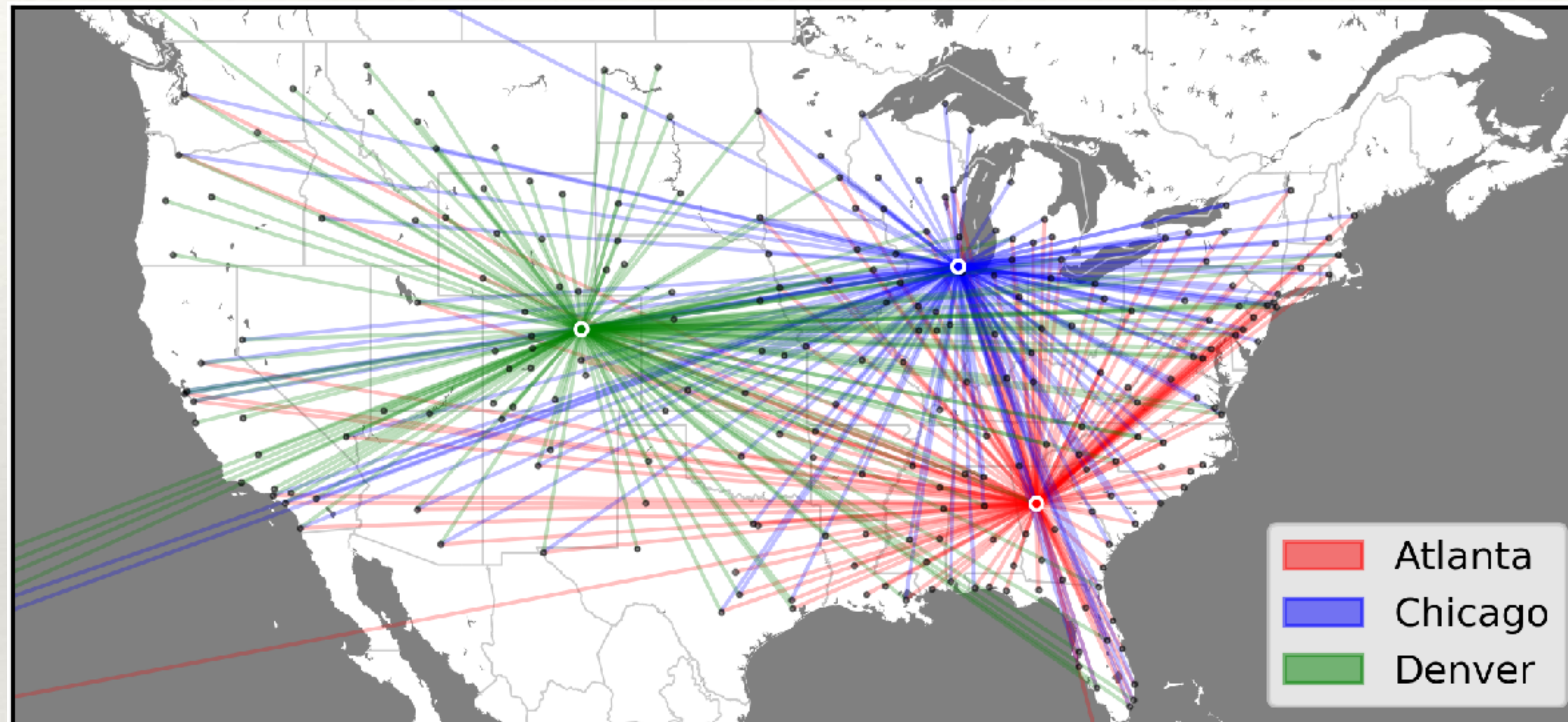


<i>b</i>	<i>c</i>	2
<i>c</i>	<i>d</i>	3
<i>d</i>	<i>e</i>	4
<i>e</i>	<i>f</i>	4
<i>e</i>	<i>g</i>	1
<i>e</i>	<i>i</i>	1
<i>e</i>	<i>j</i>	2
<i>h</i>	<i>i</i>	3
<i>i</i>	<i>j</i>	1

L



# Real networks are heterogeneous



Some nodes (and links) are much more important (**central**) than others!



---

# Centrality measures

---

- ❖ **Centrality:** measure of importance of a node
- ❖ **Measures:**
  1. Degree
  2. Closeness
  3. Betweenness



---

# Degree

---

- **Degree of a node:** number of neighbors of the node

$$k_i = \text{number of neighbors of node } i$$

- High-degree nodes are called **hubs**

- **Average degree of the network:**

$$\langle k \rangle = \frac{\sum_i k_i}{N} = \frac{2L}{N}$$

```
G.degree(2) # returns the degree of node 2
```

```
G.degree() # dict with the degree of all nodes of G
```



---

# Closeness

---

**Idea:** a node is the more central the *closer* it is to the other nodes, on average

$$g_i = \frac{1}{\sum_{j \neq i} \ell_{ij}}$$

where  $\ell_{ij}$  is the distance between nodes  $i$  and  $j$

```
nx.closeness centrality(G, node) # closeness centrality  
                                # of node
```



---

# Betweenness

---

**Idea:** a node is the more central the *more often it is crossed by paths*

$$b_i = \sum_{h \neq j \neq i} \frac{\sigma_{hj}(i)}{\sigma_{hj}}$$

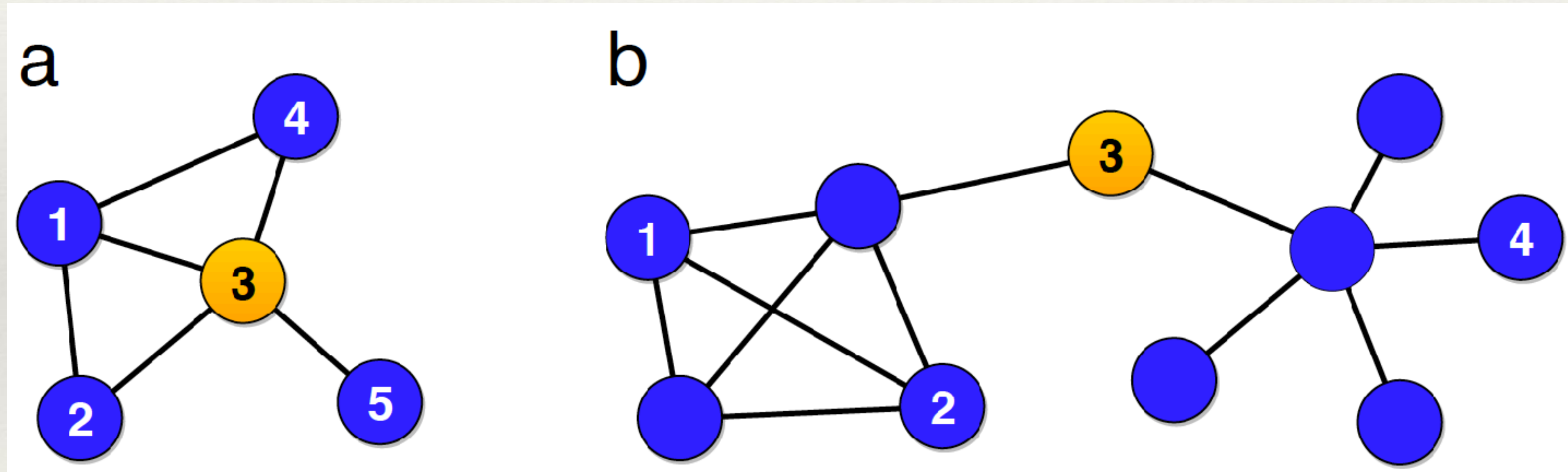
$\sigma_{hj}$  = number of shortest paths from  $h$  to  $j$

$\sigma_{hj}(i)$  = number of shortest paths from  $h$  to  $j$  running through  $i$



# Betweenness

Hubs usually have high betweenness, but there can be nodes with high betweenness **that are not hubs**





---

# Betweenness

---

- Betweenness can be easily extended to links
- **Link betweenness:** fraction of shortest paths among all possible node pairs that pass through the link



---

# Centrality distributions

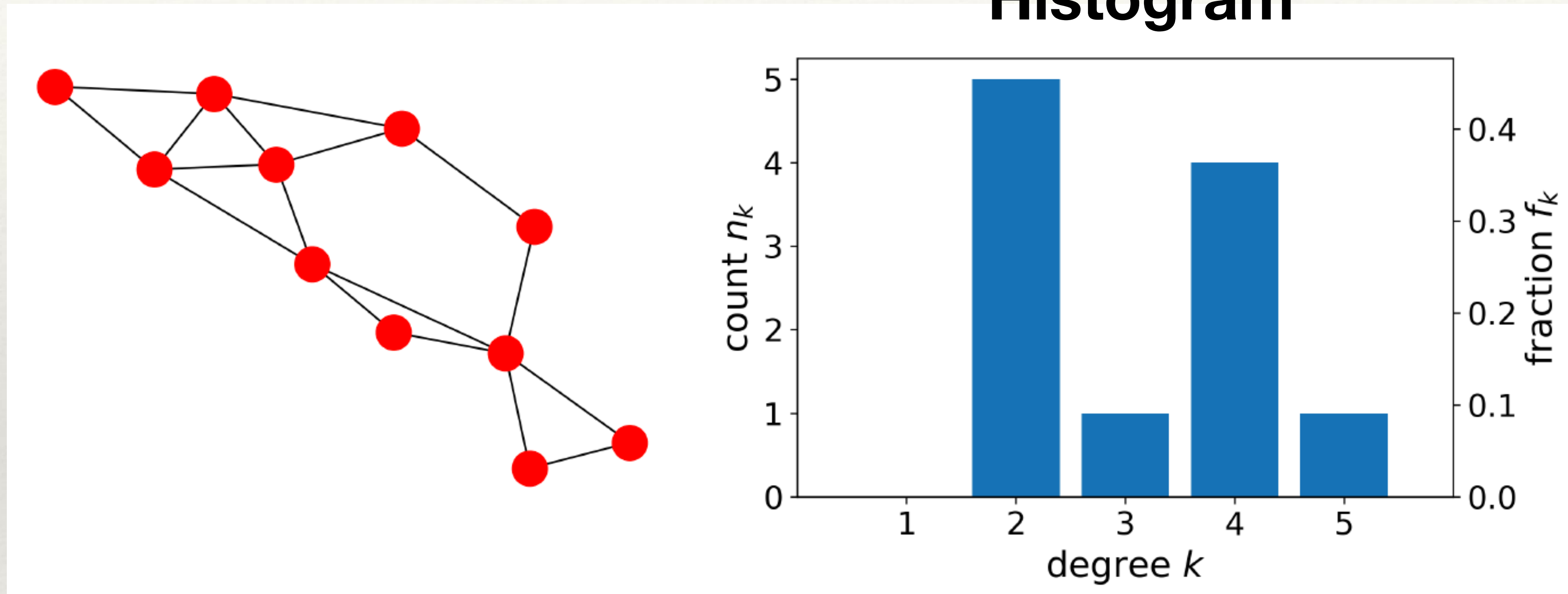
---

- On small networks it makes sense to ask which nodes or links are most important
- On large networks **it does not**
- **Solution:** statistical approach
- Instead of focusing on individual nodes and links, we consider **classes** of nodes and links with similar properties



# Centrality distributions

## Histogram

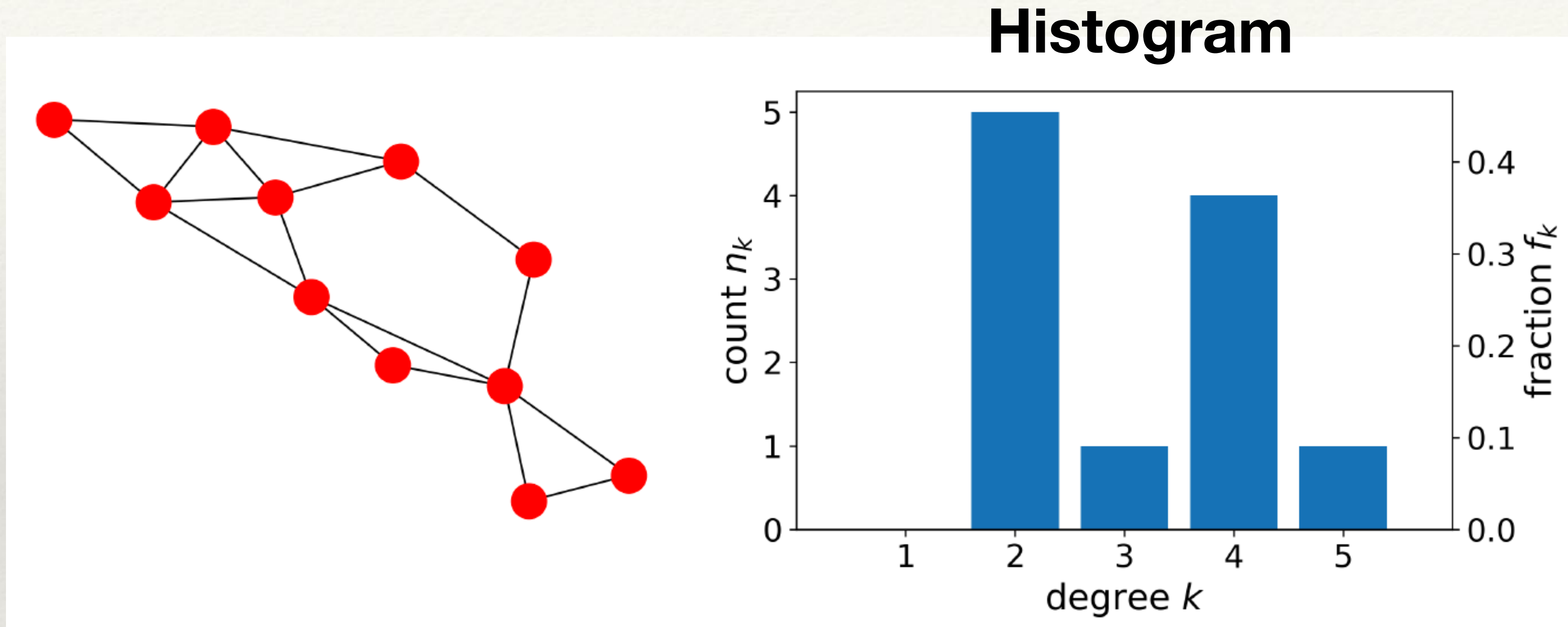


$n_k$  = number of nodes with degree  $k$

$$f_k = \frac{n_k}{N} = \text{frequency of degree } k$$



# Centrality distributions



- For large  $N$ , the frequency  $f_k$  becomes the **probability**  $p_k$  of having degree  $k$
- **Probability distribution:** plot of probability  $p_k$  versus  $k$



---

# Cumulative distributions

---

- If the variable is *not integer* (e.g., betweenness), the range of the variable is divided into intervals (bins) and we count how many values fall in each interval
- **Cumulative distribution  $P(x)$ :** probability that the variable takes values *larger* than  $x$  as a function of  $x$
- **How to compute it:** by summing the frequencies of the variable inside the intervals to the right of  $x$

$$P(x) = \sum_{v \geq x} f_v$$



---

# Logarithmic scale

---

- **Question:** how to plot a probability distribution if the variable spans a large range of values, from small to (very) large?
- **Answer:** use the **logarithmic** scale
- **How to do it:** report the logarithms of the values on the x- and y-axes

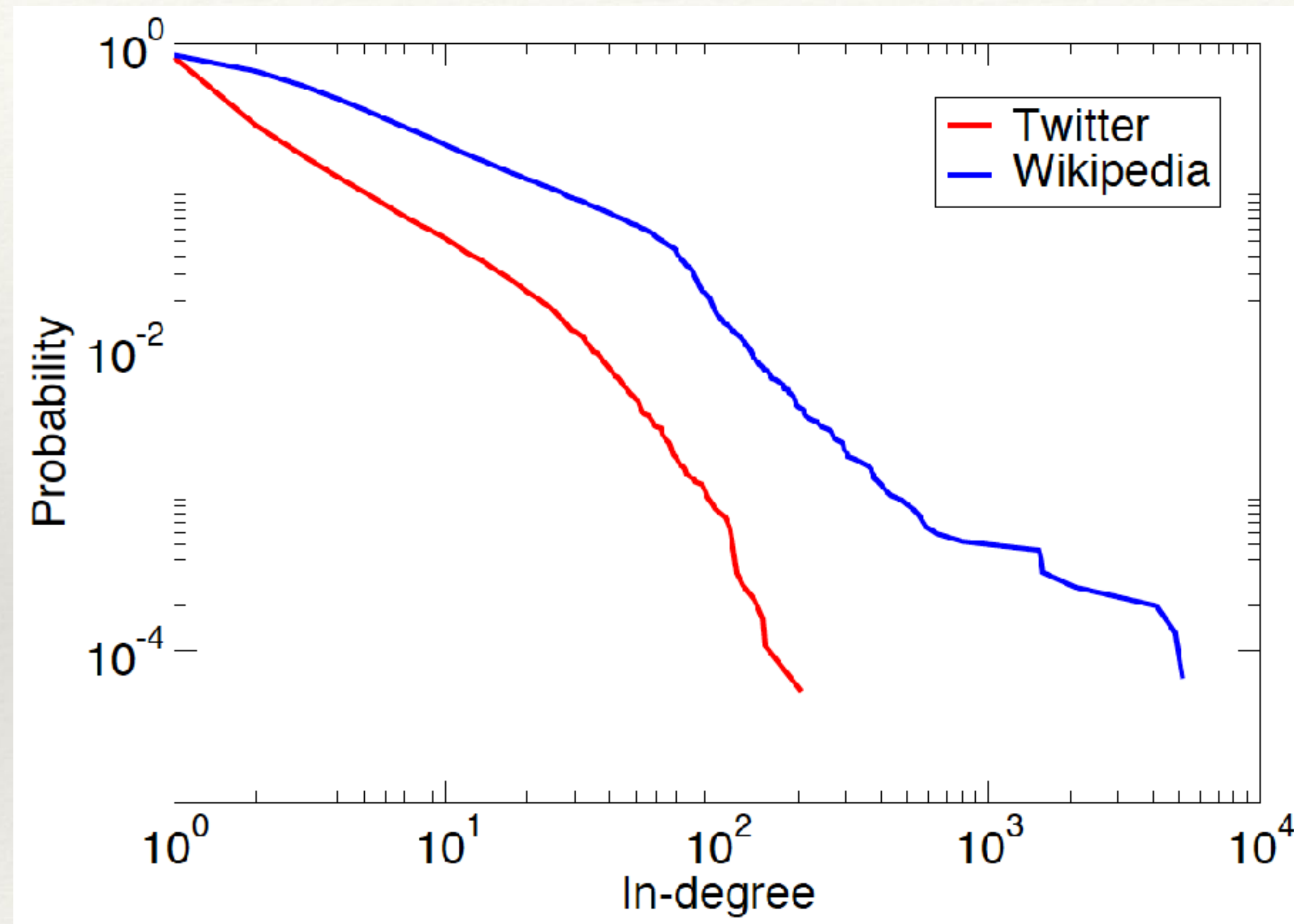
$$\log_{10} 10 = 1$$

$$\log_{10} 1,000 = \log_{10} 10^3 = 3$$

$$\log_{10} 1,000,000 = \log_{10} 10^6 = 6$$



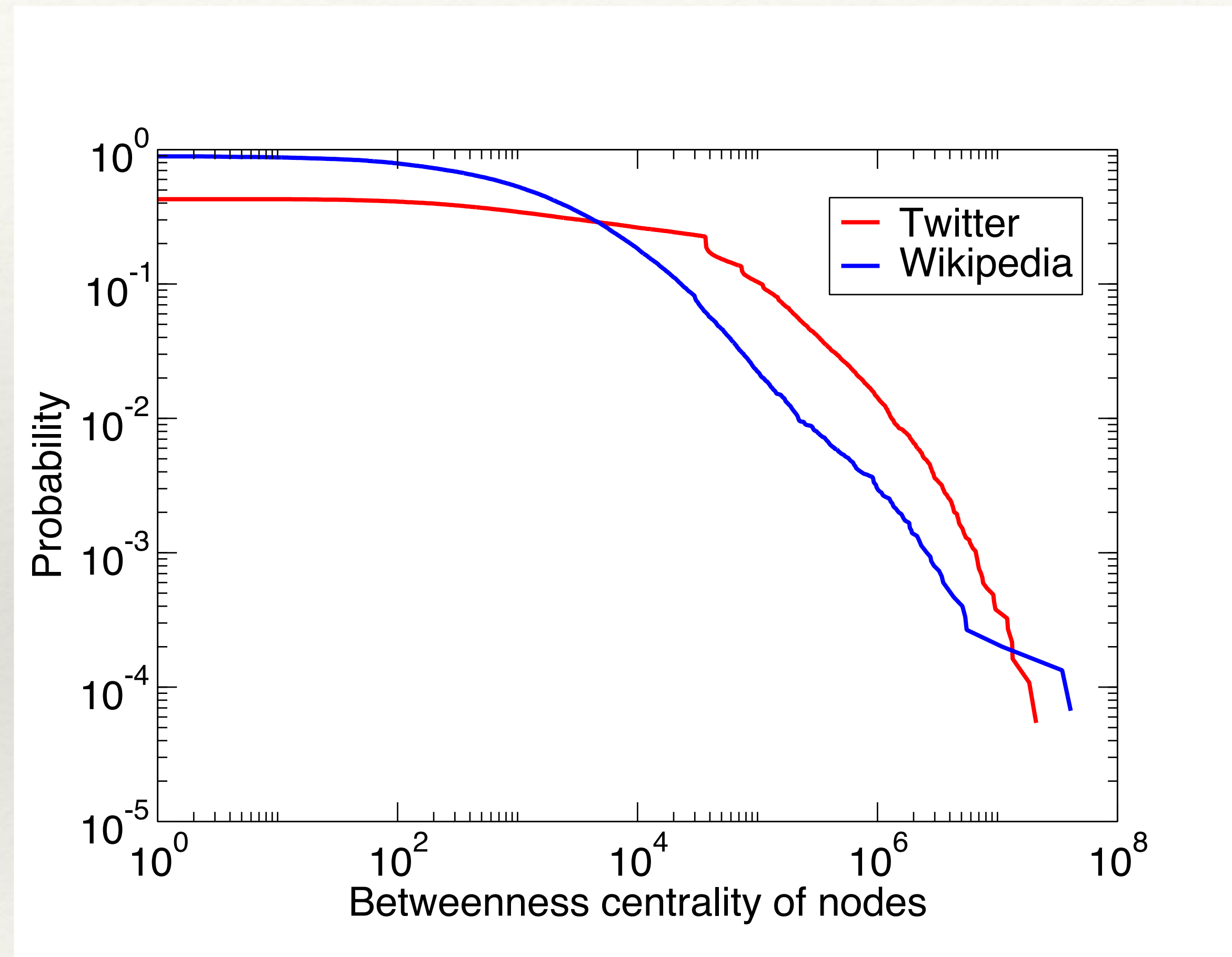
# Degree distributions



**Heavy-tail distributions:** the variable goes from small to large values



# Betweenness distributions



**Heavy-tail distribution:** the variable goes from small to large values



---

# Robustness

---

- A system is **robust** if the failure of some of its components does not affect its function
- **Question:** how can we define the robustness of a network?
- **Answer:** we remove nodes and/or links and see what happens to its structure
- **Key point:** *connectedness*
- If the Internet were not connected, it would be impossible to transmit signals (e.g., emails) between routers in different components



---

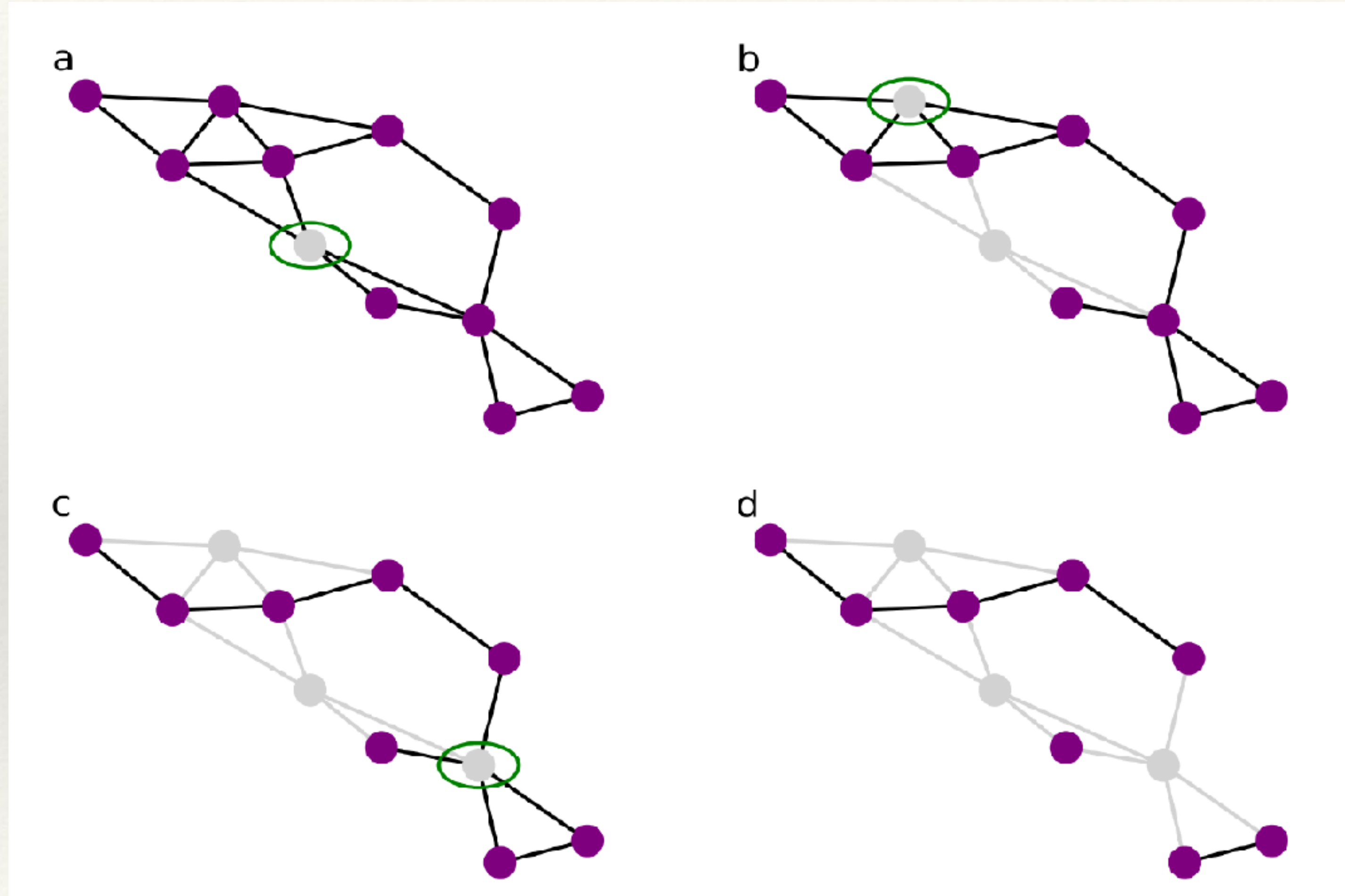
# Robustness

---

- **Robustness test:** checking how the connectedness of the network is affected as more and more nodes are removed
- **How to do it:** plot the relative size  $S$  of the largest connected component as a function of the fraction of removed nodes
- We suppose that the network is initially connected: there is only one component and  $S = 1$
- As more and more nodes (and their links) are removed, the network is progressively broken up into components and  $S$  goes down



# Robustness





---

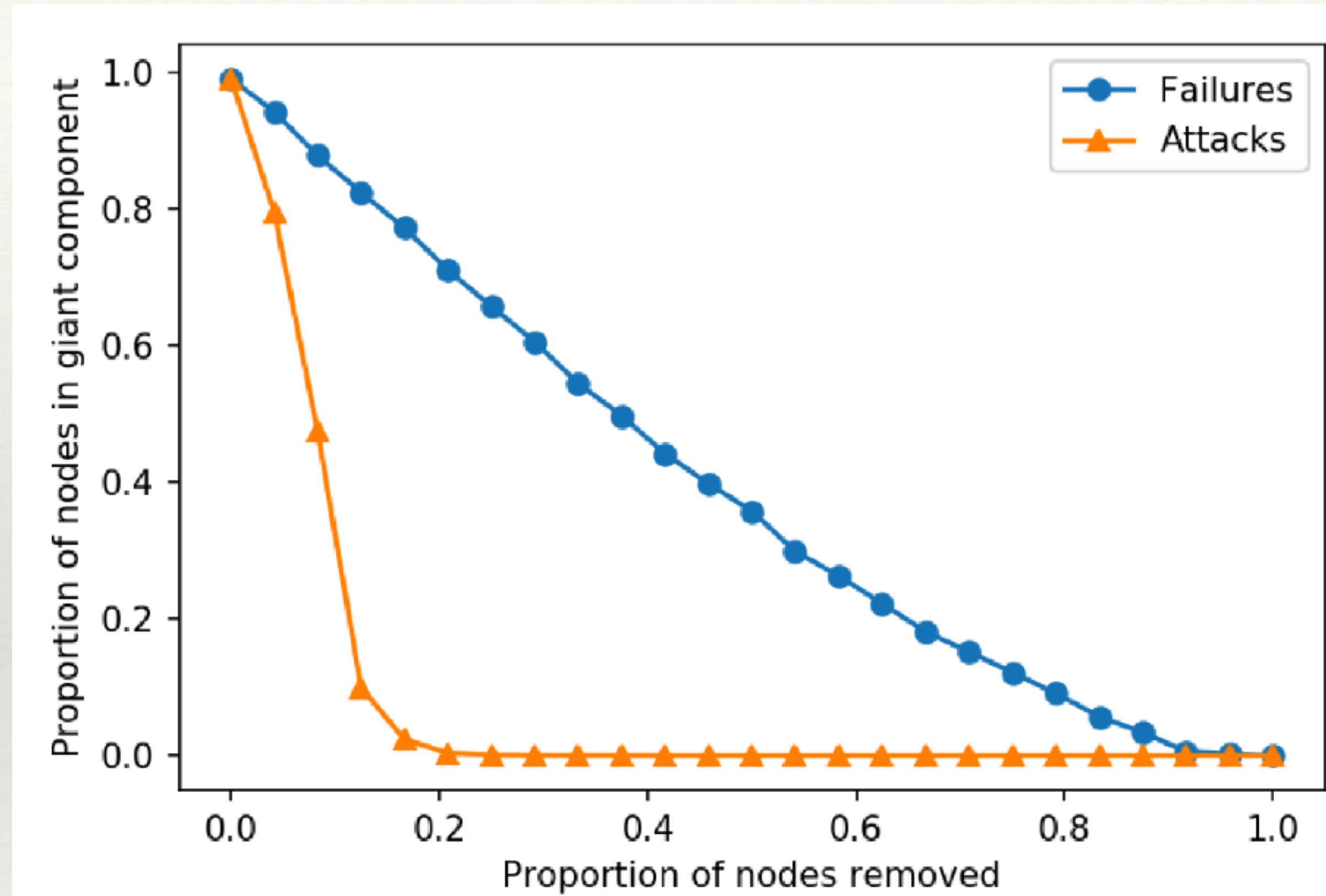
# Robustness

---

- **Two strategies:**
  1. **Random failures:** nodes break down randomly, so they are all chosen with the **same probability**
  2. **Attacks:** hubs are deliberately targeted — the larger the **degree**, the higher the probability of removing the node
- In the first approach, we remove a fraction  $f$  of nodes, chosen at random
- In the second approach, we remove the fraction  $f$  of nodes with largest degree, from the one with largest degree downwards



# Robustness



**Conclusion:** real networks are robust against random failures but fragile against targeted attacks!



---

# Pointer to epidemic modeling

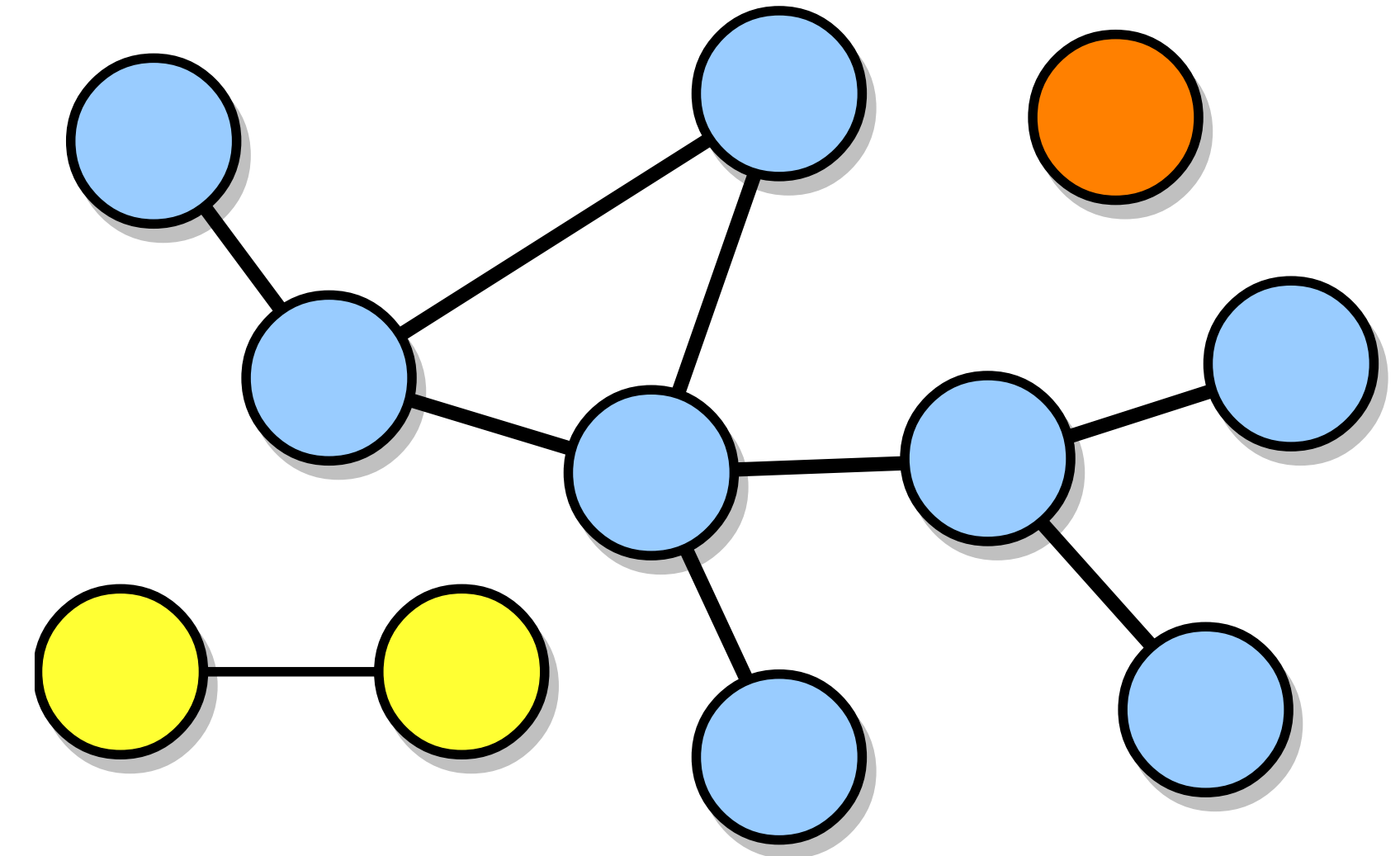
---

- ❖ Studying network robustness is a great framework for comparing different immunization (vaccination) strategies
  - ❖ this very simple idea has been applied also for mitigating the diffusion of computer viruses
- ❖ Problem: real contact network is not usually available...



# Connectedness and components

- ❖ A network is **connected** if there is a path between any two nodes
- ❖ If a network is not connected, it is **disconnected** and has multiple connected components
- ❖ A **connected component** is a connected subnetwork
  - ❖ The largest one is called **giant component**; it often includes a substantial portion of the network
  - ❖ A **singleton** is the smallest-possible connected component





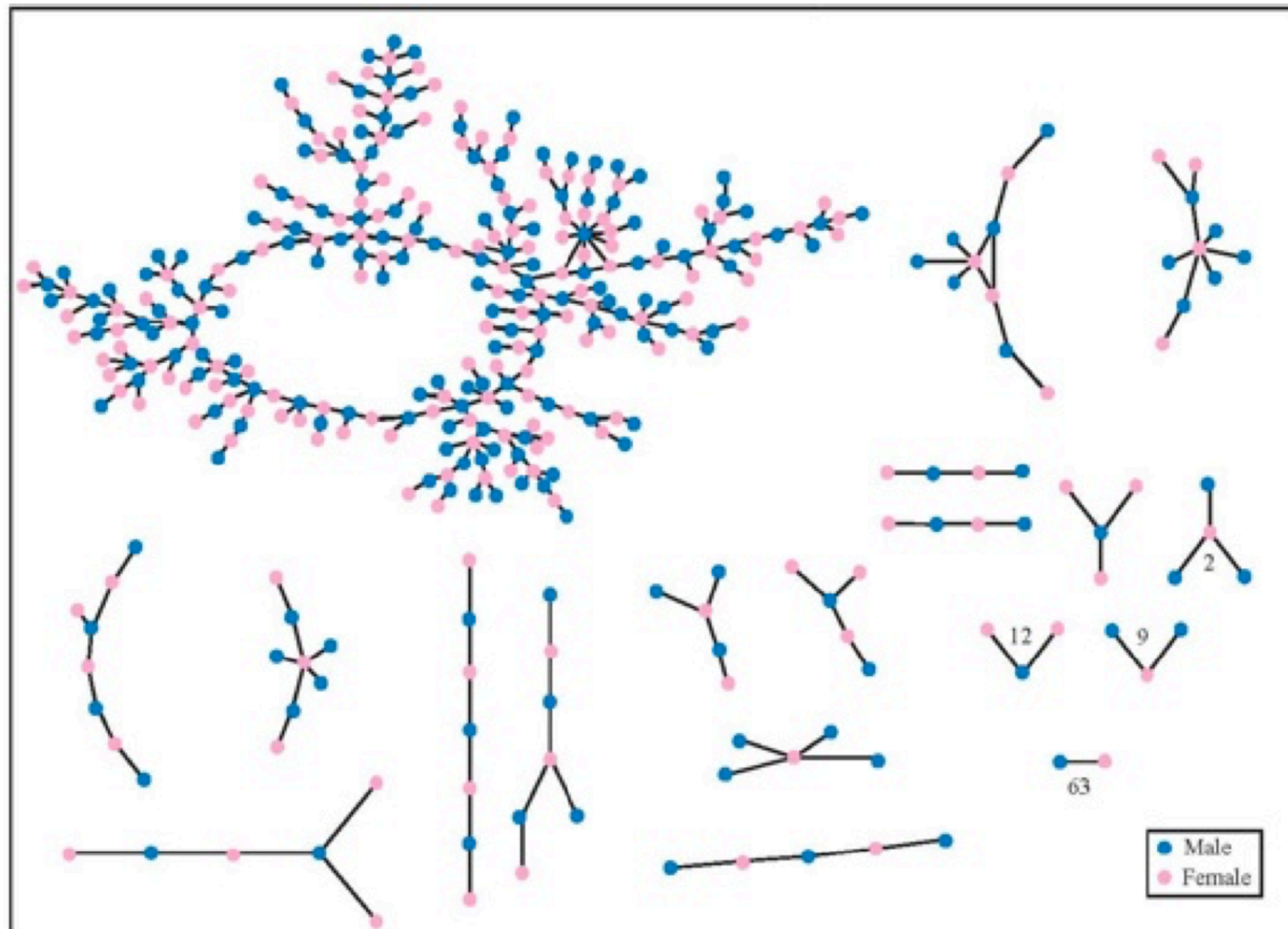


Figure 2.7: A network in which the nodes are students in a large American high school, and an edge joins two who had a romantic relationship at some point during the 18-month period in which the study was conducted [49].

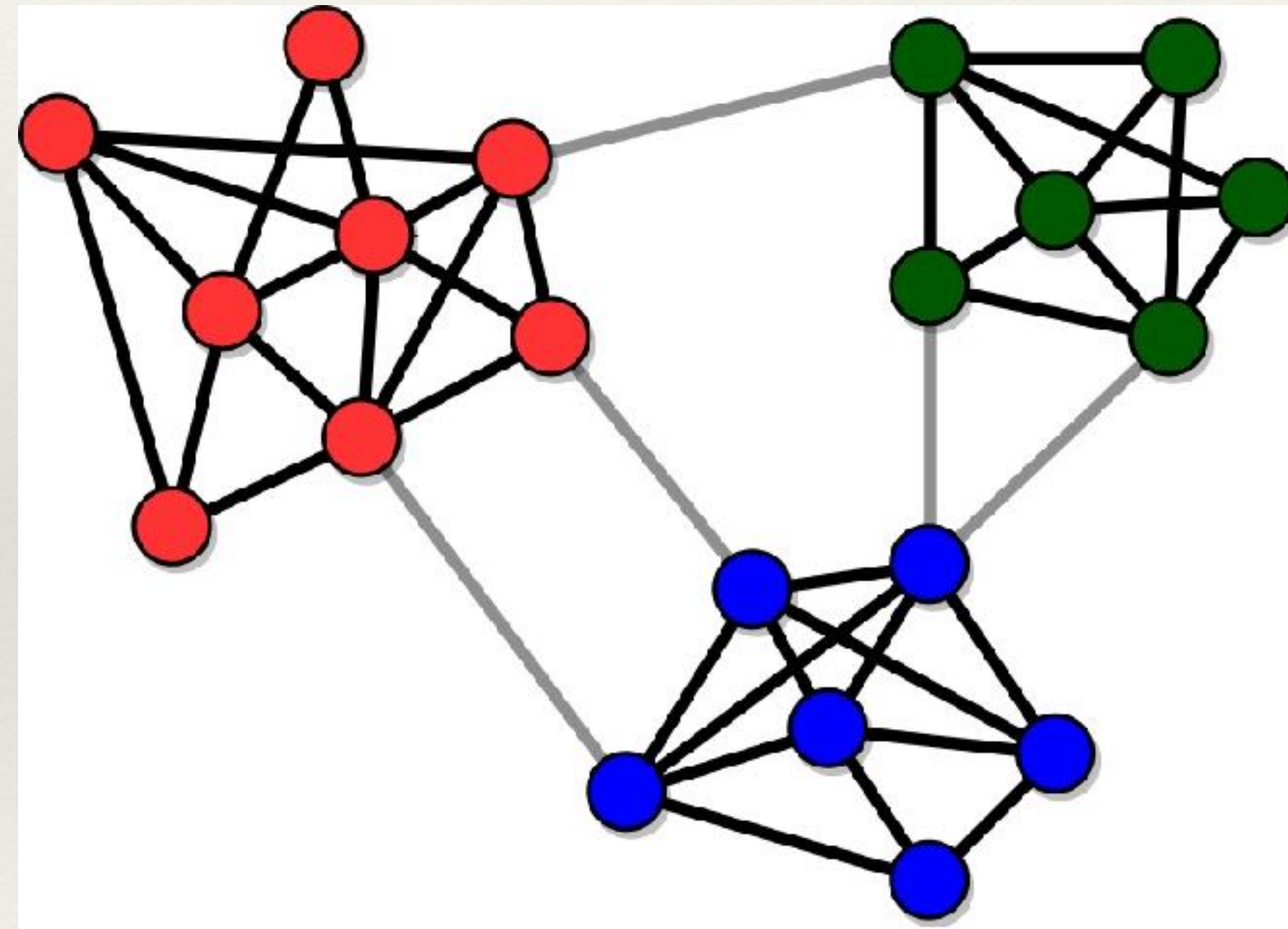


---

# Community structure

---

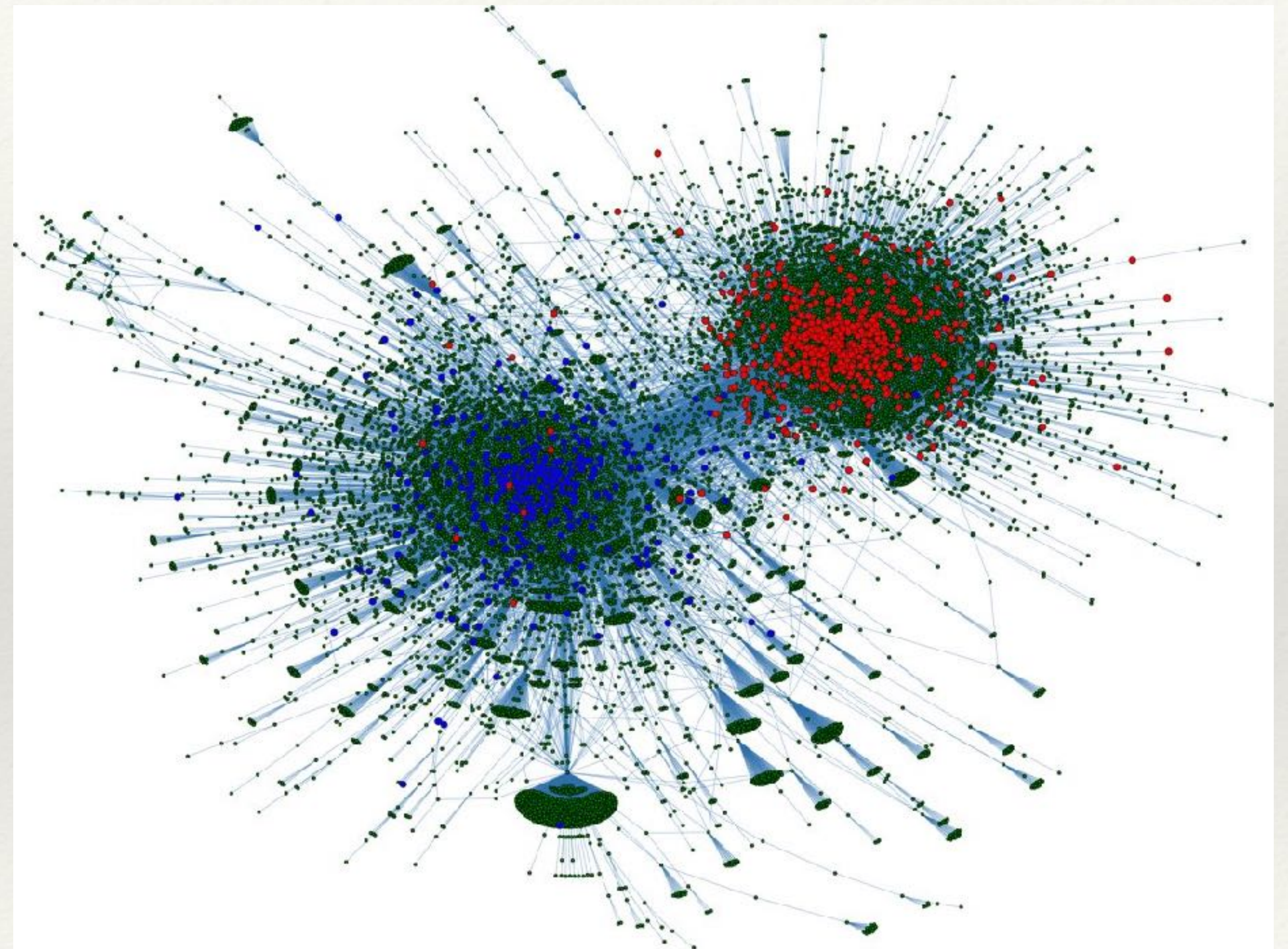
**Communities (or clusters):** sets of tightly connected nodes





# Community structure

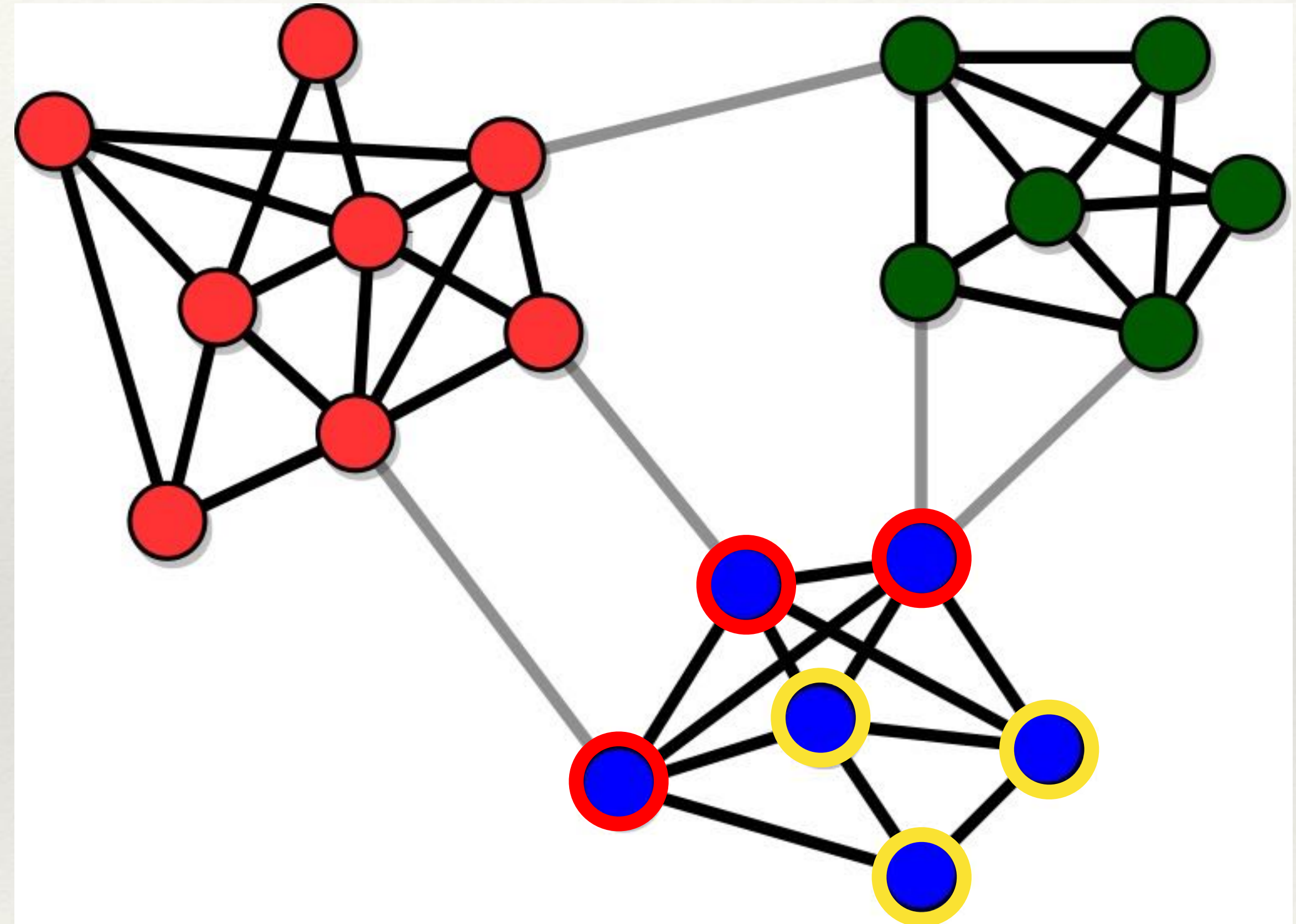
- **Example:** Twitter users with strong political preferences tend to follow those aligned with them and not to follow users with different political orientation
- **Other examples:** social circles in social networks, functional modules in protein interaction networks, groups of pages about the same topic on the Web, etc.





# Why study communities?

- Uncover the organization of the network
- Identify features of the nodes
- Classify the nodes based on their position in the clusters
- Find missing links

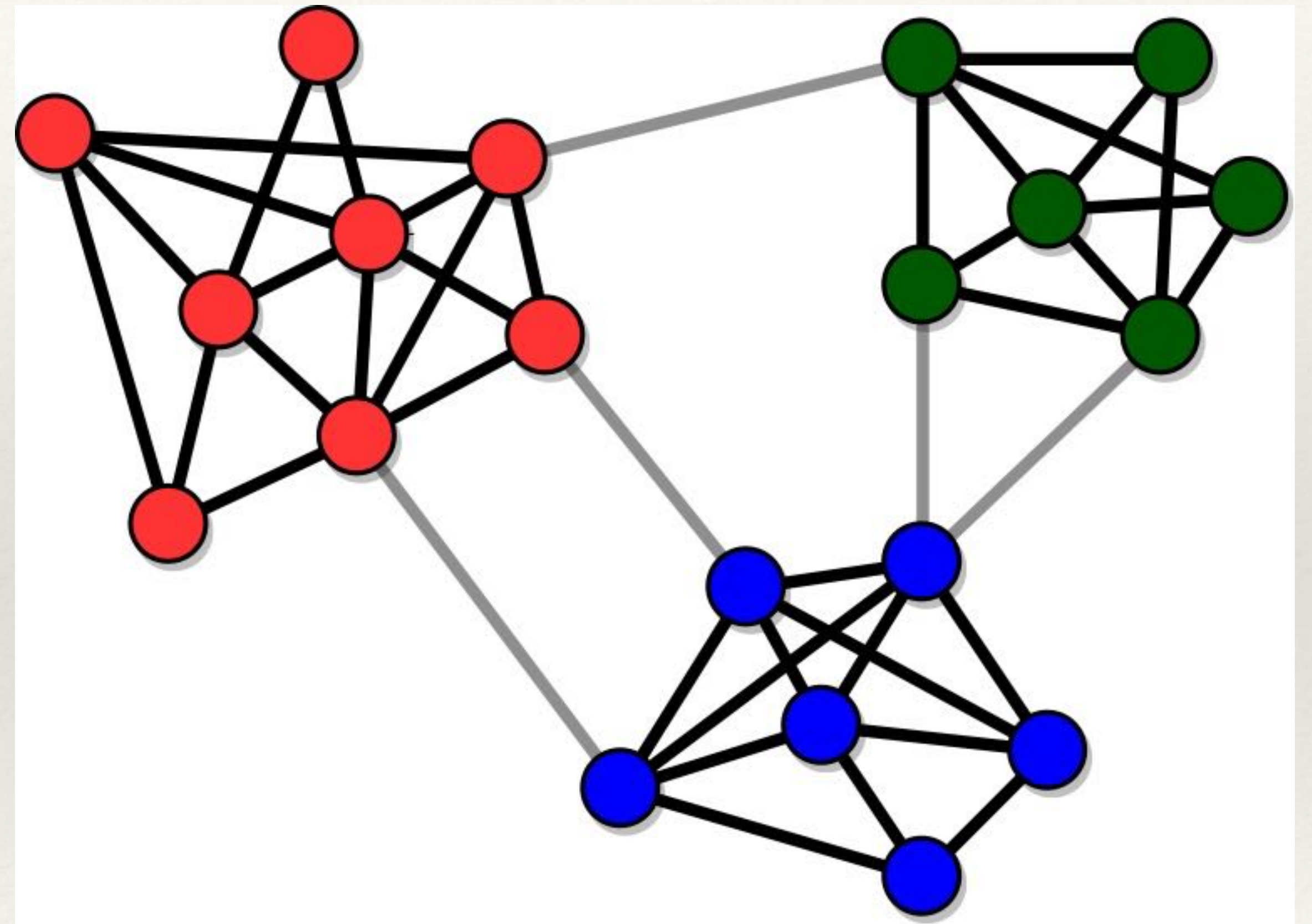




# Basic definitions: community

Two main features:

- **High cohesion:** communities have many internal links, so their nodes stick together
- **High separation:** communities are connected to each other by few links





---

# Partitions

---

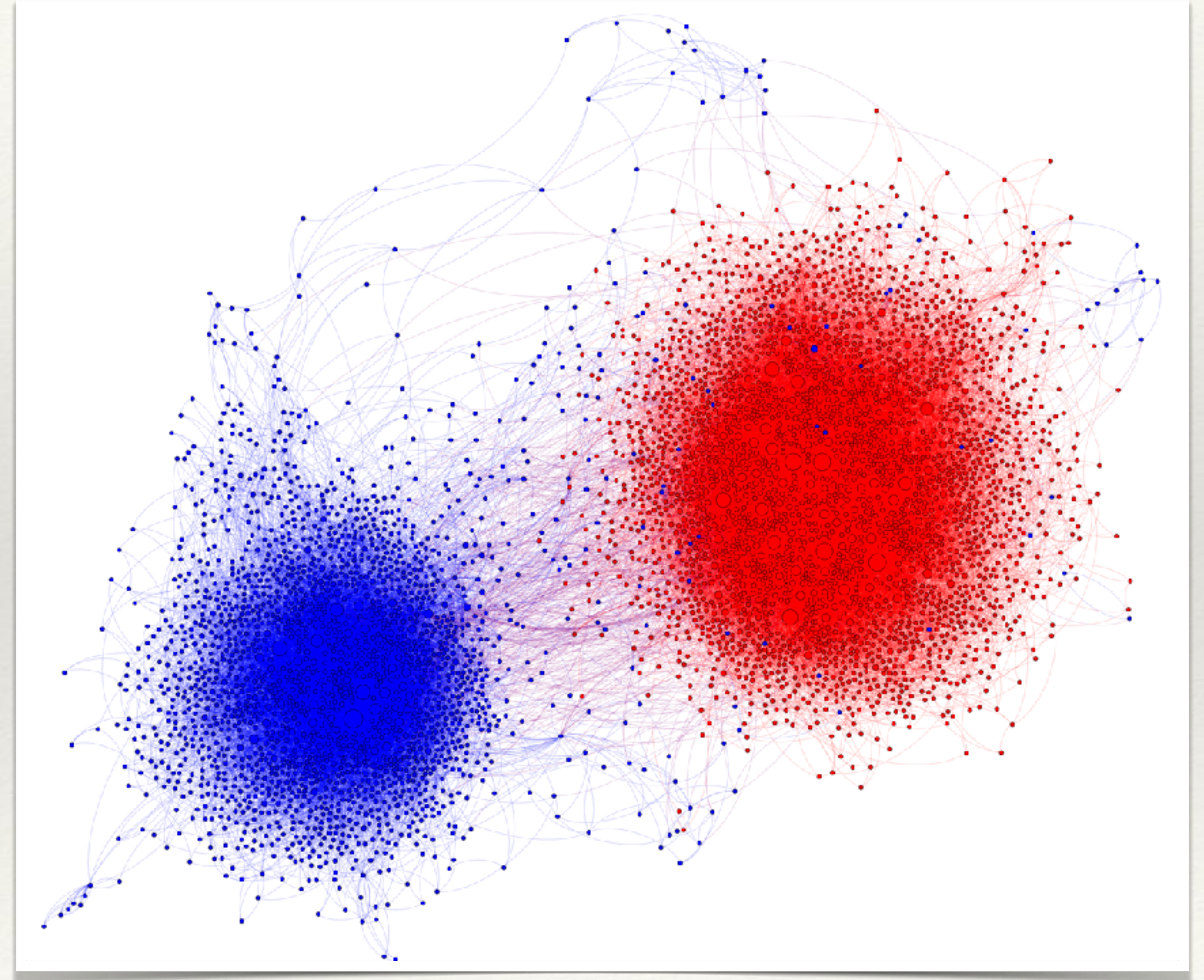
- The number of partitions of  $n$  objects is the **Bell number  $B_n$**
- The Bell number **grows faster than exponentially with  $n$**
- **Conclusion:** it makes no sense to look for interesting community structures by exploring the whole space of partitions! A smart exploration of the partition space must be performed.

$n$	$B_n$
1	1
2	2
3	5
4	15
5	52
6	203
7	877
8	4140
9	21147
10	115975
11	678570
12	4213597
13	27644437
14	190899322
15	1382958545



# example: retweet networks

- ❖ goal: detecting communities or clusters
- ❖ "echo chambers"
- ❖ **homophily**: tendency of individuals to link with similar ones
- ❖ warning: no trivial linear relationships but **interplay**





# Drawing networks

- ❖ A **network layout algorithm** places nodes on a plane to visualize the structure of the network
- ❖ There are many layout algorithms; the most commonly used are **force-directed layout** (a.k.a. **spring layout**) algorithms:
  - ❖ Connected nodes are placed near each other
  - ❖ Links have similar length
  - ❖ Link crossings are minimized
- ❖ This is done by simulating a physical systems where adjacent nodes are connected by springs and otherwise repel each other
- ❖ The community structure of the network can be revealed this way if the network is not too dense or too large

